

# An $\ell_{2,0}$ -norm constrained matrix optimization via extended discrete first-order algorithms

Ryoya Oda\* † Mineaki Ohishi‡  
Yuya Suzuki§ ¶ and Hirokazu Yanagihara||

(Last Modified: October 8, 2021)

## Abstract

The present paper is concerned with a constrained matrix optimization problem. The constraint is referred to as the  $\ell_{2,0}$ -norm of the matrix, which is defined as the number of non-zero row vectors of the matrix. We extend the discrete first-order algorithm by Bertsimas *et al.* (2016) to solve the optimization problem. The extended algorithm is useful for selecting variables in multivariate statistical models. Then, the convergence properties of the extended algorithm are established. In numerical experiments, we apply the extended algorithm to the optimization problem for the multivariate linear regression model. Furthermore, we also incorporate selecting variables using information criteria into the optimization problem.

## 1 Introduction

Optimization problems exist in multiple areas and are widely required. Among optimization problems, constrained matrix optimization problems are often used in estimating parameters with constraints for multivariate statistical models, such as the multivariate linear regression model [12, 13]. Consider the following constrained matrix optimization problem for a function  $f : \mathbb{R}^{k \times p} \rightarrow \mathbb{R}$ :

$$\min_{\Theta} f(\Theta) \text{ subject to } \|\Theta\|_{2,0} \leq q, \quad (1)$$

where  $\|\Theta\|_{2,0}$  is called the  $\ell_{2,0}$ -norm of a matrix  $\Theta \in \mathbb{R}^{k \times p}$  and is defined as

$$\|\Theta\|_{2,0} = \sum_{j=1}^k I(\theta_j \neq \mathbf{0}_p), \quad (2)$$

in which  $I(\cdot)$  denotes the indicator function,  $\theta_j$  is the  $j$ -th row vector of  $\Theta$ , i.e.,  $\Theta = (\theta_1, \dots, \theta_k)'$ , and  $\mathbf{0}_p$  is a  $p$ -dimensional vector of zeros. Note that the  $\ell_{2,0}$ -norm is not a norm in the usual sense because the  $\ell_{2,0}$ -norm lacks positive scalability:  $\|a\Theta\|_{2,0} = |a|\|\Theta\|_{2,0}$  for any  $a \in \mathbb{R}$ . For

---

\*Corresponding author. Email: ryoya-oda@hiroshima-u.ac.jp

†School of Informatics and Data Science, Hiroshima University, Higashi-Hiroshima, Japan

‡Education and Research Center for Artificial Intelligence and Data Innovation, Hiroshima University, Hiroshima, Japan

§Department of Mathematics, Graduate School of Science, Hiroshima University, Higashi-Hiroshima, Japan

¶Current address: Network Construction Department, DOCOMO CS SHIKOKU.Inc, Takamatsu, Japan.

||Graduate School of Advanced Science and Engineering, Hiroshima University, Higashi-Hiroshima, Japan

convenience, we refer to  $\|\cdot\|_{2,0}$  as the  $\ell_{2,0}$ -norm. By definition (2), the constraint in (1) explicitly restricts the number of non-zero row vectors of  $\Theta$ . From the viewpoint of statistical models, such a constraint is useful for selecting variables in estimating parameters because variables corresponding to zero-parameters can often be regarded as redundant in the model. Furthermore, since it may be desirable to find variables that affect all of the responses in multivariate statistical models (e.g., [9, 10, 14, 15]), it is important to consider the constraint in (1) because such variable selection requires the use of a vector constraint, rather than a scalar constraint.

Optimization problems with the  $\ell_{2,0}$ -norm constraint are non-convex optimization problems and are NP-hard because the  $\ell_{2,0}$ -norm is nonconvex and discontinuous. Hence, it is generally desired to obtain algorithms to achieve a sensible solution of (1) within a reasonable computation time. In multi-class classification with linear regression, Cai *et al.* [4] considered an efficient algorithm based on the general method of augmented Lagrange multipliers for solving (1) when the objective function is expressed as the  $\ell_{2,1}$ -norm of a matrix, where the  $\ell_{2,1}$ -norm is defined as  $\|\mathbf{A}\|_{2,1} = \sum_{j=1}^p (\sum_{i=1}^n a_{ij}^2)^{1/2}$  for a matrix  $\mathbf{A} = (a_{ij}) \in \mathbb{R}^{n \times p}$ . For general objective functions, Gotoh *et al.* [5] proposed a proximal algorithm for solving (1) when the objective function is represented as a difference of two convex functions. On the other hand, Bertsimas *et al.* [3] developed an algorithm for solving (1) only when  $p = 1$  with smooth convex objective functions. Note that when  $p = 1$ , the  $\ell_{2,0}$ -norm is equivalent to the  $\ell_0$ -norm of a vector, which is given by  $\|\mathbf{a}\|_0 = \sum_{j=1}^k I(a_j \neq 0)$  for a vector  $\mathbf{a} = (a_1, \dots, a_k)' \in \mathbb{R}^k$ . Their algorithm is referred to as the discrete first-order algorithm (DFA), and they derived the asymptotic convergence properties and the global convergence results of the DFA. Moreover, they confirmed that a mixed integer optimization initialized with a solution obtained from the DFA realized a near-optimal solution of (1) in numerical studies.

In the present paper, we consider the DFA proposed by Bertsimas *et al.* [3], and we extend the algorithm to solve (1) even when  $p \geq 2$ . Moreover, in the numerical experiments, we combine the extended DFA and information criteria to select variables for multivariate statistical models. In the framework of multivariate analysis, information criteria are often used to select variables. Examples of such information criteria include Akaike's information criterion (AIC) [1, 2] and the Bayesian information criterion (BIC) [11]. These criteria do not work or are not defined when the number of variables exceeds the sample size unless the candidate set of variables is narrow. However, it is expected that information criteria work because the candidate set becomes narrow by using solutions to problem (1) via the extended DFA.

The remainder of the present paper is organized as follows. In section 2, we present the optimization problem and extend the DFA. In section 3, we obtain several convergence properties of the extended DFA. In section 4, we conduct numerical experiments using the extended DFA and information criteria for selecting variables in the multivariate linear regression model. Technical details are provided in the Appendix.

## 2 Optimization problem and Algorithm

### 2.1 $\ell_{2,0}$ -norm constrained optimization problem

Suppose that the function  $g : \mathbb{R}^{k \times p} \rightarrow \mathbb{R}$  is bounded below, is convex, and has a Lipschitz continuous gradient with constant  $\ell$ , i.e., there exists a constant  $\ell > 0$  for all  $\Theta, \tilde{\Theta} \in \mathbb{R}^{k \times p}$ ,

$$\|D(\Theta) - D(\tilde{\Theta})\|_F \leq \ell \|\Theta - \tilde{\Theta}\|_F, \quad (3)$$

where  $D(\Theta) \in \mathbb{R}^{k \times p}$  is a matrix that is based on partial derivatives, i.e.,  $D(\Theta) = \partial g(\Theta) / \partial \Theta$ , and  $\|\cdot\|_F$  is the Frobenius norm of a matrix that is given as  $\|\Theta\|_F = \text{tr}(\Theta' \Theta)^{1/2}$  for a matrix  $\Theta$ . The function  $g$  defined in (3) was used by Bertsimas *et al.* [3]. Then, we consider the following  $\ell_{2,0}$ -norm constrained optimization problem for the function  $g$ :

$$\min_{\Theta} g(\Theta) \text{ subject to } \Theta = (\mathcal{B}', \Xi')', \|\Xi\|_{2,0} \leq q, \quad (4)$$

where  $\mathcal{B} \in \mathbb{R}^{k_1 \times p}$  and  $\Xi \in \mathbb{R}^{k_2 \times p}$  are the partitioned matrices of  $\Theta \in \mathbb{R}^{k \times p}$ , and  $k_1$  and  $k_2$  satisfy  $k_1 + k_2 = k$ . Problem (4) restricts the number of non-zero vectors of  $\Xi$ , but  $\mathcal{B}$  is optimized without constraints. Thus, (4) includes the following optimization problem:

$$\min_{\Theta} g(\Theta) \text{ subject to } \|\Theta\|_{2,0} \leq q, \quad (5)$$

because problem (4) can be regarded as (5) by letting  $k_1 = 0$  or  $k_2 = k$ . Problem (4) is useful for estimating parameters without constraints for a part of  $\Theta$  (e.g., the parameter corresponding to the intercept term) in multivariate statistical models. An example of (4) in the following multivariate statistical model is presented.

*Example* (Multivariate linear regression model). Suppose that  $\mathbf{Y} = (\mathbf{y}_1, \dots, \mathbf{y}_n)' \in \mathbb{R}^{n \times p}$  is an observation matrix stacking individual  $p$  response variables and  $\mathbf{X} = (\mathbf{x}_1, \dots, \mathbf{x}_n)' \in \mathbb{R}^{n \times (k-1)}$  is an observation matrix stacking individual  $k-1$  explanatory variables, where  $n$  is the sample size. We assume that the column vectors of  $\mathbf{X}$  have unit  $\ell_2$ -norm, i.e.,  $\|\mathbf{x}_j\|_2 = 1$  ( $j = 1, \dots, k-1$ ), where  $\|\cdot\|_2$  is the  $\ell_2$ -norm of vector, which is defined as  $\|\mathbf{a}\|_2 = (\mathbf{a}'\mathbf{a})^{1/2}$  for a vector  $\mathbf{a}$ . Moreover, we assume that the intercept term is included in this model. Hence, let  $\mathbf{Z} = (\mathbf{1}_n, \mathbf{X}) \in \mathbb{R}^{n \times k}$  be the matrix including the intercept term, where  $\mathbf{1}_n$  is an  $n$ -dimensional vector of ones. Then, the residual sum of squares is widely used in estimating parameters:

$$g(\Theta) = g_1(\Theta) = \frac{1}{2n} \|(\mathbf{Y} - \mathbf{Z}\Theta)\mathbf{G}^{-1/2}\|_F^2, \quad (6)$$

where  $\mathbf{G}$  is a positive definite matrix. Note that the intercept term does not vanish, because of non-centralizing of the column vectors of  $\mathbf{Y}$  and  $\mathbf{X}$ . When the constraint for the intercept term is not set, we can apply (6) to problem (4) when  $k_1 = 1$  and  $k_2 = k-1$ . Moreover, we observe that  $D(\Theta) = -n^{-1} \mathbf{Z}'(\mathbf{Y} - \mathbf{Z}\Theta)\mathbf{G}^{-1}$  and that one value of  $\ell$  is  $n^{-1} \lambda_{\max}(\mathbf{Z}'\mathbf{Z}) / \lambda_{\min}(\mathbf{G})$ , where  $\lambda_{\max}(\cdot)$  and  $\lambda_{\min}(\cdot)$  are the maximum and minimum eigenvalues, respectively, of a square matrix.

## 2.2 Extended discrete first-order algorithm

We extend the DFA proposed by Bertsimas *et al.* [3] to solve (4). First of all, for a given  $\mathbf{C} = (\mathbf{c}_1, \dots, \mathbf{c}_{k_2})' \in \mathbb{R}^{k_2 \times p}$ , we consider the following optimization problem:

$$\min_{\Xi} \|\Xi - \mathbf{C}\|_F^2 \text{ subject to } \|\Xi\|_{2,0} \leq q. \quad (7)$$

Let  $I_q(\mathbf{C})$  be the set consisting of suffixes of the  $q$  largest row vectors of  $\mathbf{C}$  in the sense of  $\ell_2$ -norm, i.e.,

$$I_q(\mathbf{C}) = \{1 \leq j \leq k_2 \mid \|\mathbf{c}_j\|_2 \text{ are among the largest } q \text{ of all } \|\mathbf{c}_j\|_2 \text{ s}\}. \quad (8)$$

Then, the optimal solutions of (7) can be derived in closed form, as is shown in the following proposition. (The proof is given in Appendix A.)

*Proposition 1.* Let  $\hat{\Xi} = (\hat{\xi}_1, \dots, \hat{\xi}_{k_2})'$  be an optimal solution to problem (7). Then,  $\hat{\Xi}$  is given by

$$\hat{\xi}_j = \begin{cases} \mathbf{c}_j & (j \in I_q(\mathbf{C})) \\ \mathbf{0}_p & (\text{otherwise}) \end{cases}, \quad (9)$$

where  $I_q(\mathbf{C})$  is defined in (8). We denote the set of optimal solutions (9) as  $\mathbf{H}_q(\mathbf{C})$ .

Note that  $\mathbf{H}_q(\mathbf{C})$  is expressed as the set of solutions because problem (7) may have some optimal solutions. The DFA is based on projected gradient decent methods in first-order convex optimization problems (see [7, 8]). The following proposition gives an upper bound of  $g$  and its minimizer with constraints. (The proof is given in Appendix B.)

*Proposition 2.* Let  $g$  be the function defined in (3). Then, for any  $L \geq \ell$ , we have

$$g(\tilde{\Theta}) \leq Q_L(\tilde{\Theta}, \Theta) = g(\Theta) + \frac{L}{2} \|\tilde{\Theta} - \Theta\|_F^2 + \text{tr} \left\{ \mathbf{D}'(\Theta)(\tilde{\Theta} - \Theta) \right\}, \quad (10)$$

for all  $\Theta = (\mathbf{B}', \Xi')' \in \mathbb{R}^{k \times p}$  and  $\tilde{\Theta} = (\tilde{\mathbf{B}}', \tilde{\Xi}')' \in \mathbb{R}^{k \times p}$  ( $\mathbf{B}, \tilde{\mathbf{B}} \in \mathbb{R}^{k_1 \times p}$ ;  $\Xi, \tilde{\Xi} \in \mathbb{R}^{k_2 \times p}$ ). Moreover, the optimal solution  $\Theta_{\dagger} = (\mathbf{B}'_{\dagger}, \Xi'_{\dagger})'$  ( $\mathbf{B}_{\dagger} \in \mathbb{R}^{k_1 \times p}$ ,  $\Xi_{\dagger} \in \mathbb{R}^{k_2 \times p}$ ) to  $\min_{\tilde{\Theta}: \|\tilde{\Xi}\|_{2,0} \leq q} Q_L(\tilde{\Theta}, \Theta)$  is given by

$$\mathbf{B}_{\dagger} = \mathbf{B} - L^{-1} \mathbf{D}_1(\Theta), \quad \Xi_{\dagger} \in \mathbf{H}_q(\Xi - L^{-1} \mathbf{D}_2(\Theta)), \quad (11)$$

where  $\mathbf{H}_q(\cdot)$  is defined in (9) and  $\mathbf{D}_1(\Theta) \in \mathbb{R}^{k_1 \times p}$  and  $\mathbf{D}_2(\Theta) \in \mathbb{R}^{k_2 \times p}$  are the partitioned matrices of  $\mathbf{D}(\Theta)$ , i.e.,  $\mathbf{D}(\Theta) = (\mathbf{D}'_1(\Theta), \mathbf{D}'_2(\Theta))'$ .

Using (11), we extend the DFA proposed by Bertsimas *et al.* [3] to solve (4), which is presented as Algorithm 1. We observe that Algorithm 1 for the parameter without constraints behaves like a vanilla gradient decent algorithm. Moreover, note that Algorithm 1 corresponds to the DFA proposed by Bertsimas *et al.* [3] when  $p = 1$  and  $k_1 = 0$ .

---

**Algorithm 1** Extended discrete first-order algorithm to solve (4)

---

**Require:** An initial value  $\Theta_1 = (\mathcal{B}'_1, \Xi'_1)'$  satisfying  $\|\Xi_1\|_{2,0} \leq q$ , a constant  $L (\geq \ell)$ , and a small value  $\varepsilon > 0$ .

$m = 1$ .

**repeat**

Obtain  $\Theta_{m+1} = (\mathcal{B}'_{m+1}, \Xi'_{m+1})'$  from (11) as follows:

$$\mathcal{B}_{m+1} = \mathcal{B}_m - L^{-1}D_1(\Theta_m), \quad \Xi_{m+1} \in \mathbf{H}_q(\Xi_m - L^{-1}D_2(\Theta_m)). \quad (12)$$

Increment  $m$  by 1.

**until**  $g(\Theta_m) - g(\Theta_{m+1}) < \varepsilon$  holds.

---

### 3 Convergence properties of Algorithm 1

We present several convergence properties of Algorithm 1. First, we define a notion of first-order optimality for problem (4).

*Definition 1.* For  $L \geq \ell$ ,  $\tilde{\Theta} = (\tilde{\mathcal{B}}', \tilde{\Xi}')'$  ( $\tilde{\mathcal{B}} \in \mathbb{R}^{k_1 \times p}$ ,  $\tilde{\Xi} \in \mathbb{R}^{k_2 \times p}$ ) is said to be an  $\ell_{2,0}$ -constrained first-order stationary point of problem (4) if  $\|\tilde{\Xi}\|_{2,0} \leq q$  holds and  $\tilde{\Theta}$  satisfies the following equation:

$$\tilde{\mathcal{B}} = \tilde{\mathcal{B}} - L^{-1}D_1(\tilde{\Theta}), \quad \tilde{\Xi} \in \mathbf{H}_q(\tilde{\Xi} - L^{-1}D_2(\tilde{\Theta})). \quad (13)$$

If  $\tilde{\Theta}$  is an  $\ell_{2,0}$ -constrained first-order stationary point, we have  $D_1(\tilde{\Theta}) = \mathbf{O}_{k_1,p}$ , where  $\mathbf{O}_{k_1,p} \in \mathbb{R}^{k_1 \times p}$  is the matrix, the elements of which are zero. Moreover, letting  $\tilde{\Xi} = (\tilde{\xi}_1, \dots, \tilde{\xi}_{k_2})'$ , it holds that  $\tilde{\xi}_j = \tilde{\xi}_j - L^{-1}d_j(\tilde{\Theta})$  for  $j \in I_q(\tilde{\Xi} - L^{-1}D_2(\tilde{\Theta}))$ , where  $d_j(\tilde{\Theta})$  is the  $j$ -th row vector of  $D_2(\tilde{\Theta})$ , i.e.,  $D_2(\tilde{\Theta}) = (d_1(\tilde{\Theta}), \dots, d_{k_2}(\tilde{\Theta}))'$ . Hence, we have  $d_j(\tilde{\Theta}) = \mathbf{0}_p$  for  $j \in I_q(\tilde{\Xi} - L^{-1}D_2(\tilde{\Theta}))$ . The following proposition is concerned with a sufficient condition for a global minimizer to the unconstrained optimization problem  $\min_{\Theta} g(\Theta)$ . (The proof is given in Appendix C.)

*Proposition 3.* If  $\tilde{\Theta} = (\tilde{\mathcal{B}}', \tilde{\Xi}')'$  satisfies (13) and  $\|\tilde{\Xi}\|_{2,0} < q$ , then we have  $\tilde{\Theta} \in \arg \min_{\Theta} g(\Theta)$ .

Next, we give several asymptotic convergence properties of Algorithm 1. To do so, we make several definitions for notational convenience. Let  $\Theta_m = (\mathcal{B}'_m, \Xi'_m)'$  be the  $m$ -iterated solution in (12) by Algorithm 1, and let  $\Xi_m = (\xi_{m,1}, \dots, \xi_{m,k_2})'$ . Moreover, let  $\mathbf{r}_m = (r_{m,1}, \dots, r_{m,k_2})'$  be the  $k_2$ -dimensional vector satisfying  $r_{m,j} = I(\xi_{m,j} \neq \mathbf{0}_p)$ . Denote as  $\alpha_{m,q} = \|\xi_{m,(q)}\|_2$  the  $\ell_2$ -norm of the  $q$ -th largest row vector of  $\Xi_m$  in the sense of  $\|\xi_{m,(1)}\|_2 \geq \dots \geq \|\xi_{m,(k_2)}\|_2$ . Using  $\alpha_{m,q}$ , we define  $\bar{\alpha}_q = \limsup_{m \rightarrow \infty} \alpha_{m,q}$  and  $\underline{\alpha}_q = \liminf_{m \rightarrow \infty} \alpha_{m,q}$ . Then, we present the several asymptotic convergence properties of Algorithm 1 as the following proposition. (The proof is given in Appendix D.)

*Proposition 4.* For problem (4), let  $\Theta_m$  be the  $m$ -iterated solution in (12) by Algorithm 1. Then, the following properties of Algorithm 1 hold:

(a) Let  $L \geq \ell$ . Then, we have

$$g(\Theta_m) - g(\Theta_{m+1}) \geq \frac{L - \ell}{2} \|\Theta_{m+1} - \Theta_m\|_F^2. \quad (14)$$

Moreover,  $g(\Theta_m)$  monotonically decreases for  $m$  and converges as  $m \rightarrow \infty$ .

- (b) For any  $L > \ell$ , it holds that  $\Theta_{m+1} - \Theta_m \rightarrow \mathbf{O}_{k,p}$  ( $m \rightarrow \infty$ ).
- (c) Let  $L > \ell$ , and  $\underline{\alpha}_q > 0$ . Then, there exists  $M > 0$  such that for all  $m \geq M$ ,  $\mathbf{r}_m = \mathbf{r}_{m+1}$ . Furthermore, the sequence  $\{\Theta_m\}$  converges to an  $\ell_{2,0}$ -constrained first-order stationary point.
- (d) Let  $L > \ell$ . Then, we have  $\lim_{m \rightarrow \infty} \mathbf{D}_1(\Theta_m) = \mathbf{O}_{k_1,p}$ . Furthermore, if  $\underline{\alpha}_q = 0$ , it holds that  $\liminf_{m \rightarrow \infty} \max_{j=1, \dots, k_2} \|\mathbf{d}_j(\Theta_m)\|_2 = 0$ .
- (e) Let  $L > \ell$ . If  $\bar{\alpha}_q = 0$  and the sequence  $\{\Theta_m\}$  has a limit point, then  $g(\Theta_m) \rightarrow \min_{\Theta} g(\Theta)$  ( $m \rightarrow \infty$ ).

The stopping rule of Algorithm 1 is based on (a) of Proposition 4. From (c), if the  $q$ -th largest vector  $\xi_{m,(q)}$  is non-zero for sufficiently large  $m$ , then the suffixes of the non-zero vectors of  $\Xi_m$  are fixed after that. Moreover, (c) ensures the global convergence to an  $\ell_{2,0}$ -constrained first-order stationary point of Algorithm 1. From (e), the objective function  $g$  converges to an optimal value for unconstrained optimization problem  $\min_{\Theta} g(\Theta)$  under minor assumptions.

Finally, we refer to the  $\ell_{2,0}$ -constrained first-order stationary point and a rate of convergence of Algorithm 1. The following proposition is concerned with some properties of the  $\ell_{2,0}$ -constrained first-order stationary point. (The proof is given in Appendix E.)

*Proposition 5.* For  $L \geq \ell$ , the following properties hold:

- (a) If  $\tilde{\Theta} = (\tilde{\mathbf{B}}', \tilde{\Xi}')'$  ( $\tilde{\mathbf{B}} \in \mathbb{R}^{k_1 \times p}$ ,  $\tilde{\Xi} \in \mathbb{R}^{k_2 \times p}$ ) is an  $\ell_{2,0}$ -constrained first-order stationary point in Definition 1, then  $\mathbf{H}_q(\tilde{\Xi} - L^{-1} \mathbf{D}_2(\tilde{\Theta}))$  has exactly one element.
- (b) Global minimizers of problem (4) are  $\ell_{2,0}$ -constrained first-order stationary points.

The following theorem presents knowledge about the rate of convergence of Algorithm 1. (The proof is given in Appendix F.)

*Theorem 1.* Let  $L \geq \ell$ . Then, Algorithm 1 iterated  $M$  times satisfies

$$\min_{m=1, \dots, M} \|\Theta_{m+1} - \Theta_m\|_F^2 \leq \frac{2\{g(\Theta_1) - g_*\}}{M(L - \ell)},$$

where  $g(\Theta_m) \downarrow g_*$  as  $m \rightarrow \infty$ .

The result of Theorem 1 is an extension of Theorem 3.1 of Bertsimas *et al.* [3] and coincides with it when  $p = 1$  and  $k_1 = 0$ .

## 4 Numerical Studies

We conduct numerical experiments based on Algorithm 1 for the  $\ell_{2,0}$ -norm constrained optimization problem (4) in terms of variable selection for the multivariate linear regression model (see Example). Denote the  $n \times p$  multivariate normal distribution with mean matrix  $\mathbf{A}$  and covariance matrix  $\mathbf{B}$  as  $N_{n \times p}(\mathbf{A}, \mathbf{B})$ . The explanatory matrix  $\mathbf{X}$ , the true parameter  $\beta_*$  corresponding to the intercept term, and  $\Xi_*$  were determined as follows:

$$\begin{aligned} \mathbf{X} &\sim N_{n \times k}(\mathbf{O}_{n,k}, \Psi \otimes \mathbf{I}_n), \quad \Theta_* = (\beta_*, \Xi_*', \mathbf{O}'_{k-k_*,p})', \\ \beta_* &\sim N_{p \times 1}(\mathbf{5}\mathbf{1}_p, \mathbf{I}_p \otimes \mathbf{1}), \quad \Xi_* \sim N_{k_* \times p}(\mathbf{5}\mathbf{1}_p \mathbf{1}'_p, \mathbf{I}_p \otimes \mathbf{I}_{k_*}), \end{aligned}$$

where  $\mathbf{I}_p \in \mathbb{R}^{p \times p}$  is the identity matrix, the  $(a, b)$ -th element of  $\Psi$  is  $(0.5)^{|a-b|}$ , and  $k_*$  is the number of non-zero row vectors of  $\Xi_*$ . Note that  $\mathbf{X}$ ,  $\beta_*$ , and  $\Xi_*$  were generated only once and were used throughout the simulation studies. Then, we made the column vectors of  $\mathbf{X}$  have unit  $\ell_2$ -norm. Using the notation in Example, the response matrix  $\mathbf{Y}$  was generated by  $\mathbf{Y} = (\mathbf{y}_{(1)}, \dots, \mathbf{y}_{(p)}) \sim N_{n \times p}(\mathbf{Z}\Theta_*, \Sigma \otimes \mathbf{I}_n)$ , where  $\mathbf{Z} = (\mathbf{1}_n, \mathbf{X})$  and  $\Sigma = 0.4\{(1 - 0.8)\mathbf{I}_p + 0.8\mathbf{1}_p\mathbf{1}_p'\}$ . Then, we made the column vectors of  $\mathbf{Y}$  have unit  $\ell_2$ -norm.

Since Algorithm 1 gives a solution for fixed  $q$ , we evaluate the best estimator by combining Algorithm 1 and information criteria, which are used to select variables. We performed the following steps:

- Step 1. Give a value  $\hat{\Theta}_{*,k_2} = (\hat{\mathbf{B}}'_{*,k_2}, \hat{\Xi}'_{*,k_2})'$  ( $\hat{\mathbf{B}}_{*,k_2} \in \mathbb{R}^{k_1 \times p}$ ,  $\hat{\Xi}_{*,k_2} \in \mathbb{R}^{k_2 \times p}$ ) satisfying  $\|\hat{\Theta}_{*,k_2}\|_{2,0} \leq k_2$  and set  $q = k_2 - 1$ .
- Step 2. Give a value  $\tilde{\Theta}_{*,q} = (\tilde{\mathbf{B}}'_{*,q}, \tilde{\Xi}'_{*,q})'$  ( $\tilde{\mathbf{B}}_{*,q} \in \mathbb{R}^{k_1 \times p}$ ,  $\tilde{\Xi}_{*,q} \in \mathbb{R}^{k_2 \times p}$ ) satisfying  $\|\tilde{\Theta}_{*,q}\|_{2,0} \leq q$ .
- Step 3. For the given  $q$ , obtain the solution  $\hat{\Theta}_{*,q} = (\hat{\mathbf{B}}'_{*,q}, \hat{\Xi}'_{*,q})'$  ( $\hat{\mathbf{B}}_{*,q} \in \mathbb{R}^{k_1 \times p}$ ,  $\hat{\Xi}_{*,q} \in \mathbb{R}^{k_2 \times p}$ ) by Algorithm 1 for the initial value  $\Theta_1 = \tilde{\Theta}_{*,q}$ . Then, decrement  $q$  by 1.
- Step 4. Repeat Steps 2 and 3 until  $q = 0$ .
- Step 5. Decide the best selection number as  $\hat{k}_* = \arg \min_{q=0, \dots, k_2} \text{IC}(\hat{\Theta}_{*,q})$  and obtain the best estimator by  $\hat{\Theta}_* = \hat{\Theta}_{*,\hat{k}_*}$ , where  $\text{IC}(\cdot)$  is an information criterion.

In the above steps and Algorithm 1,  $k_1 = 1$ ,  $\hat{\Theta}_{*,k_2} = (\mathbf{Z}'\mathbf{Z})^{-1}\mathbf{Z}'\mathbf{Y}$  and  $\varepsilon = 10^{-4}$ , and the value  $\tilde{\Theta}_{*,q}$  in Step 2 is given as follows:

- Step 2-1. Denote  $\mathcal{A}_q = \{a_1, \dots, a_{q+1}\}$  ( $a_1 < \dots < a_{q+1}$ ) as the active set of  $\hat{\Xi}_{*,q+1}$  defined by  $\{k_1 + 1 \leq j \leq k \mid \hat{\theta}_{*,q+1,j} \neq \mathbf{0}_p\}$ , where  $\hat{\theta}_{*,q+1,j}$  is the  $j$ -th row vector of  $\hat{\Theta}_{*,q+1}$  in Step 3.

- Step 2-2. Set  $\bar{\mathcal{A}}_q = \{1, \dots, k_1\} \cup \mathcal{A}_q$  and

$$(\bar{\mathbf{B}}'_q, \bar{\Xi}'_q)' = \left( \mathbf{Z}'_{\bar{\mathcal{A}}_q} \mathbf{Z}_{\bar{\mathcal{A}}_q} \right)^{-1} \mathbf{Z}'_{\bar{\mathcal{A}}_q} \mathbf{Y} \quad (\bar{\mathbf{B}}_q \in \mathbb{R}^{k_1 \times p}, \bar{\Xi}_q \in \mathbb{R}^{|\mathcal{A}_q| \times p}),$$

where  $\mathbf{Z}_{\bar{\mathcal{A}}_q}$  is the  $n \times |\bar{\mathcal{A}}_q|$  matrix consisting of columns of  $\mathbf{Z}$  indexed by the elements of  $\bar{\mathcal{A}}_q$ . Furthermore, denote the  $j$ -th row vectors of  $\bar{\mathbf{B}}_q$  and  $\bar{\Xi}_q$  as  $\bar{\beta}_{q,j}$  and  $\bar{\xi}_{q,a_j}$ , respectively.

- Step 2-3. Give the  $j$ -th row vector of  $\tilde{\Theta}_{*,q}$  used in Step 2 as follows:

$$\begin{cases} \bar{\beta}_{q,j} & (1 \leq j \leq k_1) \\ \bar{\xi}_{q,a_j} & ((j \in \mathcal{A}_q) \wedge (j \neq a_{q,\min})) \\ \mathbf{0}_p & ((j \in \{k_1 + 1, \dots, k\} \cap \mathcal{A}_q^c) \vee (j = a_{q,\min})) \end{cases},$$

where  $a_{q,\min} = \arg \min_{j \in \mathcal{A}_q} \|\bar{\xi}_{q,a_j}\|_2$ .

Furthermore, the BIC proposed by Schwarz [11] was used as an information criterion and is defined by

$$\text{IC}(\Theta) = n \log |n^{-1}(\mathbf{Y} - \mathbf{Z}\Theta)'(\mathbf{Y} - \mathbf{Z}\Theta)| + p\|\Theta\|_{2,0} \log n.$$

For problem (4), we set  $g = g_1$ ,  $\mathbf{G} = (n-k)^{-1}\mathbf{Y}'\{\mathbf{I}_n - \mathbf{Z}(\mathbf{Z}'\mathbf{Z})^{-1}\mathbf{Z}'\}\mathbf{Y}$  and  $L = \lceil n^{-1}\lambda_{\max}(\mathbf{Z}'\mathbf{Z})/\lambda_{\min}(\mathbf{G}) \rceil$ , where  $\lceil \cdot \rceil$  is the ceiling function. For these settings, the above steps were carried out for 1,000 simulation iterations.

In these numerical studies, we examine the following properties.

- The two relative mean square errors (RMSEs):

$$\text{RMSE}_{\Theta_*} = \frac{E[\|\Theta_* - \hat{\Theta}_*\|_F^2]}{E[\|\Theta_* - \hat{\Theta}_{*,k_2}\|_F^2]} \times 100 (\%),$$

$$\text{RMSE}_{\mathbf{Y}} = \frac{E[\|\mathbf{Y} - \mathbf{Z}\hat{\Theta}_*\|_F^2]}{E[\|\mathbf{Y} - \mathbf{Z}\hat{\Theta}_{*,k_2}\|_F^2]} \times 100 (\%).$$

In our numerical settings,  $E[\|\Theta_* - \hat{\Theta}_{*,k_2}\|_F^2] = \text{tr}\{(\mathbf{Z}'\mathbf{Z})^{-1}\}\text{tr}(\Sigma)$  and  $E[\|\mathbf{Y} - \mathbf{Z}\hat{\Theta}_{*,k_2}\|_F^2] = (n-k)\text{tr}(\Sigma)$ . These RMSEs are approximated by the average value of 1,000 simulation iterations. Note that the smaller the  $\text{RMSE}_{\Theta_*}$  and  $\text{RMSE}_{\mathbf{Y}}$ , the better the accuracy of the estimation of  $\Theta_*$  and the prediction accuracy of  $\mathbf{Y}$ , respectively.

- The probability (%) such that the suffixes of the non-zero vectors of  $\hat{\Theta}_*$  and  $\Theta_*$  are equivalent among 1,000 simulation iterations.
- The CPU time (s) obtained as the average value of 1,000 simulation iterations.

Table 1. Properties of the estimation results for  $\Theta_*$  by the combination of the extended DFA and the BIC in the multivariate linear regression model.

$n$	$k$	$\text{RMSE}_{\Theta_*}$	$\text{RMSE}_{\mathbf{Y}}$	Probability	CPU time
100	20	47.01	14.69	95.2	0.848
100	40	17.26	20.83	90.1	1.187
100	60	8.61	31.88	85.1	1.477
100	80	4.30	77.50	76.9	2.411
300	20	51.28	4.02	99.0	0.125
300	40	23.36	4.62	97.8	0.150
300	60	14.43	4.69	96.7	0.214
300	120	5.43	6.67	95.4	0.422
300	180	2.42	10.00	93.7	0.872
300	240	1.08	22.50	88.9	2.167
500	20	52.73	2.34	98.8	0.032
500	40	23.99	2.45	99.2	0.059
500	100	8.59	2.81	98.4	0.220
500	200	3.17	4.17	96.8	0.544
500	300	1.39	5.63	96.1	1.331
500	400	0.67	13.75	88.7	2.423

Table 1 shows the above properties when we set  $p = 3$  and  $k_* = 10$ . From Table 1, we observe that both the RMSEs ( $\text{RMSE}_{\Theta_*}$  and  $\text{RMSE}_{\mathbf{Y}}$ ) are smaller than 100. This means that



the estimator  $\hat{\Theta}_*$  by the combination of the extended DFA and the BIC is better than the least squares estimator  $\hat{\Theta}_{*,k_2}$  in terms of the accuracy of the estimation of  $\Theta_*$  and the prediction accuracy of  $\mathbf{Y}$ . Moreover, the probabilities are high and the CPU times are short. Therefore, we can confirm that the combination of the extended DFA and the BIC is valid for the estimation of  $\Theta_*$ .

## Appendix

### A Proof of Proposition 1

Since the Frobenius norm is invariant to the exchange of rows, without loss of generality, the given  $\mathbf{C} = (\mathbf{c}_1, \dots, \mathbf{c}_{k_2})'$  in (7) can be regarded as  $\|\mathbf{c}_1\|_2 \geq \dots \geq \|\mathbf{c}_{k_2}\|_2$ . Problem (7) can be rewritten as

$$\min_{\|\Xi\|_{2,0} \leq q} \sum_{j=1}^{k_2} \|\xi_j - \mathbf{c}_j\|_2^2.$$

From the above expression, we can see that the optimal solution for (7) is limited to the case in which  $q$  row vectors of  $\Xi$  become  $\mathbf{c}_j$  and the  $(k_2 - q)$  remainder becomes  $\mathbf{0}_p$ . On the other hand, let  $\mathcal{S} = \{1 \leq j \leq k_2 \mid \xi_j \neq \mathbf{0}_p\}$ . Then, we have

$$\min_{\|\Xi\|_{2,0} \leq q} \sum_{j=1}^{k_2} \|\xi_j - \mathbf{c}_j\|_2^2 = \min_{\|\Xi\|_{2,0} \leq q} \left\{ \sum_{j \in \mathcal{S}} \|\xi_j - \mathbf{c}_j\|_2^2 + \sum_{j \notin \mathcal{S}} \|\mathbf{c}_j\|_2^2 \right\}.$$

Since the above problem is optimal when  $\xi_j = \mathbf{c}_j$  for  $j \in \mathcal{S}$  and  $\sum_{j \notin \mathcal{S}} \|\mathbf{c}_j\|_2^2$  is minimum, we can see that  $\mathcal{S} = \{1, \dots, q\}$ .  $\square$

### B Proof of Proposition 2

First, we show (10). Let  $\boldsymbol{\theta} = \text{vec}(\Theta')$  and  $\tilde{\boldsymbol{\theta}} = \text{vec}(\tilde{\Theta}')$  for any  $\Theta, \tilde{\Theta} \in \mathbb{R}^{k \times p}$ . Then, by rewriting  $g(\tilde{\Theta})$  as  $h(\tilde{\boldsymbol{\theta}})$ , the following inequality can be derived (see, e.g., [7]):

$$h(\tilde{\boldsymbol{\theta}}) \leq h(\boldsymbol{\theta}) + \frac{L}{2} \|\tilde{\boldsymbol{\theta}} - \boldsymbol{\theta}\|_2^2 + (\partial h(\boldsymbol{\theta}) / \partial \boldsymbol{\theta})' (\tilde{\boldsymbol{\theta}} - \boldsymbol{\theta}).$$

From the properties of the vec operator and the Frobenius norm, the above inequality can be expressed as

$$g(\tilde{\Theta}) \leq g(\Theta) + \frac{L}{2} \|\tilde{\Theta} - \Theta\|_F^2 + \text{tr} \left\{ \mathbf{D}'(\Theta) (\tilde{\Theta} - \Theta) \right\}.$$

Next, we show (11). The following equation can be derived:

$$\begin{aligned} Q_L(\tilde{\Theta}, \Theta) &= \frac{L}{2} \left\| \tilde{\Theta} - (\Theta - L^{-1} \mathbf{D}(\Theta)) \right\|_F^2 - \frac{1}{2L} \|\mathbf{D}(\Theta)\|_F^2 + g(\Theta) \\ &= \frac{L}{2} \left\| \tilde{\mathcal{B}} - (\mathcal{B} - L^{-1} \mathbf{D}_1(\Theta)) \right\|_F^2 + \frac{L}{2} \left\| \tilde{\Xi} - (\Xi - L^{-1} \mathbf{D}_2(\Theta)) \right\|_F^2 \\ &\quad - \frac{1}{2L} \|\mathbf{D}(\Theta)\|_F^2 + g(\Theta). \end{aligned}$$

Hence, the optimal solution to  $\min_{\tilde{\Theta}: \|\tilde{\Xi}\|_{2,0} \leq q} Q_L(\tilde{\Theta}, \Theta)$ , is derived as follows:

$$\begin{aligned} \min_{\tilde{\Theta}: \|\tilde{\Xi}\|_{2,0} \leq q} Q_L(\tilde{\Theta}, \Theta) &= \frac{L}{2} \min_{\tilde{\mathcal{B}}} \left\| \tilde{\mathcal{B}} - (\mathcal{B} - L^{-1} \mathbf{D}_1(\Theta)) \right\|_F^2 \\ &\quad + \frac{L}{2} \min_{\|\tilde{\Xi}\|_{2,0} \leq q} \left\| \tilde{\Xi} - (\Xi - L^{-1} \mathbf{D}_2(\Theta)) \right\|_F^2 \\ &\quad - \frac{1}{2L} \|\mathbf{D}(\Theta)\|_F^2 + g(\Theta). \end{aligned} \tag{B.1}$$

This completes the proof of (11).  $\square$

### C Proof of Proposition 3

Let  $\mathcal{M} = (\boldsymbol{\mu}_1, \dots, \boldsymbol{\mu}_{k_2})'$ , satisfying  $\boldsymbol{\mu}_j = \tilde{\boldsymbol{\xi}}_j - L^{-1} \mathbf{d}_j(\tilde{\Theta})$  for  $j = 1, \dots, k_2$ . First, we show that  $\mathbf{d}_j(\tilde{\Theta}) = \mathbf{0}_p$  for all  $j \notin I_q(\mathcal{M})$ . From (13), it is straightforward to observe that  $\boldsymbol{\mu}_i = \tilde{\boldsymbol{\xi}}_i$  for  $i \in I_q(\mathcal{M})$ . On the other hand, since  $\tilde{\boldsymbol{\xi}}_j = \mathbf{0}_p$  for  $j \notin I_q(\mathcal{M})$ , we have  $\boldsymbol{\mu}_j = L^{-1} \mathbf{d}_j(\tilde{\Theta})$  for  $j \notin I_q(\mathcal{M})$ . These imply that  $\|\tilde{\boldsymbol{\xi}}_i\|_2 \geq \|L^{-1} \mathbf{d}_j(\tilde{\Theta})\|_2$  for any  $i \in I_q(\mathcal{M})$  and  $j \notin I_q(\mathcal{M})$  because  $\|\boldsymbol{\mu}_i\|_2 \geq \|\boldsymbol{\mu}_j\|_2$ . Note that it follows from  $\|\tilde{\Xi}\|_{2,0} < q$  that  $\min_{i \in I_q(\mathcal{M})} \|\boldsymbol{\mu}_i\|_2 = \min_{i \in I_q(\mathcal{M})} \|\tilde{\boldsymbol{\xi}}_i\|_2 = 0$ . Hence, we have  $\mathbf{d}_j(\tilde{\Theta}) = \mathbf{0}_p$  for all  $j \notin I_q(\mathcal{M})$ .

Next, we show that  $\mathbf{d}_i(\tilde{\Theta}) = \mathbf{0}_p$  for all  $i \in I_q(\mathcal{M})$ . From (13), it is straightforward to observe that  $\tilde{\boldsymbol{\xi}}_i = \tilde{\boldsymbol{\xi}}_i - L^{-1} \mathbf{d}_i(\tilde{\Theta})$  for  $i \in I_q(\mathcal{M})$ . Hence, we have  $\mathbf{d}_i(\tilde{\Theta}) = \mathbf{0}_p$  for all  $i \in I_q(\mathcal{M})$ . Therefore, it holds that  $\mathbf{D}_2(\tilde{\Theta}) = \mathbf{O}_{k_2, p}$ . Moreover, since it is straightforward to observe that  $\mathbf{D}_1(\tilde{\Theta}) = \mathbf{O}_{k_1, p}$ , we have  $\mathbf{D}(\tilde{\Theta}) = \mathbf{O}_{k, p}$ . This fact and the convexity of  $g$  lead to  $\tilde{\Theta} \in \arg \min_{\Theta} g(\Theta)$ .  $\square$

### D Proof of Proposition 4

#### D.1 Proof of (a)

Let  $\Theta_{\dagger} \in \arg \min_{\tilde{\Theta}: \|\tilde{\Xi}\|_{2,0} \leq q} Q_L(\tilde{\Theta}, \Theta)$ , where  $Q_L(\tilde{\Theta}, \Theta)$  is as defined in (10). Then, for  $\Theta = (\mathcal{B}', \Xi')'$  ( $\mathcal{B} \in \mathbb{R}^{k_1 \times p}$ ,  $\Xi \in \mathbb{R}^{k_2 \times p}$ ) such that  $\|\Xi\|_{2,0} \leq q$ , we have

$$\begin{aligned} g(\Theta) &= Q_L(\Theta, \Theta) \\ &\geq \inf_{\tilde{\Theta}: \|\tilde{\Xi}\|_{2,0} \leq q} Q_L(\tilde{\Theta}, \Theta) \\ &= Q_L(\Theta_{\dagger}, \Theta) \\ &= g(\Theta) + \frac{L}{2} \|\Theta_{\dagger} - \Theta\|_F^2 + \text{tr} \{ \mathbf{D}'(\Theta)(\Theta_{\dagger} - \Theta) \} \\ &= g(\Theta) + \frac{L-\ell}{2} \|\Theta_{\dagger} - \Theta\|_F^2 + \frac{\ell}{2} \|\Theta_{\dagger} - \Theta\|_F^2 + \text{tr} \{ \mathbf{D}'(\Theta)(\Theta_{\dagger} - \Theta) \} \\ &= \frac{L-\ell}{2} \|\Theta_{\dagger} - \Theta\|_F^2 + Q_{\ell}(\Theta_{\dagger}, \Theta) \\ &\geq \frac{L-\ell}{2} \|\Theta_{\dagger} - \Theta\|_F^2 + g(\Theta_{\dagger}). \end{aligned} \tag{D.1}$$

Since we can regard  $\Theta_{\dagger}$  as  $\Theta_{m+1}$  by letting  $\Theta = \Theta_m$ , (D.1) is expressed as

$$g(\Theta_m) - g(\Theta_{m+1}) \geq \frac{L-\ell}{2} \|\Theta_{m+1} - \Theta_m\|_F^2.$$

Hence,  $g(\Theta_m)$  monotonically decreases for  $m$ . Moreover, it is straightforward to observe that  $g(\Theta_m)$  converges as  $m \rightarrow \infty$  because it is bounded below.  $\square$

## D.2 Proof of (b)

Since  $g(\Theta_m)$  converges as  $m \rightarrow \infty$  from (a) in Proposition 4,  $g(\Theta_m) - g(\Theta_{m+1})$  converges to 0. Hence,  $\|\Theta_{m+1} - \Theta_m\|_F^2$  in (14) also converges to 0 for  $L > \ell$ . This implies that  $\Theta_{m+1} - \Theta_m \rightarrow \mathbf{O}_{k,p}$  ( $m \rightarrow \infty$ ).  $\square$

## D.3 Proof of (c)

First, we prove that there exists  $M > 0$  such that for all  $m \geq M$ ,  $\mathbf{r}_m = \mathbf{r}_{m+1}$  by contradiction. Assume that for any  $M > 0$ , there exists  $\tilde{m} \geq M$  such that  $\mathbf{r}_{\tilde{m}} \neq \mathbf{r}_{\tilde{m}+1}$ . Since  $\underline{\alpha}_q > 0$ , we observe that  $\|\Xi_m\|_{2,0} = q$  for sufficiently large  $m$ . Hence, by considering  $M > m$ , we can see that there exist  $i, j$  ( $i \neq j$ ) such that

$$\xi_{\tilde{m},i} = \mathbf{0}_p, \quad \xi_{\tilde{m},j} \neq \mathbf{0}_p, \quad \xi_{\tilde{m}+1,i} \neq \mathbf{0}_p, \quad \xi_{\tilde{m}+1,j} = \mathbf{0}_p,$$

for infinitely many  $\tilde{m} \geq M$ . Using the above equations, the following inequality can be derived:

$$\|\Xi_{\tilde{m}} - \Xi_{\tilde{m}+1}\|_F \geq \sqrt{\|\xi_{\tilde{m}+1,i}\|_2^2 + \|\xi_{\tilde{m},j}\|_2^2} \geq \frac{\|\xi_{\tilde{m}+1,i}\|_2 + \|\xi_{\tilde{m},j}\|_2}{\sqrt{2}}.$$

From the above, we observe that the  $\ell_2$ -norms of the non-zero vectors  $\|\xi_{\tilde{m}+1,i}\|_2$  and  $\|\xi_{\tilde{m},j}\|_2$  converge to 0 as  $\tilde{m} \rightarrow \infty$  because  $\|\Theta_{\tilde{m}} - \Theta_{\tilde{m}+1}\|_F \rightarrow 0$  from (b) in Proposition 4. This contradicts  $\underline{\alpha}_q > 0$ .

Next, we show that the sequence  $\{\Theta_m\}$  converges to an  $\ell_{2,0}$ -constrained first-order stationary point. Since  $\mathbf{r}_m = \mathbf{r}_{m+1}$  for sufficiently large  $m$ , we can set  $\mathcal{L} = \{1 \leq j \leq k_2 \mid r_{m,j} = 1\}$ . Note that the elements in  $\mathcal{L}$  are invariant for sufficiently large  $m$ . Hence, using (B.1), for sufficiently large  $m$  we have

$$\begin{aligned} & \min_{\tilde{\Theta}: \|\tilde{\Theta}\|_{2,0} \leq q} Q_L(\tilde{\Theta}, \Theta_m) \\ &= \frac{L}{2} \min_{\tilde{\Theta}: \|\tilde{\Theta}\|_{2,0} \leq q} \left\| \tilde{\Theta} - (\Xi_m - L^{-1}D_2(\Theta_m)) \right\|_F^2 - \frac{1}{2L} \|D(\Theta_m)\|_F^2 + g(\Theta_m) \\ &= \frac{L}{2} \min_{\tilde{\xi}_j: j \in \mathcal{L}} \sum_{j \in \mathcal{L}} \left\| \tilde{\xi}_j - (\xi_{m,j} - L^{-1}d_j(\Theta_m)) \right\|_2^2 + \frac{L}{2} \sum_{j \notin \mathcal{L}} \|\xi_{m,j} - L^{-1}d_j(\Theta_m)\|_2^2 \\ & \quad - \frac{1}{2L} \|D(\Theta_m)\|_F^2 + g(\Theta_m). \end{aligned}$$

This implies that Algorithm 1 behaves like a vanilla gradient decent algorithm for minimizing a convex function over a closed convex. Therefore, the sequence  $\{\Theta_m\}$  converges to an  $\ell_{2,0}$ -constrained first-order stationary point.  $\square$

## D.4 Proof of (d)

From (12) and (b), it is straightforward to observe that  $\lim_{m \rightarrow \infty} D_1(\Theta_m) = \mathbf{O}_{k_1,p}$ . Assume that  $\underline{\alpha}_q = 0$ . Let  $\mathcal{M}_m = (\mu_{m,1}, \dots, \mu_{m,k_2})' = \Xi_m - L^{-1}D_2(\Theta_m)$ . From (12), for any  $i \in I_q(\mathcal{M}_m)$  and  $j \notin I_q(\mathcal{M}_m)$ , we have  $\|\mu_{m,i}\|_2 \geq \|\mu_{m,j}\|_2$ . Hence, the following inequality can be derived:

$$\liminf_{m \rightarrow \infty} \min_{i \in I_q(\mathcal{M}_m)} \|\mu_{m,i}\|_2 \geq \liminf_{m \rightarrow \infty} \max_{j \notin I_q(\mathcal{M}_m)} \|\mu_{m,j}\|_2. \quad (\text{D.2})$$

On the other hand, from (12) and (b), it is straightforward to observe that

$$\boldsymbol{\xi}_{m,i} - \boldsymbol{\xi}_{m+1,i} = \begin{cases} L^{-1} \mathbf{d}_i(\boldsymbol{\Theta}_m) & \rightarrow \mathbf{0}_p \quad (m \rightarrow \infty) \quad (i \in I_q(\mathcal{M}_m)) \\ \boldsymbol{\xi}_{m,i} & \rightarrow \mathbf{0}_p \quad (m \rightarrow \infty) \quad (i \notin I_q(\mathcal{M}_m)) \end{cases}. \quad (\text{D.3})$$

Using (D.3) and the triangle inequality, we have

$$\begin{aligned} & \liminf_{m \rightarrow \infty} \min_{i \in I_q(\mathcal{M}_m)} \|\boldsymbol{\mu}_{m,i}\|_2 \\ & \leq \liminf_{m \rightarrow \infty} \min_{i \in I_q(\mathcal{M}_m)} \|\boldsymbol{\xi}_{m,i}\|_2 + \liminf_{m \rightarrow \infty} \max_{i \in I_q(\mathcal{M}_m)} \|L^{-1} \mathbf{d}_i(\boldsymbol{\Theta}_m)\|_2 \\ & = \liminf_{m \rightarrow \infty} \min_{i \in I_q(\mathcal{M}_m)} \|\boldsymbol{\xi}_{m,i}\|_2, \end{aligned} \quad (\text{D.4})$$

$$\begin{aligned} & \liminf_{m \rightarrow \infty} \max_{j \notin I_q(\mathcal{M}_m)} \|\boldsymbol{\mu}_{m,j}\|_2 \\ & \geq \limsup_{m \rightarrow \infty} \min_{i \notin I_q(\mathcal{M}_m)} \|\boldsymbol{\xi}_{m,i}\|_2 + \liminf_{m \rightarrow \infty} \max_{i \notin I_q(\mathcal{M}_m)} \|L^{-1} \mathbf{d}_i(\boldsymbol{\Theta}_m)\|_2 \\ & = \liminf_{m \rightarrow \infty} \max_{i \notin I_q(\mathcal{M}_m)} \|L^{-1} \mathbf{d}_i(\boldsymbol{\Theta}_m)\|_2. \end{aligned} \quad (\text{D.5})$$

By combining (D.2), (D.4), and (D.5), the following inequality can be derived:

$$\liminf_{m \rightarrow \infty} \min_{i \in I_q(\mathcal{M}_m)} \|\boldsymbol{\xi}_{m,i}\|_2 \geq \liminf_{m \rightarrow \infty} \max_{j \notin I_q(\mathcal{M}_m)} \|L^{-1} \mathbf{d}_i(\boldsymbol{\Theta}_m)\|_2. \quad (\text{D.6})$$

Since  $\underline{\alpha}_q = \liminf_{m \rightarrow \infty} \min_{i \in I_q(\mathcal{M}_m)} \|\boldsymbol{\xi}_{m,i}\|_2 = 0$ , the right-hand side of (D.6) becomes 0. This fact and (D.3) imply that  $\liminf_{m \rightarrow \infty} \max_{j=1, \dots, k_2} \|\mathbf{d}_j(\boldsymbol{\Theta}_m)\|_2 = 0$ .  $\square$

## D.5 Proof of (e)

By modifying  $\liminf$  to  $\limsup$  in (D.6), we can derive the following inequality:

$$\limsup_{m \rightarrow \infty} \min_{i \in I_q(\mathcal{M}_m)} \|\boldsymbol{\xi}_{m,i}\|_2 \geq \limsup_{m \rightarrow \infty} \max_{j \notin I_q(\mathcal{M}_m)} \|L^{-1} \mathbf{d}_i(\boldsymbol{\Theta}_m)\|_2, \quad (\text{D.7})$$

where  $\mathcal{M}_m = (\boldsymbol{\mu}_{m,1}, \dots, \boldsymbol{\mu}_{m,k_2})' = \boldsymbol{\Xi}_m - L^{-1} \mathbf{D}_2(\boldsymbol{\Theta}_m)$ . Since  $\bar{\alpha}_q = 0$ , the right-hand side of (D.7) becomes 0. Since this fact and (D.3) imply that  $\mathbf{D}_2(\boldsymbol{\Theta}_m) \rightarrow \mathbf{O}_{k_2,p}$  ( $m \rightarrow \infty$ ), we have  $\mathbf{D}(\boldsymbol{\Theta}_m) \rightarrow \mathbf{O}_{k,p}$  ( $m \rightarrow \infty$ ) from (d) in Proposition 4. Let  $\boldsymbol{\Theta}_\infty$  be a limit point of the sequence  $\{\boldsymbol{\Theta}_m\}$ . Then, there exists a subsequence  $\{\tilde{m}\}$  such that  $\boldsymbol{\Theta}_{\tilde{m}} \rightarrow \boldsymbol{\Theta}_\infty$  and  $g(\boldsymbol{\Theta}_{\tilde{m}}) \rightarrow g(\boldsymbol{\Theta}_\infty)$ . Since  $\mathbf{D}(\boldsymbol{\Theta})$  is Lipschitz continuous, it holds that  $\mathbf{D}(\boldsymbol{\Theta}_{\tilde{m}}) \rightarrow \mathbf{D}(\boldsymbol{\Theta}_\infty) = \mathbf{O}_{k,p}$  as  $\tilde{m} \rightarrow \infty$ . This implies that  $\boldsymbol{\Theta}_\infty$  is a solution to  $\min_{\boldsymbol{\Theta}} g(\boldsymbol{\Theta})$ . Therefore,  $g(\boldsymbol{\Theta}_m) \rightarrow \min_{\boldsymbol{\Theta}} g(\boldsymbol{\Theta})$  ( $m \rightarrow \infty$ ) holds.  $\square$

## E Proof of Proposition 5

### E.1 Proof of (a)

Suppose that  $\boldsymbol{\Theta}_\dagger = (\boldsymbol{\mathcal{B}}_\dagger', \boldsymbol{\Xi}_\dagger')$  satisfies  $\|\boldsymbol{\Xi}_\dagger\|_{2,0} \leq q$  and the following equation:

$$\boldsymbol{\mathcal{B}}_\dagger = \tilde{\boldsymbol{\mathcal{B}}} - L^{-1} \mathbf{D}_1(\tilde{\boldsymbol{\Theta}}), \quad \boldsymbol{\Xi}_\dagger \in \mathbf{H}_q \left( \tilde{\boldsymbol{\Xi}} - L^{-1} \mathbf{D}_2(\tilde{\boldsymbol{\Theta}}) \right).$$

Using (D.1), we can derive the following inequality:

$$g(\tilde{\Theta}) - g(\Theta_{\dagger}) \geq \frac{L-\ell}{2} \|\Theta_{\dagger} - \tilde{\Theta}\|_F^2. \quad (\text{E.1})$$

On the other hand, since the function  $g$  is convex and differentiable, we have

$$g(\Theta_{\dagger}) - g(\tilde{\Theta}) \geq \text{tr} \left\{ \mathbf{D}'(\tilde{\Theta})(\Theta_{\dagger} - \tilde{\Theta}) \right\}. \quad (\text{E.2})$$

Since  $\Theta_{\dagger}$  and  $\tilde{\Theta}$  are  $\ell_{2,0}$ -constrained first-order stationary points, we observe that  $\mathbf{D}_1(\tilde{\Theta}) = \mathbf{O}_{k_1,p}$ ,  $\mathbf{d}_j(\tilde{\Theta}) = \mathbf{0}_p$  for  $j \in I_q(\tilde{\Xi} - L^{-1}\mathbf{D}_2(\tilde{\Theta}))$  and  $\xi_{\dagger,j} = \tilde{\xi}_j = \mathbf{0}_p$  for  $j \notin I_q(\tilde{\Xi} - L^{-1}\mathbf{D}_2(\tilde{\Theta}))$ , where  $\xi_{\dagger,j}$  is the  $j$ -th row vector of  $\Xi_{\dagger}$ . Hence, we have  $\text{tr}\{\mathbf{D}'(\tilde{\Theta})(\Theta_{\dagger} - \tilde{\Theta})\} = 0$ . Therefore, it follows from (E.1) and (E.2) that  $\|\Theta_{\dagger} - \tilde{\Theta}\|_F^2 = 0$  holds for  $L > \ell$ . This implies that  $\mathbf{H}_q(\tilde{\Xi} - L^{-1}\mathbf{D}_2(\tilde{\Theta}))$  has exactly one element.  $\square$

## E.2 Proof of (b)

For any global minimizer  $\hat{\Theta}_{\text{glo}} = (\hat{\mathbf{B}}'_{\text{glo}}, \hat{\Xi}'_{\text{glo}})'$  of (4), let  $\tilde{\Theta} = (\tilde{\mathbf{B}}', \tilde{\Xi}')'$  be the matrix satisfying  $\|\tilde{\Xi}\|_{2,0} \leq q$  and the following equation:

$$\tilde{\mathbf{B}} = \hat{\mathbf{B}}_{\text{glo}} - L^{-1}\mathbf{D}_1(\hat{\Theta}_{\text{glo}}), \quad \tilde{\Xi} \in \mathbf{H}_q \left( \hat{\Xi}_{\text{glo}} - L^{-1}\mathbf{D}_2(\hat{\Theta}_{\text{glo}}) \right).$$

By the definition of  $\hat{\Theta}_{\text{glo}}$ , we have  $g(\tilde{\Theta}) \geq g(\hat{\Theta}_{\text{glo}})$ . Moreover, as with (E.1), the following inequality can be derived:

$$g(\hat{\Theta}_{\text{glo}}) - g(\tilde{\Theta}) \geq \frac{L-\ell}{2} \|\tilde{\Theta} - \hat{\Theta}_{\text{glo}}\|_F^2.$$

Hence,  $\|\tilde{\Theta} - \hat{\Theta}_{\text{glo}}\|_F^2 = 0$  holds for  $L > \ell$ . This implies that  $\hat{\Theta}_{\text{glo}}$  is an  $\ell_{2,0}$ -constrained first-order stationary point.  $\square$

## F Proof of Theorem 1

By summing (14) for  $m = 1, \dots, M$ , we have

$$\sum_{m=1}^M \{g(\Theta_m) - g(\Theta_{m+1})\} \geq \frac{L-\ell}{2} \sum_{m=1}^M \|\Theta_{m+1} - \Theta_m\|_F^2.$$

Since  $g(\Theta_m)$  monotonically decreases for  $m$  and converges to  $g_*$  as  $m \rightarrow \infty$ , we obtain the following inequality:

$$\begin{aligned} g(\Theta_1) - g_* &\geq \sum_{m=1}^M \{g(\Theta_m) - g(\Theta_{m+1})\} \\ &\geq \frac{L-\ell}{2} \sum_{m=1}^M \|\Theta_{m+1} - \Theta_m\|_F^2 \\ &\geq \frac{M(L-\ell)}{2} \min_{m=1, \dots, M} \|\Theta_{m+1} - \Theta_m\|_F^2. \end{aligned}$$

This completes the proof of Theorem 1.  $\square$

## Acknowledgments

Ryoya Oda is supported by JSPS KAKENHI Grant Numbers JP19K21672, JP20K14363, and JP20H04151. Mineaki Ohishi is supported by JSPS KAKENHI Grant Numbers JP20H04151 and JP21K13834. Hirokazu Yanagihara is supported by JSPS KAKENHI Grant Numbers JP16H03606, JP18K03415, and JP20H04151.

## References

- [1] 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.
- [2] AKAIKE, H. (1974). A new look at the statistical model identification. *IEEE Trans. Automatic Control*, **AC-19**, 716–723.
- [3] BERTSIMAS, D., KING, A. and MAZUMDER, R. (2016). Best subset selection via a modern optimization lens. *Ann. Statist.*, **44**, 813–852.
- [4] CAI, X., NIE, F. and HUANG, H. (2013). Exact top- $k$  feature selection via  $\ell_{2,0}$ -norm constraint. In *Proceedings of the Twenty-Third International Joint Conference on Artificial Intelligence*.
- [5] GOTOH, J., TAKEDA, A. and TONO, K. (2018). DC formulations and algorithms for sparse optimization problems. *Math. Program.*, **169**, 141–176.
- [6] HARVILLE, D. A. (1997). *Matrix Algebra from a Statistician's Perspective*. Springer-Verlag, New York.
- [7] NESTEROV, Y. (2004). *Introductory Lectures on Convex Optimization*. Kluwer Academic Publishers, Boston, MA.
- [8] NESTEROV, Y. (2013). Gradient methods for minimizing composite functions. *Math. Program.*, **140**, 125–161.
- [9] OBOZINSKI, G., WAINWRIGHT, M. J. and JORDAN, M. I. (2011). Support union recovery in high-dimensional multivariate regression. *Ann. Statist.*, **39**, 1–47.
- [10] ODA, R. and YANAGIHARA, H. (2020). A fast and consistent variable selection method for high-dimensional multivariate linear regression with a large number of explanatory variables. *Electron. J. Statist.*, **14**, 1386–1412.
- [11] SCHWARZ, G. (1978). Estimating the dimension of a model. *Ann. Statist.*, **6**, 461–464.
- [12] SRIVASTAVA, M. S. (2002). *Methods of Multivariate Statistics*. John Wiley & Sons, New York.
- [13] TIMM, N. H. (2002). *Applied Multivariate Analysis*. Springer-Verlag, New York.

- [14] YANAGIHARA, H. (2006). Corrected version of AIC for selecting multivariate normal linear regression models in a general nonnormal case. *J. Multivariate Anal.*, **97**, 1070–1089.
- [15] YANAGIHARA, H., KAMO, K., IMORI, S. and SATOH, K. (2012). Bias-corrected AIC for selecting variables in multinomial logistic regression models. *Linear Algebra Appl.*, **436**, 4329–4341.