



Contents lists available at SciVerse ScienceDirect

Linear Algebra and its Applications

journal homepage: www.elsevier.com/locate/laa



Bias-corrected AIC for selecting variables in multinomial logistic regression models[☆]

Hirokazu Yanagihara^{a,*,1}, Ken-ichi Kamo^b, Shinpei Imori^a, Kenichi Satoh^{c,2}

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

^b Department of Liberal Arts and Sciences, Sapporo Medical University, South 1, West 17, Chuo-ku, Sapporo 060-8543, Japan

^c Department of Environmetrics and Biometrics, Research Institute for Radiation Biology and Medicine, Hiroshima University, 1-2-3 Kasumi, Minami-ku, Hiroshima 734-8553, Japan

ARTICLE INFO

Article history:

Received 22 January 2012

Accepted 25 January 2012

Available online 24 February 2012

Submitted by R.A. Brualdi

AMS classification:

62H12

62J12

Keywords:

AIC

Bias correction

Multinomial logistic model

MLE

Partial differential operator

Variable selection

ABSTRACT

In this paper, we consider the bias correction of Akaike's information criterion (AIC) for selecting variables in multinomial logistic regression models. For simplifying a formula of the bias-corrected AIC, we calculate the bias of the AIC to a risk function through the expectations of partial derivatives of the negative log-likelihood function. As a result, we can express the bias correction term of the bias-corrected AIC with only three matrices consisting of the second, third, and fourth derivatives of the negative log-likelihood function. By conducting numerical studies, we verify that the proposed bias-corrected AIC performs better than the crude AIC.

© 2012 Elsevier Inc. All rights reserved.

1. Introduction

A multinomial logistic regression model is a regression model that generalizes a logistic regression by allowing more than two discrete response variables. When categories are unordered, the

[☆] The authors would like to thank the associate editor and the reviewer for valuable comments.

* Corresponding author.

E-mail address: yanagi@math.sci.hiroshima-u.ac.jp (H. Yanagihara).

¹ Supported by the Ministry of Education, Science, Sports, and Culture, Grant-in-Aid for Challenging Exploratory Research, #22650058, 2010–2012.

² Supported by the Ministry of Education, Science, Sports, and Culture, Grant-in-Aid for Young Scientists (B), #23700337, 2011–2014.

multinomial logistic model is one strategy often used. The multinomial logistic regression model has been introduced in many textbooks for applied statistical analysis (see e.g. [1, Chapter 8.1]), and even now it is widely used in biometrics, econometrics, psychometrics, sociometrics, and many other fields of applications for the prediction of probabilities of different possible outcomes of categorically distributed response variables by a set of explanatory variables (e.g. [2–4]). In addition, the multinomial logistic regression model can be easily fitted to observed data by using the “vglm” function in “R” [5] that is a free software environment for statistical computing and graphics. Since we would like to specify the factors affecting the probabilities of response variables in the regression analysis, searching for the optimal subset of explanatory variables is important.

Akaike's information criterion (AIC) proposed by Akaike [6,7] is widely used for selecting the best model among the candidate models (for details of statistical model selection, see e.g. [8–10]). The model having the smallest AIC among the candidate models is regarded as the best model. In the multinomial logistic regression model, the subset of explanatory variables in the best model is the best subset. However, the AIC may perform poorly; that is, a model having too many parameters tends to be chosen as the best model when the sample size is small or the number of unknown parameters is large. Such a problem is often resolved by using a bias-corrected AIC (see e.g. [8, Chapter 2.4]). The AIC is an estimator of the risk function consisting of predictive Kullback–Leibler (K–L) information [11], which measures the discrepancy between the true model and the candidate model. The order of the bias of the AIC is $O(n^{-1})$ when the candidate model includes the true model, where n is the sample size. Although the AIC is an asymptotic unbiased estimator of the risk function, it has a nonnegligible bias to the risk function when the sample size is small or the number of unknown parameters is large. A bias-corrected AIC called CAIC in this paper improves the bias of AIC to $O(n^{-2})$ under the assumption that the candidate model includes the true model.

The CAIC in the logistic regression models was obtained by Yanagihara et al. [12]. But the CAIC in multinomial regression models has not been derived yet, although the multinomial logistic regression model is widely used in many application fields. The CAIC can be obtained by removing the bias of the AIC to the risk function from the AIC with the use of a consistent estimator of the bias. The bias of the AIC to the risk function is then evaluated by moments of the maximum likelihood estimator (MLE) of unknown parameters. Since such moments are represented by the moments of response variables, calculating the moments of response variables is essential for evaluating the bias of the AIC in the ordinary bias correction method, which is used in [12,13], etc. However, in the case of multiple response variables, calculations and expressions of the moments of the MLE mediated by the moments of response variables become complicated. Hence, without directly calculating the moments of response variables, we derive the moments of the MLE by using expectations of the partial derivatives of the negative log-likelihood function. This different approach from the ordinary bias correction method leads to a simple expression of the bias correction term of the CAIC. In fact, the bias correction term of our CAIC is represented by only three matrices consisting of the second, third, and fourth derivatives of the negative log-likelihood function.

The present paper is organized as follows. In Section 2, we give a stochastic expansion of the MLE. In Section 3, the CAIC in multinomial logistic regression models is proposed. In Section 4, we verify that the proposed CAIC has better performance than the AIC by conducting numerical experiments. In Section 5, we conclude our discussion. Technical details are provided in Appendix.

2. Stochastic expansion of MLE

Suppose that the data consists of a sequence $\{\mathbf{y}_i, \mathbf{x}_i\}$, where $\mathbf{y}_1, \dots, \mathbf{y}_m$ are r -dimensional independent unordered discrete random vectors, and $\mathbf{x}_1, \dots, \mathbf{x}_m$ are k -dimensional vectors of known constants. Let $\boldsymbol{\beta} = (\beta_1, \dots, \beta_{kr})'$ be a kr -dimensional unknown regression coefficient vector that is partitioned as $\boldsymbol{\beta} = (\boldsymbol{\beta}'_1, \dots, \boldsymbol{\beta}'_r)'$, where $\boldsymbol{\beta}_j$ is a k -dimensional vector denoted by $\boldsymbol{\beta}_j = (\beta_{(j-1)k+1}, \dots, \beta_{jk})'$. In the multinomial logistic regression model, we assume that $(y_{i0}, \mathbf{y}'_i)' = (y_{i0}, y_{i1}, \dots, y_{ir})'$ is distributed according to the multinomial distribution with the number of events n_i ($n_i = \sum_{j=0}^r y_{ij}$, $n = \sum_{i=1}^m n_i$) and the cell probability vector $(p_{i0}(\boldsymbol{\beta}), \mathbf{p}_i(\boldsymbol{\beta})')'$, given by

$$\begin{aligned}
 p_{i0}(\boldsymbol{\beta}) &= \frac{1}{1 + \sum_{j=1}^r \exp(\mathbf{x}_i' \boldsymbol{\beta}_j)} \\
 \mathbf{p}_i(\boldsymbol{\beta}) &= (p_{i1}(\boldsymbol{\beta}), \dots, p_{ir}(\boldsymbol{\beta}))' \\
 &= \left(\frac{\exp(\mathbf{x}_i' \boldsymbol{\beta}_1)}{1 + \sum_{j=1}^r \exp(\mathbf{x}_i' \boldsymbol{\beta}_j)}, \dots, \frac{\exp(\mathbf{x}_i' \boldsymbol{\beta}_r)}{1 + \sum_{j=1}^r \exp(\mathbf{x}_i' \boldsymbol{\beta}_j)} \right)'.
 \end{aligned} \tag{1}$$

The MLE of $\boldsymbol{\beta}$ is obtained by maximizing the log-likelihood function. In this paper, we assume that the maximum exists, i.e., that the supremum does not occur at the boundary. By omitting the constant term, the log-likelihood function of the multinomial logistic regression model in (1) is expressed as

$$\ell(\boldsymbol{\beta}) = \sum_{i=1}^m \left\{ (\mathbf{y}_i \otimes \mathbf{x}_i)' \boldsymbol{\beta} - n_i \log \left(1 + \sum_{j=1}^r \exp(\mathbf{x}_i' \boldsymbol{\beta}_j) \right) \right\}.$$

Hence, the MLE of $\boldsymbol{\beta}$ is given by

$$\hat{\boldsymbol{\beta}} = \arg \max_{\boldsymbol{\beta}} \ell(\boldsymbol{\beta}).$$

To evaluate a bias of the AIC to the risk function, a stochastic expansion of $\hat{\boldsymbol{\beta}}$ is needed. The purpose of this section is to obtain the stochastic expansion $\hat{\boldsymbol{\beta}}$ up to the order $n^{-3/2}$. Two cases serve as a framework for asymptotic approximations:

Case (i): n_j 's are fixed, and $m \rightarrow \infty$,

Case (ii): m is fixed, $n_j \rightarrow \infty$ and $\rho_j^{-1} = n/n_j = O(1)$ for each j .

Although we only consider Case (i) in this paper, our formula can also be applied to Case (ii).

Suppose that $\mathbf{x}_1, \dots, \mathbf{x}_m$ are members of an admissible compact set \mathcal{F} , i.e., $\mathbf{x}_1, \dots, \mathbf{x}_m \in \mathcal{F}$. To expand the MLE, we consider the following regularity assumptions (see e.g. [14]):

C.1 : $\boldsymbol{\beta} \in \mathcal{B}$, where \mathcal{B} is a convex and open set in \mathbb{R}^k ,

C.2 : $(\mathbf{I}_r \otimes \mathbf{x}_i)' \boldsymbol{\beta} \in \Theta^0$, $i = 1, 2, \dots$, for all $\boldsymbol{\beta} \in \mathcal{B}$, where Θ^0 is the interior of the convex natural parameter space $\Theta \subset \mathbb{R}^r$,

C.3 : $\exists m_0$ s.t. $\mathbf{X}'\mathbf{X}$ has the full rank for $m \geq m_0$, where $\mathbf{X} = (\mathbf{x}_1, \dots, \mathbf{x}_m)'$.

Condition C.1 guarantees the uniqueness of the MLE if it exists. Condition C.2 is necessary to obtain the multinomial logistic regression model for all $\boldsymbol{\beta}$. Condition C.3 ensures that $\sum_{i=1}^m n_i \boldsymbol{\Sigma}_i(\boldsymbol{\beta}) \otimes \mathbf{x}_i \mathbf{x}_i'$ is positive definite for all $\boldsymbol{\beta} \in \mathcal{B}$, $m \geq m_0$, where

$$\boldsymbol{\Sigma}_i(\boldsymbol{\beta}) = \text{diag}(\mathbf{p}_i(\boldsymbol{\beta})) - \mathbf{p}_i(\boldsymbol{\beta}) \mathbf{p}_i(\boldsymbol{\beta})'. \tag{2}$$

Moreover, we prepare the following additional conditions to assure weak consistency and asymptotic normality of $\hat{\boldsymbol{\beta}}$, which can be derived by slightly modifying the results in [14]:

C.4 : sequence $\{\mathbf{x}_i\}$ lies in \mathcal{F} with $(\mathbf{I}_r \otimes \mathbf{x})' \boldsymbol{\beta} \in \Theta^0$, for all $\boldsymbol{\beta} \in \mathcal{B}$,

C.5 : $\liminf_{m \rightarrow \infty} \lambda(\sum_{i=1}^m n_i \boldsymbol{\Sigma}_i(\boldsymbol{\beta}) \otimes \mathbf{x}_i \mathbf{x}_i' / n) > 0$, where $\lambda(\mathbf{A})$ indicates the smallest eigenvalue of symmetric matrix \mathbf{A} .

According to Corollary 1 in [14], $\hat{\boldsymbol{\beta}}$ has weak consistency and asymptotic normality under these conditions. Furthermore, from Condition C.5, $\sum_{i=1}^m n_i \boldsymbol{\Sigma}_i(\boldsymbol{\beta}) \otimes \mathbf{x}_i \mathbf{x}_i' = O(n)$ is satisfied. Under the assumption that all conditions are satisfied, $\hat{\boldsymbol{\beta}}$ can be formally expanded as follows:

$$\hat{\beta} = \beta + \frac{1}{\sqrt{n}}\mathbf{b}_1 + \frac{1}{n}\mathbf{b}_2 + \frac{1}{n\sqrt{n}}\mathbf{b}_3 + O_p(n^{-2}), \quad (3)$$

where \mathbf{b}_1 , \mathbf{b}_2 , and \mathbf{b}_3 are kr -dimensional random vectors. The purpose of this section is achieved by specifying \mathbf{b}_1 , \mathbf{b}_2 , and \mathbf{b}_3 .

Since the log-likelihood function $\ell(\beta)$ is a maximum at $\beta = \hat{\beta}$, the first derivative of $\ell(\beta)$ becomes $\mathbf{0}_{kr}$ at $\beta = \hat{\beta}$, i.e.,

$$\left. \frac{\partial \ell(\beta)}{\partial \beta} \right|_{\beta=\hat{\beta}} = \sum_{i=1}^m \left\{ (\mathbf{y}_i \otimes \mathbf{x}_i) - n_i(\mathbf{p}_i(\hat{\beta}) \otimes \mathbf{x}_i) \right\} = \mathbf{0}_{kr}, \quad (4)$$

where $\mathbf{0}_{kr}$ is a kr -dimensional vector of zeros. To expand Eq. (4), we prepare the following three matrices consisting of the second, third, and fourth derivatives of $-\ell(\beta)/n$:

$$\mathbf{G}_2(\beta) = -\frac{1}{n} \frac{\partial^2 \ell(\beta)}{\partial \beta \partial \beta'}, \quad \mathbf{G}_3(\beta) = -\frac{1}{n} \left(\frac{\partial}{\partial \beta'} \otimes \frac{\partial^2}{\partial \beta \partial \beta'} \right) \ell(\beta),$$

$$\mathbf{G}_4(\beta) = -\frac{1}{n} \left(\frac{\partial^2}{\partial \beta \partial \beta'} \otimes \frac{\partial^2}{\partial \beta \partial \beta'} \right) \ell(\beta).$$

The result of the first derivative of $\ell(\beta)$ in (4) implies the following explicit forms of $\mathbf{G}_2(\beta)$, $\mathbf{G}_3(\beta)$, and $\mathbf{G}_4(\beta)$ (details of the derivations are given in [Appendix A](#)):

$$\mathbf{G}_2(\beta) = \sum_{i=1}^m \rho_i \left\{ \frac{\partial \mathbf{p}_i(\beta)}{\partial \beta'} \right\} \otimes \mathbf{x}_i = \sum_{i=1}^m \rho_i \left(\Sigma_i(\beta) \otimes \mathbf{x}_i \mathbf{x}_i' \right), \quad (5)$$

$$\begin{aligned} \mathbf{G}_3(\beta) &= \sum_{i=1}^m \rho_i \left\{ \left(\frac{\partial}{\partial \beta'} \otimes \frac{\partial}{\partial \beta'} \right) \mathbf{p}_i(\beta) \right\} \otimes \mathbf{x}_i \\ &= \sum_{i=1}^m \rho_i \left\{ \Delta_{3,i}(\beta) \otimes \mathbf{x}_i \mathbf{x}_i' \right\}, \end{aligned} \quad (6)$$

$$\begin{aligned} \mathbf{G}_4(\beta) &= \sum_{i=1}^m \rho_i \left\{ \left(\frac{\partial^2}{\partial \beta \partial \beta'} \otimes \frac{\partial}{\partial \beta'} \right) \mathbf{p}_i(\beta) \right\} \otimes \mathbf{x}_i \\ &= \sum_{i=1}^m \rho_i \left\{ \Delta_{4,i}(\beta) \otimes \mathbf{x}_i \mathbf{x}_i' \right\}, \end{aligned} \quad (7)$$

where $\Delta_{3,i}(\beta)$ and $\Delta_{4,i}(\beta)$ are $kr \times (kr)^2$ and $(kr)^2 \times (kr)^2$ matrices, respectively, which are defined by

$$\begin{aligned} \Delta_{3,i}(\beta) &= \sum_{a=1}^r p_{ia}(\beta) \mathbf{e}_a' \otimes \mathbf{x}_i' \otimes \mathbf{q}_{i,a}(\beta) \mathbf{q}_{i,a}(\beta)' - \mathbf{p}_i(\beta)' \otimes \mathbf{x}_i' \otimes \Sigma_i(\beta), \\ \Delta_{4,i}(\beta) &= \sum_{a=1}^r p_{ia}(\beta) \mathbf{q}_{i,a}(\beta) \mathbf{q}_{i,a}(\beta)' \otimes \mathbf{x}_i \mathbf{x}_i' \otimes (\mathbf{q}_{i,a}(\beta) \mathbf{q}_{i,a}(\beta)' - \mathbf{p}_i(\beta) \mathbf{p}_i(\beta)') \\ &\quad - \Sigma_i(\beta) \otimes \mathbf{x}_i \mathbf{x}_i' \otimes (\Sigma_i(\beta) - \mathbf{p}_i(\beta) \mathbf{p}_i(\beta)') \\ &\quad - \sum_{a,b}^r p_{ia}(\beta) p_{ib}(\beta) \mathbf{q}_{i,a}(\beta) \mathbf{q}_{i,b}(\beta)' \otimes \mathbf{x}_i \mathbf{x}_i' \otimes (\mathbf{e}_a \mathbf{e}_b' + \mathbf{e}_b \mathbf{e}_a'). \end{aligned} \quad (8)$$

Here, \mathbf{e}_j is an r -dimensional j th coordinate unit vector whose j th element is 1 and others are 0, $\mathbf{q}_{i,a}(\boldsymbol{\beta}) = \mathbf{e}_a - \mathbf{p}_i(\boldsymbol{\beta})$, and the notation \sum_{a_1, \dots, a_j}^r means $\sum_{a_1=1}^r \cdots \sum_{a_j=1}^r$.

Applying a Taylor expansion around $\hat{\boldsymbol{\beta}} = \boldsymbol{\beta}$ to Eq. (4) yields

$$\begin{aligned} \frac{1}{n} \sum_{i=1}^m \{(\mathbf{y}_i - n_i \mathbf{p}_i(\boldsymbol{\beta})) \otimes \mathbf{x}_i\} &= \mathbf{G}_2(\boldsymbol{\beta})(\hat{\boldsymbol{\beta}} - \boldsymbol{\beta}) + \frac{1}{2} \mathbf{G}_3(\boldsymbol{\beta}) \{(\hat{\boldsymbol{\beta}} - \boldsymbol{\beta}) \otimes (\hat{\boldsymbol{\beta}} - \boldsymbol{\beta})\} \\ &\quad + \frac{1}{6} \{\mathbf{I}_{kr} \otimes (\hat{\boldsymbol{\beta}} - \boldsymbol{\beta})'\} \mathbf{G}_4(\boldsymbol{\beta}) \{(\hat{\boldsymbol{\beta}} - \boldsymbol{\beta}) \otimes (\hat{\boldsymbol{\beta}} - \boldsymbol{\beta})\} \\ &\quad + O_p(n^{-2}). \end{aligned} \quad (9)$$

Notice that the order of the left-hand side of Eq. (9) is $O_p(n^{-1/2})$. By comparing the $O_p(n^{-1/2})$, $O_p(n^{-1})$, and $O_p(n^{-3/2})$ terms after substituting (3) into (9), \mathbf{b}_1 , \mathbf{b}_2 , and \mathbf{b}_3 in (3) are specified as

$$\begin{aligned} \mathbf{b}_1 &= \frac{1}{\sqrt{n}} \mathbf{G}_2(\boldsymbol{\beta})^{-1} \sum_{i=1}^m \{(\mathbf{y}_i - n_i \mathbf{p}_i(\boldsymbol{\beta})) \otimes \mathbf{x}_i\}, \\ \mathbf{b}_2 &= -\frac{1}{2} \mathbf{G}_2(\boldsymbol{\beta})^{-1} \mathbf{G}_3(\boldsymbol{\beta})(\mathbf{b}_1 \otimes \mathbf{b}_1), \\ \mathbf{b}_3 &= -\frac{1}{2} \mathbf{G}_2(\boldsymbol{\beta})^{-1} \left\{ \mathbf{G}_3(\boldsymbol{\beta})(\mathbf{b}_1 \otimes \mathbf{b}_2 + \mathbf{b}_2 \otimes \mathbf{b}_1) + \frac{1}{3} (\mathbf{I}_{kr} \otimes \mathbf{b}_1') \mathbf{G}_4(\boldsymbol{\beta})(\mathbf{b}_1 \otimes \mathbf{b}_1) \right\}. \end{aligned} \quad (10)$$

We use the stochastic expansion of $\hat{\boldsymbol{\beta}}$ with \mathbf{b}_1 , \mathbf{b}_2 , and \mathbf{b}_3 to evaluate the bias of the AIC to the risk function. The stochastic expansion is regarded as a special case of the general stochastic expansion of MLE, e.g., in [15].

3. Main result

Let $\mathcal{L}(\boldsymbol{\beta})$ be a loss function defined by

$$\begin{aligned} \mathcal{L}(\boldsymbol{\beta}) &= E[-2\ell(\boldsymbol{\beta})] \\ &= -2 \sum_{i=1}^m n_i \left\{ (\mathbf{p}_i^* \otimes \mathbf{x}_i)' \boldsymbol{\beta} - \log \left(1 + \sum_{j=1}^r \exp(\mathbf{x}_i' \boldsymbol{\beta}_j) \right) \right\}, \end{aligned} \quad (11)$$

where \mathbf{p}_i^* is the cell probability vector of the true model. Then, the risk function consisting of the predictive K-L information is given by

$$R = E[\mathcal{L}(\hat{\boldsymbol{\beta}})]. \quad (12)$$

In this section, we propose a CAIC that improves the bias of the AIC to $O(n^{-2})$ under the assumption that the candidate model includes the true model. Notice that the crude AIC is defined by

$$\text{AIC} = -2\ell(\hat{\boldsymbol{\beta}}) + 2kr. \quad (13)$$

Thus, it is sufficient to derive the bias of $-2\ell(\hat{\boldsymbol{\beta}})$ to R for evaluating the bias of the AIC. Also notice that $\mathbf{p}_i^* = \mathbf{p}_i(\boldsymbol{\beta})$ holds when the candidate model includes the true model. Then, the bias of $-2\ell(\hat{\boldsymbol{\beta}})$ to R under the assumption that the candidate model includes the true model is expanded as

$$\begin{aligned} B &= R - E[-2\ell(\hat{\boldsymbol{\beta}})] \\ &= 2 \sum_{i=1}^m E \left[\{(\mathbf{y}_i - n_i \mathbf{p}_i(\boldsymbol{\beta})) \otimes \mathbf{x}_i\}' \hat{\boldsymbol{\beta}} \right] \end{aligned}$$

$$\begin{aligned}
&= 2\sqrt{n}E[\mathbf{b}_1' \mathbf{G}_2(\boldsymbol{\beta}) \hat{\boldsymbol{\beta}}] \\
&= 2\left\{ \sqrt{n}E[\mathbf{b}_1' \mathbf{G}_2(\boldsymbol{\beta}) \boldsymbol{\beta}] + E[\mathbf{b}_1' \mathbf{G}_2(\boldsymbol{\beta}) \mathbf{b}_1] \right. \\
&\quad \left. + \frac{1}{\sqrt{n}}E[\mathbf{b}_1' \mathbf{G}_2(\boldsymbol{\beta}) \mathbf{b}_2] + \frac{1}{n}E[\mathbf{b}_1' \mathbf{G}_2(\boldsymbol{\beta}) \mathbf{b}_3] \right\} + O(n^{-2}), \tag{14}
\end{aligned}$$

where matrices $\mathbf{G}_2(\boldsymbol{\beta})$, $\mathbf{G}_3(\boldsymbol{\beta})$, and $\mathbf{G}_4(\boldsymbol{\beta})$ are given by (5), (6), and (7), respectively, and kr -dimensional random vectors \mathbf{b}_1 , \mathbf{b}_2 , and \mathbf{b}_3 are given by (10). In many cases of practical interest, a moment of statistic can be expanded as a power series in n^{-1} (see e.g. [16, p. 46]). Hence, the order of the remainder term of (14) is shown by $O(n^{-2})$, not $O(n^{-3/2})$. Indeed, an $n^{-3/2}$ term of the stochastic expansion of $\sum_{i=1}^m \{(\mathbf{y}_i - n_i \mathbf{p}_i(\boldsymbol{\beta})) \otimes \mathbf{x}_i\}' \hat{\boldsymbol{\beta}}$ in the bias can be expressed as a fifth-order polynomial of elements of \mathbf{b}_1 . Since \mathbf{b}_1 has an asymptotic normality, the expectation of an odd-order polynomial of elements of \mathbf{b}_1 becomes $O(n^{-1/2})$. Given this fact, the order of the remainder term of the expansion in (14) is $O(n^{-2})$.

From elementary linear algebra and the definition of \mathbf{b}_2 in (10), $\mathbf{b}_1' \mathbf{G}_2(\boldsymbol{\beta}) \mathbf{b}_2$ in (14) is expressed by the function of \mathbf{b}_1 as

$$\mathbf{b}_1' \mathbf{G}_2(\boldsymbol{\beta}) \mathbf{b}_2 = -\frac{1}{2} \mathbf{b}_1' \mathbf{G}_3(\boldsymbol{\beta}) (\mathbf{b}_1 \otimes \mathbf{b}_1) = -\frac{1}{2} \text{tr}\{\mathbf{G}_3(\boldsymbol{\beta}) (\mathbf{b}_1 \otimes \mathbf{b}_1 \mathbf{b}_1')\}. \tag{15}$$

Since the derivative is invariant to changes in the order of differentiation, we have

$$\mathbf{b}_1' \mathbf{G}_3(\boldsymbol{\beta}) (\mathbf{b}_1 \otimes \mathbf{b}_2) = \mathbf{b}_1' \mathbf{G}_3(\boldsymbol{\beta}) (\mathbf{b}_2 \otimes \mathbf{b}_1) = \mathbf{b}_2' \mathbf{G}_3(\boldsymbol{\beta}) (\mathbf{b}_1 \otimes \mathbf{b}_1) = (\mathbf{b}_1 \otimes \mathbf{b}_1)' \mathbf{G}_3(\boldsymbol{\beta})' \mathbf{b}_2.$$

It follows from the above equations and the definition of \mathbf{b}_2 in (10) that

$$\begin{aligned}
\mathbf{b}_1' \mathbf{G}_3(\boldsymbol{\beta}) (\mathbf{b}_1 \otimes \mathbf{b}_2 + \mathbf{b}_2 \otimes \mathbf{b}_1) &= 2(\mathbf{b}_1 \otimes \mathbf{b}_1)' \mathbf{G}_3(\boldsymbol{\beta})' \mathbf{b}_2 \\
&= -(\mathbf{b}_1 \otimes \mathbf{b}_1)' \mathbf{G}_3(\boldsymbol{\beta})' \mathbf{G}_2(\boldsymbol{\beta})^{-1} \mathbf{G}_3(\boldsymbol{\beta}) (\mathbf{b}_1 \otimes \mathbf{b}_1) \\
&= -\text{tr}\left\{ \mathbf{G}_3(\boldsymbol{\beta})' \mathbf{G}_2(\boldsymbol{\beta})^{-1} \mathbf{G}_3(\boldsymbol{\beta}) (\mathbf{b}_1 \mathbf{b}_1' \otimes \mathbf{b}_1 \mathbf{b}_1') \right\}.
\end{aligned}$$

Thus, from the above result and the definition of \mathbf{b}_3 in (10), $\mathbf{b}_1' \mathbf{G}_2(\boldsymbol{\beta}) \mathbf{b}_3$ in (14) is expressed by the function of \mathbf{b}_1 as

$$\begin{aligned}
\mathbf{b}_1' \mathbf{G}_2 \mathbf{b}_3 &= -\frac{1}{2} \mathbf{b}_1' \mathbf{G}_3(\boldsymbol{\beta}) (\mathbf{b}_1 \otimes \mathbf{b}_2 + \mathbf{b}_2 \otimes \mathbf{b}_1) - \frac{1}{6} \mathbf{b}_1' (\mathbf{I}_{kr} \otimes \mathbf{b}_1') \mathbf{G}_4(\boldsymbol{\beta}) (\mathbf{b}_1 \otimes \mathbf{b}_1) \\
&= \frac{1}{2} \text{tr}\left\{ \mathbf{G}_3(\boldsymbol{\beta})' \mathbf{G}_2(\boldsymbol{\beta})^{-1} \mathbf{G}_3(\boldsymbol{\beta}) (\mathbf{b}_1 \mathbf{b}_1' \otimes \mathbf{b}_1 \mathbf{b}_1') \right\} - \frac{1}{6} \text{tr}\{\mathbf{G}_4(\boldsymbol{\beta}) (\mathbf{b}_1 \mathbf{b}_1' \otimes \mathbf{b}_1 \mathbf{b}_1')\}. \tag{16}
\end{aligned}$$

Hence, Eqs. (15) and (16) indicate that the expansion of B in (14) can be calculated until the fourth moment of \mathbf{b}_1 .

Since \mathbf{b}_1 consists of a centralized \mathbf{y}_i , we can directly calculate the expectations in (14) by centralized moments of $\mathbf{y}_1, \dots, \mathbf{y}_m$. Then, all combinations of multivariate moments of $\mathbf{y}_i - n_i \mathbf{p}_i(\boldsymbol{\beta})$ are needed until the fourth-order. However, it is troublesome to calculate the third- and fourth-order multivariate moments of $\mathbf{y}_i - n_i \mathbf{p}_i(\boldsymbol{\beta})$, because we have to consider all combinations of the multivariate moments. For simplicity, the relations between the moments of \mathbf{b}_1 and the expectations of the derivatives of $-\ell(\boldsymbol{\beta})$ with respect to $\boldsymbol{\beta}$ are used instead of calculating the multivariate moments of $\mathbf{y}_i - n_i \mathbf{p}_i(\boldsymbol{\beta})$. It is easy to obtain $E[\mathbf{b}_1] = \mathbf{0}_{kr}$ because $E[\mathbf{y}_i] = n_i \mathbf{p}_i(\boldsymbol{\beta})$. From the result of the first derivative of $\ell(\boldsymbol{\beta})$ in (4) and the definition of \mathbf{b}_1 in (10), we can see that

$$\frac{\partial}{\partial \boldsymbol{\beta}} \ell(\boldsymbol{\beta}) = \sqrt{n} \mathbf{G}_2(\boldsymbol{\beta}) \mathbf{b}_1.$$

Notice that $E[\partial \ell(\boldsymbol{\beta}) / \partial \boldsymbol{\beta}] = \mathbf{0}_{kr}$ holds and $\mathbf{G}_2(\boldsymbol{\beta})$, $\mathbf{G}_3(\boldsymbol{\beta})$, and $\mathbf{G}_4(\boldsymbol{\beta})$ are constant matrices. By applying general formulas of expectations (B.6) in Appendix B to the case of the multinomial logistic regression

model, the following equations are obtained:

$$\begin{aligned} n\mathbf{G}_2(\boldsymbol{\beta}) &= n\mathbf{G}_2(\boldsymbol{\beta})E[\mathbf{b}_1\mathbf{b}_1']\mathbf{G}_2(\boldsymbol{\beta}), \\ n\mathbf{G}_3(\boldsymbol{\beta}) &= n\sqrt{n}\mathbf{G}_2(\boldsymbol{\beta})E[\mathbf{b}_1' \otimes \mathbf{b}_1\mathbf{b}_1'](\mathbf{G}_2(\boldsymbol{\beta}) \otimes \mathbf{G}_2(\boldsymbol{\beta})), \\ n\mathbf{G}_4(\boldsymbol{\beta}) &= n^2(\mathbf{G}_2(\boldsymbol{\beta}) \otimes \mathbf{G}_2(\boldsymbol{\beta}))E[\mathbf{b}_1\mathbf{b}_1' \otimes \mathbf{b}_1\mathbf{b}_1'](\mathbf{G}_2(\boldsymbol{\beta}) \otimes \mathbf{G}_2(\boldsymbol{\beta})) \\ &\quad - n^2 \left\{ (\mathbf{I}_{k^2r^2} + \mathbf{K}_{kr})(\mathbf{G}_2(\boldsymbol{\beta}) \otimes \mathbf{G}_2(\boldsymbol{\beta})) + \text{vec}(\mathbf{G}_2(\boldsymbol{\beta}))\text{vec}(\mathbf{G}_2(\boldsymbol{\beta}))' \right\}, \end{aligned}$$

where $\text{vec}(\mathbf{A})$ is an operator to transform a matrix to a vector by stacking the first to the last column of \mathbf{A} , i.e., $\text{vec}(\mathbf{A}) = (\mathbf{a}'_1, \dots, \mathbf{a}'_m)'$ when $\mathbf{A} = (\mathbf{a}_1, \dots, \mathbf{a}_m)$ (see e.g. [17, Chapter 16.2]), and \mathbf{K}_m is the $m^2 \times m^2$ vec-permutation matrix such that $\text{vec}(\mathbf{B}) = \mathbf{K}_m \text{vec}(\mathbf{B}')$ when \mathbf{B} is an $m \times m$ matrix (see e.g. [17, Chapter 16.3]). These results lead us to the simple expression of moments of \mathbf{b}_1 as

$$E[\mathbf{b}_1\mathbf{b}_1'] = \mathbf{G}_2(\boldsymbol{\beta})^{-1}, \quad (17)$$

$$E[\mathbf{b}_1' \otimes \mathbf{b}_1\mathbf{b}_1'] = \frac{1}{\sqrt{n}}\mathbf{G}_2(\boldsymbol{\beta})^{-1}\mathbf{G}_3(\boldsymbol{\beta})(\mathbf{G}_2(\boldsymbol{\beta})^{-1} \otimes \mathbf{G}_2(\boldsymbol{\beta})^{-1}), \quad (18)$$

$$\begin{aligned} E[\mathbf{b}_1\mathbf{b}_1' \otimes \mathbf{b}_1\mathbf{b}_1'] &= (\mathbf{I}_{k^2r^2} + \mathbf{K}_{kr})(\mathbf{G}_2(\boldsymbol{\beta})^{-1} \otimes \mathbf{G}_2(\boldsymbol{\beta})^{-1}) \\ &\quad + \text{vec}(\mathbf{G}_2(\boldsymbol{\beta})^{-1})\text{vec}(\mathbf{G}_2(\boldsymbol{\beta})^{-1})' + O(n^{-1}). \end{aligned} \quad (19)$$

The result in (17) implies that

$$E[\mathbf{b}_1'\mathbf{G}_2(\boldsymbol{\beta})\mathbf{b}_1] = E[\{\mathbf{G}_2(\boldsymbol{\beta})\mathbf{b}_1\mathbf{b}_1'\}] = \text{tr}\{\mathbf{G}_2(\boldsymbol{\beta})\mathbf{G}_2(\boldsymbol{\beta})^{-1}\} = kr. \quad (20)$$

Similarly, from (15) and (18), we have

$$\begin{aligned} E[\mathbf{b}_1'\mathbf{G}_2(\boldsymbol{\beta})\mathbf{b}_2] &= -\frac{1}{2}E[\text{tr}\{\mathbf{G}_3(\boldsymbol{\beta})'(\mathbf{b}_1' \otimes \mathbf{b}_1\mathbf{b}_1')\}] \\ &= -\frac{1}{2\sqrt{n}}\text{tr}\left\{\mathbf{G}_3(\boldsymbol{\beta})'\mathbf{G}_2(\boldsymbol{\beta})^{-1}\mathbf{G}_3(\boldsymbol{\beta})(\mathbf{G}_2(\boldsymbol{\beta})^{-1} \otimes \mathbf{G}_2(\boldsymbol{\beta})^{-1})\right\}. \end{aligned} \quad (21)$$

Notice that $\mathbf{G}_3(\boldsymbol{\beta})\mathbf{K}_{kr} = \mathbf{G}_3(\boldsymbol{\beta})$ holds because the derivative is invariant to changes in the order of differentiation. By using this fact and Eq. (19), the expectation of the first part in (16) is given by

$$\begin{aligned} E\left[\text{tr}\left\{\mathbf{G}_3(\boldsymbol{\beta})'\mathbf{G}_2(\boldsymbol{\beta})^{-1}\mathbf{G}_3(\boldsymbol{\beta})(\mathbf{b}_1\mathbf{b}_1' \otimes \mathbf{b}_1\mathbf{b}_1')\right\}\right] \\ &= \text{tr}\left\{\mathbf{G}_3(\boldsymbol{\beta})'\mathbf{G}_2(\boldsymbol{\beta})^{-1}\mathbf{G}_3(\boldsymbol{\beta})(\mathbf{I}_{k^2r^2} + \mathbf{K}_{kr})(\mathbf{G}_2(\boldsymbol{\beta})^{-1} \otimes \mathbf{G}_2(\boldsymbol{\beta})^{-1})\right\} \\ &\quad + \text{vec}(\mathbf{G}_2(\boldsymbol{\beta})^{-1})'\mathbf{G}_3(\boldsymbol{\beta})'\mathbf{G}_2(\boldsymbol{\beta})^{-1}\mathbf{G}_3(\boldsymbol{\beta})\text{vec}(\mathbf{G}_2(\boldsymbol{\beta})^{-1}) + O(n^{-1}) \\ &= 2\text{tr}\left\{\mathbf{G}_3(\boldsymbol{\beta})'\mathbf{G}_2(\boldsymbol{\beta})^{-1}\mathbf{G}_3(\boldsymbol{\beta})(\mathbf{G}_2(\boldsymbol{\beta})^{-1} \otimes \mathbf{G}_2(\boldsymbol{\beta})^{-1})\right\} \\ &\quad + \text{vec}(\mathbf{G}_2(\boldsymbol{\beta})^{-1})'\mathbf{G}_3(\boldsymbol{\beta})'\mathbf{G}_2(\boldsymbol{\beta})^{-1}\mathbf{G}_3(\boldsymbol{\beta})\text{vec}(\mathbf{G}_2(\boldsymbol{\beta})^{-1}) + O(n^{-1}). \end{aligned} \quad (22)$$

Moreover, since the derivative is invariant to changes in the order of differentiation, we can see that $\mathbf{G}_4(\boldsymbol{\beta})\mathbf{K}_{kr} = \mathbf{G}_4(\boldsymbol{\beta})$ and

$$\text{tr}\left\{\mathbf{G}_4(\boldsymbol{\beta})(\mathbf{G}_2(\boldsymbol{\beta})^{-1} \otimes \mathbf{G}_2(\boldsymbol{\beta})^{-1})\right\} = \text{vec}(\mathbf{G}_2(\boldsymbol{\beta})^{-1})'\mathbf{G}_4(\boldsymbol{\beta})\text{vec}(\mathbf{G}_2(\boldsymbol{\beta})^{-1}).$$

By using the above relations and Eq. (19), the expectation of the second part in (16) is given by

$$\begin{aligned} E\left[\text{tr}\{\mathbf{G}_4(\boldsymbol{\beta})(\mathbf{b}_1\mathbf{b}_1' \otimes \mathbf{b}_1\mathbf{b}_1')\}\right] &= \text{tr}\{\mathbf{G}_4(\boldsymbol{\beta})(\mathbf{I}_{k^2r^2} + \mathbf{K}_{kr})(\mathbf{G}_2(\boldsymbol{\beta})^{-1} \otimes \mathbf{G}_2(\boldsymbol{\beta})^{-1})\} \\ &\quad + \text{vec}(\mathbf{G}_2(\boldsymbol{\beta})^{-1})\mathbf{G}_4(\boldsymbol{\beta})\text{vec}(\mathbf{G}_2(\boldsymbol{\beta})^{-1}) + O(n^{-1}) \\ &= 3\text{tr}\{\mathbf{G}_4(\boldsymbol{\beta})(\mathbf{G}_2(\boldsymbol{\beta})^{-1} \otimes \mathbf{G}_2(\boldsymbol{\beta})^{-1})\} + O(n^{-1}). \end{aligned} \quad (23)$$

Hence, from Eqs. (16), (22), and (23), we can see that

$$\begin{aligned} E[\mathbf{b}'_1 \mathbf{G}_2 \mathbf{b}_3] &= \text{tr} \left\{ \mathbf{G}_3(\boldsymbol{\beta})' \mathbf{G}_2(\boldsymbol{\beta})^{-1} \mathbf{G}_3(\boldsymbol{\beta}) (\mathbf{G}_2(\boldsymbol{\beta})^{-1} \otimes \mathbf{G}_2(\boldsymbol{\beta})^{-1}) \right\} \\ &\quad + \frac{1}{2} \text{vec}(\mathbf{G}_2(\boldsymbol{\beta})^{-1})' \mathbf{G}_3(\boldsymbol{\beta})' \mathbf{G}_2(\boldsymbol{\beta})^{-1} \mathbf{G}_3(\boldsymbol{\beta}) \text{vec}(\mathbf{G}_2(\boldsymbol{\beta})^{-1}) \\ &\quad - \frac{1}{2} \text{tr} \left\{ \mathbf{G}_4(\boldsymbol{\beta}) (\mathbf{G}_2(\boldsymbol{\beta})^{-1} \otimes \mathbf{G}_2(\boldsymbol{\beta})^{-1}) \right\} + O(n^{-1}). \end{aligned} \quad (24)$$

Consequently, by substituting $E[\mathbf{b}'_1 \mathbf{G}_2(\boldsymbol{\beta}) \boldsymbol{\beta}] = 0$, and Eqs. (20), (21), and (24) into (14), the bias of $-2\ell(\hat{\boldsymbol{\beta}})$ to R is expanded as

$$B = 2kr + \frac{1}{n} \{ \alpha_1(\boldsymbol{\beta}) + \alpha_2(\boldsymbol{\beta}) - \alpha_3(\boldsymbol{\beta}) \} + O(n^{-2}),$$

where coefficients $\alpha_1(\boldsymbol{\beta})$, $\alpha_2(\boldsymbol{\beta})$, and $\alpha_3(\boldsymbol{\beta})$ are given by

$$\begin{aligned} \alpha_1(\boldsymbol{\beta}) &= \text{tr} \left\{ \mathbf{G}_3(\boldsymbol{\beta})' \mathbf{G}_2(\boldsymbol{\beta})^{-1} \mathbf{G}_3(\boldsymbol{\beta}) (\mathbf{G}_2(\boldsymbol{\beta})^{-1} \otimes \mathbf{G}_2(\boldsymbol{\beta})^{-1}) \right\}, \\ \alpha_2(\boldsymbol{\beta}) &= \text{vec}(\mathbf{G}_2(\boldsymbol{\beta})^{-1})' \mathbf{G}_3(\boldsymbol{\beta})' \mathbf{G}_2(\boldsymbol{\beta})^{-1} \mathbf{G}_3(\boldsymbol{\beta}) \text{vec}(\mathbf{G}_2(\boldsymbol{\beta})^{-1}), \\ \alpha_3(\boldsymbol{\beta}) &= \text{tr} \left\{ \mathbf{G}_4(\boldsymbol{\beta}) (\mathbf{G}_2(\boldsymbol{\beta})^{-1} \otimes \mathbf{G}_2(\boldsymbol{\beta})^{-1}) \right\}. \end{aligned} \quad (25)$$

The CAIC can then be defined by adding an estimated B to $-2\ell(\hat{\boldsymbol{\beta}})$, i.e.,

$$\text{CAIC} = -2\ell(\hat{\boldsymbol{\beta}}) + 2kr + \frac{1}{n} \{ \alpha_1(\hat{\boldsymbol{\beta}}) + \alpha_2(\hat{\boldsymbol{\beta}}) - \alpha_3(\hat{\boldsymbol{\beta}}) \}. \quad (26)$$

The CAIC improves the bias of the AIC to $O(n^{-2})$, although the order of the bias of the AIC is $O(n^{-1})$, i.e., the following equations are satisfied:

$$R - E[\text{AIC}] = O(n^{-1}), \quad R - E[\text{CAIC}] = O(n^{-2}),$$

where R is the risk function given by (12).

4. Numerical studies

In this section, we conduct numerical studies to show that the CAIC in (26) works better than the crude AIC in (13). To compare the performances of the AIC and the CAIC, the following two properties are considered:

- (I) the selection probability: the frequency of the model chosen by minimizing the information criterion.
- (II) the prediction error of the best model (PE_B): the risk function of the best model chosen by the information criterion, which is defined by

$$\text{PE}_B = E[\mathcal{L}(\hat{\boldsymbol{\beta}}_B)],$$

where $\mathcal{L}(\boldsymbol{\beta})$ is the loss function given by (11) and $\hat{\boldsymbol{\beta}}_B$ is the MLE of $\boldsymbol{\beta}$ under the best model.

These two properties were evaluated by a Monte Carlo simulation with 10,000 iterations. The information criterion with the higher selection probability of the true model and the smaller prediction error of the best model is regarded as a high-performance model selector. In the basic concept of the AIC, a good model selection method is one that chooses the best model so that the prediction is improved. Hence, PE_B is a more important property than the selection probability.

We prepared eight candidate models M_1, \dots, M_8 , with $m = 20$ and 50 , $n_i = 5$ ($i = 1, \dots, m$) and $r = 2$. An $m \times 8$ matrix of explanatory variables $\mathbf{X} = (\mathbf{x}_1, \dots, \mathbf{x}_m)'$ was constructed as follows. The first

Table 1

Selection probability of the model and the prediction error of the best model.

Case m		Criterion	Selection probability								PE _B
			M_1	M_2	M_3	M_4	M_5	M_6	M_7	M_8	
1	20	AIC	1.81	0.29	74.84	11.41	4.95	2.88	2.20	1.62	210.06
		CAIC	3.19	0.66	79.92	10.01	3.58	1.38	0.88	0.38	207.84
	50	AIC	0.01	0.00	79.15	11.16	4.77	2.21	1.56	1.14	511.32
		CAIC	0.01	0.00	81.25	10.71	4.27	1.88	1.17	0.71	511.04
2	20	AIC	77.22	10.92	4.86	2.69	1.56	1.11	0.77	0.87	202.42
		CAIC	82.63	10.06	3.81	2.07	0.76	0.38	0.19	0.10	200.41
	50	AIC	79.21	10.89	4.48	2.20	1.24	1.04	0.55	0.39	494.89
		CAIC	80.99	10.58	4.10	1.91	1.01	0.70	0.44	0.27	494.63

Note: The selection probability of the true model is marked in bold.

column of \mathbf{X} is $\mathbf{1}_m$, where $\mathbf{1}_m$ is an m -dimensional vector of ones, and the remaining seven columns of \mathbf{X} were generated randomly from the binomial distribution $B(1, 0.5)$. Simulation data were generated from the multinomial distribution with the true cell probability consisting of $\beta^* = (\beta_1^*, \beta_2^*)'$. In this simulation study, we prepared two β^* , as follows:

$$\text{Case1 : } \beta_1^* = (0, 0.2, -1.0, 0, 0, 0, 0, 0)', \quad \beta_2^* = (-0.1, -0.4, 1.2, 0, 0, 0, 0, 0)',$$

$$\text{Case2 : } \beta_1^* = (-0.5, 0, 0, 0, 0, 0, 0, 0)', \quad \beta_2^* = (0.7, 0, 0, 0, 0, 0, 0, 0)'.$$

The matrix of explanatory variables in M_j consists of the first j columns of \mathbf{X} ($j = 1, \dots, 8$). Thus, the true model in Case 1 is M_3 , and the true model in Case 2 is M_1 .

Table 1 shows the two properties (I) and (II). In the table, the selection probability of the true model is marked in bold. From this table, we can see that the selection probabilities and the prediction errors of the CAIC were improved in comparison with those of the AIC in all situations. We simulated several other models and obtained similar results.

5. Conclusion and discussion

In this paper, we proposed the CAIC for selecting variables in the multinomial logistic regression models. The proposed CAIC improves the bias of the AIC to $O(n^{-2})$, although the order of the bias of the AIC is $O(n^{-1})$. By using relations between the moments of \mathbf{b}_1 and expectations of the derivatives of $-\ell(\beta)$ instead of directly calculating the moments of \mathbf{y}_i to evaluate the moments of \mathbf{b}_1 , a simple expression of the CAIC is developed. Indeed, the bias correction term of the proposed CAIC is represented by only three matrices $\mathbf{G}_2(\hat{\beta})$, $\mathbf{G}_3(\hat{\beta})$, and $\mathbf{G}_4(\hat{\beta})$, which consist of the second, third, and fourth derivatives of $-\ell(\beta)$. Even though expressions of $\mathbf{G}_2(\hat{\beta})$, $\mathbf{G}_3(\hat{\beta})$, and $\mathbf{G}_4(\hat{\beta})$ are not simple, we can derive the bias correction term of the CAIC from linear functions of $\mathbf{G}_2(\hat{\beta})^{-1}$, $\mathbf{G}_3(\hat{\beta})$, and $\mathbf{G}_4(\hat{\beta})$. This is a desirable character of the CAIC.

In all situations of the simulation study, the CAIC improved the crude AIC in the sense of making a high selection probability of the true model and a small prediction error of the best model chosen by the information criterion. However, the improvements were smaller when the sample size was large. This is natural because the CAIC is proposed so that the bias of the AIC is corrected when the sample size is small. Needless to say, the AIC and the CAIC are asymptotical equivalents. Hence, the difference between two criteria becomes small when the sample size is increased. The sample sizes of our simulation were 100 and 250. Nevertheless, a clear difference exists in the performances of the CAIC and the AIC. This difference indicates that the CAIC is valuable even when the sample size is not so small. Consequently, we recommend using the CAIC instead of the AIC for selecting multinomial logistic regression models.

The simple expression of the proposed CAIC is based on the property that the second derivatives of $-\ell(\beta)$ do not depend on response variables. A generalized linear model (GLM) with a natural link and a known dispersion parameter, e.g., a logistic regression model or a Poisson regression model, will

have this property. Then, we can simply express the bias-corrected AIC just like the proposed CAIC in (26) in the same way presented in Section 3. Namely, the bias-corrected AIC with constant second derivatives of the negative log-likelihood function may be stated by

$$\text{CAIC} = \text{AIC} + \gamma_1(\hat{\theta}) + \gamma_2(\hat{\theta}) - \gamma_3(\hat{\theta}),$$

where $\hat{\theta}$ is the MLE of unknown parameter θ , and coefficients $\gamma_1(\theta)$, $\gamma_2(\theta)$, and $\gamma_3(\theta)$ are given by

$$\begin{aligned}\gamma_1(\theta) &= \text{tr} \left\{ \mathbf{C}(\theta)' \mathbf{H}(\theta)^{-1} \mathbf{C}(\theta) (\mathbf{H}(\theta)^{-1} \otimes \mathbf{H}(\theta)^{-1}) \right\}, \\ \gamma_2(\theta) &= \text{vec}(\mathbf{H}(\theta)^{-1})' \mathbf{C}(\theta)' \mathbf{H}(\theta)^{-1} \mathbf{C}(\theta) \text{vec}(\mathbf{H}(\theta)^{-1}), \\ \gamma_3(\theta) &= \text{tr} \left\{ \mathbf{Q}(\theta) (\mathbf{H}(\theta)^{-1} \otimes \mathbf{H}(\theta)^{-1}) \right\}.\end{aligned}$$

Here, $\mathbf{H}(\theta)$, $\mathbf{C}(\theta)$, and $\mathbf{Q}(\theta)$ are matrices consisting of the second, third, and fourth derivatives, respectively, of the negative log-likelihood function and are defined by (B.4) in Appendix B.

Before concluding this section, we consider a log-likelihood ratio statistic for testing the null hypothesis $H_0: \beta = \beta_0$. It is known that a log-likelihood ratio statistic for testing the null hypothesis H_0 is $T = 2\{\ell(\hat{\beta}) - \ell(\beta_0)\}$. By comparing Eq. (25) with the general formula of the Bartlett factor of a log-likelihood ratio statistic (see e.g. [15]), we find that the first term in the asymptotic expansion of the bias of AIC consists of the similar coefficients in a Bartlett factor of T . This is because the bias of AIC and the Bartlett factor of T are partially formed from $2E[\ell(\hat{\beta})]$. However, $B = 2E[\ell(\hat{\beta})] + R$ and $E[T] = 2E[\ell(\hat{\beta})] - 2E[\ell(\beta_0)]$ are clear different because R is not equal to $-2E[\ell(\beta_0)]$. Although both are different, we will obtain an asymptotic expansion of $E[T]$ in the same way in Section 3. Then, the Bartlett factor of T may be expressed by the linear functions of $\mathbf{G}_2(\hat{\beta})^{-1}$, $\mathbf{G}_3(\hat{\beta})$ and $\mathbf{G}_4(\hat{\beta})$ as in the bias correction term of CAIC.

Appendix A. Explicit forms of $\mathbf{G}_2(\beta)$, $\mathbf{G}_3(\beta)$, AND $\mathbf{G}_4(\beta)$

In this subsection, for simplicity, we write $\Sigma_i(\beta)$, $\mathbf{p}_i(\beta)$, and $p_{ij}(\beta)$ as Σ_i , \mathbf{p}_i , and p_{ij} , respectively. Notice that

$$\frac{\partial \mathbf{p}_i}{\partial \beta_j'} = p_{ij} (\mathbf{e}_j - \mathbf{p}_i) \mathbf{x}_i', \quad (j = 1, \dots, r),$$

where \mathbf{e}_j is the j th coordinate unit vector, which is used in Eq. (8). This result and Eq. (2) imply that

$$\frac{\partial \mathbf{p}_i}{\partial \beta'} = (p_{i1} (\mathbf{e}_1 - \mathbf{p}_i) \mathbf{x}_i', \dots, p_{ir} (\mathbf{e}_r - \mathbf{p}_i) \mathbf{x}_i') = \Sigma_i \otimes \mathbf{x}_i'.$$

Substituting the above result into the definition of $\mathbf{G}_2(\beta)$ yields Eq. (5). Furthermore, from the definitions of $\mathbf{G}_3(\beta)$ and $\mathbf{G}_4(\beta)$, we can see that $\Delta_{3,i}(\beta)$ and $\Delta_{4,i}(\beta)$ in (6) and (7), respectively, satisfy

$$\Delta_{3,i}(\beta) = \frac{\partial}{\partial \beta'} \otimes \Sigma_i, \quad \Delta_{4,i}(\beta) = \frac{\partial^2}{\partial \beta \partial \beta'} \otimes \Sigma_i.$$

Notice that the (a, b) th element of Σ_i is $p_{ia}\delta_{ab} - p_{ia}p_{ib}$, where δ_{ab} is the Kronecker delta, i.e., $\delta_{aa} = 1$ and $\delta_{ab} = 0$ for $a \neq b$. This equation leads us to other expressions of $\Delta_{3,i}(\beta)$ and $\Delta_{4,i}(\beta)$, as follows:

$$\begin{aligned}\Delta_{3,i}(\beta) &= \sum_{a,b}^r \frac{\partial}{\partial \beta'} (p_{ia}\delta_{ab} - p_{ia}p_{ib}) \otimes \mathbf{e}_a \mathbf{e}_b', \\ \Delta_{4,i}(\beta) &= \sum_{a,b}^r \frac{\partial^2}{\partial \beta \partial \beta'} (p_{ia}\delta_{ab} - p_{ia}p_{ib}) \otimes \mathbf{e}_a \mathbf{e}_b'.\end{aligned}\tag{A.1}$$

Derivatives of p_{ia} are calculated as

$$\begin{aligned}\frac{\partial p_{ia}}{\partial \boldsymbol{\beta}} &= p_{ia}(\mathbf{e}_a - \mathbf{p}_i) \otimes \mathbf{x}_i, \\ \frac{\partial^2 p_{ia}}{\partial \boldsymbol{\beta} \partial \boldsymbol{\beta}'} &= p_{ia}(\mathbf{e}_a - \mathbf{p}_i)(\mathbf{e}_a - \mathbf{p}_i)' \otimes \mathbf{x}_i \mathbf{x}_i' - p_{ia} \boldsymbol{\Sigma}_i \otimes \mathbf{x}_i \mathbf{x}_i' \\ &= p_{ia} \left\{ (\mathbf{e}_a - \mathbf{p}_i)(\mathbf{e}_a - \mathbf{p}_i)' - \boldsymbol{\Sigma}_i \right\} \otimes \mathbf{x}_i \mathbf{x}_i', \\ \frac{\partial^2 p_{ia} p_{ib}}{\partial \boldsymbol{\beta} \partial \boldsymbol{\beta}'} &= p_{ib} \frac{\partial^2 p_{ia}}{\partial \boldsymbol{\beta} \partial \boldsymbol{\beta}'} + \frac{\partial p_{ib}}{\partial \boldsymbol{\beta}} \frac{\partial p_{ia}}{\partial \boldsymbol{\beta}'} + \frac{\partial p_{ia}}{\partial \boldsymbol{\beta}} \frac{\partial p_{ib}}{\partial \boldsymbol{\beta}'} + p_{ia} \frac{\partial^2 p_{ib}}{\partial \boldsymbol{\beta} \partial \boldsymbol{\beta}'} \\ &= p_{ia} p_{ib} \left\{ (\mathbf{e}_a + \mathbf{e}_b - 2\mathbf{p}_i)(\mathbf{e}_a + \mathbf{e}_b - 2\mathbf{p}_i)' - 2\boldsymbol{\Sigma}_i \right\} \otimes \mathbf{x}_i \mathbf{x}_i' .\end{aligned}$$

By substituting the above derivatives into (A.1), we have

$$\begin{aligned}\Delta_{3,i}(\boldsymbol{\beta}) &= \sum_{a,b}^r \left(\delta_{ab} p_{ia} \mathbf{q}'_{i,a} - p_{ia} p_{ib} \mathbf{q}'_{i,a} - p_{ia} p_{ib} \mathbf{q}'_{i,b} \right) \otimes \mathbf{x}'_i \otimes \mathbf{e}_a \mathbf{e}'_b \\ &= \sum_{a,b}^r p_{ia} \left\{ (\delta_{ab} - p_{ib}) \mathbf{q}'_{i,a} - p_{ib} \mathbf{q}'_{i,b} \right\} \otimes \mathbf{x}'_i \otimes \mathbf{e}_a \mathbf{e}'_b \\ &= \sum_{a=1}^r p_{ia} (\mathbf{e}_a \otimes \mathbf{x}_i)' \otimes \mathbf{q}_{i,a} \mathbf{q}'_{i,a} - (\mathbf{p}_i \otimes \mathbf{x}_i)' \otimes \boldsymbol{\Sigma}_i\end{aligned}$$

and

$$\begin{aligned}\Delta_{4,i}(\boldsymbol{\beta}) &= \sum_{a,b}^r p_{ia} \left(\delta_{ab} (\mathbf{q}_{i,a} \mathbf{q}'_{i,a} - \boldsymbol{\Sigma}_i) - p_{ib} \left\{ (\mathbf{q}_{i,a} + \mathbf{q}_{i,b})(\mathbf{q}_{i,a} + \mathbf{q}_{i,b})' - 2\boldsymbol{\Sigma}_i \right\} \right) \otimes \mathbf{x}_i \mathbf{x}'_i \otimes \mathbf{e}_a \mathbf{e}'_b \\ &= \sum_{a=1}^r p_{ia} \mathbf{q}_{i,a} \mathbf{q}'_{i,a} \otimes \mathbf{x}_i \mathbf{x}'_i \otimes (\mathbf{q}_{i,a} \mathbf{q}'_{i,a} - \mathbf{p}_i \mathbf{p}'_i) - \boldsymbol{\Sigma}_i \otimes \mathbf{x}_i \mathbf{x}'_i \otimes (\boldsymbol{\Sigma}_i - \mathbf{p}_i \mathbf{p}'_i) \\ &\quad - \sum_{a,b}^r p_{ia} p_{ib} \mathbf{q}_{i,a} \mathbf{q}'_{i,b} \otimes \mathbf{x}_i \mathbf{x}'_i \otimes (\mathbf{e}_a \mathbf{e}'_b + \mathbf{e}_b \mathbf{e}'_a),\end{aligned}$$

where $\mathbf{q}_{i,a} = \mathbf{e}_a - \mathbf{p}_i$. The above two equations indicate that explicit forms of $\mathbf{G}_3(\boldsymbol{\beta})$ and $\mathbf{G}_4(\boldsymbol{\beta})$ are given in (6) and (7), respectively.

Appendix B. Expectations of derivatives of the negative log-likelihood function

In this subsection, we derive general formulas of the expectations of derivatives of the negative log-likelihood function. Let $f(\mathbf{u}|\boldsymbol{\theta})$ be a joint probability density function of \mathbf{u} specified by q -dimensional parameter vector $\boldsymbol{\theta}$, and $L(\boldsymbol{\theta})$ be a negative log-likelihood function defined by $L(\boldsymbol{\theta}) = -\log f(\mathbf{u}|\boldsymbol{\theta})$. Suppose that

$$\dot{f}_{a_1 \dots a_j} = \frac{\partial^j}{\partial \theta_{a_1} \dots \partial \theta_{a_j}} f(\mathbf{u}|\boldsymbol{\theta}), \quad \dot{L}_{a_1 \dots a_j} = \frac{\partial^j}{\partial \theta_{a_1} \dots \partial \theta_{a_j}} L(\boldsymbol{\theta}).$$

By carrying out tedious calculations, we have

$$\begin{aligned}\dot{L}_a &= -\frac{\dot{f}_a}{f}, & \dot{L}_{ab} &= \dot{L}_a \dot{L}_b - \frac{\dot{f}_{ab}}{f}, \\ \dot{L}_{abc} &= -\dot{L}_a \dot{L}_b \dot{L}_c + \sum_{[3]} \dot{L}_a \dot{L}_{bc} - \frac{\dot{f}_{abc}}{f}, \\ \dot{L}_{abcd} &= \dot{L}_a \dot{L}_b \dot{L}_c \dot{L}_d - \sum_{[6]} \dot{L}_a \dot{L}_b \dot{L}_{cd} + \sum_{[3]} \dot{L}_{ab} \dot{L}_{cd} + \sum_{[4]} \dot{L}_a \dot{L}_{bcd} - \frac{\dot{f}_{abcd}}{f},\end{aligned}\tag{B.1}$$

where we simplify $f(\mathbf{u}|\boldsymbol{\theta})$ as f , and $\sum_{[j]}$ is the summation of a total of j terms of different combinations, e.g., $\sum_{[3]} \dot{L}_{ab} \dot{L}_{cd} = \dot{L}_{ab} \dot{L}_{cd} + \dot{L}_{ac} \dot{L}_{bd} + \dot{L}_{ad} \dot{L}_{bc}$. It follows from $\int f d\mathbf{u} = 1$ that

$$\begin{aligned}E\left[\frac{\dot{f}_{a_1 \dots a_j}}{f}\right] &= \int \frac{\dot{f}_{a_1 \dots a_j}}{f} f d\mathbf{u} \\ &= \int \frac{\partial^j}{\partial \theta_{a_1} \dots \partial \theta_{a_j}} f d\mathbf{u} = \frac{\partial^j}{\partial \theta_{a_1} \dots \partial \theta_{a_j}} \int f d\mathbf{u} = 0.\end{aligned}\tag{B.2}$$

The above equation can be satisfied when \mathbf{u} is continuous. Even when \mathbf{u} is discrete, we can obtain the same result by replacing the integration with a summation. Eqs. (B.1) and (B.2) imply that

$$\begin{aligned}E[\dot{L}_a] &= 0, & E[\dot{L}_{ab}] &= E[\dot{L}_a \dot{L}_b], \\ E[\dot{L}_{abc}] &= -E[\dot{L}_a \dot{L}_b \dot{L}_c] + \sum_{[3]} E[\dot{L}_a \dot{L}_{bc}], \\ E[\dot{L}_{abcd}] &= E[\dot{L}_a \dot{L}_b \dot{L}_c \dot{L}_d] - \sum_{[6]} E[\dot{L}_a \dot{L}_b \dot{L}_{cd}] + \sum_{[3]} E[\dot{L}_{ab} \dot{L}_{cd}] + \sum_{[4]} E[\dot{L}_a \dot{L}_{bcd}].\end{aligned}\tag{B.3}$$

Let us consider a vector of the first derivatives, and matrices of the second, third, and fourth derivatives, which are defined as

$$\begin{aligned}\mathbf{g}(\boldsymbol{\theta}) &= -\frac{\partial}{\partial \boldsymbol{\theta}} \ell(\boldsymbol{\theta}), & \mathbf{H}(\boldsymbol{\theta}) &= -\frac{\partial^2}{\partial \boldsymbol{\theta} \partial \boldsymbol{\theta}'} \ell(\boldsymbol{\theta}), \\ \mathbf{C}(\boldsymbol{\theta}) &= -\left(\frac{\partial}{\partial \boldsymbol{\theta}'} \otimes \frac{\partial^2}{\partial \boldsymbol{\theta} \partial \boldsymbol{\theta}'}\right) \ell(\boldsymbol{\theta}), & \mathbf{Q}(\boldsymbol{\theta}) &= -\left(\frac{\partial^2}{\partial \boldsymbol{\theta} \partial \boldsymbol{\theta}'} \otimes \frac{\partial^2}{\partial \boldsymbol{\theta} \partial \boldsymbol{\theta}'}\right) \ell(\boldsymbol{\theta}).\end{aligned}\tag{B.4}$$

From the expectations in (B.3), we obtain $E[\mathbf{H}(\boldsymbol{\theta})]$, $E[\mathbf{C}(\boldsymbol{\theta})]$, and $E[\mathbf{Q}(\boldsymbol{\theta})]$ as

$$\begin{aligned}E[\mathbf{H}(\boldsymbol{\theta})] &= E[\mathbf{g}(\boldsymbol{\theta}) \mathbf{g}(\boldsymbol{\theta})'], \\ E[\mathbf{C}(\boldsymbol{\theta})] &= -E[\mathbf{g}(\boldsymbol{\theta})' \otimes \mathbf{g}(\boldsymbol{\theta}) \mathbf{g}(\boldsymbol{\theta})'] + E[\mathbf{g}(\boldsymbol{\theta})' \otimes \mathbf{H}(\boldsymbol{\theta})] \\ &\quad + E[\mathbf{H}(\boldsymbol{\theta}) \otimes \mathbf{g}(\boldsymbol{\theta})'] + E[\mathbf{g}(\boldsymbol{\theta}) \text{vec}(\mathbf{H}(\boldsymbol{\theta}))'], \\ E[\mathbf{Q}(\boldsymbol{\theta})] &= E[\mathbf{g}(\boldsymbol{\theta}) \mathbf{g}(\boldsymbol{\theta})' \otimes \mathbf{g}(\boldsymbol{\theta}) \mathbf{g}(\boldsymbol{\theta})'] \\ &\quad - (\mathbf{I}_{q^2} + \mathbf{K}_q) E[\mathbf{g}(\boldsymbol{\theta}) \mathbf{g}(\boldsymbol{\theta})' \otimes \mathbf{H}(\boldsymbol{\theta})] (\mathbf{I}_{q^2} + \mathbf{K}_q) \\ &\quad - E[\text{vec}(\mathbf{g}(\boldsymbol{\theta}) \mathbf{g}(\boldsymbol{\theta})') \text{vec}(\mathbf{H}(\boldsymbol{\theta}))'] \\ &\quad - E[\text{vec}(\mathbf{H}(\boldsymbol{\theta})) \text{vec}(\mathbf{g}(\boldsymbol{\theta}) \mathbf{g}(\boldsymbol{\theta})')'] \\ &\quad + (\mathbf{I}_{q^2} + \mathbf{K}_q) E[\{\mathbf{H}(\boldsymbol{\theta}) \otimes \mathbf{H}(\boldsymbol{\theta})\}] \\ &\quad + E[\text{vec}(\mathbf{H}(\boldsymbol{\theta})) \text{vec}(\mathbf{H}(\boldsymbol{\theta}))'] \\ &\quad + E[\{\mathbf{g}(\boldsymbol{\theta}) \otimes \mathbf{C}(\boldsymbol{\theta})\}] (\mathbf{I}_{q^2} + \mathbf{K}_q) \\ &\quad + (\mathbf{I}_{q^2} + \mathbf{K}_q) E[\{\mathbf{g}(\boldsymbol{\theta})' \otimes \mathbf{C}(\boldsymbol{\theta})'\}].\end{aligned}\tag{B.5}$$

Recall that $E[\mathbf{g}(\boldsymbol{\theta})] = \mathbf{0}_q$ holds. Furthermore, we note that $\mathbf{C}(\boldsymbol{\theta})$ and $\mathbf{Q}(\boldsymbol{\theta})$ are constant when $\mathbf{H}(\boldsymbol{\theta})$ is constant. Hence, when $\mathbf{H}(\boldsymbol{\theta})$ is constant, $\mathbf{H}(\boldsymbol{\theta})$, $\mathbf{C}(\boldsymbol{\theta})$, and $\mathbf{Q}(\boldsymbol{\theta})$ become simpler, as follows:

$$\begin{aligned}\mathbf{H}(\boldsymbol{\theta}) &= E[\mathbf{g}(\boldsymbol{\theta})\mathbf{g}(\boldsymbol{\theta})'], \\ \mathbf{C}(\boldsymbol{\theta}) &= -E[\mathbf{g}(\boldsymbol{\theta})' \otimes \mathbf{g}(\boldsymbol{\theta})\mathbf{g}(\boldsymbol{\theta})'], \\ \mathbf{Q}(\boldsymbol{\theta}) &= E[\mathbf{g}(\boldsymbol{\theta})\mathbf{g}(\boldsymbol{\theta})' \otimes \mathbf{g}(\boldsymbol{\theta})\mathbf{g}(\boldsymbol{\theta})'] - (\mathbf{I}_{q^2} + \mathbf{K}_q)\{\mathbf{H}(\boldsymbol{\theta}) \otimes \mathbf{H}(\boldsymbol{\theta})\} - \text{vec}(\mathbf{H}(\boldsymbol{\theta}))\text{vec}(\mathbf{H}(\boldsymbol{\theta}))'.\end{aligned}\tag{B.6}$$

References

- [1] D.W. Hosmer, S. Lemeshow, *Applied Logistic Regression*, second ed., John Wiley & Sons, Inc., 2000.
- [2] T. Briz, R.W. Ward, Consumer awareness of organic products in Spain. An application of multinomial logit models, *Food Policy* 34 (2009) 295–304.
- [3] S.-W. Choi, B. Sohngen, R. Alig, An assessment of the influence of bioenergy and marketed land amenity values on land uses in the Midwestern US, *Ecol. Econ.* 70 (2011) 713–720.
- [4] L. dell'Olio, A. Ibeas, P. Cecin, The quality of service desired by public transport users, *Transport Policy* 18 (2011) 217–227.
- [5] R Development Core Team, *R: A Language and Environment for Statistical Computing*, R Foundation for Statistical Computing, 2011.
- [6] H. Akaike, Information theory and an extension of the maximum likelihood principle, in: B.N. Petrov, F. Csáki (Eds.), *Second International Symposium on Information Theory*, Akadémiai Kiadó, Budapest, 1973, pp. 267–281.
- [7] H. Akaike, A new look at the statistical model identification, *IEEE Trans. Automat. Control* AC-19 (1974) 716–723.
- [8] K.P. Burnham, D.R. Anderson, *Model Selection and Multi-model Inference: A Practical Information-theoretic Approach*, second ed., Springer-Verlag, 2002.
- [9] S. Konishi, Statistical model evaluation and information criteria, in: S. Ghosh (Ed.), *Multivariate Analysis, Design of Experiments, and Survey Sampling*, Marcel Dekker, 1999, pp. 369–399.
- [10] S. Konishi, G. Kitagawa, *Information Criteria and Statistical Modeling*, Springer Science + Business Media, LLC, New York, 2008.
- [11] S. Kullback, R. Leibler, On information and sufficiency, *Ann. Math. Statist.* 22 (1951) 79–86.
- [12] H. Yanagihara, R. Sekiguchi, Y. Fujikoshi, Bias correction of AIC in logistic regression models, *J. Statist. Plann. Inference* 115 (2003) 349–360.
- [13] K. Kamo, H. Yanagihara, K. Satoh, Bias-corrected AIC for selecting variables in Poisson regression models, *Comm. Statist. Theory Methods*, in press.
- [14] L. Fahrmeir, H. Kaufmann, Consistency and asymptotic normality of the maximum likelihood estimator in generalized linear models, *Ann. Statist.* 13 (1985) 79–86.
- [15] P. McCullagh, D.R. Cox, Invariants and likelihood ratio statistics, *Ann. Statist.* 14 (1986) 1419–1430.
- [16] P. Hall, *The Bootstrap and Edgeworth Expansion*, Springer-Verlag, 1992.
- [17] D.A. Harville, *Matrix Algebra from a Statistician's Perspective*, Springer-Verlag, 1997.