

# **A Fast Algorithm for Optimizing Ridge Parameters in a Generalized Ridge Regression by Minimizing an Extended GCV Criterion**

**Mineaki Ohishi, Hirokazu Yanagihara and Yasunori Fujikoshi**

Department of Mathematics, Graduate School of Science, Hiroshima University  
1-3-1 Kagamiyama, Higashi-Hiroshima, Hiroshima 739-8626, Japan

## **Abstract**

In this paper, we deal with an optimization of ridge parameters in a generalized ridge regression by minimizing a model selection criterion. Optimization methods based on minimizations of a generalized  $C_p$  ( $GC_p$ ) criterion and GCV criterion have several important advantages and several major disadvantages. The most important advantage of two methods is that the speed of an optimization is very fast because minimizers of two criteria are given by closed forms. The serious disadvantage of the two methods is that the optimization methods do not work well when the number of explanatory variables is larger than the sample size. In order to improve the disadvantage while maintaining the advantage, at first we extend the GCV criterion, called the extended GCV (EGCV), by changing the second power of its denominator to the  $\alpha$ -power, where  $\alpha$  is some positive number larger than 2. Next, we propose a fast optimization algorithm to minimize the EGCV by using specific candidate minimizers of the EGCV of which the number is finite. By conducting numerical examinations, we verify that the optimization method based on the minimization of the EGCV performs well than those based on the minimizations of existing criteria.

(Last Modified: April 14, 2017)

**Key words:** Generalized ridge regression, Generalized cross-validation criterion, Linear regression model, Model selection criterion, Ridge parameters, Optimization of ridge parameters.

E-mail address: mineaki-ohishi@hiroshima-u.ac.jp (Mineaki Ohishi)

## **1. Introduction**

A linear regression model is one of important tools to clarify the relationship between a scalar response variable and one or more explanatory variables. Let  $\mathbf{y} = (y_1, \dots, y_n)'$  be an  $n$ -

dimensional vector of response variables and  $\mathbf{X}$  be an  $n \times k$  matrix of nonstochastic centralized explanatory variables ( $\mathbf{X}'\mathbf{1}_n = \mathbf{0}_k$ ) with  $\text{rank}(\mathbf{X}) = m \leq \min\{n - 1, k\}$ , where  $n$  is the sample size,  $k$  is the number of explanatory variables,  $\mathbf{1}_n$  is an  $n$ -dimensional vector of ones, and  $\mathbf{0}_k$  is a  $k$ -dimensional vector of zeros. An equation of the linear regression model is

$$\mathbf{y} = \mu\mathbf{1}_n + \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\varepsilon},$$

where  $\mu$  is an unknown location parameter,  $\boldsymbol{\beta}$  is a  $k$ -dimensional vector of unknown regression coefficients, and  $\boldsymbol{\varepsilon} = (\varepsilon_1, \dots, \varepsilon_n)'$  is an  $n$ -dimensional vector of independent error variables from a distribution with the mean 0 and variance  $\sigma^2$ .

One of methods to avoid a serious problem of multicollinearity among explanatory variables is the generalized ridge regression (GRR) proposed by Hoerl and Kennard (1970), which gives a shrinkage estimate of  $\boldsymbol{\beta}$  towards  $\mathbf{0}_k$  via multiple ridge parameters  $\theta_1, \dots, \theta_k$ , where  $\theta_j \in \mathbb{R}_+ = \{\theta \in \mathbb{R} \mid \theta \geq 0\}$  ( $j = 1, \dots, k$ ). Although the GRR was proposed 45 years ago, it is still used in many research fields (e.g., Liu *et al.*, 2011; Shen *et al.*, 2013). Since an estimator of the GRR estimator (GRRE) of  $\boldsymbol{\beta}$  changes depending on  $\theta_1, \dots, \theta_k$ , it is important how to optimize  $\theta_1, \dots, \theta_k$ .

An optimization method based on minimization of a model selection criterion (MSC) with respect to  $\theta_1, \dots, \theta_k$  is one of common methods to optimize  $\theta_1, \dots, \theta_k$  (see e.g., Lawless, 1981; Walker & Page, 2001; Yanagihara, 2013). Most MSCs consist of two terms; one is a discrepancy term that measures a goodness of fit of a model, and the other is a penalty term that applies a penalty to a complexity of a model. The generalized  $C_p$  ( $GC_p$ ) criterion (see Atkinson, 1980; Nagai *et al.*, 2012) and generalized cross-validation (GCV) criterion (see Craven & Wahba, 1979; Ye, 1998; Yanagihara, 2013) measure a goodness of fit of model via the residual sum of squares. An optimization method based on the minimization of  $GC_p$  criterion (called the  $GC_p$ -minimization method) has the following advantages and disadvantage:

- Advantages of the  $GC_p$ -minimization method.
  - (1) The speed of the optimization is very fast because minimizers of  $GC_p$  criterion can be derived by the closed forms (see Nagai *et al.*, 2012).
  - (2) The strength of the penalty term can be changed freely.
- A disadvantage of the  $GC_p$ -minimization method.
  - (1) This method cannot be used when  $n$  is smaller than  $k$  because the  $GC_p$  requires an asymptotically unbiased estimator of  $\sigma^2$ .

Since we encounter many opportunities to analyze data with huge explanatory variables, i.e.,

high-dimensional explanatory variables, we have to say that the  $GC_p$ -minimization method has a serious disadvantage. Since the serious problem is occurred from a requirement of an asymptotically unbiased estimator of  $\sigma^2$ , it is solvable by using a model selection criterion that does not require the asymptotically unbiased estimator. The GCV criterion is one of model selection criteria that do not require an asymptotically unbiased estimator of  $\sigma^2$ . However, an optimization method based on the minimization of GCV criterion (called the GCV-minimization method) also has the following advantages and disadvantage:

- Advantages of the GCV-minimization method.
  - (1) The speed of the optimization is very fast because minimizers of the GCV criterion can be derived by the closed forms (see Yanagihara, 2013).
  - (2) This method can be used even when  $n$  is smaller than  $k$  because the GCV does not require an asymptotically unbiased estimator of  $\sigma^2$ .
- A disadvantage of the GCV-minimization method.
  - (1) The strength of the penalty term is fixed.

Although the GCV-minimization method can be used when  $n$  is smaller than  $k$ , this method does not work well then. This is because the GRRE hardly shrink towards  $\mathbf{0}_k$  (see Yanagihara, 2013). The problem will be avoided if the strength of the penalty term in the GCV can be changed freely.

On the other hand, a goodness of fit of a model can be also measured by the Kullback-Leibler (KL) discrepancy (Kullback & Leibler, 1951). Hence, we can define a generalized information criterion (GIC; Nishii, 1984) for the GRR under a normal distribution assumption. An optimization method based on the minimization of GIC (called the GIC-minimization method) also has the following advantages and disadvantage:

- Advantages of the GIC-minimization method.
  - (1) This method can be used even when  $n$  is smaller than  $k$  because the GIC does not require an asymptotically unbiased estimator of  $\sigma^2$ .
  - (2) The strength of the penalty to a complexity of a model can be changed freely.
- A disadvantage of the GIC-minimization method.
  - (1) The speed of the optimization is slow because minimizers of the GIC are derived by an iterative calculation, e.g., the Newton-Raphson method.

In this paper, at the beginning, we propose a fast optimization algorithm to solve a minimization problem of the GIC by using specific candidate minimizers of the GIC of which the number is finite. Unfortunately, we also show that the GIC-minimization method does not work well when  $n$  is smaller than  $k$  even if the strength of the penalty term is increased. Therefore, we extend the GCV so that the strength of the penalty term can be changed freely. We call this criterion an extended GCV (EGCV), and an optimization method based on the minimization of this criterion the EGCV-minimization method. We also propose a fast optimization algorithm to solve a minimization problem of the EGCV by using specific candidate minimizers of the EGCV of which the number is finite.

This paper is organized as follows: In section 2, we describe several results for deriving main theorems. In section 3, we give a fast optimization algorithm of the GIC-minimization method, and show the reason why the GIC-minimization method does not work well when  $k$  is larger than  $n$ . In section 4, we propose the EGCV criterion and fast optimization algorithm to solve a minimization problem of EGCV. Numerical examinations are conducted in Section 5. Technical details are provided in the Appendix.

## 2. Preliminaries

It follows from the results in Yanagihara (2013) that the GRRE of  $\beta$  and a hat matrix of the GRR are invariant to any changes in  $\theta_{m+1}, \dots, \theta_k$  when  $m < k$ . From this fact, we set  $\theta_{m+1} = \dots = \theta_k = \infty$  to reduce a variance of the GRRE of  $\beta$ . Therefore, it is enough just to minimize MSC with respect to  $\theta_1, \dots, \theta_m$  for optimizing ridge parameters in the GRR.

Let  $Q$  be a  $k$ th orthogonal matrix that diagonalizes  $X'X$ , i.e.,

$$Q'X'XQ = \begin{pmatrix} D & O_{m,k-m} \\ O_{k-m,m} & O_{k-m,k-m} \end{pmatrix},$$

where  $D$  is an  $m$ th diagonal matrix defined by

$$D = \text{diag}(d_1, \dots, d_m), \quad d_1 \geq \dots \geq d_m \text{ are nonzero eigenvalues of } X'X, \quad (2.1)$$

and  $O_{k,m}$  is a  $k \times m$  matrix of zeros, and let  $\theta$  be an  $m$ -dimensional vector given by  $\theta = (\theta_1, \dots, \theta_m)' \in \mathbb{R}_+^m$ , and  $\Theta$  be a  $k$ th diagonal matrix given by  $\Theta = \text{diag}(\theta_1, \dots, \theta_m, 0, \dots, 0)$ , where  $\mathbb{R}_+^m$  is the  $m$ th Cartesian power set of  $\mathbb{R}_+$ . Notice that  $d_1, \dots, d_m$  are positive because  $X'X$  is a positive semidefinite matrix. Moreover, let  $M_\theta$  be a  $k$ th square matrix defined by  $M_\theta = X'X + Q\Theta Q'$ . In particular, we write  $M_\theta = M$  when  $\theta = \mathbf{0}_m$ . Then, the GRRE of  $\beta$  is expressed by

$$\hat{\beta}_\theta = M_\theta^+ X' y, \quad (2.2)$$

where  $A^+$  is the Moore-Penrose inverse matrix of a matrix  $A$  (for details of the Moore-Penrose inverse matrix, see e.g., Harville, 1997, chap. 20). It is well known fact that the least square estimator (LSE) of  $\beta$  is given by

$$\hat{\beta} = M^+ X' y. \quad (2.3)$$

It is clear that (2.2) coincides with (2.3) when  $\theta = \mathbf{0}_m$ . Let  $\hat{\mu}$  be the LSE of  $\mu$  given by  $\hat{\mu} = \bar{y}$ , where  $\bar{y}$  is the sample mean of  $y_1, \dots, y_n$ , i.e.,  $\bar{y} = \sum_{j=1}^n y_j/n$ . The equation (2.2) and  $\hat{\mu}$  imply a predictive value of  $y$  as

$$\hat{y}_\theta = \hat{\mu} \mathbf{1}_n + X \hat{\beta}_\theta = (J_n + X M_\theta^+ X') y,$$

where  $J_n$  is an  $n \times n$  projection matrix defined by  $J_n = \mathbf{1}_n \mathbf{1}_n' / n$ . By using  $M_\theta^+$  and  $\hat{y}_\theta$ , we define an estimator of  $\sigma^2$  and a generalized degree of freedom (df) of the GRR that are main bodies of the discrepancy and penalty terms of MSC, respectively, as

$$\hat{\sigma}^2(\theta) = \frac{1}{n} (y - \hat{y}_\theta)' (y - \hat{y}_\theta) = \frac{1}{n} y' (I_n - J_n - X M_\theta^+ X')^2 y, \quad (2.4)$$

$$\text{df}(\theta) = \text{tr}(J_n + X M_\theta^+ X') = 1 + \text{tr}(M_\theta^+ M). \quad (2.5)$$

Most MSCs for the GRR are defined by a bivariate function of  $\hat{\sigma}^2(\theta)$  and  $\text{df}(\theta)$ .

Let  $z$  be an  $m$ -dimensional vector defined by

$$z = (z_1, \dots, z_m)' = (D^{-1/2}, O_{m,k-m}) Q' X' y = D^{-1/2} Q_1' X' y, \quad (2.6)$$

where  $Q_1$  is a  $k \times m$  matrix defined by  $Q = (Q_1, Q_2)$ . Here, we assume that all  $z_1, \dots, z_m$  are not 0. If  $z_j$  is accidentally 0, then we set  $d_j = 0$  and use  $(z_1, \dots, z_{j-1}, z_{j+1}, \dots, z_m)'$  as  $z$ . By using the result in Yanagihara (2013),  $\hat{\sigma}^2(\theta)$  in (2.4) and  $\text{df}(\theta)$  in (2.5) can be rewritten as

$$\hat{\sigma}^2(\theta) = \frac{1}{n} \left\{ n \hat{\sigma}_0^2 + \sum_{j=1}^m \left( \frac{\theta_j}{d_j + \theta_j} \right)^2 z_j^2 \right\}, \quad \text{df}(\theta) = 1 + m - \sum_{j=1}^m \frac{\theta_j}{d_j + \theta_j}, \quad (2.7)$$

where  $\hat{\sigma}_0^2$  is the maximum likelihood estimator of  $\sigma^2$  under normality, which is given by

$$\hat{\sigma}_0^2 = \frac{1}{n} y' (I_n - J_n - X M^+ X') y. \quad (2.8)$$

It should be kept in mind that  $\hat{\sigma}_0^2 = 0$  holds in most cases of the linear regression with high-dimensional explanatory variables. Notice that  $\theta_j / (d_j + \theta_j)$  monotonically increases  $\theta_j$ , and  $\hat{\sigma}^2(\mathbf{0}_m) = \hat{\sigma}_0^2$ ,  $\text{df}(\mathbf{0}_m) = m + 1$ , and

$$\lim_{\theta_1 \rightarrow \infty, \dots, \theta_m \rightarrow \infty} \hat{\sigma}^2(\theta) = \hat{\sigma}_\infty^2 = \frac{1}{n} y' (I_n - J_n) y, \quad \lim_{\theta_1 \rightarrow \infty, \dots, \theta_m \rightarrow \infty} \text{df}(\theta) = 1. \quad (2.9)$$

Hence, the ranges of  $\hat{\sigma}^2(\theta)$  and  $\text{df}(\theta)$  are given by

$$\hat{\sigma}^2(\boldsymbol{\theta}) \in [\hat{\sigma}_0^2, \hat{\sigma}_\infty^2], \quad d\mathbf{f}(\boldsymbol{\theta}) \in [1, m + 1].$$

Let  $f(r, u)$  be a bivariate function which satisfies the following assumptions:

(A1)  $f(r, u)$  is a continuous function at any  $(r, u) \in (0, \hat{\sigma}_\infty^2] \times [1, n)$ .

(A2)  $f(r, u) > 0$  for any  $(r, u) \in (0, \hat{\sigma}_\infty^2] \times [1, n)$ .

(A3)  $f(r, u)$  is the first-order partially differentiable at any  $(r, u) \in (0, \hat{\sigma}_\infty^2] \times [1, n)$ , and

$$\dot{f}_r(r, u) = \frac{\partial}{\partial r} f(r, u) > 0, \quad \dot{f}_u(r, u) = \frac{\partial}{\partial u} f(r, u) > 0, \quad \forall (r, u) \in (0, \hat{\sigma}_\infty^2] \times [1, n).$$

Let  $\mathbb{D}$  be the domain of  $f$ . It is easy to see that  $m = n - 1 \Rightarrow \hat{\sigma}_0^2 = 0$ . Recall that  $m \leq n - 1$ . Hence, we have

$$\hat{\sigma}_0^2 \neq 0 \Rightarrow \mathbb{D} \subseteq (0, \hat{\sigma}_\infty^2] \times [1, n).$$

By using the bivariate function  $f$ , MSC for the GRR can be expressed as

$$\text{MSC}(\boldsymbol{\theta}) = f(\hat{\sigma}^2(\boldsymbol{\theta}), d\mathbf{f}(\boldsymbol{\theta})). \quad (2.10)$$

The relationship between an existing criterion and the bivariate function  $f$  is as follows:

$$f(r, u) = \begin{cases} f_{GC_p}(r, u) = nr/s_0^2 + \alpha u & (GC_p : \text{when } \hat{\sigma}_0^2 \neq 0) \\ f_{GCV}(r, u) = r(1 - u/n)^{-2} & (GCV : \text{when } u \neq n) \\ f_{GIC}(r, u) = r \exp(\alpha u/n) & (GIC) \end{cases}, \quad (2.11)$$

where  $\alpha$  is some positive value expressing the strength of the penalty term, and  $s_0^2$  is an ordinary unbiased estimator of  $\sigma^2$  when  $m < n - 1$ , which is given by  $s_0^2 = n\hat{\sigma}_0^2/(n - m - 1)$ . The  $GC_p$  includes several famous MSCs, e.g., the  $GC_p$  coincides with the  $C_p$  proposed by Mallows (1973) when  $\alpha = 2$ , and the modified  $C_p$  ( $MC_p$ ) proposed by Fujikoshi and Satoh (1997) and Yanagihara *et al.* (2009) when  $\alpha = 1 + 2/(n - m - 1)$ . Moreover, the original GIC under normality is expressed as  $n \log r + \alpha u$ . The GIC in this paper is defined by an exponential transformation of the original GIC divided by  $n$ . The GIC includes several famous MSCs, e.g., the GIC coincides the AIC proposed by Akaike (1973) when  $\alpha = 2$ , the HQC proposed by Hannan and Quinn (1978) when  $\alpha = 2 \log \log n$ , and the BIC proposed by Schwarz (1978) when  $\alpha = \log n$ .

In order to express an optimized GRRE of  $\boldsymbol{\beta}$  systematically, we prepare the following class of  $\boldsymbol{\theta}$  for any  $h \in \mathbb{R}_+$ :

$$\hat{\boldsymbol{\theta}}(h) = (\hat{\theta}_1(h), \dots, \hat{\theta}_m(h))', \quad \hat{\theta}_j(h) = \begin{cases} \frac{d_j h}{z_j^2 - h} & (h \leq z_j^2) \\ \infty & (h > z_j^2) \end{cases}. \quad (2.12)$$

The GRRE of  $\beta$  based on  $\hat{\theta}$  is given by

$$\hat{\beta}_{\hat{\theta}(\hat{h})} = \hat{\beta}(\hat{h}) = \mathbf{Q}_1 \mathbf{V}(\hat{h}) \mathbf{Q}'_1 \hat{\beta}, \quad (2.13)$$

where  $\hat{\beta}$  is the LSE of  $\beta$  given by (2.3), and  $\mathbf{V}(h)$  is a  $k$ th diagonal matrix of which the  $j$ th diagonal element is given by

$$v_j(h) = \frac{d_j}{d_j + \hat{\theta}_j(h)} = \begin{cases} 1 - \frac{h}{z_j^2} & (h \leq z_j^2) \\ 0 & (h > z_j^2) \end{cases}.$$

An optimized  $\theta$  by minimizing  $\text{MSC}(\theta)$  when the domain of  $f$  is a subset of  $(0, \hat{\sigma}_\infty^2] \times [1, n)$  is expressed as in the following theorem (the proof is given in Appendix A.1):

**Lemma 1.** *Let*

$$\phi(h) = \text{MSC}(\hat{\theta}(h)) \quad (h \in \mathbb{R}_+ \setminus \{0\}). \quad (2.14)$$

*Suppose that  $\mathbb{D} \subseteq (0, \hat{\sigma}_\infty^2] \times [1, n)$  and  $\exists \xi > 0$  s.t.  $\phi(\xi) < \lim_{h \rightarrow 0} \phi(h)$ . Then, optimized ridge parameters by minimizing the  $\text{MSC}(\theta)$  are given by  $\hat{\theta}(\hat{h})$ , where  $\hat{h}$  is given by*

$$\hat{h} = \arg \min_{h \in \mathbb{R}_+ \setminus \{0\}} \phi(h).$$

Let  $t_j$  ( $j = 1, \dots, m$ ) be the  $j$ th-order statistic of  $z_1^2, \dots, z_m^2$ , i.e.,

$$t_j = \begin{cases} \min\{z_1^2, \dots, z_m^2\} & (j = 1) \\ \min\{z_1^2, \dots, z_m^2\} \setminus \{t_1, \dots, t_{m-1}\} & (j = 2, \dots, m) \end{cases},$$

and let  $R_j$  ( $j = 0, \dots, m$ ) be a range defined by

$$R_j = \begin{cases} (0, t_1] & (j = 0) \\ (t_j, t_{j+1}] & (j = 1, \dots, m-1) \\ (t_m, \infty) & (j = m) \end{cases}. \quad (2.15)$$

By using  $t_j$  and  $R_j$ , we can clarify properties of  $\phi(h)$  given in (2.14) as in the following lemma (the proof is given in Appendix A.2):

**Lemma 2.** *The  $\phi(h)$  satisfies the following properties:*

(P1)  $\phi(h)$  is a continuous at any  $h \in \mathbb{R}_+ \setminus \{0\}$ .

(P2)  $\phi(h) = f(\hat{\sigma}_\infty^2, 1)$  for every  $h \geq t_m$ .

(P3)  $\phi(h)$  is a piecewise function as

$$\phi(h) = \phi_a(h) \quad (h \in R_a, a = 0, 1, \dots, m).$$

The specific form of  $\phi_a(h)$  can be expressed as

$$\phi_a(h) = f\left(\hat{\sigma}_0^2 + (c_{1,a} + h^2 c_{2,a})/n, 1 + m - a - hc_{2,a}\right),$$

where  $c_{1,a}$  and  $c_{2,a}$  are given by

$$c_{1,a} = \sum_{j=1}^a t_j, \quad c_{2,a} = \sum_{j=a+1}^m \frac{1}{t_j}. \quad (2.16)$$

Let  $s_a^2$  ( $a = 0, 1, \dots, m$ ) be a statistic defined by

$$s_a^2 = \frac{n\hat{\sigma}_0^2 + c_{1,a}}{n - m - 1 + a}, \quad (2.17)$$

and, let  $a^*$  be a positive integer defined by

$$a^* \in \{0, 1, \dots, m-1\} \text{ s.t. } s_{a^*}^2 \in R_{a^*}, \quad (2.18)$$

From Lemma 2.1 in Yanagihara (2013), we can see that  $a^*$  exists unique when  $\hat{\sigma}_0^2 \neq 0$ . Let

$$GC_p(\boldsymbol{\theta}) = f_{GC_p}(\hat{\sigma}^2(\boldsymbol{\theta}), \text{df}(\boldsymbol{\theta})), \quad \text{GCV}(\boldsymbol{\theta}) = f_{\text{GCV}}(\hat{\sigma}^2(\boldsymbol{\theta}), \text{df}(\boldsymbol{\theta})).$$

From the results in Nagai *et al.* (2012) and Yanagihara (2013), the minimizers of the  $GC_p$  and GCV can be expressed as in the following theorem (the proof is omitted because it is easy to obtain from the results in Nagai *et al.*, 2012; Yanagihara, 2013):

**Theorem 1.** *Optimized ridge parameters minimizing  $GC_p(\boldsymbol{\theta})$  and  $\text{GCV}(\boldsymbol{\theta})$  are  $\hat{\boldsymbol{\theta}}(\hat{h}_{GC_p})$  and  $\hat{\boldsymbol{\theta}}(\hat{h}_{\text{GCV}})$ , respectively, where  $\hat{h}_{GC_p}$  and  $\hat{h}_{\text{GCV}}$  are defined as*

$$\hat{h}_{GC_p} = \frac{1}{2}\alpha s_0^2 \quad (\hat{\sigma}_0^2 \neq 0), \quad \hat{h}_{\text{GCV}} = \begin{cases} s_{a^*}^2 & (\text{when } \hat{\sigma}_0^2 \neq 0) \\ t_1 & (\text{when } \hat{\sigma}_0^2 = 0 \text{ and } m = n - 1) \\ 0 & (\text{when } \hat{\sigma}_0^2 = 0 \text{ and } m < n - 1) \end{cases}$$

where  $s_a^2$  is given by (2.17) and  $a_*$  is the integer given by (2.18).

By using the results in Theorem 1, optimized GRREs of  $\boldsymbol{\beta}$  by minimizing the  $GC_p$  and GCV can be derived as  $\hat{\boldsymbol{\beta}}(\hat{h}_{GC_p})$  and  $\hat{\boldsymbol{\beta}}(\hat{h}_{\text{GCV}})$ , respectively.



### 3. Algorithm to Solve the Minimization Problem of the GIC

Theorem 1 in the previous section indicates that the  $GC_p$ -minimization method cannot be used when  $\hat{\sigma}_0^2$  in 2.8 is 0, i.e., most cases of the linear regression with high dimensional explanatory variables, the GCV-minimization method hardly shrink the GRRE of  $\beta$  towards  $0_k$  when  $\hat{\sigma}_0^2 = 0$ . On the other hand, the GIC is one of model selection criteria that can be defined even if  $\hat{\sigma}_0^2 = 0$  and of which the size of the penalty term can be changed freely. However, at the moment, there is no efficient algorithm to solve the minimization problem of the GIC. Hence, we propose the fast algorithm for the GIC-minimization method.

Let

$$\text{GIC}(\boldsymbol{\theta}) = f_{\text{GIC}}(\hat{\sigma}^2(\boldsymbol{\theta}), \text{df}(\boldsymbol{\theta})), \quad (3.1)$$

where  $f_{\text{GIC}}(r, u)$  is given by (2.11), and  $\hat{\sigma}^2(\boldsymbol{\theta})$  and  $\text{df}(\boldsymbol{\theta})$  are given by (2.7), and let

$$\phi_{\text{GIC}}(h) = f_{\text{GIC}}(\hat{\sigma}^2(\hat{\boldsymbol{\theta}}(h)), \text{df}(\hat{\boldsymbol{\theta}}(h))) \quad (h \in \mathbb{R}_+ \setminus \{0\}),$$

where  $\hat{\boldsymbol{\theta}}(h)$  is given by (2.12). From Lemma 2, we can see that  $\phi_{\text{GIC}}(h)$  is a piecewise function as

$$\phi_{\text{GIC}}(h) = \phi_{\text{GIC},a}(h) \quad (h \in R_a, a = 0, 1, \dots, m), \quad (3.2)$$

where  $R_a$  is the range given by (2.15) and  $\phi_{\text{GIC},a}(h)$  is defined by

$$\phi_{\text{GIC},a}(h) = \left\{ \hat{\sigma}_0^2 + \frac{1}{n}(c_{1,a} + h^2 c_{2,a}) \right\} \exp \left\{ \frac{1}{n} \alpha (1 + m - a - hc_{2,a}) \right\}.$$

Here  $c_{1,a}$  and  $c_{2,a}$  ( $a = 0, 1, \dots, m$ ) are given by (2.16). Furthermore, we define the function  $\psi_{\text{GIC}}(h)$  ( $h \in (0, t_m]$ ) which is a piecewise function as

$$\psi_{\text{GIC}}(h) = \psi_{\text{GIC},a}(h) \quad (h \in R_a, a = 0, 1, \dots, m-1), \quad (3.3)$$

where  $\psi_{\text{GIC},a}(h)$  is defined by

$$\begin{aligned} \psi_{\text{GIC},a}(h) &= \frac{n^2}{c_{2,a} \exp\{\alpha(1 + m - a - hc_{2,a})/n\}} \frac{\partial}{\partial h} \phi_{\text{GIC},a}(h) \\ &= -\alpha c_{2,a} h^2 + 2nh - \alpha(n\hat{\sigma}_0^2 + c_{1,a}). \end{aligned} \quad (3.4)$$

The minimizer of GIC can be expressed as in the following theorem (the proof is given in Appendix A.3):

**Theorem 2.** *Let  $\xi_{\text{GIC},a}$  be one of the two distinct roots or the double root of the quadratic equation  $\psi_{\text{GIC},a}(h) = 0$  as*

$$\xi_{\text{GIC},a} = \frac{n - \sqrt{n^2 - \alpha^2 c_{2,a}(n\hat{\sigma}_0^2 + c_{1,a})}}{\alpha c_{2,a}}, \quad (3.5)$$

and let  $\mathcal{A}_{\text{GIC}}$  and  $\mathcal{T}_{\text{GIC}}$  be sets defined by

$$\mathcal{A}_{\text{GIC}} = \{a \in \{0, 1, \dots, m-1\} \mid \xi_{\text{GIC},a} \in R_a\}, \quad \mathcal{T}_{\text{GIC}} = \begin{cases} \{t_m\} & (\text{when } \hat{\sigma}_\infty^2 > 2t_m/\alpha) \\ \{\emptyset\} & (\text{when } \hat{\sigma}_\infty^2 \leq 2t_m/\alpha) \end{cases}, \quad (3.6)$$

where  $\hat{\sigma}_\infty^2$  is given by (2.9). We define a set of candidate minimizers of  $\phi_{\text{GIC}}(h)$  as

$$\mathcal{S}_{\text{GIC}} = \left\{ \bigcup_{a \in \mathcal{A}_{\text{GIC}}} \{\xi_{\text{GIC},a}\} \right\} \bigcup \mathcal{T}_{\text{GIC}}. \quad (3.7)$$

Then, an optimized ridge parameters by minimizing  $\text{GIC}(\boldsymbol{\theta})$  are given by  $\hat{\boldsymbol{\theta}}(\hat{h}_{\text{GIC}})$ , where  $\hat{h}_{\text{GIC}}$  is given by

$$\hat{h}_{\text{GIC}} = \begin{cases} \arg \min_{h \in \mathcal{S}_{\text{GIC}}} \phi_{\text{GIC}}(h) & (\text{when } \hat{\sigma}_0^2 \neq 0) \\ 0 & (\text{when } \hat{\sigma}_0^2 = 0) \end{cases}. \quad (3.8)$$

It should be emphasized that elements of  $\mathcal{S}_{\text{GIC}}$  can be written by closed forms, and the number of elements of  $\mathcal{S}_{\text{GIC}}$  is equal to or smaller than  $m+1$ , i.e.,  $\#\mathcal{S}_{\text{GIC}} \leq m+1$ . Hence, when  $\hat{\sigma}_0^2 \neq 0$ , by using the results in Theorem 2, we can optimize  $\boldsymbol{\theta}$  quickly as follows:

### A Fast Algorithm to Derive the Optimized GRRE of $\beta$ by Minimizing the GIC

- (1) Calculate elements of  $\mathcal{S}_{\text{GIC}}$  by using  $\xi_{\text{GIC},a}$ ,  $\mathcal{A}_{\text{GIC}}$  and  $\mathcal{T}_{\text{GIC}}$ .
- (2) Determine  $\hat{h}_{\text{GIC}}$  by comparing function values of  $\phi_{\text{GIC}}(h)$  at the points in  $\mathcal{S}_{\text{GIC}}$ .
- (3) Calculate the optimized GRRE of  $\beta$  by minimizing the GIC as  $\hat{\beta}(\hat{h}_{\text{GIC}})$ , where  $\hat{\beta}(h)$  is given by (2.13).

Unfortunately, Theorem 2 also indicates that the optimized GRRE of  $\beta$  by minimizing GIC is always equal to the LSE of  $\beta$  when  $\hat{\sigma}_0^2 = 0$ . Thus, we cannot help but say that the GIC-minimization method does not work well in the case of high-dimensional explanatory variables.

When  $\hat{\sigma}_0^2 \neq 0$ , we cannot derive  $\#\mathcal{S}_{\text{GIC}}$  without comparing function values of  $\phi_{\text{GIC}}(h)$ . However,  $\#\mathcal{S}_{\text{GIC}}$  becomes 1 under the specific  $\alpha$  when  $\hat{\sigma}_0^2 \neq 0$ . The result is summarized as in the following corollary (the proof is given as Appendix A.4):

**Corollary 1.** *Suppose that  $\alpha \leq n/m$  and  $\hat{\sigma}_0^2 \neq 0$ . Let  $a^\dagger$  be a positive integer defined by*

$$a^\dagger \in \{0, 1, \dots, m-1\} \text{ s.t. } \xi_{\text{GIC},a^\dagger} \in R_{a^\dagger}.$$

The integer  $a^\dagger$  exists unique when  $\hat{\sigma}_\infty^2 \leq 2t_m/\alpha$ . Then  $\hat{h}_{\text{GIC}}$  in (3.8) can be rewritten as

$$\hat{h}_{\text{GIC}} = \begin{cases} \xi_{\text{GIC}, a^\dagger} & (\text{when } \hat{\sigma}_\infty^2 \leq 2t_m/\alpha) \\ t_m & (\text{when } \hat{\sigma}_\infty^2 > 2t_m/\alpha) \end{cases}.$$

#### 4. Algorithm to Solve the Minimization Problem of the EGCV

Theorem 2 in the previous section indicates that the GIC-minimization method does not shrink the GRRE of  $\beta$  towards  $\mathbf{0}_k$  at all when  $\hat{\sigma}_0^2$  in (2.8) is 0, i.e., most cases of the linear regression with high dimensional explanatory variables. From Theorems 1 and 2, we would have to say regrettably that there are no existing MSCs which work well when  $\hat{\sigma}_0^2 = 0$ . Hence, we have to propose new MSC which works well even when  $\hat{\sigma}_0^2 = 0$ . On the other hand, the cause of the problem in GCV is that the strength of the penalty to a complexity of a model is small. Hence, we extend the GCV so that the size of the penalty to a complexity of a model can be changed freely. It is called the extended GCV (EGCV) in this paper. Let

$$f_{\text{EGCV}}(r, u) = \frac{r}{(1 - u/n)^\alpha}, \quad \alpha > 2 \text{ and } u < n. \quad (4.1)$$

Then, the EGCV for the GRR is defined as

$$\text{EGCV}(\boldsymbol{\theta}) = f_{\text{EGCV}}(\hat{\sigma}^2(\boldsymbol{\theta}), \text{df}(\boldsymbol{\theta})), \quad (4.2)$$

where  $\hat{\sigma}^2(\boldsymbol{\theta})$  and  $\text{df}(\boldsymbol{\theta})$  are given by (2.7). It is easy to see that the EGCV coincides with the GCV if  $\alpha = 2$ .

Let

$$\phi_{\text{EGCV}}(h) = f_{\text{EGCV}}(\hat{\sigma}^2(\hat{\boldsymbol{\theta}}(h)), \text{df}(\hat{\boldsymbol{\theta}}(h))) \quad (h \in \mathbb{R}_+),$$

where  $\hat{\boldsymbol{\theta}}(h)$  is given by (2.12). From Lemma 2, we can see that  $\phi_{\text{EGCV}}(h)$  is a piecewise function as

$$\phi_{\text{EGCV}}(h) = \phi_{\text{EGCV}, a}(h) \quad (h \in R_a, a = 0, 1, \dots, m), \quad (4.3)$$

where  $R_a$  is the range given by (2.15) and  $\phi_{\text{EGCV}, a}(h)$  is defined by

$$\phi_{\text{EGCV}, a}(h) = \frac{\hat{\sigma}_0^2 + c_{1,a}/n + h^2 c_{2,a}/n}{(b + a/n + h c_{2,a}/n)^\alpha}. \quad (4.4)$$

Here  $c_{1,a}$  and  $c_{2,a}$  ( $a = 0, 1, \dots, m$ ) are given by (2.16), and  $b$  is the constant defined by

$$b = 1 - \frac{1}{n}(m + 1). \quad (4.5)$$

Furthermore, we define the function  $\psi_{\text{EGCV}}(h)$  ( $h \in (0, t_m]$ ) which is a piecewise function as

$$\psi_{\text{EGCV}}(h) = \psi_{\text{EGCV},a}(h) \quad (h \in R_a, a = 0, 1, \dots, m-1), \quad (4.6)$$

where  $\psi_{\text{EGCV},a}(h)$  is defined by

$$\begin{aligned} \psi_{\text{EGCV},a}(h) &= \frac{n^2(b + a/n + c_{2,a}h/n)^{\alpha+1}}{c_{2,a}} \frac{\partial}{\partial h} \phi_{\text{EGCV},a}(h) \\ &= -(\alpha - 2)c_{2,a}h^2 + 2(a + nb)h - \alpha(n\hat{\sigma}_0^2 + c_{1,a}). \end{aligned} \quad (4.7)$$

The minimizer of EGCV can be expressed as in the following theorem (the proof is given in Appendix A.5):

**Theorem 3.** *Let  $\xi_{\text{EGCV},a}$  be one of the two distinct roots or the double root of the quadratic equation  $\psi_{\text{EGCV},a} = 0$  as*

$$\xi_{\text{EGCV},a} = \frac{(a + nb) - \sqrt{(a + nb)^2 - \alpha(\alpha - 2)c_{2,a}(n\hat{\sigma}_0^2 + c_{1,a})}}{(\alpha - 2)c_{2,a}},$$

and let  $\mathcal{A}_{\text{EGCV}}$  and  $\mathcal{T}_{\text{EGCV}}$  be sets defined by

$$\mathcal{A}_{\text{EGCV}} = \{a \in \{0, 1, \dots, m-1\} \mid \xi_{\text{EGCV},a} \in R_a\}, \quad \mathcal{T}_{\text{EGCV}} = \begin{cases} \{t_m\} & (\text{when } \hat{\sigma}_\infty^2 > 2(1 - n^{-1})t_m/\alpha) \\ \{\emptyset\} & (\text{when } \hat{\sigma}_\infty^2 \leq 2(1 - n^{-1})t_m/\alpha) \end{cases},$$

where  $\hat{\sigma}_\infty^2$  is given by (2.9). We define a set of candidate minimizers of  $\phi_{\text{EGCV}}(h)$  as

$$\mathcal{S}_{\text{EGCV}} = \left\{ \bigcup_{a \in \mathcal{A}_{\text{EGCV}}} \{\xi_{\text{EGCV},a}\} \right\} \bigcup \mathcal{T}_{\text{EGCV}}. \quad (4.8)$$

Then, an optimized ridge parameters by minimizing EGCV( $\theta$ ) are given by  $\hat{\theta}(\hat{h}_{\text{EGCV}})$ , where  $\hat{h}_{\text{EGCV}}$  is given by

$$\hat{h}_{\text{EGCV}} = \begin{cases} \arg \min_{h \in \mathcal{S}_{\text{EGCV}}} \phi_{\text{EGCV}}(h) & (\text{when } \hat{\sigma}_0^2 \neq 0 \text{ or } b = 0) \\ 0 & (\text{when } \hat{\sigma}_0^2 = 0 \text{ and } b \neq 0) \end{cases}. \quad (4.9)$$

It should be emphasized that elements of  $\mathcal{S}_{\text{EGCV}}$  can be written by closed forms, and the number of elements of  $\mathcal{S}_{\text{EGCV}}$  is equal to or smaller than  $m+1$ , i.e.,  $\#\mathcal{S}_{\text{EGCV}} \leq m+1$ . Hence, when  $\hat{\sigma}_0^2 \neq 0$  or  $b = 0$ , by using the results in Theorem 3, we can optimize  $\beta$  quickly as follows:

### A Fast Algorithm to Derive the Optimized GRRE of $\beta$ by Minimizing the EGCV

- (1) Calculate elements of  $\mathcal{S}_{\text{EGCV}}$  by using  $\xi_{\text{EGCV},a}$ ,  $\mathcal{A}_{\text{EGCV}}$  and  $\mathcal{T}_{\text{EGCV}}$ .
- (2) Determine  $\hat{h}_{\text{EGCV}}$  by comparing function values of  $\phi_{\text{EGCV}}(h)$  at the points in  $\mathcal{S}_{\text{EGCV}}$ .

- (3) Calculate the optimized GRRE of  $\beta$  by minimizing the EGCV as  $\hat{\beta}(\hat{h}_{\text{EGCV}})$ , where  $\hat{\beta}(h)$  is given by (2.13).

When  $\hat{\sigma}_0^2 \neq 0$  or  $m = n - 1$ , we cannot derive  $\#(\mathcal{S}_{\text{EGCV}})$  without comparing function values of  $\phi_{\text{EGCV}}(h)$ . However,  $\#(\mathcal{S}_{\text{EGCV}})$  becomes 1 under the specific  $\alpha$  when  $\hat{\sigma}_0^2 \neq 0$  or  $m = n - 1$ . The result is summarized as in the following corollary (the proof is given as Appendix A.6):

**Corollary 2.** *Suppose that  $\alpha \leq (n + m - 1)/m$ , and  $\{\hat{\sigma}_0^2 \neq 0\} \vee \{b = 0\}$ . Let  $a^\ddagger$  be a positive integer defined by*

$$a^\ddagger \in \{0, 1, \dots, m - 1\} \text{ s.t. } \xi_{\text{EGCV}, a^\ddagger} \in R_{a^\ddagger}.$$

The integer  $a^\ddagger$  exists unique when  $\hat{\sigma}_\infty^2 \leq 2(1 - n^{-1})t_m/\alpha$ . Then  $\hat{h}_{\text{EGCV}}$  in (4.9) can be rewritten as

$$\hat{h}_{\text{EGCV}} = \begin{cases} \xi_{\text{EGCV}, a^\ddagger} & (\text{when } \hat{\sigma}_\infty^2 \leq 2(1 - n^{-1})t_m/\alpha) \\ t_m & (\text{when } \hat{\sigma}_\infty^2 > 2(1 - n^{-1})t_m/\alpha) \end{cases}.$$

## 5. Numerical Study

In this section, we compared performances of the  $GC_p$ -, GIC- and EGCV-minimization methods with  $\alpha = 2$ ,  $2 \log \log n$  and  $\log n$  by conducting numerical examinations. We generated data from the simulation model  $N_n(\mathbf{X}\beta, \mathbf{I}_n)$ , where  $\mathbf{X} = (\mathbf{I}_n - \mathbf{J}_n)\mathbf{X}_0\Phi(\rho)^{1/2}$  and  $\beta = \mathbf{M}^+\mathbf{X}'\eta$ . Here  $\mathbf{X}_0$  is an  $n \times k$  matrix of which elements are identically and independently distributed according to  $U(-1, 1)$ ,  $\Phi(\rho)$  is a  $k \times k$  symmetric matrix of which the  $(i, j)$ th element is  $\rho^{|i-j|}$  and  $\eta$  is an  $n$ -dimensional vector of which the  $j$ th element is given by

$$\sqrt{\frac{12n(n-1)}{4n^2 + 6n - 1}} \left\{ (-1)^{j-1} \left( 1 - \frac{j-1}{n} \right) - \frac{1}{2n} \right\}.$$

The simulation model is the same as that in Yanagihara (2013).

Let  $\hat{\theta}$  be an optimized ridge parameters by minimizing MSC,  $\hat{y}$  be the predictive value of  $y$  based on the LSE of  $\beta$  and  $\hat{y}_\theta$  be the predictive value of  $y$  in the GRR based on  $\theta$ , i.e.,

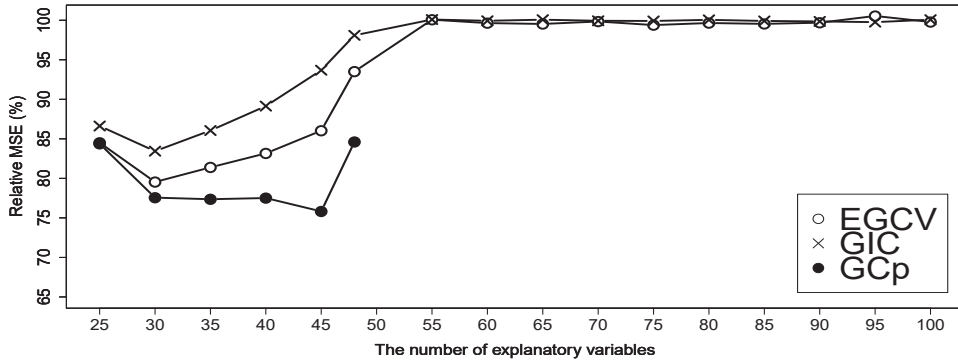
$$\hat{y} = (\mathbf{J}_n + \mathbf{X}\mathbf{M}^+\mathbf{X}')y, \quad \hat{y}_\theta = (\mathbf{J}_n + \mathbf{X}\mathbf{M}_\theta^+\mathbf{X}')y.$$

We used the following relative mean square error (MSE) as a measurement of performance of each method.

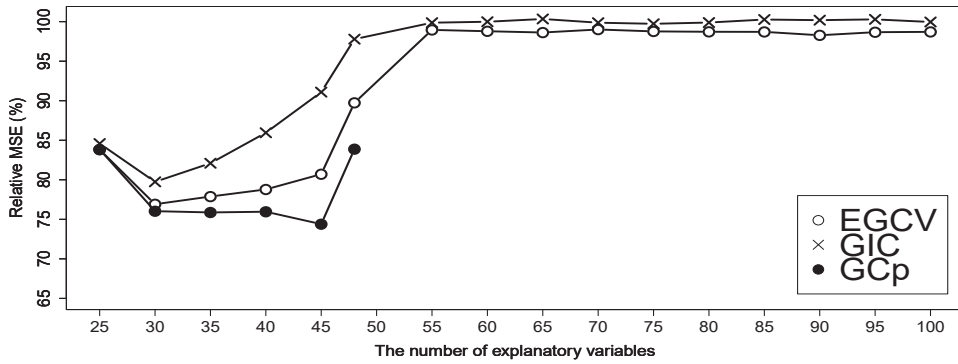
$$\frac{\text{E}[(\hat{y}_\theta - \mathbf{X}\beta)'(\hat{y}_\theta - \mathbf{X}\beta)]}{\text{E}[(\hat{y} - \mathbf{X}\beta)'(\hat{y} - \mathbf{X}\beta)]} \times 100(\%)$$

Notice that  $\text{E}[(\hat{y} - \mathbf{X}\beta)'(\hat{y} - \mathbf{X}\beta)] = m + 1$  in this case. The above expectation was evaluated by Monte Carlo simulation with 10,000 iterations.

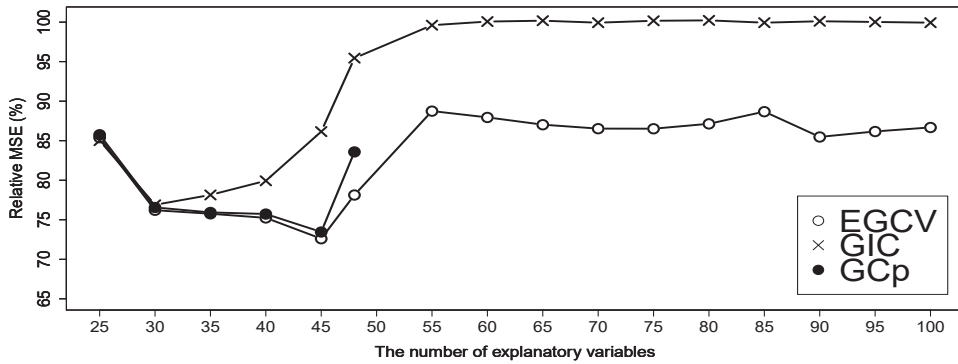
A Fast Algorithm for Minimizing EGCV in GRR



(a) Case of  $\alpha = 2$

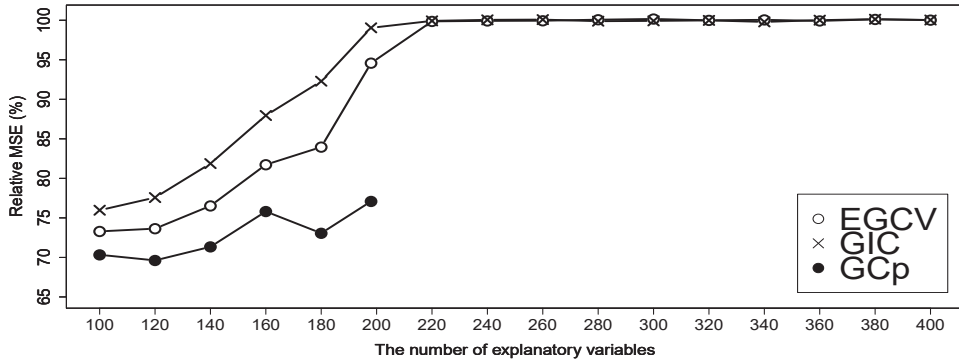


(b) Case of  $\alpha = 2 \log \log n$

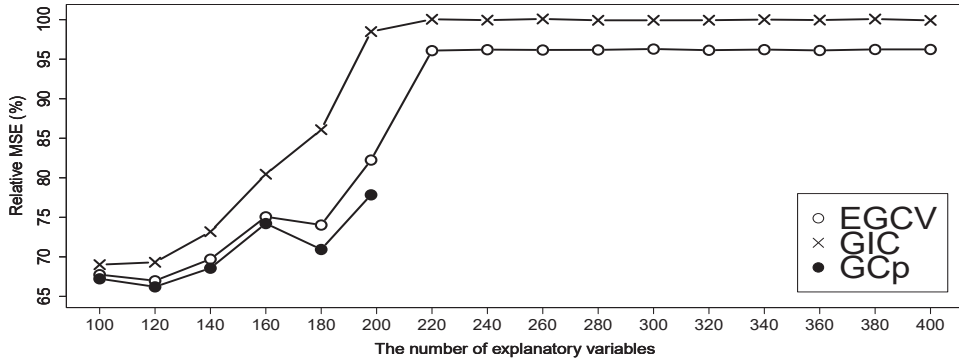


(c) Case of  $\alpha = \log n$

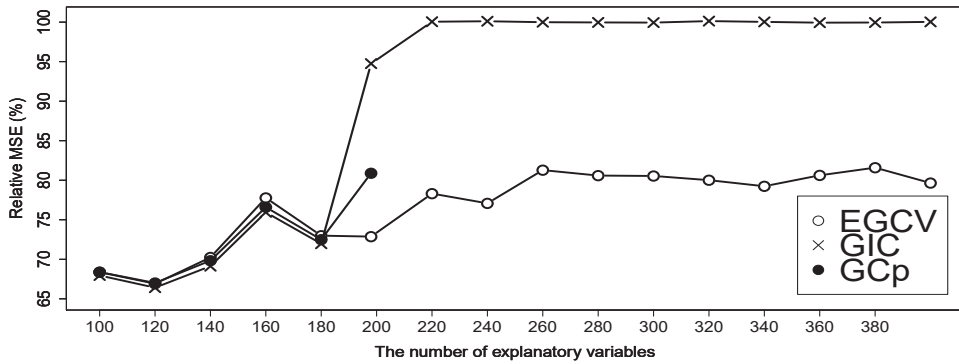
Figure 1. Case of  $n = 50$



(a) Case of  $\alpha = 2$



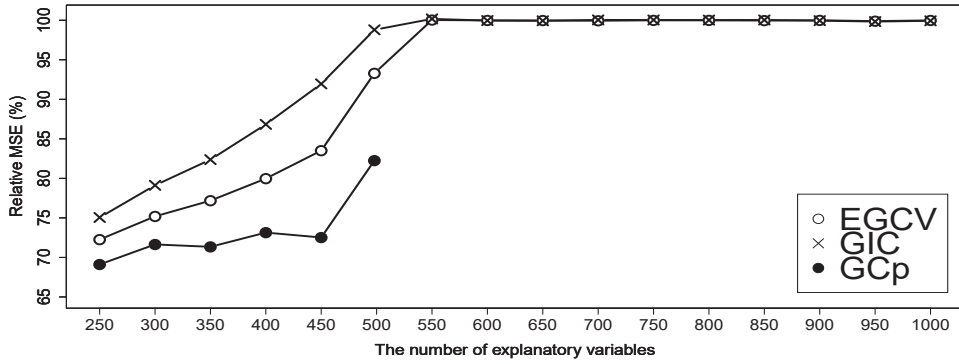
(b) Case of  $\alpha = 2 \log \log n$



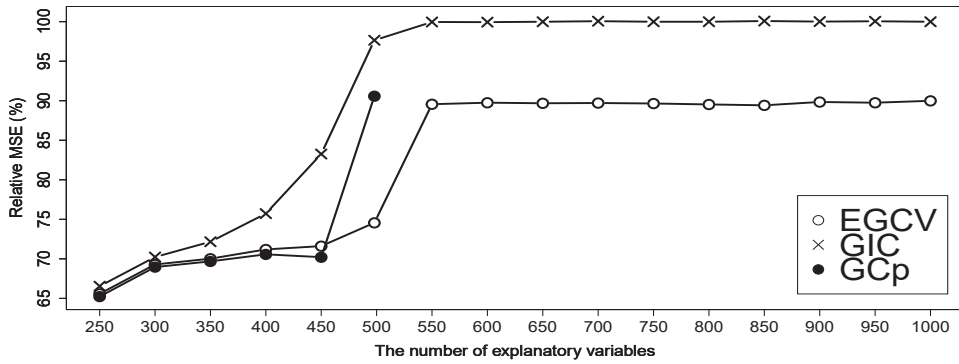
(c) Case of  $\alpha = \log n$

Figure 2. Case of  $n = 200$

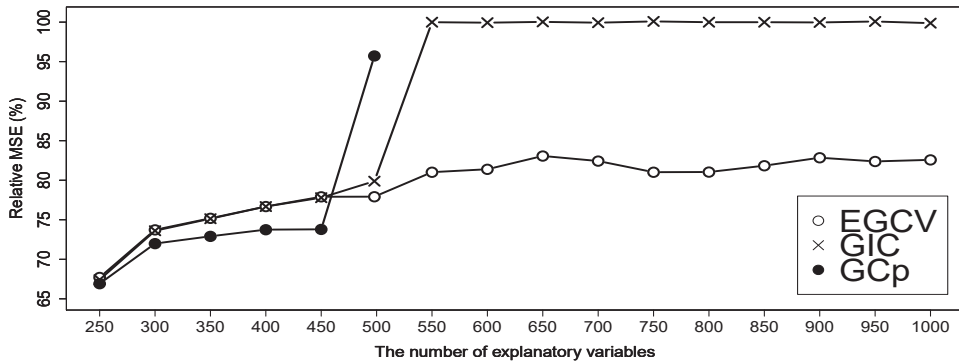
A Fast Algorithm for Minimizing EGCV in GRR



(a) Case of  $\alpha = 2$



(b) Case of  $\alpha = 2 \log \log n$



(c) Case of  $\alpha = \log n$

Figure 3. Case of  $n = 500$



Figures 1, 2 and 3 shows the results when  $n = 50, 200$  and  $500$ , respectively, with  $\rho = 0.99$  and  $k = k_1, \dots, k_{16}$ , where

$$k_j = \begin{cases} n(4 + j)/10 & (j = 1, \dots, 5) \\ n - 2 & (j = 6) \\ n\{1 + (j - 6)/10\} & (j = 7, \dots, 16) \end{cases} .$$

In the figures, there were no results of  $GC_p$ -minimization methods when  $k = k_7, \dots, k_{16}$  because the  $GC_p$  was not definable when  $k \geq k_7$ . From figures, we can see that the performances of the EGCV-minimization method with  $\alpha = \log n$  were well in most cases. On the other hand, the performances of the GIC-minimization method were very bad when  $k \geq k_7$  because the GIC-minimization method dose not shrink the GRRE of  $\beta$  towards  $\mathbf{0}_k$  at all when  $k \geq k_7$ . Also, the performances of the GCV-minimization method were also bad when  $k \geq k_7$  because the GCV-minimization method hardly shrink the GRRE of  $\beta$  towards  $\mathbf{0}_k$  when  $k \geq k_7$ . Moreover, as for the number of candidate minimizers, the following results were derived:

$$\text{mode } \{\#(\mathcal{S}_{\text{GIC}})\} = 2, \quad \min \{\#(\mathcal{S}_{\text{GIC}})\} = 2, \quad \max \{\#(\mathcal{S}_{\text{GIC}})\} = 7,$$

and

$$\text{mode } \{\#(\mathcal{S}_{\text{EGCV}})\} = \begin{cases} 3 & (k \geq n - 1 \ \& \ \alpha = 2 \log \log n) \\ 2 & (\text{otherwise}) \end{cases} ,$$

$$\min \{\#(\mathcal{S}_{\text{GIC}})\} = 2, \quad \max \{\#(\mathcal{S}_{\text{EGCV}})\} = \begin{cases} 3 & (k < n - 1) \\ 12 & (k \geq n - 1) \end{cases} .$$

Hence our algorithm will be efficient because the number of candidate minimizers was small.

**Acknowledgment** The authors thank Dr. Shintaro Hashimoto, Mr. Tomoyuki Nakagawa, Mr. Ryoya Oda and Mr. Shota Ochiai, of Hiroshima University, for helpful comments.

## References

- Akaike, H. (1973). Information theory and an extension of the maximum likelihood principle. In *2nd International Symposium on Information Theory* (eds. B. N. Petrov & F. Csáki), pp. 267–281. Akadémiai Kiadó, Budapest.
- Atkinson, A. C. (1980). A note on the generalized information criterion for choice of a model. *Biometrika*, **67**, 413–418.
- Craven, P. & Wahba, G. (1979). Smoothing noisy data with spline functions: estimating the correct degree of smoothing by the method of generalized cross-validation. *Numer. Math.*, **31**, 377–403.
- Fujikoshi, Y. and Satoh, K. (1997). Modified AIC and  $C_p$  in multivariate linear regression. *Biometrika*, **84**, 707–716.
- Hannan, E. J. & Quinn, B. G. (1979). The determination of the order of an autoregression. *J. Roy. Stat. Soc. Ser. B*, **41**, 190–195.

- Harville, D. A. (1997). *Matrix Algebra from a Statistician's Perspective*. Springer-Verlag, New York.
- Hoerl, A. E. & Kennard, R. W. (1970). Ridge regression: biased estimation for nonorthogonal problems. *Technometrics*, **12**, 55–67.
- Kullback, S. & Leibler, R. A. (1951). On information and sufficiency. *Ann. Math. Statistics*, **22**, 79–86.
- Lawless, J. F. (1981). Mean squared error properties of generalized ridge regression. *J. Amer. Statist. Assoc.*, **76**, 462–466.
- Liu, Z., Shen, Y. & Ott, J. (2011). Multilocus association mapping using generalized ridge logistic regression. *BMC Bioinformatics*, **12**, 384–391.
- Mallows, C. L. (1973). Some comments on  $C_p$ . *Technometrics*, **15**, 661–675.
- Nagai, I., Yanagihara, H. & Satoh, K. (2012). Optimization of ridge parameters in multivariate generalized ridge regression by plug-in methods. *Hiroshima Math. J.*, **42**, 301–324.
- Nishii, R. (1984). Asymptotic properties of criteria for selection of variables in multiple regression. *Ann. Statist.*, **12**, 758–765.
- Schwarz, G. (1978). Estimating the dimension of a model. *Ann. Statist.*, **6**, 461–464.
- Shen, X., Alam, M., Fikse, F. & Rönnegård, L. (2013). A novel generalized ridge regression method for quantitative genetics. *Genetics*, **193**, 1255–1268.
- Walker, S. G. & Page, C. J. (2001). Generalized ridge regression and a generalization of the  $C_p$  statistic. *J. Appl. Statist.*, **28**, 911–922.
- Yanagihara, H. (2013). Explicit solution to the minimization problem of generalized cross-validation criterion for selecting ridge parameters in generalized ridge regression. *TR-No. 13-07, Hiroshima Statistical Research Group, Hiroshima University*.
- Yanagihara, H., Nagai, I. & Satoh, K. (2009). A bias-corrected  $C_p$  criterion for optimizing ridge parameters in multivariate generalized ridge regression. *Japanese J. Appl. Statist.*, **38**, 151–172 (in Japanese).
- Ye, J. (1998). On measuring and correcting the effects of data mining and model selection. *J. Amer. Statist. Assoc.*, **93**, 120–131.

## Appendix

### A.1. Proof of Lemma 1

Let  $\boldsymbol{\delta} = (\delta_1, \dots, \delta_m)'$  be an  $m$ -dimensional vector of which the  $j$ th element  $\delta_j \in [0, 1]$  ( $j = 1, \dots, m$ ) is defined by

$$\delta_j = \frac{\theta_j}{d_j + \theta_j},$$

where  $\theta_1, \dots, \theta_m$  are ridge parameters, and  $d_j$  is given by (2.1). Since  $\boldsymbol{\theta} \rightarrow \boldsymbol{\delta}$  is one to one, we minimize MSC via  $\boldsymbol{\delta}$  instead of  $\boldsymbol{\theta}$ . By using  $\boldsymbol{\delta}$ , we rewrite  $\hat{\sigma}^2(\boldsymbol{\theta})$  and  $df(\boldsymbol{\theta})$  in (2.7) as a function with respect to  $\boldsymbol{\delta}$  as

$$\hat{\sigma}^2(\boldsymbol{\theta}) = r(\boldsymbol{\delta}) = \frac{1}{n} \left( n\hat{\sigma}_0^2 + \sum_{j=1}^m \delta_j^2 z_j^2 \right), \quad df(\boldsymbol{\theta}) = u(\boldsymbol{\delta}) = 1 + m - \sum_{j=1}^m \delta_j, \quad (\text{A.1})$$

where  $z_j^2$  and  $\hat{\sigma}_0^2$  are given by (2.6) and (2.8), respectively. By using  $r(\boldsymbol{\delta})$  and  $u(\boldsymbol{\delta})$ , we rewrite  $MSC(\boldsymbol{\theta})$  in (2.10) as a function with respect to  $\boldsymbol{\delta}$  as

$$MSC(\boldsymbol{\theta}) = g(\boldsymbol{\delta}) = f(r(\boldsymbol{\delta}), u(\boldsymbol{\delta})). \quad (\text{A.2})$$

Since  $\mathbb{D} \subseteq (0, \hat{\sigma}_\infty^2] \times [1, n)$ , we know that  $\delta \neq \mathbf{0}_m$ . From the assumption of  $f$ , we have  $\dot{f}_r(r(\delta), u(\delta)) > 0$  and  $\dot{f}_u(r(\delta), u(\delta)) > 0$ , where  $\dot{f}_r(r, u)$  and  $\dot{f}_u(r, u)$  are the first partial derivatives of  $f$  with respect to  $r$  and  $u$ , respectively. Let

$$\tau(\delta) = \frac{n\dot{f}_u(r(\delta), u(\delta))}{2\dot{f}_r(r(\delta), u(\delta))}.$$

It is clear that  $\tau(\delta) > 0$  when  $\mathbb{D} \subseteq (0, \hat{\sigma}_\infty^2] \times [1, n)$ . Notice that

$$\begin{aligned} \frac{\partial}{\partial \delta_j} g(\delta) &= \frac{\partial r(\delta)}{\partial \delta_j} \cdot \frac{\partial f(r, u)}{\partial r} \Big|_{(r, u) = (r(\delta), u(\delta))} + \frac{\partial u(\delta)}{\partial \delta_j} \cdot \frac{\partial f(r, u)}{\partial u} \Big|_{(r, u) = (r(\delta), u(\delta))} \\ &= \frac{1}{n} 2z_j^2 \delta_j \dot{f}_r(r(\delta), u(\delta)) - \dot{f}_u(r(\delta), u(\delta)) = \frac{2}{n} z_j^2 \dot{f}_r(r(\delta), u(\delta)) \left\{ \delta_j - \frac{\tau(\delta)}{z_j^2} \right\}. \end{aligned} \quad (\text{A.3})$$

Let  $\delta^* = (\delta_1^*, \dots, \delta_m^*)'$  be the minimizer of  $g(\delta)$ , i.e.,

$$\delta^* = \arg \min_{\delta \in [0, 1]^m \setminus \{\mathbf{0}_m\}} g(\delta),$$

where  $[0, 1]^m$  is the  $m$ th Cartesian power of the set  $[0, 1]$ . From (A.3), necessary condition of  $\delta^*$  is given by

$$\delta_j^* = \begin{cases} \frac{\tau(\delta^*)}{z_j^2} & (\tau(\delta^*) \leq z_j^2) \\ 1 & (\tau(\delta^*) > z_j^2) \end{cases} \quad (j = 1, \dots, m).$$

Let  $\mathcal{G}$  be a set defined by

$$\mathcal{G} = \left\{ \delta = (\delta_1, \dots, \delta_m)' \in [0, 1]^m \mid \delta = \hat{\delta}(h), \forall h \in \mathbb{R}_+ \setminus \{0\} \right\},$$

where  $\hat{\delta}(h)$  is the  $m$ -dimensional vector of which  $j$ th element is defined by

$$\hat{\delta}_j(h) = \begin{cases} \frac{h}{z_j^2} & (h \leq z_j^2) \\ 1 & (h > z_j^2) \end{cases} \quad (j = 1, \dots, m). \quad (\text{A.4})$$

It follows from the fact  $\mathcal{G} \subseteq [0, 1]^m \setminus \{\mathbf{0}_m\}$  that

$$g(\delta^*) = \min_{\delta \in [0, 1]^m \setminus \{\mathbf{0}_m\}} g(\delta) \leq \min_{\delta \in \mathcal{G}} g(\delta) = \min_{h \in \mathbb{R}_+ \setminus \{0\}} g(\hat{\delta}(h)).$$

Moreover, it is clear that  $\delta^* \in \mathcal{G}$ . This implies that

$$g(\delta^*) \geq \min_{h \in \mathbb{R}_+ \setminus \{0\}} g(\hat{\delta}(h)).$$

Therefore, we have  $\boldsymbol{\delta}^* = \hat{\boldsymbol{\delta}}(\hat{h})$ , where

$$\hat{h} = \arg \min_{h \in \mathbb{R}_+ \setminus \{0\}} g(\hat{\boldsymbol{\delta}}(h)).$$

Recall that  $g(\boldsymbol{\delta}) = \text{MSC}(\boldsymbol{\theta})$  and

$$\theta_j = \begin{cases} \frac{d_j \delta_j}{1 - \delta_j} & (\text{if } \delta_j \neq 1) \\ \infty & (\text{if } \delta_j = 1) \end{cases}.$$

Moreover,  $\hat{\boldsymbol{\beta}}_{\boldsymbol{\theta}}$  is rewritten as

$$\begin{aligned} \hat{\boldsymbol{\beta}}_{\boldsymbol{\theta}} &= (\mathbf{X}'\mathbf{X} + \mathbf{Q}\boldsymbol{\Theta}\mathbf{Q}')^+ \mathbf{X}'\mathbf{y} = \mathbf{Q}(\mathbf{D} + \boldsymbol{\Theta})^+ \mathbf{Q}'\mathbf{X}\mathbf{y} \\ &= \mathbf{Q}(\mathbf{D} + \boldsymbol{\Theta})^+ \mathbf{D}\mathbf{Q}'\mathbf{Q}\mathbf{D}^+ \mathbf{Q}'\mathbf{X}'\mathbf{y} = \mathbf{Q}\mathbf{V}\mathbf{Q}'\hat{\boldsymbol{\beta}}, \end{aligned}$$

where  $\mathbf{V} = (\mathbf{D} + \boldsymbol{\Theta})^+ \mathbf{D}$  and  $\hat{\boldsymbol{\beta}}$  is the LSE of  $\boldsymbol{\beta}$  given by (2.3). When  $\boldsymbol{\theta} \neq \mathbf{0}_m$ , the  $j$ th diagonal element of  $\mathbf{V}$  is given by

$$v_j = \begin{cases} \frac{d_j}{d_j + \theta_j} & (\text{when } \theta_j < \infty) \\ 0 & (\text{when } \theta_j = \infty) \end{cases}.$$

It follows from the simple calculation that  $d_j / \{d_j + \hat{\theta}_j(h)\} = 1 - h/z_j^2$  when  $h \leq z_j^2$ . Consequently, Lemma 1 is proved.

## A.2. Proof of Lemma 2

From (A.2), we can see that

$$\phi(h) = \text{MSC}(\hat{\boldsymbol{\theta}}(h)) = f(r(\hat{\boldsymbol{\delta}}(h)), u(\hat{\boldsymbol{\delta}}(h))),$$

where  $\hat{\boldsymbol{\theta}}(h)$  is given by (2.12),  $\hat{\boldsymbol{\delta}}(h)$  are given by (A.4), and  $r(\boldsymbol{\delta})$  and  $u(\boldsymbol{\delta})$  are given by (A.1), respectively. Since  $t_1, \dots, t_m$  are order statistics of  $z_1^2, \dots, z_m^2$ , we have

$$\begin{aligned} r(\hat{\boldsymbol{\delta}}(h)) &= \hat{\sigma}_0^2 + \frac{1}{n} \sum_{j=1}^m \left\{ \mathbb{I}(z_j^2 \geq h) \left( \frac{h}{z_j^2} - 1 \right) + 1 \right\}^2 z_j^2 = \hat{\sigma}_0^2 + \frac{1}{n} \sum_{j=1}^m \left\{ \mathbb{I}(t_j \geq h) \left( \frac{h}{t_j} - 1 \right) + 1 \right\}^2 t_j, \\ u(\hat{\boldsymbol{\delta}}(h)) &= 1 + m - \sum_{j=1}^m \left\{ \mathbb{I}(z_j^2 \geq h) \left( \frac{h}{z_j^2} - 1 \right) + 1 \right\} = 1 + m - \sum_{j=1}^m \left\{ \mathbb{I}(t_j \geq h) \left( \frac{h}{t_j} - 1 \right) + 1 \right\}, \end{aligned}$$

where  $\mathbb{I}(x \geq h)$  is an indicator function, i.e.,  $\mathbb{I}(x \geq h) = 1$  if  $x \geq h$  and  $\mathbb{I}(x \geq h) = 0$  if  $x < h$ . It is clear that  $r(\hat{\boldsymbol{\delta}}(h))$  and  $u(\hat{\boldsymbol{\delta}}(h))$  are continuous functions on  $h \in \mathbb{R}_+ \setminus \{0\}$ . Therefore, the property (P1) is proved.

Notice that when  $h \in R_a$ , we have

$$r(\hat{\delta}(h)) = r_a(\hat{\delta}(h)) = \hat{\sigma}_0^2 + \frac{1}{n} \left( \sum_{j=1}^a t_j + \sum_{j=a+1}^m \frac{h^2}{t_j} \right) = \hat{\sigma}_0^2 + \frac{1}{n} (c_{1,a} + h^2 c_{2,a}),$$

$$u(\hat{\delta}(h)) = u_a(\hat{\delta}(h)) = 1 + m - \left( \sum_{j=1}^a 1 + \sum_{j=a+1}^m \frac{1}{t_j} \right) = 1 + m - a - h c_{2,a},$$

where  $c_{1,a}$  and  $c_{2,a}$  are given by (2.16). The property (P3) is proved by this equation. Notice that  $c_{2,m} = 0$  and

$$n\hat{\sigma}_0^2 + c_{1,m} = n\hat{\sigma}_0^2 + \mathbf{z}'\mathbf{z} = \mathbf{y}'(\mathbf{I}_n - \mathbf{J}_n - \mathbf{X}\mathbf{M}^+\mathbf{X}')\mathbf{y} + \mathbf{y}'\mathbf{X}\mathbf{M}^+\mathbf{X}'\mathbf{y} = n\hat{\sigma}_\infty^2.$$

These imply that  $r_m(\hat{\delta}(h)) = \hat{\sigma}_\infty^2$  and  $u_m(\hat{\delta}(h)) = 1$  when  $h \in R_m$ . Hence, the property (P2) is proved.

### A.3. Proof of Theorem 2

At first, we derive the minimizer of  $\text{GIC}(\boldsymbol{\theta})$  in (3.1) when  $\hat{\sigma}_0^2$  in (2.8) is 0. It is easy to see that  $\hat{\sigma}^2(\mathbf{0}_m) = 0$  when  $\hat{\sigma}_0^2 = 0$ , and  $f_{\text{GIC}}(0, u) = 0$  holds for any  $u \in [1, n]$ , where  $\hat{\sigma}(\boldsymbol{\theta})$  is given by (2.4) and  $f_{\text{GIC}}(r, u)$  is given by (2.11). Moreover,  $f_{\text{GIC}}(r, u) > 0$  holds for any  $(r, u) \in (0, \hat{\sigma}_\infty^2) \times [1, n]$ , where  $\hat{\sigma}_\infty^2$  is given by (2.9). Notice that  $\hat{\sigma}^2(\boldsymbol{\theta}) = 0$  holds if and only if  $\boldsymbol{\theta} = \mathbf{0}_m$  and  $\hat{\sigma}_0^2 = 0$ , and  $\hat{\boldsymbol{\theta}}(0) = \mathbf{0}_m$  holds, where  $\hat{\boldsymbol{\theta}}(h)$  is given by (2.12). These results imply that the minimizer of  $\text{GIC}(\boldsymbol{\theta})$  is  $\hat{\boldsymbol{\theta}}(0)$  when  $\hat{\sigma}_0^2 = 0$ .

Next, we derive the minimizer of  $\text{GIC}(\boldsymbol{\theta})$  when  $\hat{\sigma}_0^2 \neq 0$ . Then the domain of  $f_{\text{GIC}}$  is included in  $(0, \hat{\sigma}_\infty^2) \times [1, n)$ . Moreover,

$$\lim_{h \rightarrow 0} \psi_{\text{GIC},0}(h) = -\alpha \hat{\sigma}_0^2 < 0,$$

where  $\psi_{\text{GIC},a}(h)$  is given by (3.4). This indicates that some positive number  $\xi$  exists such that  $\phi_{\text{GIC}}(\xi) < \lim_{h \rightarrow 0} \phi_{\text{GIC}}(h)$ , where  $\phi_{\text{GIC}}(h)$  is given by (3.2). Therefore, by using Lemma 1, the minimizer of  $\text{GIC}(\boldsymbol{\theta})$  is given by  $\hat{\boldsymbol{\theta}}(\hat{h}_{\text{GIC}})$ , where  $\hat{h}_{\text{GIC}}$  is the minimizer of  $\phi_{\text{GIC}}(h)$  as

$$\hat{h}_{\text{GIC}} = \arg \min_{h \in \mathbb{R}_+ \setminus \{0\}} \phi_{\text{GIC}}(h).$$

Hence, in order to prove Theorem 2 when  $\hat{\sigma}_0^2 \neq 0$ , it is enough just to show  $\hat{h}_{\text{GIC}} \in \mathcal{S}_{\text{GIC}}$ , where  $\mathcal{S}_{\text{GIC}}$  is a set of candidate minimizers of  $\phi_{\text{GIC}}(h)$  given by (3.7). From (P1) and (P2) in Lemma 2, we can see that  $\phi_{\text{GIC}}(h)$  is continuous at any  $h \in \mathbb{R}_+ \setminus \{0\}$  and  $\phi_{\text{GIC}}(h) = \hat{\sigma}_\infty^2 \exp(\alpha/n)$  when  $h \in R_m$ . The  $\psi_{\text{GIC}}(h)$  in (3.3) is also continuous at any  $h \in (0, t_m]$  because of

$$\begin{aligned} \psi_{\text{GIC},a}(t_{a+1}) &= -\alpha c_{2,a} t_{a+1}^2 + 2nt_{a+1} - \alpha(n\hat{\sigma}_0^2 + c_{1,a}) \\ &= -\alpha \left( c_{2,a+1} + \frac{1}{t_{a+1}} \right) t_{a+1}^2 + 2nt_{a+1} - \alpha(n\hat{\sigma}_0^2 + c_{1,a+1} - t_{a+1}) \\ &= -\alpha c_{2,a+1} t_{a+1}^2 + 2nt_{a+1} - \alpha(n\hat{\sigma}_0^2 + c_{1,a+1}) = \psi_{\text{GIC},a+1}(t_{a+1}) \quad (a = 0, \dots, m-2), \end{aligned}$$

where  $c_{1,a}$  and  $c_{2,a}$  are given by (2.16). It should be emphasized that the sign of  $\partial\psi_{\text{GIC},a}(h)/\partial h$  is equal to that of  $\psi_{\text{GIC},a}(h)$ . Hence, we have the necessary condition of  $\hat{h}_{\text{GIC}}$  as

$$\left\{ \left\{ \psi_{\text{GIC}}(\hat{h}_{\text{GIC}}) = 0 \right\} \bigwedge \left\{ \exists \epsilon_0 > 0 \text{ s.t. } \forall \epsilon \in (0, \epsilon_0), \psi_{\text{GIC}}(\hat{h}_{\text{GIC}} - \epsilon) < 0 \right\} \right\} \bigvee \left\{ \hat{h}_{\text{GIC}} = t_m \text{ s.t. } \psi_{\text{GIC},m-1}(t_m) < 0 \right\}. \quad (\text{A.5})$$

Notice that

$$\psi_{\text{GIC},m-1}(t_m) = -\frac{\alpha}{t_m}t_m^2 + 2nt_m - \alpha(n\hat{\sigma}_\infty^2 - t_m) = 2nt_m - \alpha n\hat{\sigma}_\infty^2.$$

Hence, we derive

$$\psi_{\text{GIC},m-1}(t_m) < 0 \iff \hat{\sigma}_\infty^2 < \frac{2}{\alpha}t_m. \quad (\text{A.6})$$

Recall that  $\psi_{\text{GIC},a}(h)$  is the concave quadratic function. Hence, we have

$$\begin{aligned} & \left\{ \psi_{\text{GIC}}(\hat{h}_{\text{GIC}}) = 0 \right\} \bigwedge \left\{ \exists \epsilon_0 > 0 \text{ s.t. } \forall \epsilon \in (0, \epsilon_0), \psi_{\text{GIC}}(\hat{h}_{\text{GIC}} - \epsilon) < 0 \right\} \\ \iff & \left\{ \hat{h}_{\text{GIC}} \text{ is the smaller of the two real distinct roots or the double} \right. \\ & \left. \text{root of } \psi_{\text{GIC},a}(h) = 0 \text{ which is included in } R_a \text{ (} a = 0, 1, \dots, m-1 \text{)} \right\}. \end{aligned} \quad (\text{A.7})$$

Notice that  $\xi_{\text{GIC},a}$  in (3.5) is the smaller of the two real distinct roots or the double root of  $\psi_{\text{GIC},a}(h) = 0$  when  $\alpha^2 c_{2,a}(n\hat{\sigma}_0^2 + c_{1,a}) \leq n^2$ . Therefore, from (A.5), (A.6) and (A.7),  $\hat{h}_{\text{GIC}} \in \mathcal{S}_{\text{GIC}}$  can be shown. Consequently, Theorem 2 is proved.

#### A.4. Proof of Corollary 1

It follows from (3.4) that

$$\psi_{\text{GIC},a}(h) = -\alpha c_{2,a} \left( h - \frac{n}{\alpha c_{2,a}} \right)^2 + \frac{n^2}{\alpha c_{2,a}} - \alpha(n\hat{\sigma}_0^2 + c_{1,a}),$$

where  $\hat{\sigma}_0^2$  is given by (2.8), and  $c_{1,a}$  and  $c_{2,a}$  are given by (2.16). This indicates that the  $h$ -coordinate of the vertex of  $\psi_{\text{GIC},a}(h)$  is  $n/(\alpha c_{2,a})$ . It follows from the equation  $t_1 \leq \dots \leq t_m$  that for any  $a \in \{0, 1, \dots, m-1\}$ ,

$$c_{a,2} = \sum_{j=a+1}^m \frac{1}{t_j} \leq \sum_{j=a+1}^m \frac{1}{t_{a+1}} = \frac{m-a}{t_{a+1}} \leq \frac{m}{t_{a+1}}. \quad (\text{A.8})$$

Hence, we have

$$\frac{n}{\alpha c_{2,a}} \geq \frac{n}{\alpha m} t_{a+1} \quad (a = 0, 1, \dots, m-1).$$

From this result, we can see that  $n/(\alpha c_{2,a}) \geq t_{a+1}$  if  $\alpha \leq n/m$ . This indicates that  $\psi_{\text{GIC},a}(h)$  is a strictly increasing function at  $h \in R_a$  in (2.15) if  $\alpha \leq n/m$  because the  $h$ -coordinate of the vertex of  $\psi_{\text{GIC},a}(h)$  is larger or equal to  $t_{a+1}$ , which is the right end point of  $R_a$ . Hence, it should

be emphasized that  $\psi_{\text{GIC}}(h)$  is a strictly increasing function at  $h \in (0, t_m]$  when  $\alpha \leq n/m$ . Recall that  $\psi_{\text{GIC}}(h)$  is continuous at any  $h \in (0, t_m]$  and  $\lim_{h \rightarrow 0} \psi_{\text{GIC},0}(h) < 0$  when  $\hat{\sigma}_0^2 \neq 0$ , where  $\hat{\sigma}_0^2$  is given by (2.8). Therefore, when  $\hat{\sigma}_0^2 \neq 0$  and  $\alpha \leq n/m$ , we have

$$\begin{aligned} \psi_{\text{GIC},m-1}(t_m) \geq 0 &\implies \text{the unique solution of } \psi_{\text{GIC}}(h) = 0 \text{ exists in } h \in (0, t_m] \\ &\iff \{\#\mathcal{A}_{\text{GIC}} = 1\} \wedge \{\mathcal{T}_{\text{GIC}} = \emptyset\}, \\ \psi_{\text{GIC},m-1}(t_m) < 0 &\implies \text{there is no solution of } \psi_{\text{GIC}}(h) = 0 \text{ in } h \in (0, t_m] \\ &\iff \{\mathcal{A}_{\text{GIC}} = \emptyset\} \wedge \{\mathcal{T}_{\text{GIC}} = \{t_m\}\}, \end{aligned}$$

where  $\mathcal{A}_{\text{GIC}}$  and  $\mathcal{T}_{\text{GIC}}$  are given by (3.6). From the above results and (A.6), Corollary 1 can be proved.

### A.5. Proof of Theorem 3

At first, we derive the minimizer of  $\text{EGCV}(\boldsymbol{\theta})$  in (4.2) when  $\hat{\sigma}_0^2 = 0$  and  $b \neq 0$ , where  $\hat{\sigma}_0^2$  and  $b$  are given by (2.8) and (4.5), respectively. It is easy to see that  $\hat{\sigma}^2(\mathbf{0}_m) = 0$  when  $\hat{\sigma}_0^2 = 0$ ,  $\text{df}(\boldsymbol{\theta}) < n$  when  $b \neq 0$ , and  $f_{\text{EGCV}}(0, u) = 0$  holds for any  $u \in [1, n)$ , where  $\hat{\sigma}^2(\boldsymbol{\theta})$  is given by (2.4) and  $f_{\text{EGCV}}(r, u)$  is given by (4.1). Moreover,  $f_{\text{EGCV}}(r, u) > 0$  holds for any  $(r, u) \in (0, \hat{\sigma}_\infty^2] \times [1, n)$ , where  $\hat{\sigma}_\infty^2$  is given by (2.9). Notice that  $\hat{\sigma}^2(\boldsymbol{\theta}) = 0$  holds if and only if  $\boldsymbol{\theta} = \mathbf{0}_m$  and  $\hat{\sigma}_0^2 = 0$ , and  $\hat{\boldsymbol{\theta}}(0) = \mathbf{0}_m$  holds, where  $\hat{\boldsymbol{\theta}}(h)$  is given by (2.12). These results imply that the minimizer of  $\text{EGCV}(\boldsymbol{\theta})$  is  $\hat{\boldsymbol{\theta}}(0)$  when  $\hat{\sigma}_0^2 = 0$  and  $b \neq 0$ .

Next, we consider when  $\hat{\sigma}_0^2 \neq 0$  or  $b = 0$ . Since  $r < n$  in  $f_{\text{EGCV}}(r, u)$ , it is necessary to satisfy the inequality  $\text{df}(\boldsymbol{\theta}) < n$ , where  $\text{df}(\boldsymbol{\theta})$  is given by (2.5). From (2.7), we can see that  $\text{df}(\boldsymbol{\theta}) < n \iff \boldsymbol{\theta} \neq \mathbf{0}_m \implies \hat{\sigma}^2(\boldsymbol{\theta}) \neq 0$ , where  $\hat{\sigma}^2(\boldsymbol{\theta})$  is given by (2.4). This indicates that the domain of  $f_{\text{EGCV}}$  is included in  $(0, \hat{\sigma}_\infty^2] \times [1, n)$ . Notice that

$$\phi_{\text{EGCV},0}(h) = \frac{\hat{\sigma}_0^2 + h^2 c_{2,0}/n}{(b + hc_{2,0}/n)^\alpha}, \quad \psi_{\text{EGCV},0}(h) = -(\alpha - 2)c_{2,0}h^2 + 2nbh - \alpha n \hat{\sigma}_0^2,$$

where  $\phi_{\text{EGCV},a}(h)$ ,  $\psi_{\text{EGCV},a}(h)$  and  $c_{2,a}$  are given by (4.4), (4.7) and (2.16), respectively. Recall that  $b = 0 \implies \hat{\sigma}_0^2 = 0$ . It follows from the above equations and the assumption  $\alpha > 2$  that

$$\lim_{h \rightarrow 0} \phi_{\text{EGCV},0}(h) = \begin{cases} \hat{\sigma}_0^2/b^\alpha & (\hat{\sigma}_0^2 \neq 0) \\ \infty & (b = 0) \end{cases}, \quad \lim_{h \rightarrow 0} \psi_{\text{EGCV},0}(h) = \begin{cases} -\alpha n \hat{\sigma}_0^2 < 0 & (\hat{\sigma}_0^2 \neq 0) \\ 0 & (b = 0) \end{cases}.$$

These results imply that some positive number  $\xi$  exists such that  $\phi_{\text{EGCV}}(\xi) < \lim_{h \rightarrow 0} \phi_{\text{EGCV}}(h)$  when  $\hat{\sigma}_0^2 \neq 0$  or  $b = 0$ , where  $\phi_{\text{EGCV}}(h)$  is given by (4.3). Therefore, by using Lemma 1, the minimizer of  $\text{EGCV}(\boldsymbol{\theta})$  is given by  $\hat{\boldsymbol{\theta}}(\hat{h}_{\text{EGCV}})$ , where  $\hat{h}_{\text{EGCV}}$  is the minimizer of  $\phi_{\text{EGCV}}(h)$  as

$$\hat{h}_{\text{EGCV}} = \arg \min_{h \in \mathbb{R}_+ \setminus \{0\}} \phi_{\text{EGCV}}(h).$$

From (P1) and (P2) in Lemma 2, we can see that  $\phi_{\text{EGCV}}(h)$  is continuous at any  $h \in \mathbb{R}_+ \setminus \{0\}$  and  $\phi_{\text{EGCV}}(h) = \hat{\sigma}_\infty^2 / (1 - n^{-1})^\alpha$  when  $h \in R_m$ , where  $\hat{\sigma}_\infty^2$  is given by (2.9). The  $\psi_{\text{EGCV}}(h)$  which is given by (4.6) is also continuous at any  $h \in (0, t_m]$  because of

$$\begin{aligned} \psi_{\text{EGCV},a}(t_{a+1}) &= -(\alpha - 2)c_{2,a}t_{a+1}^2 + 2(a + nb)t_{a+1} - \alpha(n\hat{\sigma}_0^2 + c_{1,a}) \\ &= -(\alpha - 2)\left(c_{2,a+1} + \frac{1}{t_{a+1}}\right)t_{a+1}^2 + 2(a + nb)t_{a+1} - \alpha\left(n\hat{\sigma}_0^2 + c_{1,a+1} - t_{a+1}\right) \\ &= -(\alpha - 2)c_{2,a+1}t_{a+1}^2 + 2(a + 1 + nb)t_{a+1} - \alpha(n\hat{\sigma}_0^2 + c_{1,a+1}) \\ &= \psi_{\text{EGCV},a+1}(t_{a+1}) \quad (a = 0, \dots, m - 2), \end{aligned}$$

where  $c_{1,a}$  is given by (2.16). Notice that

$$\psi_{\text{EGCV},m-1}(t_m) = -\frac{\alpha - 2}{t_m}t_m^2 + 2(n - 2)t_m - \alpha(n\hat{\sigma}_\infty^2 - t_m) = 2(n - 1)t_m - \alpha n\hat{\sigma}_\infty^2.$$

Hence, by the same way as in the proof of Theorem 2, we can show that a set of candidate minimizers of  $\phi_{\text{GIC}}(h)$  is given by  $\mathcal{S}_{\text{EGCV}}$  in (4.8). Consequently, Theorem 3 is proved.

## A.6. Proof of Corollary 2

It follows from (4.7) that

$$\psi_{\text{EGCV},a}(h) = -(\alpha - 2)c_{2,a} \left\{ h - \frac{a + nb}{(\alpha - 2)c_{2,a}} \right\}^2 + \frac{(a + nb)^2}{(\alpha - 2)c_{a,2}} - \alpha(n\hat{\sigma}_0^2 + c_{1,a}),$$

where  $\hat{\sigma}_0^2$  and  $b$  are given by (2.8) and (4.5), respectively, and  $c_{a,1}$  and  $c_{a,2}$  are given by (2.16). This indicates that the  $h$ -coordinate of the vertex of  $\psi_{\text{EGCV},a}(h)$  is  $(a + nb)/\{(\alpha - 2)c_{2,a}\}$ . From (A.8), we derive

$$\frac{a + nb}{(\alpha - 2)c_{2,a}} \geq \frac{a + nb}{(\alpha - 2)m}t_{a+1} \geq \frac{n - m - 1}{(\alpha - 2)m}t_{a+1} \quad (a = 0, 1, \dots, m - 1).$$

From this result, we can see that  $(a + nb)/\{(\alpha - 2)c_{2,a}\} \geq t_{a+1}$  if  $\alpha \leq (n + m - 1)/m$ . This indicates that  $\psi_{\text{EGCV},a}(h)$  is a strictly increasing function at  $h \in R_a$  in (2.15) if  $\alpha \leq (n + m - 1)/m$  because the  $h$ -coordinate of the vertex of  $\psi_{\text{EGCV},a}(h)$  is larger or equal to  $t_{a+1}$ , which is the right end point of  $R_a$ . It should be emphasized that  $\psi_{\text{GIC}}(h)$  is a strictly increasing function at  $h \in (0, t_m]$  when  $\alpha \leq (n + m - 1)/m$ . Hence, by the same way as in the proof of Corollary 1, we can prove Corollary 2.