

CIRJE-F-575

**Conditional Information Criteria for Selecting
Variables in Linear Mixed Models**

Muni S. Srivastava
University of Toronto

Tatsuya Kubokawa
University of Tokyo

July 2008; Revised in December 2008

CIRJE Discussion Papers can be downloaded without charge from:

<http://www.e.u-tokyo.ac.jp/cirje/research/03research02dp.html>

Discussion Papers are a series of manuscripts in their draft form. They are not intended for circulation or distribution except as indicated by the author. For that reason Discussion Papers may not be reproduced or distributed without the written consent of the author.

Conditional Information Criteria for Selecting Variables in Linear Mixed Models

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

December 24, 2008

Abstract

In this paper, we consider the problem of selecting the variables of the fixed effects in the linear mixed models where the random effects are present and the observation vectors have been obtained from many clusters. As the variable selection procedure, we here use the Akaike Information Criterion, AIC. In the context of the mixed linear models, two kinds of AIC have been proposed: marginal AIC and conditional AIC. In this paper, we derive three versions of conditional AIC depending upon different estimators of the regression coefficients and the random effects. Through the simulation studies, it is shown that the proposed conditional AIC's are superior to the marginal and conditional AIC's proposed in the literature in the sense of selecting the true model. Finally, the results are extended to the case when the random effects in all the clusters are of the same dimension but have a common unknown covariance matrix.

Key words and phrases: Akaike information criterion, analysis of variance, linear mixed model, nested error regression model, random effect, selection of variables.

1 Introduction

Consider the model in which the n_i -vector of response variables \mathbf{y}_i in the i -th cluster is related by

$$\mathbf{y}_i = \mathbf{X}_i\boldsymbol{\beta} + \mathbf{Z}_i\mathbf{v}_i + \boldsymbol{\epsilon}_i, \quad i = 1, \dots, k, \quad (1)$$

where the k observation vectors $\mathbf{y}_1, \dots, \mathbf{y}_k$ are independently distributed, \mathbf{X}_i is a known $n_i \times p$ matrix, $\boldsymbol{\beta} = (\beta_0, \dots, \beta_{p-1})'$ is a p -vector of unknown parameters, \mathbf{Z}_i is an $n_i \times r_i$,

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

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

$n_i \geq r_i$, matrix of known covariables and \mathbf{v}_i is an r_i -vector of random effects. The error n_i -vector $\boldsymbol{\epsilon}_i$ is distributed independently of the random effects vector \mathbf{v}_i , both are assumed to be normally distributed; $\mathbf{v}_i \sim \mathcal{N}_{r_i}(\mathbf{0}, \sigma^2 \mathbf{G}_i)$ and $\boldsymbol{\epsilon}_i \sim \mathcal{N}_{n_i}(\mathbf{0}, \sigma^2 \mathbf{I}_{n_i})$. It is assumed that $\sum_{i=1}^k (n_i - r_i) > p$, $\sum_{i=1}^k r_i \geq p$, and \mathbf{Z}_i are of full rank r_i . Writing $\mathbf{y} = (\mathbf{y}'_1, \dots, \mathbf{y}'_k)'$, $\mathbf{X} = (\mathbf{X}'_1, \dots, \mathbf{X}'_k)'$, $\mathbf{Z} = \text{block diagonal}(\mathbf{Z}_1, \dots, \mathbf{Z}_k)$, $\mathbf{v} = (\mathbf{v}'_1, \dots, \mathbf{v}'_k)'$ and $\boldsymbol{\epsilon} = (\boldsymbol{\epsilon}'_1, \dots, \boldsymbol{\epsilon}'_k)'$, we can express the model in (1) in matrix notation as

$$\mathbf{y} = \mathbf{X}\boldsymbol{\beta} + \mathbf{Z}\mathbf{v} + \boldsymbol{\epsilon}, \quad (2)$$

where $\boldsymbol{\epsilon}$ and \mathbf{v} are independently distributed as $\mathcal{N}_N(\mathbf{0}, \sigma^2 \mathbf{I}_N)$ and $\mathcal{N}_R(\mathbf{0}, \sigma^2 \mathbf{G})$ for

$$R = \sum_{i=1}^k r_i, \quad N = \sum_{i=1}^k n_i, \quad \mathbf{G} = \text{block diag}(\mathbf{G}_1, \dots, \mathbf{G}_k).$$

The usual goal of the model (2) is to provide a good predicted value for a future observation. Often, it is achieved by reducing the dimension of the parameters, or equivalently by using fewer members of fixed variables than p , either by testing the hypothesis that some specified β_i 's are zero or by model selection method such as Akaike Information Criterion (AIC) of Akaike (1973, 1974). Let $f(\mathbf{y}|\mathbf{v}, \boldsymbol{\beta}, \sigma^2)$ and $f(\mathbf{v}|\sigma^2)$ be the conditional density of \mathbf{y} given \mathbf{v} and the marginal density of \mathbf{v} , respectively. Then, the marginal density of \mathbf{y} is written by $f_m(\mathbf{y}|\boldsymbol{\beta}, \sigma^2) = \int f(\mathbf{y}|\mathbf{v}, \boldsymbol{\beta}, \sigma^2)f(\mathbf{v}|\sigma^2)d\mathbf{v}$, which has

$$\mathcal{N}_N(\mathbf{X}\boldsymbol{\beta}, \sigma^2 \boldsymbol{\Lambda}(\mathbf{G})) \quad \text{for} \quad \boldsymbol{\Lambda}(\mathbf{G}) = \mathbf{I} + \mathbf{Z}\mathbf{G}\mathbf{Z}'.$$

When \mathbf{G} is known, the Akaike information based on the marginal density $f_m(\mathbf{y}|\boldsymbol{\beta}, \sigma^2)$ is

$$AI = -2 \int \int \{\log f_m(\mathbf{y}^*|\hat{\boldsymbol{\beta}}, \hat{\sigma}^2)\} f_m(\mathbf{y}|\boldsymbol{\beta}, \sigma^2) d\mathbf{y}^* d\mathbf{y},$$

where \mathbf{y}^* is a future observation having the same distribution as \mathbf{y} but independently distributed of \mathbf{y} , $\hat{\boldsymbol{\beta}}$ and $\hat{\sigma}^2$ are the maximum likelihood estimators (MLE) of $\boldsymbol{\beta}$ and σ^2 based on the observation \mathbf{y} where the marginal distribution is given above. The expressions for the MLE $\hat{\boldsymbol{\beta}}$ and $\hat{\sigma}^2$ are given in (5) and (6), respectively. Selection of the variables of the fixed effects \mathbf{X} is based on the minimum value of an unbiased estimator of AI . In (7), an exact unbiased estimate of AI , which we denote by AIC is given.

Another interesting approach, proposed by Vaida and Blanchard (2005), is based on the so-called conditional Akaike information given by

$$cAI = -2 \int \int \int \log \{f(\mathbf{y}^*|\hat{\mathbf{v}}, \hat{\boldsymbol{\beta}}, \hat{\sigma}^2)\} f(\mathbf{y}^*|\mathbf{v}, \boldsymbol{\beta}, \sigma^2) f(\mathbf{y}|\mathbf{v}, \boldsymbol{\beta}, \sigma^2) f(\mathbf{v}|\sigma^2) d\mathbf{y}^* d\mathbf{y} d\mathbf{v}, \quad (3)$$

where $\hat{\mathbf{v}}$ is the empirical Bayes estimator of \mathbf{v} given in (8) and $\hat{\boldsymbol{\beta}}$ and $\hat{\sigma}^2$ are given in (5) and (6), respectively. Vaida and Blanchard (2005) derived an unbiased estimator of cAI given by (9). We will denote this estimator by $cAIC$. Vaida and Blanchard have observed that

when \mathbf{G} is known $cAIC$ takes the same value as DIC , the deviance information criterion proposed by Spiegelhalter, Best, Carlin and van der Linde (2002) for Bayesian inference.

In our simulation results, we, however, find that the performance of $cAIC$ is no better than the marginal $mAIC$. Thus, the performance of both of them is not good. This may be due to the fact that the estimate $\hat{\sigma}^2(\mathbf{G})$ ignores the existence of the random effects although it provides the largest degrees of freedom available to estimate σ^2 . But it also makes it necessary to obtain cAI marginally, that is averaging with respect to the density of \mathbf{v} . Thus, we consider estimating σ^2 from the conditional model, where

$$\mathbf{y}|\mathbf{v} \sim \mathcal{N}_N(\mathbf{W}\boldsymbol{\gamma}, \sigma^2\mathbf{I}),$$

where $\mathbf{W} = (\mathbf{X}, \mathbf{Z})$ and $\boldsymbol{\gamma} = (\boldsymbol{\beta}', \mathbf{v}')'$. Hence, we need to obtain the bias term only conditionally. That is, the conditional Akaike information we consider is

$$CAI(\mathbf{v}) = -2 \int \int \log\{f(\mathbf{y}^*|\tilde{\mathbf{v}}, \tilde{\boldsymbol{\beta}}, \hat{\sigma}_0^2)\} f(\mathbf{y}^*|\mathbf{v}, \boldsymbol{\beta}, \sigma^2) f(\mathbf{y}|\mathbf{v}, \boldsymbol{\beta}, \sigma^2) d\mathbf{y}^* d\mathbf{y}, \quad (4)$$

where $\hat{\sigma}_0 = (\mathbf{y} - \mathbf{W}\hat{\boldsymbol{\gamma}}_0)'(\mathbf{y} - \mathbf{W}\hat{\boldsymbol{\gamma}}_0)/N$ for the least squares estimator $\hat{\boldsymbol{\gamma}}_0$ of $\boldsymbol{\gamma}$, and $\tilde{\mathbf{v}}$ and $\tilde{\boldsymbol{\beta}}$ are linear estimators of \mathbf{y} for \mathbf{v} and $\boldsymbol{\beta}$. It is noted that $\hat{\sigma}_0^2$ is not affected by the random effect term \mathbf{v} and $N\hat{\sigma}_0^2/\sigma^2$ has a chisquare distribution with $(N - p - R)$ -degrees of freedom. The proposed new conditional AIC is an unbiased estimator of CAI given by

$$CAIC(\mathbf{G}) = -2 \log f(\mathbf{y}|\tilde{\mathbf{v}}, \tilde{\boldsymbol{\beta}}, \hat{\sigma}_0^2) + 2 \frac{N(\text{tr}[\mathbf{W}\mathbf{L}] + 1)}{N - r_w - 2},$$

where $r_w = \text{rank}(\mathbf{W})$ and \mathbf{L} is a $(p + R) \times N$ matrix such that $(\tilde{\boldsymbol{\beta}}', \tilde{\mathbf{v}}')' = \mathbf{L}\mathbf{y}$. It is interesting to note that although this is obtained conditionally given \mathbf{v} , the bias term $-2N(\text{tr}[\mathbf{W}\mathbf{L}] + 1)/(N - r_w - 2)$ does not depend on \mathbf{v} , and thus, it is also an unbiased estimator of cAI defined by Vaida and Blanchard (2005). Although $\hat{\sigma}_0^2$ is based on less degrees of freedom than $\hat{\sigma}^2$, we argue that it is the most appropriate estimate of σ^2 to use in $-2 \log f(\mathbf{y}|\tilde{\mathbf{v}}, \tilde{\boldsymbol{\beta}}, \hat{\sigma}_0^2)$ as it takes into account that not only $\boldsymbol{\beta}$ but \mathbf{v} has also been estimated while $\hat{\sigma}^2$ does not take this into account; in fact in the conditional framework $\hat{\sigma}^2$ is a biased estimator of σ^2 . Simulation experiments carried out in Section 4 show that the performance of $CAIC(\mathbf{G})$ is significantly better than $mAIC(\mathbf{G})$ and $cAIC(\mathbf{G})$ in the sense of selecting the true model. This method also facilitates to consider the case when \mathbf{G} is a function of unknown parameters.

The organization of the paper is as follows. In Section 2, we introduce the concept of AIC and marginal and conditional AIC's. In Section 3, we derive three conditional AIC using three different kinds of estimators available for $(\boldsymbol{\beta}, \mathbf{v})$, assuming that \mathbf{G} is known. The case of unknown \mathbf{G} is considered in Section 4. We give simulation results in Section 5. The paper concludes in Section 6. The proof is given in the Appendix.

2 Marginal and Conditional AIC in the linear mixed model

2.1 Concept of AIC

We now introduce the marginal and conditional AIC's in the linear mixed model in the case of known \mathbf{G} . Before describing them, we first explain the concept of AIC briefly. The AIC is based on Kullback-Leibler distance. For a true density f and an approximating one g_ω , this distance is defined as

$$I(f, g_\omega) = E_{y^*} \log f(y^*) - E_f \log g_\omega(y^*),$$

where E_{y^*} denotes the expectation with respect to the true density $f(y^*)$. Let $\mathcal{G} = \{g_\omega : \omega \in \Omega\}$ be the class of approximating densities. If $f \in \mathcal{G}$, then there exist a $g_{\omega_0} \in \mathcal{G}$ such that $I(f, g_{\omega_0}) = 0$, otherwise $I(f, g_{\omega_0}) \geq 0$. Thus a g is chosen for which $I(f, g)$ is minimum. Usually ω is not known and estimated from the data \mathbf{y} by $\hat{\omega} = \hat{\omega}(\mathbf{y})$. Thus, $I(f, g_\omega)$ is approximated by $I(f, g_{\hat{\omega}})$ and the quality of approximation is judged by

$$E_y f[I(f, g_{\hat{\omega}})] = E_{y^*} [\log f(y^*)] - E_y E_{y^*} [\log g_{\hat{\omega}(\mathbf{y})}(y^*)],$$

where E_y denotes the expectation with respect to the true density $f(\mathbf{y})$, which is independent of y^* . Akaike information is defined by

$$AI = -2E_y E_{y^*} [\log g_{\hat{\omega}(\mathbf{y})}(y^*)].$$

An unbiased estimator of AI is given by

$$AIC = -2 \log g_{\hat{\omega}(\mathbf{y})}(\mathbf{y}) + \Delta,$$

where Δ is the bias

$$\Delta = E_y [-2 \log g_{\hat{\omega}(\mathbf{y})}(\mathbf{y})] - AI.$$

Akaike (1973, 1974) used an approximate value of the bias given by the number of free parameters. Thus, Akaike used the number of free parameters in place of Δ . It is noted that AIC is a criterion for selecting a good model in terms of minimizing the prediction error. It may be noted that the estimator $\hat{\omega}$ of ω need not be an MLE as any consistent estimator of ω may perform as good, see Konishi and Kitagawa (2007).

2.2 Marginal AIC

The marginal AIC in the linear mixed model is AIC based on the marginal distribution

$$\mathbf{y} \sim \mathcal{N}_N(\mathbf{X}\boldsymbol{\beta}, \sigma^2\boldsymbol{\Lambda}),$$

for $\boldsymbol{\Lambda} = \boldsymbol{\Lambda}(\mathbf{G}) = \mathbf{I}_N + \mathbf{Z}\mathbf{G}\mathbf{Z}'$, where the marginal density is given by

$$f_m(\mathbf{y}|\boldsymbol{\beta}, \sigma^2) = (2\pi\sigma^2)^{-N/2} |\boldsymbol{\Lambda}|^{-1/2} \exp\{-\frac{1}{2\sigma^2}(\mathbf{y} - \mathbf{X}\boldsymbol{\beta})' \boldsymbol{\Lambda}^{-1}(\mathbf{y} - \mathbf{X}\boldsymbol{\beta})\}.$$

When \mathbf{G} does not include unknown parameters, the maximum likelihood estimators (MLE) of $\boldsymbol{\beta}$ and σ^2 are given by

$$\widehat{\boldsymbol{\beta}} = \widehat{\boldsymbol{\beta}}(\mathbf{G}) = (\mathbf{X}'\boldsymbol{\Lambda}^{-1}\mathbf{X})^+ \mathbf{X}'\boldsymbol{\Lambda}^{-1}\mathbf{y}, \quad (5)$$

$$\begin{aligned} \widehat{\sigma}^2 &= \widehat{\sigma}^2(\mathbf{G}) = (\mathbf{y} - \mathbf{X}\widehat{\boldsymbol{\beta}})' \boldsymbol{\Lambda}^{-1} (\mathbf{y} - \mathbf{X}\widehat{\boldsymbol{\beta}}) / N \\ &= \mathbf{y}' \{ \boldsymbol{\Lambda}^{-1} - \boldsymbol{\Lambda}^{-1} (\mathbf{X}'\boldsymbol{\Lambda}^{-1}\mathbf{X})^{-} \boldsymbol{\Lambda}^{-1} \} \mathbf{y}, \end{aligned} \quad (6)$$

where \mathbf{A}^+ denotes the Moore-Penrose inverse of \mathbf{A} and satisfies (i) $\mathbf{A}\mathbf{A}^+\mathbf{A} = \mathbf{A}$, (ii) $\mathbf{A}^+\mathbf{A}\mathbf{A}^+ = \mathbf{A}^+$, (iii) $(\mathbf{A}^+\mathbf{A})' = (\mathbf{A}^+\mathbf{A})$, (iv) $(\mathbf{A}\mathbf{A}^+)' = \mathbf{A}\mathbf{A}^+$. Then, $-2\log f_m(\mathbf{y}|\widehat{\boldsymbol{\beta}}, \widehat{\sigma}^2)$ is expressed as

$$-2\log f_m(\mathbf{y}|\widehat{\boldsymbol{\beta}}, \widehat{\sigma}^2) = N[\log(2\pi\widehat{\sigma}^2(\mathbf{G})) + 1] + \log|\boldsymbol{\Lambda}|,$$

and the exact bias correction AIC based on the marginal likelihood is given by

$$mAIC(\mathbf{G}) = N[\log(2\pi\widehat{\sigma}^2(\mathbf{G})) + 1] + \log|\boldsymbol{\Lambda}| + 2N(r_x + 1)/(N - r_x - 2), \quad (7)$$

where $r_x = \text{rank}(\mathbf{X})$. When \mathbf{G} includes unknown parameters, we can use the criteria $AIC(\widehat{\mathbf{G}})$ and $mAIC(\widehat{\mathbf{G}})$ when a consistent estimator $\widehat{\mathbf{G}}$ of \mathbf{G} is available.

2.3 Conditional AIC

The conditional AIC in the linear mixed model was proposed by Vaida and Blanchard (2005), who considered estimating the random effects \mathbf{v} as well as $\boldsymbol{\beta}$ by the mixed model equation

$$\begin{pmatrix} \mathbf{X}'\mathbf{X} & \mathbf{X}'\mathbf{Z} \\ \mathbf{Z}'\mathbf{X} & \mathbf{Z}'\mathbf{Z} + \mathbf{G}^{-1} \end{pmatrix} \begin{pmatrix} \widehat{\boldsymbol{\beta}} \\ \widehat{\mathbf{v}} \end{pmatrix} = \begin{pmatrix} \mathbf{X}' \\ \mathbf{Z}' \end{pmatrix} \mathbf{y},$$

which was given by Henderson (1950). Then the solution $\widehat{\boldsymbol{\beta}}$ is the generalized least squares estimator (GLS) given in (5), and $\widehat{\mathbf{v}}$ is given by

$$\widehat{\mathbf{v}} = \mathbf{G}\mathbf{Z}'\boldsymbol{\Lambda}^{-1}(\mathbf{y} - \mathbf{X}\widehat{\boldsymbol{\beta}}), \quad (8)$$

which can also be derived as an empirical Bayes estimator by considering the conditional distribution of \mathbf{v} given \mathbf{y} . Thus, using the estimator $\widehat{\sigma}^2(\mathbf{G})$ defined in (6), they define the conditional Akaike information by (3). The conditional AIC, denoted by $cAIC(\mathbf{G})$, is given by

$$cAIC(\mathbf{G}) = -2\log f(\mathbf{y}|\widehat{\mathbf{v}}, \widehat{\boldsymbol{\beta}}, \widehat{\sigma}^2) - \Delta_c, \quad (9)$$

where

$$\begin{aligned} \Delta_c &= E_{\mathbf{y}}[-2\log f(\mathbf{y}|\widehat{\mathbf{v}}, \widehat{\boldsymbol{\beta}}, \widehat{\sigma}^2)] - cAI \\ &= -2 \frac{N(N - r_x - 1)}{(N - r_x)(N - r_x - 2)} (\rho + 1) + \frac{N(r_x + 1)}{(N - r_x)(N - r_x - 2)}. \end{aligned}$$

Here, ρ is defined by $\rho = \text{tr}(\mathbf{H}_1)$, where

$$\mathbf{H}_1 = (\mathbf{X}, \mathbf{Z}) \begin{pmatrix} \mathbf{X}'\mathbf{X} & \mathbf{X}'\mathbf{Z} \\ \mathbf{Z}'\mathbf{X} & \mathbf{Z}'\mathbf{Z} + \mathbf{G}^{-1} \end{pmatrix}^{-1} \begin{pmatrix} \mathbf{X}' \\ \mathbf{Z}' \end{pmatrix}.$$

which can be also expressed as

$$\begin{aligned} \mathbf{H}_1 &= \mathbf{X}(\mathbf{X}'\mathbf{\Lambda}^{-1}\mathbf{X})^+ \mathbf{X}'\mathbf{\Lambda}^{-1} + \mathbf{Z}\mathbf{G}\mathbf{Z}'\mathbf{\Lambda}^{-1} \{ \mathbf{I} - \mathbf{X}(\mathbf{X}'\mathbf{\Lambda}^{-1}\mathbf{X})^+ \mathbf{X}'\mathbf{\Lambda}^{-1} \} \\ &= \mathbf{\Lambda}^{-1} \mathbf{X}(\mathbf{X}'\mathbf{\Lambda}^{-1}\mathbf{X})^+ \mathbf{X}'\mathbf{\Lambda}^{-1} + \mathbf{I} - \mathbf{\Lambda}^{-1}. \end{aligned} \quad (10)$$

3 Proposed conditional AIC

In our simulation results, we find that the performance of $cAIC(\mathbf{G})$ is only slightly better than the marginal $mAIC(\mathbf{G})$ in the sense of selecting true models. To improve the performance, we here propose another type of conditional AIC's. In this section, we assume that \mathbf{G} is known, and we begin with rewriting the model (2) as

$$\mathbf{y} = \mathbf{X}\boldsymbol{\beta} + \mathbf{Z}\mathbf{v} + \boldsymbol{\epsilon} \equiv \mathbf{W}\boldsymbol{\gamma} + \boldsymbol{\epsilon},$$

where $\mathbf{W} = (\mathbf{X}, \mathbf{Z})$ is an $N \times (p + R)$ matrix and $\boldsymbol{\gamma} = (\boldsymbol{\beta}', \mathbf{v}')'$ is a $(p + R)$ -dimensional vector. Given \mathbf{v} , the model $\mathbf{y} = \mathbf{W}\boldsymbol{\gamma} + \boldsymbol{\epsilon}$ can be conditionally regarded as a usual linear regression model, and the conditional MLE of $\boldsymbol{\gamma}$ and σ^2 are given by

$$\hat{\boldsymbol{\gamma}}_0 = (\mathbf{W}'\mathbf{W})^+ \mathbf{W}'\mathbf{y}, \quad (11)$$

$$\hat{\sigma}_0^2 = (\mathbf{y} - \mathbf{W}\hat{\boldsymbol{\gamma}}_0)'(\mathbf{y} - \mathbf{W}\hat{\boldsymbol{\gamma}}_0)/N. \quad (12)$$

It can be shown that the MLE of $\boldsymbol{\gamma}$ can be expressed as

$$\begin{aligned} \hat{\boldsymbol{\gamma}}_0 &= \begin{pmatrix} \hat{\boldsymbol{\beta}}_0 \\ \hat{\mathbf{v}}_0 \end{pmatrix} = \begin{pmatrix} (\widetilde{\mathbf{X}}'\widetilde{\mathbf{X}})^+ \widetilde{\mathbf{X}}' \\ (\mathbf{Z}'\mathbf{Z})^{-1} \mathbf{Z}'(\mathbf{I} - \widetilde{\mathbf{X}}(\widetilde{\mathbf{X}}'\widetilde{\mathbf{X}})^+ \widetilde{\mathbf{X}}') \end{pmatrix} \mathbf{y} \\ &\equiv \mathbf{L}_0 \mathbf{y}, \end{aligned}$$

where $\widetilde{\mathbf{X}} = (\mathbf{I} - \mathbf{H}_z)\mathbf{X}$ for $\mathbf{H}_z = \mathbf{Z}(\mathbf{Z}'\mathbf{Z})^{-1}\mathbf{Z}'$. It is noted that the MLE of $\boldsymbol{\gamma}$ is a linear function of \mathbf{y} . Thus, we may consider a general estimator of $\boldsymbol{\gamma}$ as a linear function of \mathbf{y} , namely,

$$\tilde{\boldsymbol{\gamma}} = \begin{pmatrix} \tilde{\boldsymbol{\beta}} \\ \tilde{\mathbf{v}} \end{pmatrix} = \mathbf{L} \mathbf{y},$$

where \mathbf{L} is a $(p + R) \times N$ matrix which will be specified later. It is noted that $N\hat{\sigma}_0^2/\sigma^2$ has a chisquare distribution with $(N - r_w)$ degrees of freedom for $r_w = \text{rank}(\mathbf{W})$. Thus, we shall obtain our $CAIC$ using $\tilde{\boldsymbol{\gamma}}$ and $\hat{\sigma}_0^2$ as estimators of $\boldsymbol{\gamma}$ and σ^2 , respectively.

The conditional Akaike information considered here is given in (4), which is different from cAI given in (3) in that $cAI = \int CAI(\mathbf{v})f(\mathbf{v}|\sigma^2)d\mathbf{v}$. Since $-2\log f(\mathbf{y}|\tilde{\mathbf{v}}, \tilde{\boldsymbol{\beta}}, \hat{\sigma}_0^2)$ is written as

$$-2\log f(\mathbf{y}|\tilde{\mathbf{v}}, \tilde{\boldsymbol{\beta}}, \hat{\sigma}_0^2) = N \log(2\pi\hat{\sigma}_0^2) + (\mathbf{y} - \mathbf{W}\tilde{\boldsymbol{\gamma}})'(\mathbf{y} - \mathbf{W}\tilde{\boldsymbol{\gamma}})/\hat{\sigma}_0^2.$$

We now consider another random vector \mathbf{y}^* distributed independently of \mathbf{y} but having the same distribution as \mathbf{y} . Hence,

$$\mathbf{y}^*|\mathbf{v} \sim \mathcal{N}_N(\mathbf{W}\boldsymbol{\gamma}, \sigma^2\mathbf{I}).$$

Thus,

$$\begin{aligned} E_{\mathbf{y}^*}[-2\log f(\mathbf{y}^*|\tilde{\mathbf{v}}, \tilde{\boldsymbol{\beta}}, \hat{\sigma}_0^2)|\mathbf{v}] &= N\log(2\pi\hat{\sigma}_0^2) + E_{\mathbf{y}^*}[(\mathbf{y}^* - \mathbf{W}\tilde{\boldsymbol{\gamma}})'(\mathbf{y}^* - \mathbf{W}\tilde{\boldsymbol{\gamma}})/\hat{\sigma}_0^2|\mathbf{v}] \\ &= N\log(2\pi\hat{\sigma}_0^2) + N\frac{\sigma^2}{\hat{\sigma}_0^2} + \frac{(\tilde{\boldsymbol{\gamma}} - \boldsymbol{\gamma})'\mathbf{W}'\mathbf{W}(\tilde{\boldsymbol{\gamma}} - \boldsymbol{\gamma})}{\hat{\sigma}_0^2}. \end{aligned}$$

Hence, the conditional Akaike information given \mathbf{v} is given by

$$CAI(\mathbf{v}) = E_{\mathbf{y}} \left[N\log(2\pi\hat{\sigma}_0^2) + N\frac{\sigma^2}{\hat{\sigma}_0^2} + \frac{(\tilde{\boldsymbol{\gamma}} - \boldsymbol{\gamma})'\mathbf{W}'\mathbf{W}(\tilde{\boldsymbol{\gamma}} - \boldsymbol{\gamma})}{\hat{\sigma}_0^2} \middle| \mathbf{v} \right],$$

and the bias is expressed as

$$\begin{aligned} \Delta_C(\mathbf{v}) &= E_{\mathbf{y}} \left[-2\log f(\mathbf{y}|\tilde{\mathbf{v}}, \tilde{\boldsymbol{\beta}}, \hat{\sigma}_0^2) \middle| \mathbf{v} \right] - CAI(\mathbf{v}) \\ &= E_{\mathbf{y}} \left[\frac{(\mathbf{y} - \mathbf{W}\boldsymbol{\gamma})'(\mathbf{y} - \mathbf{W}\boldsymbol{\gamma})}{\hat{\sigma}_0^2} - 2\frac{(\mathbf{y} - \mathbf{W}\boldsymbol{\gamma})'\mathbf{W}(\tilde{\boldsymbol{\gamma}} - \boldsymbol{\gamma})}{\hat{\sigma}_0^2} - N\frac{\sigma^2}{\hat{\sigma}_0^2} \middle| \mathbf{v} \right]. \end{aligned} \quad (13)$$

It is here noted that $E_{\mathbf{y}}[\sigma^2/\hat{\sigma}_0^2|\mathbf{v}] = N/(N - r_w - 2)$ and $E_{\mathbf{y}}[(\mathbf{y} - \mathbf{W}\boldsymbol{\gamma})'(\mathbf{y} - \mathbf{W}\boldsymbol{\gamma})/\hat{\sigma}_0^2|\mathbf{v}] = N + Nr_w/(N - r_w - 2)$. Also for $\tilde{\boldsymbol{\gamma}} = \mathbf{L}\mathbf{y}$, it is observed that

$$\begin{aligned} E_{\mathbf{y}} \left[\frac{(\mathbf{y} - \mathbf{W}\boldsymbol{\gamma})'\mathbf{W}(\tilde{\boldsymbol{\gamma}} - \boldsymbol{\gamma})}{\hat{\sigma}_0^2} \middle| \mathbf{v} \right] &= E_{\mathbf{y}} \left[\frac{(\mathbf{y} - \mathbf{W}\boldsymbol{\gamma})'\mathbf{W}\mathbf{L}\mathbf{y}}{\hat{\sigma}_0^2} \middle| \mathbf{v} \right] \\ &= \frac{1}{2}E_{\mathbf{y}} \left[\frac{(\mathbf{y} - \mathbf{W}\boldsymbol{\gamma})'(\mathbf{W}\mathbf{L} + \mathbf{L}'\mathbf{W}')(\mathbf{y} - \mathbf{W}\boldsymbol{\gamma})}{\hat{\sigma}_0^2} \middle| \mathbf{v} \right] \\ &= \frac{N}{2}E_{\mathbf{y}} \left[\frac{\mathbf{u}'(\mathbf{W}\mathbf{L} + \mathbf{L}'\mathbf{W}')\mathbf{u}}{\mathbf{u}'(\mathbf{I} - \mathbf{H}_w)\mathbf{u}} \middle| \mathbf{v} \right], \end{aligned}$$

where $\mathbf{u} = (\mathbf{y} - \mathbf{W}\boldsymbol{\gamma})/\sigma$ and $\mathbf{H}_w = \mathbf{W}(\mathbf{W}'\mathbf{W})^{-1}\mathbf{W}'$. Note that given \mathbf{v} , the conditional distribution of \mathbf{u} has $\mathcal{N}_N(\mathbf{0}, \mathbf{I})$, so that \mathbf{u} is independent of \mathbf{v} . Using Lemma A.1 in the Appendix, we can see that

$$\begin{aligned} \frac{N}{2}E \left[\frac{\mathbf{u}'(\mathbf{W}\mathbf{L} + \mathbf{L}'\mathbf{W}')\mathbf{u}}{\mathbf{u}'(\mathbf{I} - \mathbf{H}_w)\mathbf{u}} \middle| \mathbf{v} \right] &= N \left\{ \frac{\text{tr}(\mathbf{W}\mathbf{L})}{N - r_w - 2} - \frac{2\text{tr}[\mathbf{W}\mathbf{L}(\mathbf{I} - \mathbf{H}_w)]}{(N - r_w)(N - r_w - 2)} \right\} \\ &= \frac{N\text{tr}(\mathbf{W}\mathbf{L})}{N - r_w - 2}, \end{aligned}$$

since $\text{tr}[\mathbf{W}\mathbf{L}(\mathbf{I} - \mathbf{H}_w)] = 0$. Combining these evaluations gives the following expression for the bias term $\Delta_C(\mathbf{v})$:

$$\Delta_C(\mathbf{v}) = -2\frac{N(\text{tr}[\mathbf{W}\mathbf{L}] + 1)}{N - r_w - 2}, \quad (14)$$

which yields the conditional AIC given by

$$CAIC(\mathbf{G}) = -2 \log f(\mathbf{y}|\tilde{\mathbf{v}}, \tilde{\boldsymbol{\beta}}, \hat{\sigma}_0^2) + 2 \frac{N(\text{tr}[\mathbf{W}\mathbf{L}] + 1)}{N - r_w - 2}. \quad (15)$$

It is interesting to note that the bias term $\Delta_C(\mathbf{v})$ does not depend on the given \mathbf{v} . Given \mathbf{v} , $CAIC$ is an unbiased estimator of CAI given in (4), which turns out to be an unbiased estimator of the conditional Akaike information cAI defined by Vaida and Blanchard (2005).

The matrix \mathbf{L} given in (15) depends on the choice of $\hat{\boldsymbol{\gamma}}$ as an estimator of $\boldsymbol{\gamma}$. We consider three kinds of estimators of $\boldsymbol{\gamma}$ as described below.

[1] **Using maximum likelihood estimator $\hat{\boldsymbol{\gamma}}_0$ for $\boldsymbol{\gamma}$.** Since the maximum likelihood or least squares estimator of $\boldsymbol{\gamma}$ is given in (11) as $\hat{\boldsymbol{\gamma}}_0 = (\mathbf{W}'\mathbf{W})^+\mathbf{W}'\mathbf{y}$, the matrix \mathbf{L} corresponds to $\mathbf{L} = (\mathbf{W}'\mathbf{W})^+\mathbf{W}'$. Hence, it is seen that $\text{tr}[\mathbf{W}\mathbf{L}] = \text{tr}[(\mathbf{W}'\mathbf{W})^+(\mathbf{W}'\mathbf{W})] = \text{rank}(\mathbf{W}) = r_w$, and we get the conditional AIC based on MLE $\hat{\boldsymbol{\gamma}}_0 = (\hat{\boldsymbol{\beta}}_0', \hat{\mathbf{v}}_0)'$ as

$$\begin{aligned} CAIC_{ML}(\mathbf{G}) &= -2 \log f(\mathbf{y}|\hat{\mathbf{v}}_0, \hat{\boldsymbol{\beta}}_0, \hat{\sigma}_0^2) + 2N(r_w + 1)/(N - r_w - 1) \\ &= N[\log(2\pi\hat{\sigma}_0^2) + 1] + 2N(r_w + 1)/(N - r_w - 1). \end{aligned} \quad (16)$$

[2] **Using empirical Bayes (EB) estimator of \mathbf{v} and simple estimator of $\boldsymbol{\beta}$.** Consider the case that the empirical Bayes estimator is used for \mathbf{v} and the simple estimator $\hat{\boldsymbol{\beta}}_1 = (\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\mathbf{y}$ is utilized for $\boldsymbol{\beta}$. The empirical Bayes estimator of \mathbf{v} is given in (8) as $\hat{\mathbf{v}} = \mathbf{G}\mathbf{Z}'\boldsymbol{\Lambda}^{-1}(\mathbf{y} - \mathbf{X}\hat{\boldsymbol{\beta}})$, and the matrix $\text{tr}[\mathbf{W}\mathbf{L}]$ corresponds to

$$\text{tr}[\mathbf{W}\mathbf{L}] = \text{tr}[\mathbf{X}(\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}' + \mathbf{Z}\mathbf{G}\mathbf{Z}'\boldsymbol{\Lambda}^{-1}\{\mathbf{I} - \mathbf{X}(\mathbf{X}'\boldsymbol{\Lambda}^{-1}\mathbf{X})^+\mathbf{X}'\boldsymbol{\Lambda}^{-1}\}].$$

Note that $\text{tr}[\mathbf{X}(\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'] = \text{tr}[\mathbf{X}(\mathbf{X}'\boldsymbol{\Lambda}^{-1}\mathbf{X})^{-1}\mathbf{X}'\boldsymbol{\Lambda}^{-1}] = \text{rank}(\mathbf{X})$. From the argument around (10), it is seen that $\text{tr}[\mathbf{W}\mathbf{L}] = \text{tr}[\mathbf{H}_1] = \rho$. Hence, we get the conditional AIC based on the estimator $(\hat{\boldsymbol{\beta}}_1', \hat{\mathbf{v}})'$ of $\boldsymbol{\gamma}$ as

$$\begin{aligned} CAIC_{SL}(\mathbf{G}) &= -2 \log f(\mathbf{y}|\hat{\mathbf{v}}, \hat{\boldsymbol{\beta}}_1, \hat{\sigma}_0^2) + 2N(\rho + 1)/(N - r_w - 2) \\ &= N \log(2\pi\hat{\sigma}_0^2) + \frac{(\mathbf{y} - \mathbf{X}\hat{\boldsymbol{\beta}}_1 - \mathbf{Z}\hat{\mathbf{v}})'(\mathbf{y} - \mathbf{X}\hat{\boldsymbol{\beta}}_1 - \mathbf{Z}\hat{\mathbf{v}})}{\hat{\sigma}_0^2} + \frac{2N(\rho + 1)}{N - r_w - 2}. \end{aligned} \quad (17)$$

[3] **Using generalized least squares estimator of $\boldsymbol{\beta}$ and EB estimator of \mathbf{v} .** Consider the generalized least squares estimator $\hat{\boldsymbol{\beta}}$ and the EB estimator $\hat{\mathbf{v}}$ given in (5) and (8), respectively. These are estimators treated by Vaida and Blanchard (2005), and it is seen that $\text{tr}[\mathbf{W}\mathbf{L}] = \rho$ since $\mathbf{W}\mathbf{L} = \mathbf{H}_1$. Hence, we get the conditional AIC based on the estimator $(\hat{\boldsymbol{\beta}}', \hat{\mathbf{v}})'$ of $\boldsymbol{\gamma}$ as

$$\begin{aligned} CAIC_{GL}(\mathbf{G}) &= -2 \log f(\mathbf{y}|\hat{\mathbf{v}}, \hat{\boldsymbol{\beta}}, \hat{\sigma}_0^2) + 2N(\rho + 1)/(N - r_w - 2) \\ &= N \log(2\pi\hat{\sigma}_0^2) + \frac{(\mathbf{y} - \mathbf{X}\hat{\boldsymbol{\beta}} - \mathbf{Z}\hat{\mathbf{v}})'(\mathbf{y} - \mathbf{X}\hat{\boldsymbol{\beta}} - \mathbf{Z}\hat{\mathbf{v}})}{\hat{\sigma}_0^2} + \frac{2N(\rho + 1)}{N - r_w - 2}. \end{aligned} \quad (18)$$

4 Extensions to the case of unknown variance components

We have explained the marginal and conditional AIC's under the assumption that \mathbf{G} is known. In most applications, however, \mathbf{G} depends on unknown parameters. In this section, we handle the case of \mathbf{G} including unknown parameters. If a consistent estimator $\widehat{\mathbf{G}}$ is available for \mathbf{G} , then it can be substituted into $mAIC(\mathbf{G})$, $cAIC(\mathbf{G})$ and $CAIC(\mathbf{G})$ to get the marginal AIC $mAIC(\widehat{\mathbf{G}})$, the conditional AIC's $cAIC(\widehat{\mathbf{G}})$ and $CAIC(\widehat{\mathbf{G}})$, which will be suggested in this paper. In this case, the problem is how to estimate \mathbf{G} . The maximum likelihood method is an approach, but we need heavy computation as well as convergence of numerical iterations. Thus, in this section, we provide estimators of \mathbf{G} in explicit forms in some specific models.

For the model given in (2), we begin with making the transformation

$$\begin{pmatrix} \tilde{\mathbf{y}}_{1i} \\ \tilde{\mathbf{y}}_{2i} \end{pmatrix} = \begin{pmatrix} \mathbf{\Gamma}_i \\ (\mathbf{Z}'_i \mathbf{Z}_i)^{-1} \mathbf{Z}'_i \end{pmatrix} \mathbf{y}_i, \quad (19)$$

where $\mathbf{\Gamma}_i$ is an $(n_i - r_i) \times n_i$ matrix such that $\mathbf{\Gamma}_i \mathbf{Z}_i = \mathbf{0}$ and $\mathbf{\Gamma}_i \mathbf{\Gamma}'_i = \mathbf{I}_{n_i - r_i}$. It may be noted that $\mathbf{I} - \mathbf{Z}_i (\mathbf{Z}'_i \mathbf{Z}_i)^{-1} \mathbf{Z}'_i = \mathbf{\Gamma}'_i \mathbf{\Gamma}_i$ and thus $\mathbf{\Gamma}_i$ can easily be computed. According to this transformation, let $\widetilde{\mathbf{X}}_{1i} = \mathbf{\Gamma}_i \mathbf{X}_i$, $\tilde{\boldsymbol{\epsilon}}_{1i} = \mathbf{\Gamma}_i \boldsymbol{\epsilon}_i$, $\widetilde{\mathbf{X}}_{2i} = (\mathbf{Z}'_i \mathbf{Z}_i)^{-1} \mathbf{Z}'_i \mathbf{X}_i$ and $\tilde{\boldsymbol{\epsilon}}_{2i} = (\mathbf{Z}'_i \mathbf{Z}_i)^{-1} \mathbf{Z}'_i \boldsymbol{\epsilon}_i$ for $i = 1, \dots, k$. For $j = 1, 2$, let $\tilde{\mathbf{y}}_j = (\tilde{\mathbf{y}}'_{j1}, \dots, \tilde{\mathbf{y}}'_{jk})'$, $\widetilde{\mathbf{X}}_j = (\widetilde{\mathbf{X}}'_{j1}, \dots, \widetilde{\mathbf{X}}'_{jk})'$ and $\tilde{\boldsymbol{\epsilon}}_j = (\tilde{\boldsymbol{\epsilon}}'_{j1}, \dots, \tilde{\boldsymbol{\epsilon}}'_{jk})'$. Then, the model is decomposed as

$$\begin{aligned} \tilde{\mathbf{y}}_1 &= \widetilde{\mathbf{X}}_1 \boldsymbol{\beta} + \tilde{\boldsymbol{\epsilon}}_1, \\ \tilde{\mathbf{y}}_2 &= \widetilde{\mathbf{X}}_2 \boldsymbol{\beta} + \mathbf{v} + \tilde{\boldsymbol{\epsilon}}_2, \end{aligned} \quad (20)$$

where $\tilde{\mathbf{y}}_1 : (N - R) \times 1$, $\tilde{\mathbf{y}}_2 : R \times 1$, $\widetilde{\mathbf{X}}_1 : (N - R) \times p$, $\widetilde{\mathbf{X}}_2 : R \times p$, $N = \sum_{i=1}^k n_i$ and $R = \sum_{i=1}^k r_i$. Let $\tilde{\mathbf{y}} = (\tilde{\mathbf{y}}'_1, \tilde{\mathbf{y}}'_2)'$ and $\mathbf{W}_2 = \text{block diag}(\mathbf{G}_i + (\mathbf{Z}'_i \mathbf{Z}_i)^{-1}; i = 1, \dots, k)$. Then, $\tilde{\mathbf{y}}_1$ and $\tilde{\mathbf{y}}_2$ are mutually independently distributed as

$$\begin{aligned} \tilde{\mathbf{y}}_1 &\sim \mathcal{N}_{N-R}(\widetilde{\mathbf{X}}_1 \boldsymbol{\beta}, \sigma^2 \mathbf{I}_{N-R}), \\ \tilde{\mathbf{y}}_2 &\sim \mathcal{N}_R(\widetilde{\mathbf{X}}_2 \boldsymbol{\beta}, \sigma^2 \mathbf{W}_2). \end{aligned}$$

From the marginal likelihood of $\tilde{\mathbf{y}}_1$, we get an unbiased estimator of σ^2 given by

$$\tilde{\sigma}_1^2 = \tilde{\mathbf{y}}'_1 \left\{ \mathbf{I}_M - \widetilde{\mathbf{X}}_1 (\widetilde{\mathbf{X}}'_1 \widetilde{\mathbf{X}}_1)^+ \widetilde{\mathbf{X}}'_1 \right\} \tilde{\mathbf{y}}_1 / (N - R - r_{(\widetilde{\mathbf{X}}_1)}), \quad (21)$$

where $r_{(\widetilde{\mathbf{X}}_1)} = \text{rank}(\widetilde{\mathbf{X}}_1)$. It can be seen that $\tilde{\sigma}_1^2$ is a consistent estimator of σ^2 .

We now consider the estimation of \mathbf{G} . For this purpose, we handle two specific cases: (1) $\mathbf{G}_i = \psi \mathbf{D}_i$ for unknown scalar ψ and known matrix \mathbf{D}_i and (2) $\mathbf{G}_i = \boldsymbol{\Psi}$ and $r_i = r$ for $i = 1, \dots, k$.

[1] **Case of $\mathbf{G}_i = \psi \mathbf{D}_i$ for unknown scalar ψ and known matrix \mathbf{D}_i .** In this case, the marginal distribution of \mathbf{y} has $\mathcal{N}_N(\mathbf{X}\boldsymbol{\beta}, \sigma^2\boldsymbol{\Lambda})$, where $\boldsymbol{\Lambda} = \mathbf{I} + \psi \mathbf{Z}\mathbf{D}\mathbf{Z}'$ for $\mathbf{D} = \text{blockdiag}(\mathbf{D}_1, \dots, \mathbf{D}_k)$. Let $S = \mathbf{y}'(\mathbf{I} - \mathbf{H}_x)\mathbf{y}$ for $\mathbf{H}_x = \mathbf{X}(\mathbf{X}'\mathbf{X})^+\mathbf{X}'$, and the expectation is written as

$$E[S] = \sigma^2 \text{tr}[(\mathbf{I} - \mathbf{H}_x)(\mathbf{I} + \psi \mathbf{Z}\mathbf{D}\mathbf{Z}')] = \sigma^2\{N - r_x + \psi \text{tr}[\mathbf{Z}'(\mathbf{I} - \mathbf{H}_x)\mathbf{Z}\mathbf{D}]\},$$

for $r_x = \text{rank}(\mathbf{X})$. Thus, we get an estimator of ψ given by

$$\widehat{\psi} = \{S/\widehat{\sigma}_1^2 - (N - r_x)\}/\text{tr}[\mathbf{Z}'(\mathbf{I} - \mathbf{H}_x)\mathbf{Z}\mathbf{D}],$$

for $\widehat{\sigma}_1^2$ given in (21).

For the consistency of $\widehat{\psi}$, we assume that $\text{tr}[\mathbf{Z}'(\mathbf{I} - \mathbf{H}_x)\mathbf{Z}\mathbf{D}] = O(k)$ for large k . From the consistency of $\widehat{\sigma}_1^2$ and the fact that $S/\text{tr}[\mathbf{Z}'(\mathbf{I} - \mathbf{H}_x)\mathbf{Z}\mathbf{D}] = O_p(1)$, it follows that $\widehat{\psi} - \psi = \{S/\sigma^2 - (N - r_x)\}/\text{tr}[\mathbf{Z}'(\mathbf{I} - \mathbf{H}_x)\mathbf{Z}\mathbf{D}] - \psi + O_p(k^{-1/2})$. It can be seen that $\text{Var}[\widehat{\psi}] = \text{Var}[S]/\{\sigma^4(\text{tr}[\mathbf{Z}'(\mathbf{I} - \mathbf{H}_x)\mathbf{Z}\mathbf{D}])^2\} + O(k^{-1})$ and that $\text{Var}[S] = \sigma^4 2\text{tr}\{[(\mathbf{I} - \mathbf{H}_x)(\mathbf{I} + \psi \mathbf{Z}\mathbf{D}\mathbf{Z}')]^2\} = O(k)$. Hence, $\text{Var}[\widehat{\psi}] = O(k^{-1})$ under the assumption that $\text{tr}[\mathbf{Z}'(\mathbf{I} - \mathbf{H}_x)\mathbf{Z}\mathbf{D}] = O(k)$, which implies the consistency of $\widehat{\psi}$.

The estimator $\widehat{\psi}$ can take negative values with positive probability. Thus, we consider a truncated estimator

$$\widehat{\psi}^{TR} = \max\{\widehat{\psi}, R^{-2/3}\}, \quad (22)$$

which can be shown to be positive and consistent as $k \rightarrow \infty$.

[2] **Case of $\mathbf{G}_i = \boldsymbol{\Psi}$, an unknown matrix, and $r_i = r$ for $i = 1, \dots, k$.** In this case, we recall the transformed model given in (20). Let $\mathbf{S}_i = (\widetilde{\mathbf{y}}_{2i} - \widetilde{\mathbf{X}}_{2i}\widetilde{\boldsymbol{\beta}}_1)(\widetilde{\mathbf{y}}_{2i} - \widetilde{\mathbf{X}}_{2i}\widetilde{\boldsymbol{\beta}}_1)'$ for $\widetilde{\boldsymbol{\beta}}_1 = (\widetilde{\mathbf{X}}_1'\widetilde{\mathbf{X}}_1)^+\widetilde{\mathbf{X}}_1'\mathbf{y}$. It is noted that $E[\mathbf{S}_i] = \sigma^2\{(\mathbf{Z}'_i\mathbf{Z}_i)^{-1} + \boldsymbol{\Psi} + \widetilde{\mathbf{X}}_{2i}(\widetilde{\mathbf{X}}_1'\widetilde{\mathbf{X}}_1)^+\widetilde{\mathbf{X}}_{2i}'\}$, which yields that $\sum_{i=1}^k E[\mathbf{S}_i] = \sigma^2 \sum_{i=1}^k \{(\mathbf{Z}'_i\mathbf{Z}_i)^{-1} + \widetilde{\mathbf{X}}_{2i}(\widetilde{\mathbf{X}}_1'\widetilde{\mathbf{X}}_1)^+\widetilde{\mathbf{X}}_{2i}'\} + \sigma^2 k \boldsymbol{\Psi}$. Thus, $\boldsymbol{\Psi}$ can be estimated by

$$\widehat{\boldsymbol{\Psi}}^U = \frac{1}{k\widehat{\sigma}_1^2} \sum_{i=1}^k \mathbf{S}_i - \frac{1}{k} \sum_{i=1}^k \left\{ (\mathbf{Z}'_i\mathbf{Z}_i)^{-1} + \widetilde{\mathbf{X}}_{2i}(\widetilde{\mathbf{X}}_1'\widetilde{\mathbf{X}}_1)^+\widetilde{\mathbf{X}}_{2i}' \right\}, \quad (23)$$

for $\widehat{\sigma}_1^2$ given in (21). It can be verified that $\widehat{\boldsymbol{\Psi}}^U$ is consistent. For the proof, see the Appendix.

For the covariance matrix $\text{Cov}(\mathbf{v}_i) = \sigma^2\boldsymbol{\Psi}$, we consider the two cases of $\boldsymbol{\Psi}$; $\boldsymbol{\Psi}$ is fully unknown, and $\boldsymbol{\Psi} = \text{diag}(\psi_1, \dots, \psi_r)$. When $\boldsymbol{\Psi}$ is a fully unknown covariance matrix, $\widehat{\boldsymbol{\Psi}}$ should be estimated by a positive definite matrix. Let \mathbf{P} be an orthogonal matrix such that $\widehat{\boldsymbol{\Psi}}^U = \mathbf{P}'\text{diag}(\omega_1, \dots, \omega_m)\mathbf{P}$ where $\omega_1, \dots, \omega_r$ are eigenvalues of $\widehat{\boldsymbol{\Psi}}^U$. Then, we can use the truncated estimator

$$\widehat{\boldsymbol{\Psi}}^{TR} = \mathbf{P}' \text{diag}(\max\{\omega_1, R^{-2/3}\}, \dots, \max\{\omega_r, R^{-2/3}\}) \mathbf{P},$$

for $R = kr$. It can be shown that $\widehat{\Psi}^{TR}$ is consistent with order $\widehat{\Psi}^{TR} - \Psi = O_p(k^{-1/2})$.

When Ψ has a covariance structure, we can use the structure to construct an appropriate estimator based on $\widehat{\Psi}^U$. For instance, assume that $\Psi = \text{diag}(\psi_1, \dots, \psi_r)$. This case implies that $\mathbf{Z}_i \mathbf{v}_i$ in (1) is expressed as $\mathbf{Z}_i \mathbf{v}_i = \mathbf{z}_{i1} v_{i1} + \dots + \mathbf{z}_{ir} v_{ir}$ where $v_{ij} \sim \mathcal{N}(0, \sigma^2 \psi_j)$ and \mathbf{z}_{ij} is an $n_i \times 1$ vector for $j = 1, \dots, r$. This model may be useful when several factors of random effects are considered in practical situations. Then, we can estimate each ψ_i by

$$\widehat{\psi}_i^{TR} = \max\{(\widehat{\Psi}^U)_{ii}, R^{-2/3}\}, \quad (24)$$

where $(\widehat{\Psi}^U)_{ii}$ denotes the (i, i) diagonal element of $\widehat{\Psi}^U$. The resulting estimator of Ψ is given by $\widehat{\Psi}^{TR} = \text{diag}(\widehat{\psi}_1^{TR}, \dots, \widehat{\psi}_r^{TR})$.

5 Simulation Studies

We now investigate the numerical performances of the marginal and the conditional AIC's derived in the previous sections through simulation and compare them in terms of the frequencies of selecting the true model.

The simulation experiments have been carried out for $k = 20$, $r_1 = \dots = r_{20} = r = 1, 2, 3$ and $p = 7$. For the sample sizes n_i 's, n_i 's are generated as $n_i = 1 + \text{Bin}(8, 1/2)$ for $i = 1, \dots, k$, where $\text{Bin}(8, 1/2)$ is a random variable distributed as a binomial distribution with mean 4 and success probability 1/2. For the $N \times p$ matrix \mathbf{X} of the regressor variables in the model (2), the row vectors $\mathbf{x}_1, \dots, \mathbf{x}_N$ for $\mathbf{X}' = (\mathbf{x}_1, \dots, \mathbf{x}_N)$ are generated as mutually independent random variables distributed as $\mathcal{N}_p(\mathbf{0}, \Sigma_x)$ where $\Sigma_x = (1 - \rho_x)\mathbf{I}_p + \rho_x \mathbf{J}_p$ for $\rho_x = 0.3$, where $\mathbf{J}_p = \mathbf{j}_p \mathbf{j}_p'$ for $\mathbf{j}_p = (1, \dots, 1)'$, a p -vector of ones. For the $n \times r$ matrix \mathbf{Z}_i 's in the model (1), the row vectors in \mathbf{Z}_i are generated as mutually independent random variables distributed as $\mathcal{N}_r(\mathbf{0}, \Sigma_z)$ where $\Sigma_z = (1 - \rho_z)\mathbf{I}_r + \rho_z \mathbf{J}_r$ for $\rho_z = 0.3$. In this experiment, we assume that the true model is given by

$$(p^*) \quad \mathbf{y} = \mathbf{X} \boldsymbol{\beta}^* + \mathbf{Z} \mathbf{v} + \boldsymbol{\epsilon},$$

where $1 \leq p^* \leq 7$, $\boldsymbol{\beta}^* = (\beta_1, \dots, \beta_{p^*}, 0, \dots, 0)'$, and \mathbf{v} and $\boldsymbol{\epsilon}$ are mutually independent random variables having $\mathbf{v} \sim \mathcal{N}_R(\mathbf{0}, \sigma^2 \psi \mathbf{I}_R)$ and $\boldsymbol{\epsilon} \sim \mathcal{N}_N(\mathbf{0}, \sigma^2 \mathbf{I}_N)$; it may be noted that when $\text{Cov}(\mathbf{v}) = \sigma^2 \psi \mathbf{D}$, where \mathbf{D} is known, we may assume without any loss of generality that $\mathbf{D} = \mathbf{I}_R$, as it can be absorbed with the matrix \mathbf{Z} by defining $\mathbf{Z}^* = \mathbf{Z} \mathbf{D}^{-1/2}$. Here, we handle the cases that $\sigma^2 = 1$ and $\psi = 0.01, 0.5, 1.0$. Also, β_ℓ for $1 \leq \ell \leq p^*$ is generated as a random variable distributed as $\beta_\ell = 2(-1)^{\ell+1} \{1 + U(0, 1)\}$ for a uniform random variable $U(0, 1)$ on the interval $(0, 1)$. Let (m) be the set $\{1, \dots, m\}$, and we write the model using the first m regressor variables β_1, \dots, β_m by M_m or simply (m) . Then, the full model is (7) and the true model is (p^*) . As candidate models, we consider the nested subsets (1), \dots , (7) of $\{1, \dots, 7\}$, namely,

$$(m) \quad \mathbf{y} = \mathbf{X} \boldsymbol{\beta}^{(m)} + \mathbf{Z} \mathbf{v} + \boldsymbol{\epsilon},$$

where $\boldsymbol{\beta}^{(m)} = (\beta_1, \dots, \beta_m, 0, \dots, 0)'$.

In the simulation experiments, 10 observations of the regressor variables \mathbf{X} and \mathbf{Z} are generated, and for each observation of \mathbf{X} and \mathbf{Z} , 30 observations of the response variable \mathbf{y} are generated from the true model (p^*) for $p^* = 2, 4, 6$. Thus, we have $10 \times 30 (= 300)$ total data sets. For each data set, we calculate the values of $mAIC$ given in (7), $cAIC$ given in (9) and $CAIC_{ML}$, $CAIC_{SL}$ and $CAIC_{GL}$ given (16), (17) and (18), respectively, for the eight candidate models (1), \dots , (8), and we select the models minimizing the values of the information criteria. For each criterion and each candidate model (m), the number of selecting the model (m) is counted for 300 data set. We thus obtain the frequencies of the model (m) selected by the criteria by dividing the number by 300. These frequencies are reported in Table 1, where standard deviations in selecting the true model are less than 0.02. From the table, we can see that the proposed conditional AIC's $CAIC_{ML}$, $CAIC_{SL}$ and $CAIC_{GL}$ are superior to $mAIC$ and $cAIC$ for most of the cases.

When the random effects \mathbf{v}_i has the covariance matrix such that $\boldsymbol{\Psi} = \text{diag}(\psi_1, \dots, \psi_r)$, we next investigate the similar performances of the five criteria $mAIC$, $cAIC$, $CAIC_{ML}$, $CAIC_{SL}$ and $CAIC_{GL}$ with known and unknown $\boldsymbol{\Psi}$, where ψ_i is estimated by (24) in the unknown case of $\boldsymbol{\Psi}$. The frequencies of selecting the true model are reported in Table 2 for the balanced case of $n_1 = \dots = n_{20} = 10$, $k = 20$ and $r = 2, 3$. This numerical results show that the conditional AIC's $CAIC_{ML}$, $CAIC_{SL}$ and $CAIC_{GL}$ are better than $mAIC$ and $cAIC$. It is interesting to note that the performance of $CAIC_{ML}$ does not depend on whether $\boldsymbol{\Psi}$ is known or unknown, since $CAIC_{ML}$ does not include $\boldsymbol{\Psi}$ or its estimator. This means that $CAIC_{ML}$ can be used even if $\boldsymbol{\Psi}$ cannot be estimated appropriately, or n_i 's are small.

6 Concluding Remarks

In this paper, we have considered linear mixed models. To select the fixed-effects variables, we have derived three conditional Akaike information criteria $CAIC_{ML}$, $CAIC_{SL}$ and $CAIC_{GL}$, and have shown that these $CAIC$'s perform better than $mAIC$ as well as better than $cAIC$, proposed by Vaida and Blanchard (2005). We have also considered the case when $\text{Cov}(\mathbf{v}_i) = \sigma^2\psi\mathbf{D}_i$ as well as when $\text{Cov}(\mathbf{v}_i) = \sigma^2\boldsymbol{\Psi}$, but $r_i = r$, where ψ is a scalar unknown parameter, and $\boldsymbol{\Psi}$ is an unknown covariance matrix. The proposed $CAIC$'s perform better than $mAIC$ and $cAIC$. However, when n_i 's are small, it is recommended to use $CAIC_{ML}$ when the matrix $\boldsymbol{\Psi}$ is completely unknown as it does not depend on the unknown parameters.

Table 1: Frequencies selected by the five criteria $mAIC$, $cAIC$, $CAIC_{ML}$, $CAIC_{SL}$ and $CAIC_{GL}$, abbreviated by AIC , C_{VB} , C_{ML} , C_{EB} and C_{GL} , in 300 replications for the unbalanced case of n_i and $k = 20$: the dimension of a full model is $p = 7$ and the true model is $(p^*) = \{1, \dots, p^*\}$

M_k	known ψ					unknown ψ				
	AIC	C_{VB}	C_{ML}	C_{EB}	C_{GL}	AIC	C_{VB}	C_{ML}	C_{EB}	C_{GL}
$p^* = 2, \psi = 0.01, r = 1$										
(1)	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0
(2)	76.3	76.3	84.0	84.7	84.7	75.3	74.7	84.0	85.0	84.7
(3)	11.0	11.0	9.7	9.0	9.0	11.3	12.7	9.7	8.7	9.0
(4)	4.3	4.3	3.3	2.3	2.3	5.0	4.7	3.3	2.7	2.7
(5)	4.3	4.3	1.0	2.3	2.3	3.7	3.7	1.0	2.0	2.0
(6)	2.0	2.0	1.0	1.3	1.3	2.3	2.0	1.0	1.0	1.0
(7)	2.0	2.0	1.0	0.3	0.3	2.3	2.3	1.0	0.7	0.7
$p^* = 4, \psi = 0.5, r = 2$										
(1)	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0
(2)	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0
(3)	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0
(4)	76.3	75.0	94.3	98.3	96.3	77.0	72.0	94.3	98.7	97.3
(5)	12.7	12.0	3.3	1.7	2.7	12.0	11.7	3.3	1.3	2.0
(6)	8.0	9.7	2.0	0.0	1.0	8.0	9.3	2.0	0.0	0.7
(7)	3.0	3.3	0.3	0.0	0.0	3.0	7.0	0.3	0.0	0.0
$p^* = 6, \psi = 1.0, r = 3$										
(1)	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0
(2)	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0
(3)	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0
(4)	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	1.0	0.0
(5)	0.0	0.0	0.0	0.0	0.03	0.0	0.0	0.0	9.7	0.0
(6)	88.0	87.3	98.7	96.0	99.7	87.0	77.3	98.7	89.0	100.0
(7)	12.0	12.7	1.3	0.0	0.3	13.0	22.7	1.3	0.3	0.0

Table 2: Frequencies selected by the five criteria $mAIC$, $cAIC$, $CAIC_{ML}$, $CAIC_{SL}$ and $CAIC_{GL}$, abbreviated by AIC , C_{VB} , C_{ML} , C_{EB} and C_{GL} , in 300 replications for known and unknown Ψ and the unbalanced case of $n_1 = \dots = n_{20} = 10$ and $k = 20$: the dimension of a full model is $p = 7$ and the true model is $(p^*) = \{1, \dots, p^*\}$

M_k	known Ψ					unknown Ψ				
	AIC	C_{VB}	C_{ML}	C_{EB}	C_{GL}	AIC	C_{VB}	C_{ML}	C_{EB}	C_{GL}
$p^* = 2, (\psi_1, \psi_2) = (0.01, 1.0), r = 2$										
(1)	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0
(2)	78.0	74.7	87.7	94.7	86.3	75.7	78.0	87.7	95.0	89.7
(3)	10.7	12.0	8.0	3.7	8.3	12.3	10.3	8.0	2.3	5.7
(4)	3.7	4.0	1.3	0.7	1.3	3.0	4.0	1.3	0.7	1.7
(5)	3.0	4.3	1.3	0.3	2.0	4.0	3.0	1.3	1.0	2.0
(6)	2.7	2.3	1.0	0.3	1.0	2.3	2.3	1.0	0.7	0.7
(7)	2.0	2.7	0.7	0.3	1.0	2.7	2.3	0.7	0.3	0.3
$p^* = 4, (\psi_1, \psi_2, \psi_3) = (0.01, 0.5, 1.0), r = 3$										
(1)	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0
(2)	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0
(3)	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0
(4)	78.0	76.7	86.0	95.7	91.0	74.7	73.0	86.0	93.3	87.3
(5)	10.3	11.7	8.7	3.7	5.3	11.3	13.0	8.7	5.0	7.0
(6)	6.7	7.3	3.7	0.3	2.0	8.3	8.0	3.7	1.3	3.0
(7)	5.0	4.3	1.7	0.3	1.7	5.7	6.0	1.7	0.3	2.7
$p^* = 6, (\psi_1, \psi_2, \psi_3) = (1.0, 1.0, 1.0), r = 3$										
(1)	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0
(2)	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0
(3)	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0
(4)	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0
(5)	0.0	0.0	0.0	0.0	0.0	0.7	6.3	0.0	0.0	0.0
(6)	85.3	85.3	91.3	95.3	91.0	85.7	85.3	91.3	97.0	93.7
(7)	14.7	14.7	8.7	4.7	9.0	14.3	14.7	8.7	3.0	6.3

A Appendix

A.1 Expectation of a ratio of quadratic forms

We here evaluate the expected value of the ratio of two quadratic forms which have been used to derive the bias terms of the conditional AIC.

Lemma A.1 *Assume that \mathbf{A} is an $N \times N$ symmetric matrix, and that \mathbf{H} is an idempotent matrix of rank q . Then,*

$$E \left[\frac{\mathbf{u}' \mathbf{A} \mathbf{u}}{\mathbf{u}' (\mathbf{I}_N - \mathbf{H}) \mathbf{u}} \right] = \frac{\text{tr } \mathbf{A}}{N - q - 2} - \frac{2 \text{tr} [A(\mathbf{I}_N - \mathbf{H})]}{(N - q)(N - q - 2)},$$

where $\mathbf{u} \sim \mathcal{N}_N(\mathbf{0}, \mathbf{I}_N)$.

Proof. Since \mathbf{H} is an idempotent matrix of rank q , we can write $\mathbf{H} = \mathbf{\Gamma}_2 \mathbf{\Gamma}_2'$, $\mathbf{\Gamma}_2 = N \times q$, and $\mathbf{\Gamma}_2' \mathbf{\Gamma}_2 = \mathbf{I}_q$. Let $\mathbf{\Gamma}_1$ be $N \times (N - q)$ matrix such that $\mathbf{\Gamma}_1' \mathbf{\Gamma}_2 = \mathbf{0}$ and $\mathbf{\Gamma}_1' \mathbf{\Gamma}_1 = \mathbf{I}_{N - q}$. Then, $\mathbf{\Gamma} = (\mathbf{\Gamma}_1, \mathbf{\Gamma}_2)$ is an $N \times N$ orthogonal matrix and hence $\mathbf{\Gamma}_1 \mathbf{\Gamma}_1' + \mathbf{\Gamma}_2 \mathbf{\Gamma}_2' = \mathbf{I}_N$, and $\mathbf{I}_N - \mathbf{H} = \mathbf{\Gamma}_1 \mathbf{\Gamma}_1'$. Let $\mathbf{v} = \mathbf{\Gamma}' \mathbf{u}$, $\mathbf{v}_1 = \mathbf{\Gamma}_1' \mathbf{u}$, $\mathbf{v}_2 = \mathbf{\Gamma}_2' \mathbf{u}$ and

$$\mathbf{B} = \mathbf{\Gamma} \mathbf{A} \mathbf{\Gamma}' = \begin{pmatrix} \mathbf{B}_{11} & \mathbf{B}_{12} \\ \mathbf{B}'_{12} & \mathbf{B}_{22} \end{pmatrix}.$$

Then,

$$Q = \frac{\mathbf{u}' \mathbf{A} \mathbf{u}}{\mathbf{u}' (\mathbf{I}_N - \mathbf{H}) \mathbf{u}} = \frac{\mathbf{v}' \mathbf{B} \mathbf{v}}{\mathbf{v}'_1 \mathbf{v}_1} = \frac{\mathbf{v}'_1 \mathbf{B}_{11} \mathbf{v}_1 + 2 \mathbf{v}'_1 \mathbf{B}_{12} \mathbf{v}_2 + \mathbf{v}'_2 \mathbf{B}_{22} \mathbf{v}_2}{\mathbf{v}'_1 \mathbf{v}_1}.$$

Noting that $\mathbf{v} \sim \mathcal{N}_N(\mathbf{0}, \mathbf{I}_N)$, we can see that

$$\begin{aligned} E[Q] &= \frac{\text{tr } \mathbf{B}_{11}}{N - q} + \frac{\text{tr } \mathbf{B}_{22}}{N - q - 2} = \frac{\text{tr } \mathbf{A} \mathbf{\Gamma}_1 \mathbf{\Gamma}_1'}{N - q} + \frac{\text{tr } \mathbf{A} \mathbf{\Gamma}_2 \mathbf{\Gamma}_2'}{N - q - 2} \\ &= \frac{\text{tr } \mathbf{A} (\mathbf{I}_N - \mathbf{H})}{N - q} + \frac{\text{tr } \mathbf{A} \mathbf{H}}{N - q - 2} \\ &= \frac{\text{tr } \mathbf{A}}{N - q - 2} - \frac{2 \text{tr } \mathbf{A} (\mathbf{I}_N - \mathbf{H})}{(N - q)(N - q - 2)}, \end{aligned}$$

which proves Lemma A.1. ■

A.2 Consistency of $\widehat{\Psi}^U$

From the consistency of $\tilde{\sigma}_1^2$ and the fact that $\sum_{i=1}^k \mathbf{S}_i / k = O_p(1)$, it follows that $\widehat{\Psi}^U = \Psi^* / \sigma^2 + O_p(k^{-1/2})$ for $\Psi^* = \sum_{i=1}^k \mathbf{S}_i / k - \sigma^2 \sum_{i=1}^k \{(\mathbf{Z}'_i \mathbf{Z}_i)^{-1} + \widetilde{\mathbf{X}}_{2i} (\widetilde{\mathbf{X}}'_1 \widetilde{\mathbf{X}}_1)^+ \widetilde{\mathbf{X}}'_{2i}\} / k$. Note that $\mathbf{S}_i = (\tilde{\mathbf{y}}_{2i} - \widetilde{\mathbf{X}}_{2i} \widetilde{\beta}_1) (\tilde{\mathbf{y}}_{2i} - \widetilde{\mathbf{X}}_{2i} \widetilde{\beta}_1)'$ and that $\tilde{\mathbf{y}}_{2i}$ and $\widetilde{\beta}_1$ are mutually independently

distributed as $\tilde{\mathbf{y}}_{2i} \sim \mathcal{N}(\widetilde{\mathbf{X}}_{2i}\boldsymbol{\beta}, \boldsymbol{\Sigma}_i)$ and $\tilde{\boldsymbol{\beta}}_1 \sim \mathcal{N}(\boldsymbol{\beta}, \sigma^2(\widetilde{\mathbf{X}}_1'\widetilde{\mathbf{X}}_1)^+)$, where $\boldsymbol{\Sigma}_i = \sigma^2\{\boldsymbol{\Psi} + (\mathbf{Z}'_i\mathbf{Z}_i)^{-1}\}$. Then, $\widehat{\boldsymbol{\Psi}}^* - \sigma^2\boldsymbol{\Psi}$ is expressed as

$$\begin{aligned}\widehat{\boldsymbol{\Psi}}^* - \sigma^2\boldsymbol{\Psi} &= \frac{1}{k} \sum_{i=1}^k \left\{ (\tilde{\mathbf{y}}_{2i} - \widetilde{\mathbf{X}}_{2i}\boldsymbol{\beta})(\tilde{\mathbf{y}}_{2i} - \widetilde{\mathbf{X}}_{2i}\boldsymbol{\beta})' - \boldsymbol{\Sigma}_i \right\} \\ &\quad + \frac{1}{k} \sum_{i=1}^k \left\{ \widetilde{\mathbf{X}}_{2i}(\tilde{\boldsymbol{\beta}}_1 - \boldsymbol{\beta})(\tilde{\boldsymbol{\beta}}_1 - \boldsymbol{\beta})'\widetilde{\mathbf{X}}_{2i}' - \sigma^2\widetilde{\mathbf{X}}_{2i}(\widetilde{\mathbf{X}}_1'\widetilde{\mathbf{X}}_1)^+\widetilde{\mathbf{X}}_{2i}' \right\} \\ &\quad - \frac{1}{k} \sum_{i=1}^k \left\{ (\tilde{\mathbf{y}}_{2i} - \widetilde{\mathbf{X}}_{2i}\boldsymbol{\beta})(\tilde{\boldsymbol{\beta}}_1 - \boldsymbol{\beta})'\widetilde{\mathbf{X}}_{2i}' + \widetilde{\mathbf{X}}_{2i}(\tilde{\boldsymbol{\beta}}_1 - \boldsymbol{\beta})(\tilde{\mathbf{y}}_{2i} - \widetilde{\mathbf{X}}_{2i}\boldsymbol{\beta})' \right\} \\ &= \mathbf{K}_1 + \mathbf{K}_2 - \mathbf{K}_3. \quad (\text{say})\end{aligned}$$

It is easy to see that $E[\mathbf{K}_j] = 0$ for $j = 1, 2, 3$. For the consistency of $\boldsymbol{\Psi}^U$ or $\boldsymbol{\Psi}^*$, it is sufficient to show that $E[\text{tr}[\mathbf{K}_j^2]] = O(k^{-1})$ for $j = 1, 2, 3$. Let $\mathbf{u}_i = \tilde{\mathbf{y}}_{2i} - \widetilde{\mathbf{X}}_{2i}\boldsymbol{\beta}$, and it has $\mathcal{N}(\mathbf{0}, \boldsymbol{\Sigma}_i)$. Then,

$$E[\text{tr}[\mathbf{K}_1^2]] = \frac{1}{k^2} E\left[\text{tr}\left[\left(\sum_{i=1}^k \mathbf{u}_i\mathbf{u}_i'\right)^2\right]\right] - \frac{1}{k^2} \text{tr}\left[\left(\sum_{i=1}^k \boldsymbol{\Sigma}_i\right)^2\right]$$

and

$$\begin{aligned}E\left[\text{tr}\left[\sum_{i=1}^k \mathbf{u}_i\mathbf{u}_i'\right]^2\right] &= E\left[\sum_{i=1}^k (\mathbf{u}_i'\mathbf{u}_i)^2 + \sum_{i=1}^k \sum_{j \neq i} \text{tr}[\mathbf{u}_i\mathbf{u}_i'\mathbf{u}_j\mathbf{u}_j']\right] \\ &= \sum_{i=1}^k \left\{ 2\text{tr}[\boldsymbol{\Sigma}_i^2] + (\text{tr}[\boldsymbol{\Sigma}_i])^2 \right\} + \sum_{i=1}^k \sum_{j \neq i} \text{tr}[\boldsymbol{\Sigma}_i\boldsymbol{\Sigma}_j] \\ &= 2 \sum_{i=1}^k \text{tr}[\boldsymbol{\Sigma}_i^2] + \text{tr}\left[\left(\sum_{i=1}^k \boldsymbol{\Sigma}_i\right)^2\right].\end{aligned}$$

Hence,

$$E[\text{tr}[\mathbf{K}_1^2]] = 2 \sum_{i=1}^k \text{tr}[\boldsymbol{\Sigma}_i^2]/k^2,$$

which has the order $O(k^{-1})$. Similarly, we can see that

$$\begin{aligned}E[\text{tr}[\mathbf{K}_2^2]] &= 2\sigma^4 \text{tr}\left[\left(\sum_{i=1}^k \widetilde{\mathbf{X}}_{2i}(\widetilde{\mathbf{X}}_1'\widetilde{\mathbf{X}}_1)^+\widetilde{\mathbf{X}}_{2i}'\right)^2\right]/k^2, \\ E[\text{tr}[\mathbf{K}_3^2]] &= 4\sigma^2 \sum_{i=1}^k \text{tr}[\boldsymbol{\Sigma}_i] \text{tr}\left[\widetilde{\mathbf{X}}_{2i}(\widetilde{\mathbf{X}}_1'\widetilde{\mathbf{X}}_1)^+\widetilde{\mathbf{X}}_{2i}'\right]/k^2,\end{aligned}$$

both of which have the order $O(k^{-1})$ under the assumption that Σ_i is bounded and $\sum_{i=1}^k \widetilde{\mathbf{X}}_{2i}(\widetilde{\mathbf{X}}_1' \widetilde{\mathbf{X}}_1)^+ \widetilde{\mathbf{X}}_{2i}' = O(1)$. Thus, the consistency of $\widehat{\Psi}^U$ is verified. ■

Acknowledgments. The research of the first author was supported by NSERC. The research of the second author was supported in part by a grant from the Ministry of Education, Japan, No. 19200020.

References

- [1] Akaike, H. (1973). Information theory and an extension of the maximum likelihood principle. In *2nd International Symposium on Information Theory*, (B.N. Petrov and Csaki, F, eds.), 267-281, Akademia Kiado, Budapest.
- [2] Akaike, H. (1974). A new look at the statistical model identification. System identification and time-series analysis. *IEEE Trans. Autom. Contr.*, **AC-19**, 716-723.
- [3] Henderson, C.R. (1950). Estimation of genetic parameters. *Ann. Math. Statist.*, **21**, 309-310.
- [4] Konishi, S. and Kitagawa, G. (2007). *Information Criteria and Statistical Modeling*. Springer.
- [5] Spiegelhalter, D.J., Best, N.G., Carlin, B.P., and van der Linde, A. (2002). Bayesian measures of model complexity and fit (with Discussion). *J. Royal Statist. Soc.*, **B 64**, 583-639.
- [6] Vaida, F., and Blanchard, S. (2005). Conditional Akaike information for mixed-effects models. *Biometrika*, **92**, 351-370.