# Modified $C_p$ Criterion for Optimizing Ridge and Smooth Parameters in the MGR Estimator for the Nonparametric GMANOVA Model

**Isamu Nagai**

*Department of Mathematics, Graduate School of Science, Hiroshima University*
*1-3-1 Kagamiyama, Higashi-Hiroshima, Hiroshima 739-8626, Japan*
*E-mail: d093481@hiroshima-u.ac.jp*

## Abstract

Longitudinal trends of observations can be estimated using the generalized multivariate analysis of variance (GMANOVA) model proposed by [10]. In the present paper, we consider estimating the trends nonparametrically using known basis functions. Then, as in nonparametric regression, an overfitting problem occurs. [13] showed that the GMANOVA model is equivalent to the varying coefficient model with non-longitudinal covariates. Hence, as in the case of the ordinary linear regression model, when the number of covariates becomes large, the estimator of the varying coefficient becomes unstable. In the present paper, we avoid the overfitting problem and the instability problem by applying the concept behind penalized smoothing spline regression and multivariate generalized ridge regression. In addition, we propose two criteria to optimize hyper parameters, namely, a smoothing parameter and ridge parameters. Finally, we compare the ordinary least square estimator and the new estimator.

## 1. Introduction

We consider the generalized multivariate analysis of variance (GMANOVA) model with $n$ observations of $p$-dimensional vectors of response variables. This model was proposed by [10]. Let $Y = (y_1, \cdots, y_n)'$, $A$, $X$ and $E = (\varepsilon_1, \cdots, \varepsilon_n)'$ be an $n \times p$ matrix of response variables, an $n \times k$ matrix of non-stochastic centerized between-individual explanatory variables (i.e., $A'\mathbf{1}_n = \mathbf{0}_k$) of $\mathrm{rank}(A) = k$ $(<n)$, a $p \times q$ matrix of non-stochastic within-individual explanatory variables of $\mathrm{rank}(X) = q$ $(q \le p)$, and an $n \times p$ matrix of error variables, respectively, where $n$ is the sample size, $\mathbf{1}_n$ is an $n$-dimensional vector of ones and $\mathbf{0}_k$ is a $k$-dimensional vector of zeros. Then, the GMANOVA model is expressed as

$$Y = \mathbf{1}_n \mu' X' + A\Xi X' + E$$

where $\Xi = (\xi_1, \cdots, \xi_k)'$ is a $k \times q$ unknown regression coefficient matrix an $\mu$ is $q$-dimensional unknown vector. We assume that $\varepsilon_1, \cdots, \varepsilon_n \sim_{\text{i.i.d.}} N_p(\mathbf{0}_p, \Sigma)$ where $\Sigma$

is a $p \times p$ unknown covariance matrix of $\mathrm{rank}(\Sigma) = p$. Then we can express the GMANOVA model as

$$Y \sim N_{n \times p}(\mathbf{1}_n \mu' X' + A\Xi X', \Sigma \otimes I_n)$$

Let $S$ be an unbiased estimator of the unknown covariance matrix $\Sigma$ that is given by

$$S = Y'\left\{I_n - n^{-1}\mathbf{1}_n\mathbf{1}_n' - A(A'A)A'\right\}Y \big/ (n - k - 1)$$

Then, the maximum likelihood (ML) estimators of $\mu$ and $\Xi$ are given by $n^{-1}(X'S^{-1}X)^{-1}X'S^{-1}Y'\mathbf{1}_n$ and $(A'A)^{-1}A'YS^{-1}X(X'S^{-1}X)^{-1}$, respectively. The ML estimators are the unbiased and asymptotically efficiency estimators of $\mu$ and $\Xi$.

In the GMANOVA model, $x(t) = (1, t, \cdots, t^{q-1})'$, $(t = t_1, \cdots, t_p)$ is often used as the $j$ th row vector of $X$. Then, we estimate the longitudinal trends of $Y$ using $(q-1)$-polynomial curves. However, occasionally, the

Using these criteria, we show that the optimized ridge parameters are obtained in closed form under the fixed $\lambda$. We also show the magnitude relationship between the optimized ridge parameters. In Section 4, we compare the LS estimator in (3) with the proposed estimator through numerical studies. In Section 5, we present our conclusions.

## 2. The New Estimators

In the model (1), we consider estimating the longitudinal trends nonparametrically by using basis functions $X$. Then, we consider the following estimators in order to avoid the overfitting problem in the nonparametric GMANOVA model, $\hat{\mu} = n^{-1}(X'X + \lambda K)^{-1} X'Y'\mathbf{1}_n$ and

$$\hat{\Xi}_\lambda = \left(\hat{\xi}_{\lambda,1}, \cdots, \hat{\xi}_{\lambda,k}\right)'$$
$$= \left(A'A\right)^{-1} A'YX\left(X'X + \lambda K\right)^{-1} \quad (5)$$

where $\lambda (\geq 0)$ is a smoothing parameter and $K$ is a $q \times q$ known penalty matrix. In this estimator, we must determine $K$ before using this estimator. Since $K$ is usually set as some nonnegative definite matrix, we assume that $K$ is a nonnegative definite matrix. If $\lambda K = O_{q,q}$, where $O_{q,q}$ is a $q \times q$ matrix of zeros, then this estimator corresponds to the LS estimators $\hat{\mu}$ and $\hat{\Xi}$ in (1). Note that this estimator controls the smoothness of each estimated curve $\hat{\phi}_0(t) = x(t)'\hat{\mu}_\lambda$ and $\hat{\phi}_j(t) = x(t)'\hat{\xi}_{\lambda,j}$, $(j = 1, \cdots, k)$ through only one parameter $\lambda$. When we use this estimator, we need to optimize the parameter $\lambda$ because this estimator changes with $\lambda$.

If multicollinearity occurs in $A$, then the LS estimator $\hat{\Xi}$ in (1) and the proposed estimator $\hat{\Xi}_\lambda$ in (5) are not good estimators in the sense of having large variance. Note that neither the LS estimator $\hat{\mu}$ nor the proposed estimator $\hat{\mu}_\lambda$ depend on $A$. Hence, we avoid the multicollinearity problem for estimating $\Xi$. Multicollinearity often occurs when $k$ becomes large. Using the following estimator, the multicollinearity problem in $A$ can be avoided,

$$\hat{\Xi} = \left(\hat{\xi}_{\theta\lambda,1}, \cdots, \hat{\xi}_{\theta\lambda,k}\right)'$$
$$= \left(A'A + \theta I_k\right)^{-1} A'YX(X'X + \lambda K)^{-1} \quad (6)$$

where $\theta \geq 0$ is a ridge parameter. This estimator with $K = I_q$ corresponds to the estimator of [16]. If $\theta = 0$, then this estimator corresponds to the estimator in (5). Note that $\left(A'A + \theta I_k\right)^{-1} A'Y$ in this estimator corresponds to the ridge estimator for a multivariate linear model [17]. In this estimator, we need to optimize $\theta$ and $\lambda$ because this estimator changes with these pa-

rameters. However, we cannot obtain the optimized $\theta$ and $\lambda$ in closed form. Thus, we need to use an iterative computational algorithm to optimize two parameters. From another point of view, this estimator controls the smoothness of each estimated curve $\hat{\phi}_j = x(t)'\hat{\xi}_{\theta\lambda,j}$, $(j = 1, \cdots, k)$ through only one parameter $\lambda$. Hence, this estimator is not a well fitting curve when the smoothnesses of the true curves differ.

Hence, we apply the concept of the MGR estimator [18] to $\left(A'A + \theta I_k\right)^{-1} A'Y$ in order to obtain the optimized ridge parameter in closed form. Here, we derive the MGR estimator for the nonparametric GMANOVA model as follows:

$$\hat{\Xi}_{\theta,\lambda} = \left(A'A + Q\Theta Q'\right)^{-1} A'YX\left(X'X + \lambda K\right)^{-1} \quad (7)$$

where $\theta = \left(\theta_1, \cdots, \theta_k\right)'$, $(\theta_i \geq 0, i = 1, \ldots, k)$ is also a ridge parameter, $\Theta = \mathrm{diag}(\theta)$, and $Q$ is the $k \times k$ orthogonal matrix that diagonalizes $A'A$, i.e., $Q'A'AQ = D$ where $D = \mathrm{diag}(d_1, \cdots, d_k)$ and $d_1, \cdots, d_k$ are eigenvalues of $A'A$. It is clearly that $d_i \geq 0$, $(i = 1, \cdots, k)$. In this estimator, since $\theta$ shrinks the estimators of $\phi_j(t)$, $(j = 1, \cdots, k)$ to 0, we can regard $\theta$ as controlling the smoothness of $\phi_1(t), \cdots, \phi_k(t)$. Therefore, in this estimator, rough smoothness of the estimated curves is controlled by $\lambda$, and each smoothness of $\phi_j(t)$, $(j = 1, \cdots, k)$ is controlled by $\theta$.

Clearly, $\hat{\Xi}_{0_k,0} = \hat{\Xi}$ and $\hat{\Xi}_{0_k,\lambda} = \hat{\Xi}_\lambda$. The $\hat{\Xi}_{\theta,\lambda}$ with $\theta = \mathbf{1}_k \theta$ for some $\theta (\geq 0)$ corresponds to $\hat{\Xi}_{\theta,\lambda}$ in (6). Thus, the estimator $\hat{\Xi}_{\theta,\lambda}$ includes these estimators. The estimator $\hat{\Xi}_{\theta,\lambda}$ is more flexible than these estimators $\hat{\Xi}_\lambda$ and $\hat{\Xi}_{\theta,\lambda}$ because $\hat{\Xi}_{\theta,\lambda}$ has $k + 1$ parameters and $\hat{\Xi}_\lambda$ or $\hat{\Xi}_{\theta,\lambda}$ has only one or two parameters. Hence, we consider $\hat{\mu}_\lambda$ and $\hat{\Xi}_{\theta,\lambda}$ in estimating the longitudinal trends or the varying coefficient curve, while avoiding the overfitting and multicollinearity problems in the nonparametric GMANOVA model. When $X = I_p$ and $\lambda K = O_{q,q}$, $\hat{\Xi}_{\theta,\lambda}$ corresponds to the MGR estimator in [18].

## 3. Main Results

### 3.1. Target MSE

In order to define the MSE of the predicted value of

$Y$, we prepare the following discrepancy function for measuring the distance between $n \times p$ matrices $F_1$ and

$$r(F_1, F_2) = \operatorname{tr}\left\{(F_1 - F_2)\Sigma^{-1}(F_1 - F_2)'\right\}:$$

Since $\Sigma$ is an unknown covariance matrix, we use the unbiased estimator $S$ in (2) instead of $\Sigma$ to estimate $r(F_1, F_2)$. Hence, we estimate $r(F_1, F_2)$ using the following sample discrepancy function:

$$\hat{r}(F_1, F_2) = \operatorname{tr}\left\{(F_1 - F_2)S^{-1}(F_1 - F_2)'\right\} \qquad (8)$$

These two functions, $r(F_1, F_2)$ and $\hat{r}(F_1, F_2)$, correspond to the summation of the Mahalanobis distance and the sample Mahalanobis distance between the rows of $F_1$ and $F_2$, respectively. Clearly, $r(F_1, F_2) = r(F_2, F_1)$ and $\hat{r}(F_1, F_2) = \hat{r}(F_2, F_1)$. Through simple calculation, we obtain the following properties:

$$r(F_1 + F_2, F_3) = r(F_1, F_3) + r(F_2, O_{n,p}) \\ + 2\operatorname{tr}\left\{(F_1 - F_3)\Sigma^{-1}F_2'\right\},$$

$$\hat{r}(F_1 + F_2, F_3) = \hat{r}(F_1, F_3) + \hat{r}(F_2, O_{n,p}) \\ + 2\operatorname{tr}\{(F_1 - F_3)S^{-1}F_2'\},$$

for any $n \times p$ matrices $F_1$, $F_2$ and $F_3$, Using the discrepancy function $r$, the MSE of the predicted value of $Y$ is defined as

$$\operatorname{MSE}\left[\hat{Y}_{\theta,\lambda}\right] = E\left[r\left(E[Y], \hat{Y}_{\theta,\lambda}\right)\right], \qquad (9)$$

where $\hat{Y}_{\theta,\lambda} = 1_n \hat{\mu}_\lambda' X' + A\hat{\Xi}_{\theta,\lambda} X'$, which is the predicted value of $Y$ when we use $\hat{\mu}_\lambda$ and $\hat{\Xi}_{\theta,\lambda}$ in (7). In the present paper, we regard $\theta$ and $\lambda$ making the MSE the smallest as the principle optimum. However, we cannot use the MSE in (9) in actual application because this MSE includes unknown parameters. Hence, we must estimate (9) in order to estimate the optimum $\theta$ and $\lambda$.

## 3.2. The $C_p$ and $MC_p$ Criteria

Let $H_\theta = A(A'A + Q\Theta Q')^{-1} A'$ and $G_\lambda = X(X'X + \lambda K)^{-1} X'$. Note that $Y = E[Y] + E$. Hence, we obtain

$$\operatorname{MSE}\left[\hat{Y}_{\theta,\lambda}\right] = E\left[r\left(Y - E, \hat{Y}_{\theta,\lambda}\right)\right].$$

From the properties of the function $r$ and using $E\left[\operatorname{tr}\left(Y\Sigma^{-1}E'\right)\right] = E\left[\operatorname{tr}\left(E\Sigma^{-1}E'\right)\right]$, since $E[Y]$ is a nonstochastic variable and $E[E] = O_{n,p}$, **and**

$E\left[\operatorname{tr}(F_4)\right] = \operatorname{tr}\left(E[F_4]\right)$ for any square matrix $F_4$, we obtain

$$\begin{aligned}\operatorname{MSE}\left[\hat{Y}_{\theta,\lambda}\right] &= E\left[r\left(Y, \hat{Y}_{\theta,\lambda}\right)\right] + E\left[r\left(E, O_{n,p}\right)\right] \\ &\quad - 2E\left[\operatorname{tr}\left\{\left(Y - \hat{Y}_{\theta,\lambda}\right)\Sigma^{-1}E'\right\}\right] \\ &= E\left[r\left(Y, \hat{Y}_{\theta,\lambda}\right)\right] + E\left[\operatorname{tr}\left(E\Sigma^{-1}E'\right)\right] \\ &\quad - 2E\left[\operatorname{tr}\left(Y\Sigma^{-1}E'\right)\right] \\ &\quad + 2E\left[\operatorname{tr}\left(\hat{Y}_{\theta,\lambda}\Sigma^{-1}E'\right)\right] \\ &= E\left[r\left(Y, \hat{Y}_{\theta,\lambda}\right)\right] - E\left[\operatorname{tr}\left(E\Sigma^{-1}E'\right)\right] \\ &\quad + 2E\left[\operatorname{tr}\left(\hat{Y}_{\theta,\lambda}\Sigma^{-1}E'\right)\right].\end{aligned}$$

Note that $\hat{Y}_{\theta,\lambda} = \left(n^{-1}1_n 1_n' + H_\theta\right)\left(E[Y] + E\right)G_\lambda$. Thus, we can calculate $E\left[\operatorname{tr}\left(\hat{Y}_{\theta,\lambda}\Sigma^{-1}E'\right)\right]$ as follows:

$$\begin{aligned}&E\left[\operatorname{tr}\left(\hat{Y}_{\theta,\lambda}\Sigma^{-1}E'\right)\right] \\ &= E\left[\operatorname{tr}\left\{\left(n^{-1}1_n 1_n' + H_\theta\right)\left(E[Y] + E\right)G_\lambda\Sigma^{-1}E'\right\}\right] \\ &= \operatorname{tr}\left\{\left(n^{-1}1_n 1_n' + H_\theta\right)E\left[EG_\lambda\Sigma^{-1}E'\right]\right\},\end{aligned}$$

because $E[Y]$, $G_\lambda$ and $\Sigma^{-1}$ are non-stochastic variables. For calculating the expectations in the MSE, we prove the following lemma.

**Lemma 3.1.** For any $p \times p$ non-stochastic matrix $J$, we obtain $E\left[EJ\Sigma^{-1}E'\right] = \operatorname{tr}(J)I_n$.

*proof.* Since $E = (\varepsilon_1, \cdots, \varepsilon_n)'$, we obtain the $(i, j)$th element of $E\left[EJ\Sigma^{-1}E'\right]$ as $E\left[\varepsilon_i'J\Sigma^{-1}\varepsilon_j\right]$, $(i = 1, \cdots, n; j = 1, \cdots, n)$. We obtain $E\left[\varepsilon_i \varepsilon_i'\right] = \delta_{i,j}\Sigma$ because $\varepsilon_i$   $\varepsilon_j$ for any $i \neq j$ and $\operatorname{Cov}(\varepsilon_i) = \Sigma$ for any $i$, where $\delta_{i,j}$ is defined as $\delta_{i,j} = 1$ if $i = j$ and $\delta_{i,j} = 0$ if $i \neq j$. Hence we obtain $E\left[\varepsilon_i'J\Sigma^{-1}\varepsilon_j\right] = \operatorname{tr}\left(\delta_{i,j}J\Sigma^{-1}\Sigma\right) = \delta_{i,j}\operatorname{tr}(J)$. This result means that $E\left[\varepsilon_i'J\Sigma^{-1}\varepsilon_j\right] = \operatorname{tr}(J)$ if $i = j$ and $E\left[\varepsilon_i'J\Sigma^{-1}\varepsilon_j\right] = 0$ if $i \neq j$. Thus, the lemma is proven.

Using this lemma, we obtain $E\left[\operatorname{tr}\left(E\Sigma^{-1}E'\right)\right] = np$ and $E\left[\operatorname{tr}\left(EG_\lambda\Sigma^{-1}E'\right)\right] = \operatorname{tr}(G_\lambda)I_n$. Hence, we obtain

$$\begin{aligned}\operatorname{MSE}\left[\hat{Y}_{\theta,\lambda}\right] &= E\left[r\left(Y, \hat{Y}_{\theta,\lambda}\right)\right] - np \\ &\quad + 2\operatorname{tr}\left\{\left(n^{-1}1_n 1_n' + H_\theta\right)\operatorname{tr}(G_\lambda)I_n\right\} \\ &= E\left[r\left(Y, \hat{Y}_{\theta,\lambda}\right)\right] - np \\ &\quad + 2\operatorname{tr}(G_\lambda)\operatorname{tr}\left(n^{-1}1_n 1_n' + H_\theta\right) \\ &= E\left[r\left(Y, \hat{Y}_{\theta,\lambda}\right)\right] - np \quad + 2\operatorname{tr}(G_\lambda)\left\{1 + \operatorname{tr}(H_\theta)\right\}\end{aligned}$$

By replacing $E\left[r\left(\boldsymbol{Y}, \hat{\boldsymbol{Y}}_{\theta,\lambda}\right)\right]$ with $\hat{r}\left(\boldsymbol{Y}, \hat{\boldsymbol{Y}}_{\theta,\lambda}\right)$, we can propose the instinctive estimator of MSE, referred to as the $C_p$ criterion, as follows:

$$C_p\left(\boldsymbol{\theta},\lambda\right)=\hat{r}\left(\boldsymbol{Y}, \hat{\boldsymbol{Y}}_{\theta,\lambda}\right)-np+2\mathrm{tr}\left(\boldsymbol{G}_\lambda\right)\left\{\mathrm{tr}\left(\boldsymbol{H}_\theta\right)+1\right\}. \quad (10)$$

When we use this criterion, we optimize the ridge parameter $\boldsymbol{\theta}$ and the smoothing parameter $\lambda$ by the following algorithm:

1) We obtain $\hat{\boldsymbol{\theta}}^{(\mathrm{C})}\left(\lambda\right)=\arg\min_\theta C_p\left(\boldsymbol{\theta},\lambda\right)$ , where $\hat{\boldsymbol{\theta}}^{(\mathrm{C})}\left(\lambda\right)=\left(\hat{\theta}_1^{(\mathrm{C})}\left(\lambda\right),\cdots,\hat{\theta}_k^{(\mathrm{C})}\left(\lambda\right)\right)'$ $\left(\hat{\theta}_i^{(\mathrm{C})}\left(\lambda\right)\geq 0, i=1,\cdots,k\right)$ if $\lambda$ is given.

2) We obtain $\hat{\lambda}^{(\mathrm{C})}=\arg\min_{\lambda\geq 0} C_p\left(\hat{\boldsymbol{\theta}}^{(\mathrm{C})}\left(\lambda\right),\lambda\right)$.

3) We obtain $\hat{\boldsymbol{\theta}}^{(\mathrm{C})}\left(\hat{\lambda}^{(\mathrm{C})}\right)=\arg\min_\theta C_p\left(\boldsymbol{\theta},\hat{\lambda}^{(\mathrm{C})}\right)$ , where $\hat{\boldsymbol{\theta}}^{(\mathrm{C})}\left(\hat{\lambda}^{(\mathrm{C})}\right)=\left(\hat{\theta}_1^{(\mathrm{C})}\left(\hat{\lambda}^{(\mathrm{C})}\right),\cdots,\hat{\theta}_k^{(\mathrm{C})}\left(\hat{\lambda}^{(\mathrm{C})}\right)\right)'$ , $\left(\hat{\theta}_i^{(\mathrm{C})}\left(\hat{\lambda}^{(\mathrm{C})}\right)\geq 0, i=1,\cdots,k\right)$ under fixed $\hat{\lambda}^{(\mathrm{C})}$.

4) We optimize the ridge parameter and the smoothing parameter as $\hat{\boldsymbol{\theta}}^{(\mathrm{C})}\left(\hat{\lambda}^{(\mathrm{C})}\right)$ and $\hat{\lambda}^{(\mathrm{C})}$, respectively.

Note that this $C_p$ criterion corresponds to that in [18] when $\boldsymbol{X}=\boldsymbol{I}_p$ and $\lambda\boldsymbol{K}=\boldsymbol{O}_{q,q}$.

There is some bias between the MSE in (9) and the $C_p$ criterion in (10) because the $C_p$ criterion is obtained by replacing $E\left[r\left(\boldsymbol{Y}, \hat{\boldsymbol{Y}}_{\theta,\lambda}\right)\right]$ in the MSE with $\hat{r}\left(\boldsymbol{Y}, \hat{\boldsymbol{Y}}_{\theta,\lambda}\right)$. Generally, when the sample size $n$ is small or the number of explanatory variables $k$ is large, this bias becomes large. Then, we cannot obtain the higher-accuracy estimation of the optimum parameters because we cannot obtain the higher-accuracy estimation of MSE of $\hat{\boldsymbol{Y}}_{\theta,\lambda}$ in (9). Hence, we correct the bias between $\mathrm{MSE}\left[\hat{\boldsymbol{Y}}_{\theta,\lambda}\right]$ and the $C_p$ criterion. To correct the bias, we assume $n-k-p-2>0$.

Let $\boldsymbol{W}=\left(n-k-1\right)\boldsymbol{S}$ and $\boldsymbol{W}_{\theta,\lambda}=\left(\boldsymbol{Y}-\hat{\boldsymbol{Y}}_{\theta,\lambda}\right)'\left(\boldsymbol{Y}-\hat{\boldsymbol{Y}}_{\theta,\lambda}\right)$.

$$E\left[r\left(\boldsymbol{Y}, \hat{\boldsymbol{Y}}_{\theta,\lambda}\right)\right]=\left(n-k-1\right)E[\mathrm{tr}\left\{\left(\boldsymbol{W}_{\theta,\lambda}-\boldsymbol{W}\right)\boldsymbol{W}^{-1}+\boldsymbol{I}_p\right\}.$$

Note that $\boldsymbol{W}\sim W_p\left(n-k-1,\boldsymbol{\Sigma}\right)$ and $\boldsymbol{W}_{\theta,\lambda}-\boldsymbol{W}$ $\boldsymbol{W}^{-1}$ because $\boldsymbol{A}'\boldsymbol{1}_n=\boldsymbol{0}_k$ and $\boldsymbol{A}'\left\{\boldsymbol{I}_n-\boldsymbol{A}(\boldsymbol{A}'\boldsymbol{A})^{-1}\boldsymbol{A}'\right\}=\boldsymbol{O}_{k,p}$ . Then, we obtain

Since $E\left[\boldsymbol{W}^{-1}\right]=\boldsymbol{\Sigma}^{-1}/\left(n-k-p-2\right)$ , $E\left[\boldsymbol{W}\right]=$

$\left(n-k-1\right)\boldsymbol{\Sigma}$ (see, e.g., [14]) and $\mathrm{tr}\left\{E\left[\boldsymbol{W}_{\theta,\lambda}\boldsymbol{\Sigma}^{-1}\right]\right\}=E\left[r\left(\boldsymbol{Y}, \hat{\boldsymbol{Y}}_{\theta,\lambda}\right)\right]$, we obtain

$$E\left[\hat{r}\left(\boldsymbol{Y}, \hat{\boldsymbol{Y}}_{\theta,\lambda}\right)\right]$$
$$=\frac{n-k-1}{n-k-p-2}\mathrm{tr}\left\{E\left[\boldsymbol{W}_{\theta,\lambda}-\boldsymbol{W}\right]\boldsymbol{\Sigma}^{-1}+\left(n-k-p-2\right)\boldsymbol{I}_p\right\}$$
$$=\frac{n-k-1}{n-k-p-2}\left\{E\left[r\left(\boldsymbol{Y}, \hat{\boldsymbol{Y}}_{\theta,\lambda}\right)\right]-p\left(p+1\right)\right\}$$

Therefore, we obtain the unbiased estimator for $E\left[r\left(\boldsymbol{Y}, \hat{\boldsymbol{Y}}_{\theta,\lambda}\right)\right]$ as $c_{\mathrm{M}}\hat{r}\left(\boldsymbol{Y}, \hat{\boldsymbol{Y}}_{\theta,\lambda}\right)+p\left(p+1\right)$ , where $c_{\mathrm{M}}=1-\left(p+1\right)/\left(n-k-1\right)$. This implies that the bias corrected $C_p$ criterion, denoted as $MC_p$ (modified $C_p$) criterion, is obtained by

$$\begin{aligned}MC_p\left(\boldsymbol{\theta},\lambda\right)&=c_{\mathrm{M}}\hat{r}\left(\boldsymbol{Y}, \hat{\boldsymbol{Y}}_{\theta,\lambda}\right)+p\left(p+1-n\right)\\&+2\mathrm{tr}\left(\boldsymbol{G}_\lambda\right)\left\{\mathrm{tr}\left(\boldsymbol{H}_\theta\right)+1\right\}.\end{aligned} \quad (11)$$

As in the case of using the $C_p$, we optimize $\boldsymbol{\theta}$ and $\lambda$ using this criterion as follows:

1) We obtain $\hat{\boldsymbol{\theta}}^{(\mathrm{M})}\left(\lambda\right)=\arg\min_\theta MC_p\left(\boldsymbol{\theta},\lambda\right)$, where $\hat{\boldsymbol{\theta}}^{(\mathrm{M})}\left(\lambda\right)=\left(\hat{\theta}_1^{(\mathrm{M})}\left(\lambda\right),\cdots,\hat{\theta}_k^{(\mathrm{M})}\left(\lambda\right)\right)'$, $\left(\hat{\theta}_i^{(\mathrm{M})}\left(\lambda\right)\geq 0, i=1,\cdots,k\right)$ if $\lambda$ is given.

2) We obtain $\hat{\lambda}^{(\mathrm{M})}=\arg\min_{\lambda\geq 0} MC_p\left(\hat{\boldsymbol{\theta}}^{(\mathrm{M})}\left(\lambda\right),\lambda\right)$.

3) We obtain $\hat{\boldsymbol{\theta}}^{(\mathrm{M})}\left(\hat{\lambda}^{(\mathrm{M})}\right)=\arg\min_\theta MC_p\left(\boldsymbol{\theta},\hat{\lambda}^{(\mathrm{M})}\right)$ , where $\hat{\boldsymbol{\theta}}^{(\mathrm{M})}\left(\hat{\lambda}^{(\mathrm{M})}\right)=\left(\hat{\theta}_1^{(\mathrm{M})}\left(\hat{\lambda}^{(\mathrm{M})}\right),\cdots,\hat{\theta}_k^{(\mathrm{M})}\left(\hat{\lambda}^{(\mathrm{M})}\right)\right)'$ , $\left(\hat{\theta}_i^{(\mathrm{M})}\left(\hat{\lambda}^{(\mathrm{M})}\right)\geq 0, i=1,\cdots,k\right)$ under fixed $\hat{\lambda}^{(\mathrm{M})}$.

4) We optimize the ridge parameter and the smoothing parameter as $\hat{\boldsymbol{\theta}}^{(\mathrm{M})}\left(\hat{\lambda}^{(\mathrm{M})}\right)$ and $\hat{\lambda}^{(\mathrm{M})}$, respectively.

Note that the $MC_p$ criterion corresponds to that in [18] when $\boldsymbol{X}=\boldsymbol{I}_p$ and $\lambda\boldsymbol{K}=\boldsymbol{O}_{q,q}$. The $MC_p$ criterion completely omits the bias between the MSE of $\hat{\boldsymbol{Y}}_{\theta,\lambda}$ in (9) and the $C_p$ criterion in (10) by using a number of constant terms $c_{\mathrm{M}}$ and $p\left(p+1\right)$ . If $\hat{\boldsymbol{\theta}}^{(\mathrm{C})}\left(\lambda\right)$ and $\hat{\boldsymbol{\theta}}^{(\mathrm{M})}\left(\lambda\right)$ can be expressed in closed form for any fixed $\lambda\geq 0$, we do not need the above iterative computational algorithm.

## 3.3. Optimizations using the $C_p$ and $MC_p$ Criteria

Using the generalized $C_p$ $\left(GC_p\right)$ criterion, which is

given in (14), we can express the $C_p$ and $MC_p$ criteria as follows:

$$C_p(\boldsymbol{\theta}, \lambda) = GC_p(\boldsymbol{\theta}, \lambda \mid 1) + 2\mathrm{tr}(\boldsymbol{G}_\lambda) - np,$$

$$MC_p(\boldsymbol{\theta}, \lambda) = GC_p(\boldsymbol{\theta}, \lambda \mid c_\mathrm{M}) + 2\mathrm{tr}(\boldsymbol{G}_\lambda) + p(p+1-n)$$

Note that the terms with respect to $\boldsymbol{\theta}$ in the $C_p$ and $MC_p$ criteria correspond to $GC_p(\boldsymbol{\theta}, \lambda \mid 1)$ and $GC_p(\boldsymbol{\theta}, \lambda \mid c_\mathrm{M})$, respectively. Hence, we consider obtaining the optimum $\boldsymbol{\theta}$ by minimizing the $GC_p$ criterion. From Theorem A, the optimum $\boldsymbol{\theta}$ is obtained in closed form as (15). Using the closed form in (15), we obtain $\hat{\theta}_i^{(\mathrm{C})}(\lambda)$ and $\hat{\theta}_i^{(\mathrm{M})}(\lambda)$ for each $i = 1, \cdots, k$ and

any fixed $\lambda \geq 0$ as follows:

$$\hat{\theta}_i^{(\mathrm{C})}(\lambda) = \hat{\theta}_i^{(\mathrm{G})}(\lambda \mid 1)$$

$$= \begin{cases} 0 & \left(0 \leq t_i^{(\mathrm{C})}(\lambda)\right) \\ \dfrac{-d_i t_i^{(\mathrm{C})}(\lambda)}{t_i^{(\mathrm{C})}(\lambda) + u_{ii}} & \left(t_i^{(\mathrm{C})}(\lambda) < 0 < t_i^{(\mathrm{C})}(\lambda) + u_{ii}\right), \\ \infty & (\text{otherwise}) \end{cases} \quad (12)$$

$$\hat{\theta}_i^{(\mathrm{M})}(\lambda) = \hat{\theta}_i^{(\mathrm{G})}(\lambda \mid c_\mathrm{M})$$

$$= \begin{cases} 0 & \left(0 \leq t_i^{(\mathrm{M})}(\lambda)\right) \\ \dfrac{-d_i t_i^{(\mathrm{M})}(\lambda)}{t_i^{(\mathrm{M})}(\lambda) + c_\mathrm{M} u_{ii}} & \left(t_i^{(\mathrm{M})}(\lambda) < 0 < t_i^{(\mathrm{M})}(\lambda) + c_\mathrm{M} u_{ii}\right), \\ \infty & (\text{otherwise}) \end{cases}$$

$$(13)$$

where $u_{ii}$ and $v_{ii}$ are the $(i, i)$ th elements of $\boldsymbol{Q'A'Y G_\lambda S^{-1} G_\lambda Y'AQ}$ and $\boldsymbol{Q'A'Y S^{-1} G_\lambda Y'AQ}$, respectively, $t_i^{(\mathrm{C})}(\lambda) = t_i(\lambda \mid 1) = v_{ii} - u_{ii} - d_i \mathrm{tr}(\boldsymbol{G}_\lambda)$ and $t_i^{(\mathrm{M})}(\lambda) = t_i(\lambda \mid c_\mathrm{M}) = c_\mathrm{M}(v_{ii} - u_{ii}) - d_i \mathrm{tr}(\boldsymbol{G}_\lambda)$. Note that $u_{ii}$ and $v_{ii}$ vary with $\lambda$. Since $\hat{\boldsymbol{\theta}}^{(\mathrm{C})}(\lambda)$ and $\hat{\boldsymbol{\theta}}^{(\mathrm{M})}(\lambda)$ are regarded as a function of $\lambda$, we can regard the $C_p$ and $MC_p$ criteria for optimizing $\boldsymbol{\theta}$ and $\lambda$ in (10) and (11) as a function of $\lambda$. This means that we can use these criteria to optimize $\lambda$.

Then, we can rewrite the optimization algorithms to optimize the ridge parameter $\boldsymbol{\theta}$ and the smoothing parameter $\lambda$ by minimizing the $C_p$ and $MC_p$ criteria in (10) and (11) as follows:

1) We obtain $\hat{\lambda}^{(\mathrm{C})} = \arg\min_{\lambda \geq 0} C_p\left(\hat{\boldsymbol{\theta}}^{(\mathrm{C})}(\lambda), \lambda\right)$ and $\hat{\lambda}^{(\mathrm{M})} = \arg\min_{\lambda \geq 0} MC_p\left(\hat{\boldsymbol{\theta}}^{(\mathrm{M})}(\lambda), \lambda\right)$.

2) We optimize the ridge parameter and the smoothing parameter as $\hat{\boldsymbol{\theta}}^{(\mathrm{C})}\left(\hat{\lambda}^{(\mathrm{C})}\right)$ and $\hat{\boldsymbol{\theta}}^{(\mathrm{M})}\left(\hat{\lambda}^{(\mathrm{M})}\right)$, respectively, by using $\hat{\lambda}^{(\mathrm{C})}$, $\hat{\lambda}^{(\mathrm{M})}$ and the closed forms in (12) and (13).

This means that we can reduce the processing time to optimize the parameters, and we need to use the optimization algorithm for only one parameter, $\lambda$, for any $k$.

### 3.4. Magnitude Relationships between Optimized Ridge Parameters

In this subsection, we prove the magnitude relationships between $\hat{\theta}_i^{(\mathrm{C})}\left(\hat{\lambda}^{(\mathrm{C})}\right)$ and $\hat{\theta}_i^{(\mathrm{M})}\left(\hat{\lambda}^{(\mathrm{M})}\right)$, $(i = 1, \cdots, k)$.

**Lemma 3.2.** For any $\lambda \geq 0$, we obtain $\mathrm{tr}(\boldsymbol{G}_\lambda) \geq 0$.

*proof.* Since we assume $\boldsymbol{K}$ as a nonnegative definite matrix, there exists $\boldsymbol{L}$ that satisfies $\boldsymbol{K} = \boldsymbol{L'L}$ (see, e.g., [3]). Then, since $\lambda \geq 0$, we have $\boldsymbol{X'X} + \lambda \boldsymbol{K} = \left(\boldsymbol{X'}, \lambda^{1/2} \boldsymbol{L'}\right)\left(\boldsymbol{X'}, \lambda^{1/2} \boldsymbol{L'}\right)'$. Hence, $\boldsymbol{X'X} + \lambda \boldsymbol{K}$ is a nonnegative definite matrix. This means that all of the eigenvalues of $\boldsymbol{X'X} + \lambda \boldsymbol{K}$ are nonnegative. Hence, all of the eigenvalues of $\left(\boldsymbol{X'X} + \lambda \boldsymbol{K}\right)^{-1}$ are nonnegative. Thus, $\left(\boldsymbol{X'X} + \lambda \boldsymbol{K}\right)^{-1}$ is also a nonnegative definite matrix for any $\lambda \geq 0$. Since $\boldsymbol{G}_\lambda = \boldsymbol{X}\left(\boldsymbol{X'X} + \lambda \boldsymbol{K}\right)^{-1} \boldsymbol{X'}$, we obtain $\boldsymbol{G}_\lambda$ as a nonnegative definite matrix for any $\lambda \geq 0$. Thus, the lemma is proven.

Using the same idea, we have $\mathrm{tr}(\boldsymbol{H}_\theta) \geq 0$ for any $\boldsymbol{\theta}$, $(\theta_i \geq 0, i = 1, ..., k)$. Therefore, the final terms of the $C_p$ and $MC_p$ criteria in (10) and (11) are always greater than $\mathrm{tr}(\boldsymbol{G}_\lambda) \geq 0$. In order to prove the magnitude relationship between $\hat{\theta}_i^{(\mathrm{C})}\left(\hat{\lambda}^{(\mathrm{C})}\right)$ and $\hat{\theta}_i^{(\mathrm{M})}\left(\hat{\lambda}^{(\mathrm{M})}\right)$, we consider two situations in which $\hat{\lambda}^{(\mathrm{C})} = \hat{\lambda}^{(\mathrm{M})}$ is satisfied and $\hat{\lambda}^{(\mathrm{C})} \neq \hat{\lambda}^{(\mathrm{M})}$ is satisfied.

First, we consider $\hat{\lambda}^{(\mathrm{C})} = \hat{\lambda}^{(\mathrm{M})}$ to be satisfied. Let $\hat{\lambda} = \hat{\lambda}^{(\mathrm{C})} = \hat{\lambda}^{(\mathrm{M})}$ $\left(\hat{\lambda} \geq 0\right)$. Using $\hat{\lambda}$, we obtain the following corollary:

**Corollary 3.1.** For any $\hat{\lambda} \geq 0$, we obtain $c_\mathrm{M} t_i^{(\mathrm{C})}\left(\hat{\lambda}\right) \geq t_i^{(\mathrm{M})}\left(\hat{\lambda}\right)$.

*proof.* Through simple calculation, we obtain

$$c_\mathrm{M} t_i^{(\mathrm{C})}\left(\hat{\lambda}\right) - t_i^{(\mathrm{M})}\left(\hat{\lambda}\right) = d_i \mathrm{tr}\left(\boldsymbol{G}_{\hat{\lambda}}\right)(1 - c_\mathrm{M}).$$

Since $d_i > 0$, $0 < c_\mathrm{M} < 1$ and $\mathrm{tr}\left(\boldsymbol{G}_{\hat{\lambda}}\right) \geq 0$ from

lemma 3.1, the corollary is proven.

This corollary indicates that $t_i^{(C)}(\hat{\lambda}) \geq 0$ is satisfied when $t_i^{(M)}(\hat{\lambda}) \geq 0$ is satisfied because $c_M > 0$, **and** $t_i^{(C)}(\hat{\lambda}) + u_{ii} > 0$ is satisfied when $t_i^{(M)}(\hat{\lambda}) + c_M u_{ii} > 0$ is satisfied because $c_M \{ t_i^{(C)}(\hat{\lambda}) + u_{ii} \} > t_i^{(M)}(\hat{\lambda}) + c_M u_{ii}$ and $c_M > 0$. Using these relationships, we obtain the following theorem.

**Theorem 3.1**. For any $\hat{\lambda} \geq 0$, we obtain $\hat{\theta}_i^{(M)}(\hat{\lambda}) \geq \hat{\theta}_i^{(C)}(\hat{\lambda})$.

*proof*. We consider the following situations:

1) $t_i^{(M)}(\hat{\lambda}) \geq 0$ is satisfied,

2) $t_i^{(M)}(\hat{\lambda}) < 0 < t_i^{(M)}(\hat{\lambda}) + c_M u_{ii}$ is satisfied,

3) $t_i^{(M)}(\hat{\lambda}) + c_M u_{ii} \leq 0$ is satisfied.

In (1), $\hat{\theta}_i^{(M)}(\hat{\lambda}) = \hat{\theta}_i^{(C)}(\hat{\lambda}) = 0$, because $t_i^{(C)}(\hat{\lambda}) \geq 0$. In (3), $\hat{\theta}_i^{(M)}(\hat{\lambda}) \geq \hat{\theta}_i^{(C)}(\hat{\lambda})$, because $\hat{\theta}_i^{(M)}(\hat{\lambda})$ becomes $\infty$. Hence, we only consider situation (2). Note that $t_i^{(C)}(\hat{\lambda}) + u_{ii} > 0$, because $c_M \{ t_i^{(C)}(\hat{\lambda}) + u_{ii} \} > 0$ and $c_M > 0$. This means that $\hat{\theta}_i^{(C)}(\hat{\lambda})$ does not become $\infty$. This theorem holds when $t_i^{(C)}(\hat{\lambda}) \geq 0$, because, in this case, $\hat{\theta}_i^{(C)}(\hat{\lambda}) = 0$ and $\hat{\theta}_i^{(M)}(\hat{\lambda}) \geq 0$. We also consider $t_i^{(C)}(\hat{\lambda}) < 0 < t_i^{(C)}(\hat{\lambda}) + u_{ii}$ to be satisfied. Then, we obtain

$$\hat{\theta}_i^{(M)}(\hat{\lambda}) - \hat{\theta}_i^{(C)}(\hat{\lambda}) = \frac{d_i u_{ii} \{ c_M t_i^{(C)}(\hat{\lambda}) - t_i^{(M)}(\hat{\lambda}) \}}{\{ t_i^{(M)}(\hat{\lambda}) + c_M u_{ii} \} \{ t_i^{(C)}(\hat{\lambda}) + c_M u_{ii} \}}$$

Since $S^{-1}$ is a positive definite matrix, $u_{ii} \geq 0$ for any $\hat{\lambda} \geq 0$. From corollary 3.1, we have $c_M t_i^{(C)}(\hat{\lambda}) \geq t_i^{(M)}(\hat{\lambda})$ for any $\hat{\lambda} \geq 0$. Hence we obtain $\hat{\theta}_i^{(M)}(\hat{\lambda}) \geq \hat{\theta}_i^{(C)}(\hat{\lambda})$ for any $\hat{\lambda} \geq 0$ since $d_i > 0$, $t_i^{(M)}(\hat{\lambda}) + c_M u_{ii} > 0$ and $t_i^{(C)}(\hat{\lambda}) + u_{ii} > 0$. Thus, this theorem is proven.

This theorem corresponds to that in [9] when $X = I_p$ and $\lambda K = O_{q,q}$.

From Theorem 3.1, we obtained the relationships be-tween $\hat{\theta}_i^{(C)}(\hat{\lambda}^{(C)})$ and $\hat{\theta}_i^{(M)}(\hat{\lambda}^{(M)})$ for the case in which the optimized smoothing parameters $\hat{\lambda}^{(C)}$ and $\hat{\lambda}^{(M)}$ are the same. However, $\hat{\lambda}^{(C)}$ and $\hat{\lambda}^{(M)}$ are op-timized by minimizing the $C_p$ and $MC_p$ criteria in (10) and (11). Hence, $\hat{\lambda}^{(C)}$ and $\hat{\lambda}^{(M)}$ are generally different. Thus, we consider the relationship be-tween $\hat{\theta}_i^{(C)}(\hat{\lambda}^{(C)})$ and $\hat{\theta}_i^{(M)}(\hat{\lambda}^{(M)})$ when $\hat{\lambda}^{(C)} \neq \hat{\lambda}^{(M)}$. Since $u_{ii}$ is regarded as a function of $\lambda$, we write $u_{ii}$ as $u_{ii}(\hat{\lambda}^{(C)})$ and $u_{ii}(\hat{\lambda}^{(M)})$ for each optimized smoothing parameter.

**Theorem 3.2**. We consider the following situations:

1) $t_i^{(C)}(\hat{\lambda}^{(C)}) + u_{ii}(\hat{\lambda}^{(C)}) \leq 0$ or $t_i^{(M)}(\hat{\lambda}^{(M)}) \geq 0$ is satisfied,

2) $t_i^{(C)}(\hat{\lambda}^{(C)}) < 0 \leq t_i^{(C)}(\hat{\lambda}^{(C)}) + u_{ii}(\hat{\lambda}^{(C)})$ and $t_i^{(M)}(\hat{\lambda}^{(M)}) < 0 \leq t_i^{(M)}(\hat{\lambda}^{(M)}) + c_M u_{ii}(\hat{\lambda}^{(M)})$ are satisfied,

3) $c_M t_i^{(C)}(\hat{\lambda}^{(C)}) u_{ii}(\hat{\lambda}^{(C)}) \leq t_i^{(M)}(\hat{\lambda}^{(M)}) u_{ii}(\hat{\lambda}^{(C)})$ is sat-isfied,

4) $t_i^{(M)}(\hat{\lambda}^{(M)}) u_{ii}(\hat{\lambda}^{(C)}) \leq c_M t_i^{(C)}(\hat{\lambda}^{(C)}) u_{ii}(\hat{\lambda}^{(C)})$ is sat-isfied,

5) $t_i^{(C)}(\hat{\lambda}^{(C)}) \geq 0$ or $t_i^{(M)}(\hat{\lambda}^{(M)}) + c_M u_{ii}(\hat{\lambda}^{(M)}) \leq 0$ is satisfied.

For any $\hat{\lambda}^{(C)} \geq 0$ and $\hat{\lambda}^{(M)} \geq 0$, we obtain the fol-lowing relationships based on the above situations:

1) If (1), then $\hat{\theta}_i^{(M)}(\hat{\lambda}^{(M)}) \leq \hat{\theta}_i^{(C)}(\hat{\lambda}^{(C)})$,

2) If (2) and (3), then $\hat{\theta}_i^{(M)}(\hat{\lambda}^{(M)}) \leq \hat{\theta}_i^{(C)}(\hat{\lambda}^{(C)})$,

3) If (2) and (4), then $\hat{\theta}_i^{(C)}(\hat{\lambda}^{(C)}) \leq \hat{\theta}_i^{(M)}(\hat{\lambda}^{(M)})$,

4) If (5), then $\hat{\theta}_i^{(C)}(\hat{\lambda}^{(C)}) \leq \hat{\theta}_i^{(M)}(\hat{\lambda}^{(M)})$.

*proof*. In (1) and (5), the relationships (i) and (iv) are true. Hence we need only prove relationships (ii) and (iii). Then we obtain $\hat{\theta}_i^{(C)}(\hat{\lambda}^{(C)})$ and $\hat{\theta}_i^{(M)}(\hat{\lambda}^{(M)})$ using the closed forms of (12) and (13). Through simple calcula-tion, we obtain

Since $d_i > 0$ **and the denominator is positive**, the sign of $\hat{\theta}_i^{(M)}(\hat{\lambda}^{(M)}) - \hat{\theta}_i^{(C)}(\hat{\lambda}^{(C)})$ is the same as the sign of $c_M t_i^{(C)}(\hat{\lambda}^{(C)}) u_{ii}(\hat{\lambda}^{(C)}) - t_i^{(M)}(\hat{\lambda}^{(M)}) u_{ii}(\hat{\lambda}^{(C)})$. Hence we obtain relationships (ii) and (iii). Thus, the theorem is proven.

## 4. Numerical Studies

In this section, we compare the LS estimator $\hat{\boldsymbol{\mu}}$ and $\hat{\boldsymbol{\Xi}}$ in (3) with the proposed estimator $\hat{\boldsymbol{\mu}}_\lambda$ and $\hat{\boldsymbol{\Xi}}_{\boldsymbol{\theta},\lambda}$ in (7) through a numerical study. Let $R_r = \mathrm{diag}(1,\cdots,r)$, and let $\boldsymbol{\Delta}_r(\rho)$ be an $r \times r$ matrix as follows:

$$\boldsymbol{\Delta}_r(\rho) = \begin{pmatrix} 1 & \rho & \rho^2 & \cdots \rho^{r-1} \\ \rho & 1 & \rho & \cdots \rho^{r-2} \\ \rho^2 & \rho & 1 & \cdots \rho^{r-3} \\ \vdots & \vdots & \vdots & \ddots \vdots \\ \rho^{r-1} & \rho^{r-2} & \rho^{r-3} & \cdots 1 \end{pmatrix}.$$

The explanatory matrix $A$ is given by $A = N\boldsymbol{\Psi}^{1/2}$ where $\boldsymbol{\Psi} = R_k^{1/2}\boldsymbol{\Delta}_k(\rho_a)R_k^{1/2}$, $N$ is an $n \times k$ matrix and each row vector of $N$ is generated from the independent $k$-dimensional normal distribution with mean $\mathbf{0}_k$ and covariance matrix $I_k$. Let $\boldsymbol{m}_i$, $(i=1,\cdots,12)$ be a $p$-dimensional vector. We set each $\boldsymbol{m}_i$ as follows:

$$\boldsymbol{m}_1 = \boldsymbol{h}(t;e^2,e^{-1.5},e^1), \boldsymbol{m}_2 = \boldsymbol{h}(t;e^2,e^{-1.5},e^2),$$
$$\boldsymbol{m}_3 = \boldsymbol{h}(t;e^2,e^{-2.0},e^1), \boldsymbol{m}_4 = \boldsymbol{h}(t;e^2,e^{-2.0},e^2),$$
$$\boldsymbol{m}_5 = \boldsymbol{h}(t;e^2,e^{-2.5},e^1), \boldsymbol{m}_6 = \boldsymbol{h}(t;e^2,e^{-2.5},e^2),$$
$$\boldsymbol{m}_7 = \boldsymbol{h}(t;e^3,e^{-1.5},e^1), \boldsymbol{m}_8 = \boldsymbol{h}(t;e^3,e^{-1.5},e^2),$$
$$\boldsymbol{m}_9 = \boldsymbol{h}(t;e^3,e^{-2.0},e^1), \boldsymbol{m}_{10} = \boldsymbol{h}(t;e^3,e^{-2.0},e^2),$$
$$\boldsymbol{m}_{11} = \boldsymbol{h}(t;e^3,e^{-2.5},e^1), \boldsymbol{m}_{12} = \boldsymbol{h}(t;e^3,e^{-2.5},e^2),$$

where $t = (1,\cdots,p)'$ and the $i$th element of $\boldsymbol{h}(t;z_1,z_2,z_3)$ is $z_1\{1 - \exp(-z_2 t_i)\}^{z_3}$. Each element of $\boldsymbol{h}(t;z_1,z_2,z_3)$ is Richard's growth curve model [12]. We set the longitudinal trends using these $\boldsymbol{m}_i$ as $\boldsymbol{M}_k(t) = (\boldsymbol{m}_1,\cdots,\boldsymbol{m}_k)'$. Note that $\boldsymbol{m}_{i+6} = e\boldsymbol{m}_i$, $(i=1,\cdots,6)$, which indicates that the last six rows of $\boldsymbol{M}_{12}(t)$ are obtained by changing the scale of $\boldsymbol{M}_6(t)$. The response matrix $Y$ is generated by $N_{n \times p}(A\boldsymbol{M}_k(t), \Sigma \otimes I_n)$ where $\Sigma = R_p^{1/2}\boldsymbol{\Delta}_p(\rho_y)R_p^{1/2}$. Then, we standardized $A$. Let $\boldsymbol{k}_i = (\mathbf{0}'_{i-1}, 1, -2, 1, \mathbf{0}'_{q-2-i})'$, $(i=1,\cdots,q-2)$ and $K = (\boldsymbol{k}_1,\cdots,\boldsymbol{k}_{q-2})(\boldsymbol{k}_1,\cdots,\boldsymbol{k}_{q-2})'$. We set each element of $X$ as a cubic $B$-spline basis function. Since $X$ is set using the cubic $B$-spline, we note that $3 \le q \le p$. Additional details concerning $K$ and $X$ are reported in

[2]. We simulate $10,000$ repetitions for each $n$, $p$, $k$, $\rho_a$ and $\rho_y$. In each repetition, we fixed $A$, but $Y$ varies. We search $\hat{\lambda}^{(\mathrm{C})}$ and $\hat{\lambda}^{(\mathrm{M})}$ using fminsearch, which is a program in the software Matlab used to search for a minimum value, because $\hat{\lambda}^{(\mathrm{C})}$ and $\hat{\lambda}^{(\mathrm{M})}$ cannot be obtained in closed form. In searching $\hat{\lambda}^{(\mathrm{C})}$ and $\hat{\lambda}^{(\mathrm{M})}$, we transform $\lambda' = \exp(\lambda)$ and search optimized $\lambda'$ by each criterion because $\hat{\lambda}^{(\mathrm{C})} \ge 0$ and $\hat{\lambda}^{(\mathrm{M})} \ge 0$. In the search algorithm, the starting point for the search is set as $\lambda = 0$. Then, we obtain the optimized ridge parameters $\hat{\theta}_i^{(\mathrm{C})}(\hat{\lambda}^{(\mathrm{C})})$ and $\hat{\theta}_i^{(\mathrm{M})}(\hat{\lambda}^{(\mathrm{M})})$ using the closed forms of (12) and (13) in each repetition. In each repetition, we need to optimize $q$ because $X$ and $K$ vary with $q$. We calculate $C_p(\hat{\boldsymbol{\theta}}^{(\mathrm{C})}(\hat{\lambda}^{(\mathrm{C})}), \hat{\lambda}^{(\mathrm{C})})$ and $MC_p(\hat{\boldsymbol{\theta}}^{(\mathrm{M})}(\hat{\lambda}^{(\mathrm{M})}), \hat{\lambda}^{(\mathrm{M})})$ for each $q = 3,...,p$ in each repetition. Then, we adopt the optimized $q$ by minimizing each criterion in each repetition. After that, we calculate for $r(E[Y], \hat{Y}_{\hat{\theta}(\hat{\lambda}),\hat{\lambda}})/(np)$ each criterion, where $\hat{Y}_{\hat{\theta}(\hat{\lambda}),\hat{\lambda}} = \mathbf{1}_n\hat{\boldsymbol{\mu}}'_{\hat{\lambda}}X' + A\hat{\boldsymbol{\Xi}}_{\hat{\theta}(\hat{\lambda}),\hat{\lambda}}X'$, which is obtained using $\hat{\lambda}$ and $\hat{\boldsymbol{\theta}}(\hat{\lambda})$ for each criterion and the optimized $q$ in each repetition. The average of $r(E[Y], \hat{Y}_{\hat{\theta}(\hat{\lambda}),\hat{\lambda}})$ over $10,000$ repetitions is regarded as the MSE of $\hat{Y}_{\hat{\theta}(\hat{\lambda}),\hat{\lambda}}$. We compare the values predicted using the estimators $\hat{\boldsymbol{\mu}}_{\hat{\lambda}}$ and $\hat{\boldsymbol{\Xi}}_{\hat{\theta}(\hat{\lambda}),\hat{\lambda}}$ with those using the LS estimators $\hat{\boldsymbol{\mu}}$ and $\hat{\boldsymbol{\Xi}}$, **and** the estimators $\hat{\boldsymbol{\mu}}_{\hat{\lambda}}$ and $\hat{\boldsymbol{\Xi}}_{\hat{\lambda}}$ in (5). When we use $\hat{\boldsymbol{\Xi}}_{\hat{\lambda}}$, we obtain $\hat{\lambda}$ by minimizing $C_p(\mathbf{0}_k, \lambda)$ and $MC_p(\mathbf{0}_k, \lambda)$. As in the case of using $\hat{\boldsymbol{\Xi}}_{\hat{\theta}(\hat{\lambda}),\hat{\lambda}}$, we adopt $q$ by using each criterion in each repetition for $\hat{\boldsymbol{\Xi}}$ and $\hat{\boldsymbol{\Xi}}_{\hat{\lambda}}$. Some of the results are shown in **Tables 1 and 2.** The values in the tables are obtained by $\mathrm{MSE}[\hat{\boldsymbol{Y}}_{\hat{\theta}(\hat{\lambda}),\hat{\lambda}}]/(np)$,

$\mathrm{MSE}[\hat{\boldsymbol{Y}}_{\hat{\lambda}}]/(np)$ where $\hat{Y}_{\hat{\lambda}} = \mathbf{1}_n\hat{\boldsymbol{\mu}}'_{\hat{\lambda}}X' + A\hat{\boldsymbol{\Xi}}_{\hat{\lambda}}X'$, and $\mathrm{MSE}[\hat{\boldsymbol{Y}}]/(np)$ where $\hat{Y} = \mathbf{1}_n\hat{\boldsymbol{\mu}}'X' + A\hat{\boldsymbol{\Xi}}X'$.

Each estimator optimized by using the $MC_p$ criterion for $\lambda$, $\hat{\boldsymbol{\theta}}$, and $q$ is more improve than that by using the $C_p$ criterion for each estimator in almost all situations. This indicates that the $MC_p$ criterion is a better estimator of the MSE of each predicted value of $Y$ than the $C_p$ criterion. The reasons for this are that the $MC_p$ criterion is an **unbiased** estimator of MSE

**Table 1. MSE when $q$ is selected using each criterion for each method in each repetition $(k = 6)$.**

| $\rho_y$ | $\rho_a$ | $p$ | $n$ | Using $\hat{Y}_{\hat{\theta}(\hat{\lambda}),\hat{\lambda}}$ | | Using $\hat{Y}_{\hat{\lambda}}$ | | Using $\hat{Y}$ | |
|---|---|---|---|---|---|---|---|---|---|
| | | | | $C_p$ | $MC_p$ | $C_p$ | $MC_p$ | $C_p$ | $MC_p$ |
| 0.2 | 0.2 | 5 | 30 | 0.127 | **0.123** | 0.133 | 0.125 | 0.206 | 0.199 |
| | | | 50 | 0.080 | **0.079** | 0.082 | 0.080 | 0.121 | 0.119 |
| | | 10 | 30 | 0.119 | 0.098 | 0.121 | **0.090** | 0.168 | 0.119 |
| | | | 50 | 0.063 | 0.058 | 0.062 | **0.056** | 0.080 | 0.070 |
| | 0.8 | 5 | 30 | 0.110 | **0.101** | 0.143 | 0.135 | 0.206 | 0.199 |
| | | | 50 | 0.067 | **0.065** | 0.088 | 0.086 | 0.122 | 0.119 |
| | | 10 | 30 | 0.111 | **0.080** | 0.128 | 0.093 | 0.170 | 0.119 |
| | | | 50 | 0.056 | **0.049** | 0.063 | 0.057 | 0.080 | 0.070 |
| | 0.99 | 5 | 30 | 0.090 | **0.078** | 0.147 | 0.140 | 0.207 | 0.199 |
| | | | 50 | 0.054 | **0.050** | 0.090 | 0.088 | 0.122 | 0.120 |
| | | 10 | 30 | 0.095 | **0.060** | 0.129 | 0.094 | 0.169 | 0.118 |
| | | | 50 | 0.045 | **0.036** | 0.064 | 0.058 | 0.079 | 0.069 |
| 0.8 | 0.2 | 5 | 30 | 0.133 | **0.131** | 0.154 | 0.147 | 0.208 | 0.201 |
| | | | 50 | 0.087 | **0.086** | 0.093 | 0.092 | 0.122 | 0.120 |
| | | 10 | 30 | 0.123 | **0.101** | 0.136 | 0.106 | 0.179 | 0.133 |
| | | | 50 | 0.069 | **0.065** | 0.070 | 0.065 | 0.089 | 0.080 |
| | 0.8 | 5 | 30 | 0.113 | **0.103** | 0.159 | 0.153 | 0.207 | 0.200 |
| | | | 50 | 0.066 | **0.063** | 0.094 | 0.092 | 0.122 | 0.120 |
| | | 10 | 30 | 0.108 | **0.074** | 0.140 | 0.107 | 0.178 | 0.131 |
| | | | 50 | 0.055 | **0.047** | 0.072 | 0.065 | 0.088 | 0.078 |
| | 0.99 | 5 | 30 | 0.092 | **0.078** | 0.162 | 0.156 | 0.208 | 0.201 |
| | | | 50 | 0.053 | **0.049** | 0.095 | 0.094 | 0.122 | 0.120 |
| | | 10 | 30 | 0.096 | **0.059** | 0.142 | 0.108 | 0.178 | 0.131 |
| | | | 50 | 0.046 | **0.037** | 0.073 | 0.066 | 0.087 | 0.078 |
| | Average | | | 0.086 | **0.074** | 0.110 | 0.098 | 0.147 | 0.130 |

**Table 2. MSE when $q$ is selected using each criterion for each method in each repetition $(k = 12)$.**

| $\rho_y$ | $\rho_a$ | $p$ | $n$ | Using $\hat{Y}_{\hat{\theta}(\hat{\lambda}),\hat{\lambda}}$ | | Using $\hat{Y}_{\hat{\lambda}}$ | | Using $\hat{Y}$ | |
|---|---|---|---|---|---|---|---|---|---|
| | | | | $C_p$ | $MC_p$ | $C_p$ | $MC_p$ | $C_p$ | $MC_p$ |
| 0.2 | 0.2 | 5 | 30 | 0.299 | **0.292** | 0.312 | 0.296 | 0.383 | 0.364 |
| | | | 50 | 0.184 | 0.183 | 0.183 | **0.180** | 0.222 | 0.217 |
| | | 10 | 30 | 0.317 | 0.247 | 0.326 | **0.226** | 0.382 | 0.248 |
| | | | 50 | 0.146 | 0.137 | 0.146 | **0.134** | 0.165 | 0.150 |
| | 0.8 | 5 | 30 | 0.285 | **0.279** | 0.313 | 0.295 | 0.384 | 0.365 |
| | | | 50 | 0.175 | **0.173** | 0.182 | 0.179 | 0.223 | 0.218 |
| | | 10 | 30 | 0.305 | 0.223 | 0.329 | **0.216** | 0.378 | 0.226 |
| | | | 50 | 0.145 | 0.132 | 0.145 | **0.129** | 0.155 | 0.135 |
| | 0.99 | 5 | 30 | 0.224 | **0.204** | 0.314 | 0.296 | 0.383 | 0.364 |
| | | | 50 | 0.142 | **0.138** | 0.183 | 0.180 | 0.222 | 0.218 |
| | | 10 | 30 | 0.270 | **0.173** | 0.330 | 0.211 | 0.377 | 0.221 |
| | | | 50 | 0.134 | **0.119** | 0.143 | 0.123 | 0.148 | 0.127 |
| 0.8 | 0.2 | 5 | 30 | 0.323 | **0.321** | 0.342 | 0.331 | 0.387 | 0.368 |
| | | | 50 | 0.204 | 0.204 | 0.205 | **0.203** | 0.224 | 0.219 |
| | | 10 | 30 | 0.330 | 0.277 | 0.344 | **0.256** | 0.389 | 0.282 |
| | | | 50 | 0.165 | 0.153 | 0.167 | **0.152** | 0.178 | 0.159 |
| | 0.8 | 5 | 30 | 0.298 | **0.294** | 0.337 | 0.321 | 0.385 | 0.367 |
| | | | 50 | 0.191 | **0.191** | 0.200 | 0.197 | 0.224 | 0.220 |
| | | 10 | 30 | 0.309 | **0.244** | 0.346 | 0.251 | 0.386 | 0.265 |
| | | | 50 | 0.161 | **0.150** | 0.166 | 0.151 | 0.175 | 0.159 |
| | 0.99 | 5 | 30 | 0.228 | **0.208** | 0.338 | 0.322 | 0.386 | 0.368 |
| | | | 50 | 0.142 | **0.137** | 0.199 | 0.196 | 0.223 | 0.219 |
| | | 10 | 30 | 0.263 | **0.170** | 0.347 | 0.236 | 0.384 | 0.247 |
| | | | 50 | 0.126 | **0.106** | 0.161 | 0.139 | 0.166 | 0.145 |
| | Average | | | 0.086 | **0.074** | 0.110 | 0.098 | 0.289 | 0.245 |

and each of the parameters in each estimator is optimized by minimizing the $MC_p$ criterion. When $k = 6$, $\hat{\Xi}_{\theta,\lambda}$ provides a greater improvement than either $\hat{\Xi}_\lambda$ or $\hat{\Xi}$ in all situations. The estimator $\hat{\Xi}_{\theta,\lambda}$, which is optimized using the $MC_p$ criterion, has the smallest MSE among these estimators for almost situations when $k = 6$. Here, $\hat{\Xi}_\lambda$ provides a greater improvement than $\hat{\Xi}$ when $k = 6$ in all situations. When $\rho_a$ is large, the estimator $\hat{\Xi}_{\theta,\lambda}$ provides a greater improvement than $\hat{\Xi}_\lambda$ in most situations when $k = 12$. On the other hand, $\hat{\Xi}_\lambda$ provides a greater improvement than $\hat{\Xi}_{\theta,\lambda}$ in most situations when $\rho_a$ is small, $k = 12$ and $p = 10$. If $k = 12$, then $\hat{\Xi}_{\theta,\lambda}$ and $\hat{\Xi}_\lambda$ improve the LS estimator. Comparing the results for $k = 6$ with the results for $k = 12$ reveals that these estimators become poor estimators when $k$ becomes large. The reasons for this are thought to be that $\mathbf{S}^{-1}$ and $\mathbf{A}$ become unstable and the $\mathbf{M}_{12}(t)$ has some curves that are in a different scale. Each MSE using each method and the $C_p$ criterion is similar to that using the $MC_p$ criterion if $n$ becomes large because $c_\mathrm{M}$ is close to 1. When $\rho_a$ becomes large, $\hat{\Xi}_{\theta,\lambda}$ improves the LS estimator more than when $\rho_a$ is small. Since $\rho_a$ controls the correlation in $\mathbf{A}$, the multicollinearity in $\mathbf{A}$ becomes large when $\rho_a$ becomes large. Then, $\hat{\Xi}_\lambda$ is not a good estimator because $(\mathbf{A}'\mathbf{A})^{-1}$ is unstable. Hence, we can avoid the multicollinearity problem in $\mathbf{A}$ by using $\hat{\Xi}_{\theta,\lambda}$, which is one of the purposes of the present study. In all situations, the new estimators improve the LS estimator $\hat{\Xi}$. In addition, $\hat{\Xi}_{\theta,\lambda}$ is better than $\hat{\Xi}_\lambda$ in most situations, especially when $k$ is small or $\rho_a$ is large. In general, $\hat{\Xi}_{\theta,\lambda}$ optimized using $MC_p$ is the best method.

## 5. Conclusions

In the present paper, we estimate the longitudinal trends nonparametrically by using the nonparametric GMANOVA model in (1), which is defined using basis functions as $\mathbf{X}$ in the GMANOVA model. When we use basis functions as $\mathbf{X}$, the LS estimators $\hat{\boldsymbol{\mu}}$ and $\hat{\Xi}$ incur overfitting. In order to avoid this problem, we proposed $\hat{\boldsymbol{\mu}}_\lambda$ and $\hat{\Xi}_\lambda$ in (5) using the smoothing parameter $\lambda (\geq 0)$ and the $q \times q$ known penalty non-negative definite matrix $\mathbf{K}$. However, if multicollinearity occurs in $\mathbf{A}$, $\hat{\Xi}$ and $\hat{\Xi}_\lambda$ are not good estimators due to large variance. In the present paper, we also proposed $\hat{\Xi}_{\theta,\lambda}$ in (7) in order to avoid the multicollinearity problem that occurs in $\mathbf{A}$ and the overfitting problem by using basis functions as $\mathbf{X}$. The estimator $\hat{\Xi}_\lambda$ controls the smoothness of each estimated longitudinal curve using only one parameter $\lambda$. On the other hand, in the estimator $\hat{\Xi}_{\theta,\lambda}$, the rough smoothness of estimated longitudinal curves is controlled using $\lambda$, and each smoothness of $\phi_1(t), \cdots,$

$\phi_k(t)$ in the varying coefficient model (4) is controlled by $\boldsymbol{\theta}$.

We also proposed the $C_p$ and $MC_p$ criteria in (10) and (11) for optimizing the ridge parameter $\boldsymbol{\theta}$ and the smoothing parameter $\lambda$. Then, using the $GC_p$ criterion in (14) and minimizing this criterion in Theorem A, we obtain the optimized $\boldsymbol{\theta}$ using the $C_p$ and $MC_p$ criteria in closed form as (12) and (13) for any $\lambda$. Thus, we can regard the $C_p$ and $MC_p$ criteria as a function of $\lambda$.

Hence, we need to optimize only one parameter $\lambda$ in order to optimize $k + 1$ parameters in $\hat{\Xi}_{\theta,\lambda}$ using these criteria. On the other hand, we must optimize two parameters when we use $\hat{\Xi}_{\theta,\lambda}$ in (6). This optimization is difficult and requires a complicated program and a long processing time for simulation or analysis of real data because the optimized $\boldsymbol{\theta}$ cannot be obtained in closed form even if $\lambda$ is fixed. This is the advantage of using $\hat{\Xi}_{\theta,\lambda}$. This advantage does not appear to be important because of the high calculation power of CPUs. However, this advantage is made clear when we use $\hat{\Xi}_{\theta,\lambda}$ together with variable selection. Even if $k$ becomes large, then this advantage remains when $\hat{\Xi}_{\theta,\lambda}$ is used because the optimized $\boldsymbol{\theta}$ obtained using each criterion is always obtained as (12) and (13) for any $k$. Furthermore, we must optimize $q$ if we use model (1) to estimate the longitudinal trends. This means that we optimize the parameters in the estimators and calculate the valuation of the estimator for each $q$, and then we compare these values in order to optimize $q$. Since this optimization requires an iterative computational algorithm, we must reduce the processing time for estimating the parameters in the estimator. Hence, the advantage of using $\hat{\Xi}_{\theta,\lambda}$ is very important. This optimized ridge parameter in (12) and (13) corresponds to that in [18] when $\mathbf{X} = \mathbf{I}_p$ and $\lambda \mathbf{K} = \mathbf{O}_{q,q}$.

Using some matrix properties, we showed that $\mathrm{tr}(\mathbf{G}_\lambda)$ and $\mathrm{tr}(\mathbf{H}_\theta)$ in the $C_p$ and $MC_p$ criteria are always nonnegative. From $\mathrm{tr}(\mathbf{G}_\lambda) \geq 0$ for any $\lambda \geq 0$ in lemma 3.1, we also established the relationship between $t_i^{(\mathrm{C})}(\lambda)$ and $t_i^{(\mathrm{M})}(\lambda)$ for any $\lambda \geq 0$ in corollary 3.1. Then, in Theorem 3.1, we established the relationship between $\hat{\theta}_i^{(\mathrm{C})}(\hat{\lambda}^{(\mathrm{C})})$ and $\hat{\theta}_i^{(\mathrm{M})}(\hat{\lambda}^{(\mathrm{M})})$ if $\hat{\lambda}^{(\mathrm{C})}$ and $\hat{\lambda}^{(\mathrm{M})}$ are the same, where $\hat{\lambda}^{(\mathrm{C})}$ and $\hat{\lambda}^{(\mathrm{M})}$ are obtained by minimizing the $C_p$ and $MC_p$ criteria. Note that this relationship corresponds to that in [9] when $\mathbf{X} = \mathbf{I}_p$ and $\lambda \mathbf{K} = \mathbf{O}_{q,q}$. In Theorem 3.2, we also established the relationships between $\hat{\theta}_i^{(\mathrm{C})}(\hat{\lambda}^{(\mathrm{C})})$ and

$\hat{\theta}_i^{(M)}\left(\hat{\lambda}^{(M)}\right)$ for the more general case, in which $\hat{\lambda}^{(C)}$ and $\hat{\lambda}^{(M)}$ are different. The reason of the relationship in Theorem 3.2 is occurred is that $\hat{\theta}_i^{(C)}(\lambda)$ and $\hat{\theta}_i^{(M)}(\lambda)$ for each $i = 1, \cdots, k$ can be regarded as a function of $\lambda$.

The numerical results reveal that $\hat{\Xi}_\lambda$ and $\hat{\Xi}_{\theta,\lambda}$ have some following properties. These estimation methods $\hat{\Xi}_\lambda$ and $\hat{\Xi}_{\theta,\lambda}$ improve the LS estimator in all situations, especially when $\rho_a$ is large. This indicates that the proposed estimators are better than the LS estimator. Even if $\rho_a$ becomes large, we note that $\hat{\Xi}_{\theta,\lambda}$ is stable because we add the ridge parameter to $A'A$ in the LS estimator. This result indicates that the multicollinearity problem in $A$ can be avoided by using the estimator in (7). These estimators can be used to estimate the true longitudinal trends nonparametrically using basis functions as $X$ without overfitting. The LS estimator and the proposed estimators $\hat{\Xi}_\lambda$ and $\hat{\Xi}_{\theta,\lambda}$ optimized using the $MC_p$ criterion provide a greater improvement than the estimators optimized using the $C_p$ criterion in most situations. The reason for this is that the $MC_p$ criterion is the unbiased estimator of MSE of the predicted value of $Y$. Based on the present numerical study, $\hat{\mu}_\lambda$ and $\hat{\Xi}_{\theta,\lambda}$ can be used to estimate the longitudinal trends in most situations. In addition, the $MC_p$ can be used to optimize the smoothing parameter $\lambda$ and the number of basis functions $q$. Hence, we can use $\hat{\mu}_\lambda$ and $\hat{\Xi}_{\theta,\lambda}$, the parameters $\theta$, $\lambda$, and $q$ of which are optimized by the $MC_p$ criterion for estimating the longitudinal trends.

## 6. Acknowledgments

## 7. Appendix

### 7.1. Minimization of the $GC_p$ Criterion

In this appendix, we show that the optimizations using the $C_p$ and $MC_p$ criteria in (10) and (11) are obtained in closed form as (12) and (13) for any $\lambda$ $(\geq 0)$. [9]

proposed the generalized $C_p$ $\left(GC_p\right)$ criterion for the MGR regression (originally the $GC_p$ criterion for selection variables in the univariate regression was proposed by [1]). Similar to their idea, we proposed the $GC_p$ criterion for the nonparametric GMANOVA model.

By omitting constant terms and some terms with respect to $\lambda$ in the $C_p$ and $MC_p$ criteria in (10) and (11), these criteria are included in a class of criteria specified by $\alpha$ $(>0)$. This class is expressed by the $GC_p$ criterion as

$$GC_p\left(\theta, \lambda \mid \alpha\right) = \alpha r\left(Y, \hat{Y}_{\theta,\lambda}\right) + 2\text{tr}\left(G_\lambda\right)\text{tr}\left(H_\theta\right), \quad (14)$$

where the function $\hat{r}$ is given by (8). Note that $GC_p\left(\theta, \lambda \mid 1\right)$ and $GC_p\left(\theta, \lambda \mid c_M\right)$ correspond to the terms with respect to $\theta$ in the $C_p$ and $MC_p$ criteria. Using this $GC_p$ criterion, we can deal systematically with the $C_p$ and $MC_p$ criteria for optimizing $\theta$. Let $\hat{\theta}^{(G)}\left(\lambda \mid \alpha\right) = \left(\hat{\theta}_1^{(G)}\left(\lambda \mid \alpha\right), \cdots, \hat{\theta}_k^{(G)}\left(\lambda \mid \alpha\right)\right)'$, $\left(\hat{\theta}_i^{(G)}\left(\lambda \mid \alpha\right) \geq 0, i = 1, \cdots, k\right)$ **which minimize** the $GC_p$ criterion for any $\lambda$ $(\geq 0)$. Then, $\hat{\theta}^{(C)}\left(\lambda\right)$ and $\hat{\theta}^{(M)}\left(\lambda\right)$ are obtained as $\hat{\theta}^{(C)}\left(\lambda\right) = \hat{\theta}^{(G)}\left(\lambda \mid 1\right)$ and $\hat{\theta}^{(M)}\left(\lambda\right) = \hat{\theta}^{(G)}\left(\lambda \mid c_M\right)$, respectively. Thus, we can deal systematically with the optimizations of $\theta$ when we use the $C_p$ and $MC_p$ criteria. This means that we need only obtain $\hat{\theta}^{(G)}\left(\lambda \mid \alpha\right)$ in order to obtain $\hat{\theta}^{(C)}\left(\lambda\right)$ and $\hat{\theta}^{(M)}\left(\lambda\right)$ for any $\lambda$ and some $\alpha$. If $\hat{\theta}^{(G)}\left(\lambda \mid \alpha\right)$ is obtained in closed form for any fixed $\lambda$, we do not need to use the iterative computational algorithm for optimizing the ridge parameter $\theta$. In order to obtain $\hat{\theta}^{(G)}\left(\lambda \mid \alpha\right)$, we obtain $\hat{\theta}_i^{(G)}\left(\lambda \mid \alpha\right)$, $\left(i = 1, \cdots, k\right)$ in closed form, as shown in the following theorem.

**Theorem A.** For any $i$ and $\lambda$ $(\geq 0)$, $\hat{\theta}_i^{(G)}\left(\lambda \mid \alpha\right)$ is obtained as

$$\hat{\theta}_i^{(G)}\left(\lambda \mid \alpha\right) = \begin{cases} 0 & \left(0 \leq t_i\left(\lambda \mid \alpha\right)\right) \\ \dfrac{-d_i t_i\left(\lambda \mid \alpha\right)}{t_i\left(\lambda \mid \alpha\right) + \alpha u_{ii}} & \left(t_i\left(\lambda \mid \alpha\right) < 0 < t_i\left(\lambda \mid \alpha\right) + \alpha u_{ii}\right), \\ \infty & \left(\text{otherwise}\right) \end{cases} \quad (15)$$

where $t_i(\lambda \mid \alpha) = \alpha(v_{ii} - u_{ii}) - d_i \mathrm{tr}(G_\lambda)$.

*proof.* Since $\hat{Y}_{\theta,\lambda} = 1_n \hat{\mu}'_\lambda X' + H_\theta Y G_\lambda$ and we use the

$$\mathrm{tr}\left\{(1_n \hat{\mu}_\lambda'X' - Y)S^{-1}(H_\theta Y G_\lambda)'\right\} = -\mathrm{tr}\left(YS^{-1}G_\lambda Y'H_\theta\right).$$

properties of the function $\hat{r}$ in Section 3.1, we can calculate $\hat{r}(Y, \hat{Y}_{\theta,\lambda})$ in the $GC_p$ criterion in (14) as follows:

$$\hat{r}(Y, \hat{Y}_{\theta,\lambda}) = \hat{r}(Y, 1_n \hat{\mu}'_\lambda X') + \hat{r}(H_\theta Y G_\lambda, O_{n,p})$$
$$+ 2\mathrm{tr}\left\{(1_n \hat{\mu}'_\lambda X' - Y)S^{-1}(H_\theta Y G_\lambda)'\right\}.$$

Since $G_\lambda = G'_\lambda$ for any $\lambda$, $H_\theta = H'_\theta$ and $H_\theta 1_n = 0_k$ for any $\theta$, the second term in the right-hand side of the above equation can be calculated as

Note that $H_\theta = A(A'A + Q\Theta Q')^{-1}A' = AQ(D+\Theta))^{-1}Q'A'$ because $Q$ is an orthogonal matrix and $Q'A'AQ = D$. Hence, we obtain the following results:

$$\mathrm{tr}\left(YS^{-1}G_\lambda Y'H_\theta\right) = \mathrm{tr}\left\{Q'A'YS^{-1}G_\lambda Y'AQ(D+\Theta)^{-1}\right\},$$

$$\hat{r}(H_\theta Y G_\lambda, O_{n,p})$$
$$= \mathrm{tr}\left(H_\theta Y G_\lambda S^{-1} G_\lambda Y'H_\theta\right)$$
$$= \mathrm{tr}\left\{Q'A'Y G_\lambda S^{-1} G_\lambda Y'AQ(D+\Theta)^{-1}D(D+\Theta)^{-1}\right\}.$$

Since $D$ and $(D+\Theta)^{-1}$ are diagonal matrices, we obtain $(D+\Theta)^{-1}D(D+\Theta)^{-1} = D(D+\Theta)^{-2}$. Hence $\hat{r}(Y, \hat{Y}_{\theta,\lambda})$ is calculated as

$$\hat{r}(Y, \hat{Y}_{\theta,\lambda}) = \hat{r}(Y, 1_n \hat{\mu}'_\lambda X')$$
$$- 2\mathrm{tr}\left\{V(D+\Theta)^{-1}\right\} + \mathrm{tr}\left\{UD(D+\Theta)^{-2}\right\},$$

where $U = Q'A'Y G_\lambda S^{-1} G_\lambda Y'AQ$ and $V = Q'A'YS^{-1}G_\lambda Y'AQ$. Clearly, $U$ and $V$ change with $\lambda$. Based on this result and $\mathrm{tr}(H_\theta) = \mathrm{tr}\left\{D(D+\Theta)^{-1}\right\}$, we can calculate the $GC_p$ criterion in (14) as follows:

$$GC_p(\theta, \lambda \mid \alpha) = \alpha\hat{r}(Y, 1_n \hat{\mu}'_\lambda X') + \alpha\mathrm{tr}\left\{UD(D+\Theta)^{-2}\right\}$$
$$- 2\mathrm{tr}\left\{(\alpha V - \mathrm{tr}(G_\lambda)D)(D+\Theta)^{-1}\right\}.$$

Then, we calculate the second and third terms in the right-hand side of the above equation as follows:

$$\alpha\mathrm{tr}\left\{UD(D+\Theta)^{-2}\right\} - 2\mathrm{tr}\left\{\alpha V - \mathrm{tr}(G_\lambda D)(D+\Theta)^{-1}\right\}$$
$$= \sum_{i=1}^{k}\left\{\frac{\alpha d_i u_{ii}}{(d_i + \theta_i)^2} - 2\frac{\alpha v_{ii} - d_i \mathrm{tr}(G_\lambda)}{d_i + \theta_i}\right\},$$

where $u_{ij}$ and $v_{ij}$ are the $(i, j)$th element of $U$ and $V$, respectively. Clearly, $u_{ij}$ and $v_{ij}$ also vary with $\lambda$. Note that $u_{ii} \geq 0$, $(i = 1, \cdots, k)$ for any $\lambda \geq 0$ because $S^{-1}$ is a positive definite matrix (see, e.g., [3]). Let $\varphi_i(\theta_i, \lambda \mid \alpha)$, $(i = 1, \cdots, k)$ be as follows:

$$\varphi_i(\theta_i, \lambda \mid \alpha) = \frac{\alpha d_i u_{ii}}{(d_i + \theta_i)^2} - 2\frac{\alpha v_{ii} - d_i \mathrm{tr}(G_\lambda)}{d_i + \theta_i}. \quad (16)$$

Using $\varphi_i(\theta_i, \lambda \mid \alpha)$, we can express

$$GC_p(\theta, \lambda \mid \alpha) = \alpha\hat{r}(Y, 1_n \hat{\mu}'_\lambda X') + \sum_{i=1}^{k}\varphi_i(\theta_i, \lambda \mid \alpha).$$

Since $\alpha\hat{r}(Y, 1_n \hat{\mu}'_\lambda X')$ does not depend on $\theta$, we can obtain $\hat{\theta}_i^{(G)}(\lambda \mid \alpha)$ by minimizing $\varphi_i(\theta_i, \lambda \mid \alpha)$ for each $i = 1, \cdots, k$ and any $\lambda$ $(\geq 0)$. In order to obtain $\hat{\theta}_i^{(G)}(\lambda \mid \alpha)$, we consider the following function for $w \in \mathbb{R}$:

$$\varphi_i(w) = \frac{\alpha d_i u_{ii}}{(d_i + w)^2} - 2\frac{\alpha v_{ii} - d_i \mathrm{tr}(G_\lambda)}{d_i + w}. \quad (17)$$

If we restrict $w$ to be greater than or equal to 0, then this function is equivalent to the function $\varphi_i(\theta_i, \lambda \mid \alpha)$ in (16), which must be minimized. Note that $\lim_{w \to \pm\infty} \varphi_i(w) = 0$ and $\lim_{w \to -d_i \pm 0} \varphi_i(w) = \infty$. Letting $\dot{\varphi}_i(w) = \partial\varphi_i(w)/\partial w$, we obtain

$$\dot{\varphi}_i(w) = -\frac{2}{(d_i + w)^3}\left\{\alpha d_i u_{ii} - (\alpha v_{ii} - d_i \mathrm{tr}(G_\lambda))(d_i + w)\right\}$$

Let $\hat{w}$ satisfy $\dot{\varphi}_i(w)|_{w=\hat{w}} = 0$ and $\hat{w} \neq \pm\infty$, then $\hat{w}$ is obtained by

$$\hat{w} = \frac{-d_i t_i(\lambda \mid \alpha)}{t_i(\lambda \mid \alpha) + \alpha u_{ii}}, \quad (\text{if } t_i(\lambda \mid \alpha) + \alpha u_{ii} \neq 0),$$

where $t_i(\lambda \mid \alpha) = \alpha(v_{ii} - u_{ii}) - d_i \mathrm{tr}(G_\lambda)$. Note that $\varphi_i(w)$ in (17) has a minimum value at $\hat{w}$, which is $\dot{\varphi}_i(w)|_{w<\hat{w}} < 0$ and $\dot{\varphi}_i(w)|_{w>\hat{w}} > 0$. Note that the sign of $t_i(\lambda \mid \alpha)$ is the same as the sign of $\dot{\varphi}_i(w)|_{w=0}$. In order to obtain $\hat{\theta}_i^{(G)}(\lambda \mid \alpha)$ $(\geq 0)$, we consider the following situations:

1) $t_i(\lambda \mid \alpha) \geq 0$ is satisfied,

2) $t_i(\lambda \mid \alpha) < 0$ and $t_i(\lambda \mid \alpha) + \alpha u_{ii} > 0$ are satisfied,

3) $t_i(\lambda \mid \alpha) < 0$ and $t_i(\lambda \mid \alpha) + \alpha u_{ii} < 0$ are satisfied.

      

In (1), $-d_i < \hat{w} < 0$, because $u_{ii} \geq 0$ and $\alpha > 0$. In addition, $\varphi_i(w) \geq \varphi_i(0)$ for any $w \geq 0$, because $\hat{w} < 0$, and $t_i(\lambda \mid \alpha) \geq 0$ indicates that the sign of $\dot{\varphi}_i(w)|_{w=0}$ is nonnegative. This means that the minimum value of $\varphi_i(w)$ in $w \geq 0$ is obtained when $w = 0$ in situation (1). In (2), $\hat{w} > 0$, and then the minimum value of $\varphi_i(w)$ in $w \geq 0$ is obtained when $w = \hat{w}$. In (3), since $\hat{w} < -d_i$ and $\dot{\varphi}_i(w)|_{w=0} < 0$, we obtain $\varphi_i(0) > \varphi_i(w_1) > \varphi_i(w_2)$ for any $w_2 > w_1 > 0$. Hence, $\varphi_i(w)$ is minimized when $w = \infty$ in $w \geq 0$.

From the above results, we obtain $\hat{\theta}_i^{(G)}(\lambda \mid \alpha) \ (\geq 0)$ as follows:

$$\hat{\theta}_i^{(G)}(\lambda \mid \alpha)$$
$$= \begin{cases} 0 & \left(0 \leq t_i(\lambda \mid \alpha)\right) \\ \dfrac{-d_i t_i(\lambda \mid \alpha)}{t_i(\lambda \mid \alpha) + \alpha u_{ii}} & \left(t_i(\lambda \mid \alpha) < 0 < t_i(\lambda \mid \alpha) + \alpha u_{ii}\right), \\ \infty & \left(\text{otherwise}\right) \end{cases}$$
$$(i = 1, \cdots, k).$$

Thus, the theorem is proven.

Note that $\hat{\boldsymbol{\theta}}^{(G)}(\lambda \mid \alpha)$ corresponds to that in [9] when $\boldsymbol{X} = \boldsymbol{I}_p$ and $\lambda \boldsymbol{K} = \boldsymbol{O}_{q,q}$. Since we obtain $\hat{\boldsymbol{\theta}}^{(C)}(\lambda)$ and $\hat{\boldsymbol{\theta}}^{(M)}(\lambda)$ in closed form as (15) for any $\lambda$, we must optimize only one parameter $\lambda$ in order to optimize $k + 1$ parameters. The use of $\hat{\boldsymbol{\Xi}}_{\boldsymbol{\theta},\lambda}$ is advantageous because only an iterative computational algorithm is required for optimizing only one parameter $\lambda$ for any $k$. This means that we can reduce the processing time required to optimize the parameters in the estimator $\hat{\boldsymbol{\Xi}}_{\boldsymbol{\theta},\lambda}$ which is defined by (7). When we use $\hat{\boldsymbol{\Xi}}_\lambda$ in (5), we also need the same iterative computational algorithm to optimize only one parameter $\lambda$.

On the other hand, when we use $\hat{\boldsymbol{\Xi}}_{\boldsymbol{\theta},\lambda}$ in (6), the $GC_p$ criterion for optimizing $\boldsymbol{\theta}$ for any fixed $\lambda$ is obtained as

$$GC_p(\boldsymbol{\theta}, \lambda \mid \alpha) = \alpha \hat{r}(\boldsymbol{Y}, \boldsymbol{1}_n \hat{\boldsymbol{\mu}}_\lambda' \boldsymbol{X}') + \sum_{i=1}^{k} \varphi_i(\boldsymbol{\theta}, \lambda \mid \alpha).$$

Since we need to minimize $\sum_{i=1}^{k} \varphi_i(\boldsymbol{\theta}, \lambda \mid \alpha)$ in order to optimize $\boldsymbol{\theta}$, we cannot obtain $\hat{\boldsymbol{\theta}}^{(G)}(\lambda \mid \alpha)$ that minimizes this $GC_p$ criterion for $\hat{\boldsymbol{\Xi}}_{\boldsymbol{\theta},\lambda}$ in closed form, even if $\lambda$ is fixed. Thus, we use an iterative computa-tional algorithm to optimize the parameters $\lambda$ and $\boldsymbol{\theta}$ simultaneously. This iterative computational algorithm for optimizing two parameters is difficult and requires a longer processing time than the optimization of a single parameter

# 8. References

[1] A. C. Atkinson, "A note on the generalized information criterion for choice of a model," Biometrika, vol. 67, no. 2, March 1980, pp. 413-418., pp. 291-293.

[2] P. J. Green and B. W. Silverman, "Nonparametric Regression and Generalized Linear Models," Chapman & Hall/CRC, 1994.

[3] D. A. Harville, "Matrix Algebra from a Statistician's Perspective," New York Springer, 1997.

[4] A. E. Hoerl and R. W. Kennard, "Ridge regression: biased estimation for nonorthogonal problems," Technometrics, vol. 12, No. 1, February 1970, pp. 55-67.

[5] A. M. Kshirsagar and W. B. Smith, "Growth Curves," Marcel Dekker, 1995.

[6] J. F. Lawless, "Mean squared error properties of generalized ridge regression," Journal of the American Statistical Association, vol. 76, no. 374, 1981, pp. 462-466.

[7] C. L. Mallows, "Some comments on Cp," Technometrics, vol. 15, no. 1, November 1973, pp. 661-675.

[8] C. L. Mallows, "More comments on Cp," Technometrics, vol. 37, no. 4, November 1995, pp. 362-372.

[9] I. Nagai, H. Yangihara and K. Satoh, "Optimization of Ridge Parameters in Multivariate Generalized Ridge Regression by Plug-in Methods," TR 10-03, Statistical Research Group, Hiroshima University, 2010.

[10] R. F. Potthoff and S. N. Roy, "A generalized multivariate analysis of variance model useful especially for growth curve problems," Biometrika, vol. 51, no. 3–4, December 1964, pp. 313-326.

[11] K. S. Riedel and K. Imre, "Smoothing spline growth curves with covariates," Communications in Statistics – Theory and Methods, vol. 22, no. 7, 1993, pp. 1795-1818.

[12] F. J. Richard, "A flexible growth function for empirical use," Journal of Experimental Botany, vol. 10, no. 2, 1959, pp. 290–301.

[13] K. Satoh and H. Yanagihara, "Estimation of varying coefficients for a growth curve model," American Journal of Mathematical and Management Sciences, 2010 (in press).

[14] M. Siotani, T. Hayakawa and Y. Fujikoshi, "Modern Multivariate Statistical Analysis: A Graduate Course and Handbook," American Sciences Press, Columbus, Ohio, 1985.

[15] R. S. Sparks, D. Coutsourides and L. Troskie, "The multivariate $C_p$," Communications in Statistics - Theory and Methods, vol. 12, no. 15, 1983, pp. 1775-1793.

[16] Y. Takane, K. Jung and H. Hwang, "Regularized reduced rank growth curve models," Computational Statistics and

Data Analysis, vol. 55, no. 2, February 2011, pp. 1041-1052.

[17] H. Yanagihara and K. Satoh, "An unbiased Cp criterion for multivariate ridge regression," Journal of Multivariate Analysis, vol. 101, no. 5, May 2010, pp. 1226-1238.

[18] H. Yanagihara, I. Nagai and K. Satoh, "A bias-corrected Cp criterion for optimizing ridge parameters in multivariate generalized ridge regression," Japanese Journal of Applied Statistics, vol. 38, no. 3, October 2009, pp. 151-172 (in Japanese).