

information criterion (TIC; Takeuchi, 1976), the regularization information criterion (RIC; Shibata, 1989), the generalized information criterion (GIC; Konishi & Kitagawa, 1996, 2003), the extended information criterion (EIC; Ishiguro et al., 1997) and AIC_b (Cavanaugh and Shumway, 1997) among others.

In each of these information criteria, $D+a$ is an estimator of -2 times the predictive expected $\log-L$, which is essentially equivalent to the Kullback-Leibler (1951) information or distance. The penalty term a is also seen as that of bias correction up to order $O(1)$ for AIC when the model holds (correctly specified model) or includes a correct model (overspecified model) and for the other information criteria given above when the model is possibly misspecified or underspecified. The stochastic a of TIC is given by observed information matrices or its expectations followed by substituting the maximum likelihood estimators (MLEs) for parameters. In EIC and AIC_b the term a is obtained using the bootstrap (see also Kishino & Hasegawa, 1989; Shibata, 1997; Shimodaira & Hasegawa, 1999; Konishi & Kitagawa, 2008).

While D is typically given by the MLEs of parameters in these criteria, RIC and GIC extended it to include robust and ridge-type estimators of parameters. Ogasawara (2015a) gives the asymptotic cumulants of $n^{-1}AIC_W$ and $n^{-1}TIC_W$, where the subscript W indicates that the weighted score estimators (WSEs) of parameters are used in place of the MLEs in $n^{-1}AIC$ and $n^{-1}TIC$. Note that the MLE is a special case of the WSE when the weight is null. The factor n^{-1} is not for essential transformation but for tractability so that $n^{-1}AIC_W$ and $n^{-1}TIC_W$ are $O_p(1)$ as for the usual MLEs of parameters.

A well-known exception of (1.1) is the Bayes information criterion (BIC; Schwarz, 1978), where a is replaced by $a_n = q \log n = O(\log n)$. The original definition of Mallows' (1973) $C_P (= O_p(1))$ for multiple linear regression is also an exception of (1.1). However, C_P can be transformed to the form of (1.1) using the predictive mean square errors (MSEs) of observations (Fujikoshi and Satoh, 1997). Further, it can be shown that C_P is asymptotically equivalent to AIC under the assumption of normality with a correct model or overspecified model.

Assume that we have n independent and identically distributed observations. Divide (1.1) by n when $D = -2 \log \hat{L}_W \equiv -2 \hat{l}_W$ and $a = -\hat{b}$, where $\hat{l}_W = \log \hat{L}_W = \log L(\mathbf{X}, \hat{\boldsymbol{\theta}}_W) = l(\mathbf{X}, \boldsymbol{\theta}_W)$; \mathbf{X} is the $n \times p$ matrix of n observations for p variables; $\boldsymbol{\theta}_W$ is the $q \times 1$ vector of the WSEs of parameters defined explicitly in the next section; $n^{-1}b$ is the bias of $-n^{-1}2\hat{l}_W = O_p(1)$ up to order $O(n^{-1})$ under possible model misspecification;

and $\widehat{b} = \widehat{b}_W^{(1)} = b(\mathbf{X}, \widehat{\boldsymbol{\theta}}_W)$, $\widehat{b} = \widehat{b}_W^{(2)} = b(\widehat{\boldsymbol{\theta}}_W)$, $\widehat{b} = \widehat{b}_{ML}^{(1)} = b(\mathbf{X}, \widehat{\boldsymbol{\theta}}_{ML})$ or $\widehat{b} = \widehat{b}_{ML}^{(2)} = b(\widehat{\boldsymbol{\theta}}_{ML})$ is a consistent estimator of b . Then, (1.1) becomes

$$-n^{-1}2\widehat{l}_W - n^{-1}\widehat{b} \equiv -2\widehat{l}_W - n^{-1}\widehat{b}, \quad (1.2)$$

where $\widehat{l}_W = n^{-1}\widehat{l}_W = \widehat{O}_p(1)$ is the sample log- L given by \mathbf{X} and $\widehat{\boldsymbol{\theta}}_W (= \boldsymbol{\theta}_W(\mathbf{X}))$ averaged over observations. The population counterpart of $\widehat{\boldsymbol{\theta}}_W$ denoted by $\boldsymbol{\theta}_0$ under model misspecification is defined as the solution of $\boldsymbol{\theta}$ given by the expectation of the first order condition of $\widehat{\boldsymbol{\theta}}_{ML}$, where the expectation is that of the true distribution rather the fitted model.

Note that (1.2) is seen as a case of usual bias correction for e.g., a MLE of order $O_p(1)$ when we see $D = -2\widehat{l}_W$ as an estimator of a “parameter”. Let $\widehat{\theta}_W$ be an element of $\widehat{\boldsymbol{\theta}}_W$. Ogasawara (2013, 2014, 2015b) dealt with the bias adjusted estimator $\widehat{\theta}_{W(k)} \equiv \widehat{\theta}_W - n^{-1}k\widehat{\alpha}_1$, where $n^{-1}\alpha_1$ is the bias of $\widehat{\theta}_W$ up to order $O(n^{-1})$, $\widehat{\alpha}_1$ is a consistent estimator of α_1 and k is a constant. In Ogasawara (2013, 2014, 2015b), the optimal value k was derived such that k minimizes the asymptotic mean square error (AMSE) of $\widehat{\theta}_W$ up to order $O(n^{-2})$.

Recall that $-2\widehat{l}_W$ is an estimator of $-2n^{-1}$ times the predictive expected log- L . In this paper, Ogasawara’s method of bias adjustment is applied to information criteria. Further, another type of adjustment for $-2\widehat{l}_W$ by shrinkage is considered in combination with bias correction. In the following sections, it is shown that the AMSEs of $-2\widehat{l}_W$ by the adjustments are smaller than that of (1.2). Three examples using basic distributions in statistics are shown under model misspecification or correct model specification. A simulation and a real example for model selection in linear regression are also given.

2 Optimal Bias Adjustment of TIC

Define $f(\mathbf{X}^* = \mathbf{X} | \boldsymbol{\theta}_0) = f(\mathbf{X} | \boldsymbol{\theta}_0)$ as a probability density or mass of \mathbf{X}^* at \mathbf{X} when the model holds while $g(\mathbf{X}^* = \mathbf{X} | \boldsymbol{\zeta}_0) = g(\mathbf{X} | \boldsymbol{\zeta}_0) = g(\mathbf{X})$ is the corresponding true density or mass under model misspecification with $\boldsymbol{\zeta}_0$ being a vector of the population parameters of an appropriate size. Let \mathbf{x}_j^{*j} be the j -th row of \mathbf{X}^* ($j = 1, \dots, n$). Then, the log- L for $\boldsymbol{\theta}$ averaged over observations is

$$\bar{l} \equiv \bar{l}(\boldsymbol{\theta} | \mathbf{X}^*) \equiv n^{-1}l \equiv n^{-1} \sum_{j=1}^n l_j \equiv n^{-1} \sum_{j=1}^n \log f(\mathbf{x}_j^* | \boldsymbol{\theta}). \quad (2.1)$$

Let \mathbf{Z}^* be an independent copy of \mathbf{X}^* to be obtained in the future. Denote n^{-1} times the expected log- L by \bar{l}_0^* when $\boldsymbol{\theta} = \boldsymbol{\theta}_0$ from a viewpoint of prediction:

$$\bar{l}_0^* \equiv E_g\{\bar{l}(\boldsymbol{\theta}_0 | \mathbf{Z}^*)\} \equiv \int_{R(\mathbf{Z})} \bar{l}(\boldsymbol{\theta}_0 | \mathbf{Z})g(\mathbf{Z})d\mathbf{Z} (= E_g\{\bar{l}(\boldsymbol{\theta}_0 | \mathbf{X}^*)\}), \quad (2.2)$$

where the summation is to be used when $g(\cdot)$ is a probability function. The sample counterpart of (2.2) using $\hat{\boldsymbol{\theta}}_W$ is

$$\begin{aligned} \hat{\bar{l}}_W^* &\equiv E_g[\bar{l}\{\boldsymbol{\theta}_W(\mathbf{X}^*) | \mathbf{Z}^*\}] = \int_{R(\mathbf{Z})} \bar{l}\{\boldsymbol{\theta}_W(\mathbf{X}^*) | \mathbf{Z}\}g(\mathbf{Z})d\mathbf{Z} \\ &= \int_{R(\mathbf{Z})} \bar{l}(\hat{\boldsymbol{\theta}}_W | \mathbf{Z})g(\mathbf{Z})d\mathbf{Z}, \end{aligned} \quad (2.3)$$

where the vector $\hat{\boldsymbol{\theta}}_W$ of WSEs addressed earlier is defined as the solution of

$$\frac{\partial \bar{l}(\boldsymbol{\theta} | \mathbf{X}^*)}{\partial \boldsymbol{\theta}} + n^{-1} \mathbf{q}^* = \mathbf{0}, \quad (2.4)$$

\mathbf{q}^* is the $q \times 1$ vector of weights for the WSEs, which correspond to the log-prior derivatives in the case of Bayesian estimation but can be other weights.

The expected predictive log- L averaged over observations when $\boldsymbol{\theta} = \hat{\boldsymbol{\theta}}_W = \boldsymbol{\theta}_W(\mathbf{X}^*)$ is

$$E_g(\hat{\bar{l}}_W^*) = \int_{R(\mathbf{X})} \int_{R(\mathbf{Z})} \bar{l}\{\boldsymbol{\theta}_W(\mathbf{X}) | \mathbf{Z}\}g(\mathbf{Z})g(\mathbf{X})d\mathbf{Z}d\mathbf{X}. \quad (2.5)$$

On the other hand, the usual sample log- L averaged over observations denoted by $\hat{\bar{l}}_W$ using $\hat{\boldsymbol{\theta}}_W$ is

$$\hat{\bar{l}}_W = \bar{l}(\hat{\boldsymbol{\theta}}_W | \mathbf{X}^*) = \bar{l}\{\boldsymbol{\theta}_W(\mathbf{X}^*) | \mathbf{X}^*\}. \quad (2.6)$$

When we use information criteria, $E_g(\hat{\bar{l}}_W^*)$ of (2.5) is of primary interest although \bar{l}_0^* of (2.2) is also of interest since \bar{l}_0^* is the limiting value of $\hat{\bar{l}}_W^*$ when n goes to infinity. The bias of $-2\hat{\bar{l}}_W$ as an estimator of $-2E_g(\hat{\bar{l}}_W^*)$ is assumed to be of order $O_p(n^{-1})$ i.e.,

$$-2E_g(\hat{\bar{l}}_W) + 2E_g(\hat{\bar{l}}_W^*) = n^{-1}b + n^{-2}c + O(n^{-3}), \quad (2.7)$$

where b , which is unchanged when $W = ML$, is well-known:

$$\begin{aligned} b &= 2 \operatorname{tr}(\mathbf{\Lambda}^{-1}\mathbf{\Gamma}), \mathbf{\Lambda} = E_g \left(\left. \frac{\partial^2 \bar{l}}{\partial \boldsymbol{\theta} \partial \boldsymbol{\theta}'} \right|_{\boldsymbol{\theta}=\boldsymbol{\theta}_0} \right) \equiv E_g \left(\frac{\partial^2 \bar{l}}{\partial \boldsymbol{\theta}_0 \partial \boldsymbol{\theta}'_0} \right), \\ \mathbf{\Gamma} &\equiv n E_g \left(\frac{\partial \bar{l}}{\partial \boldsymbol{\theta}_0} \frac{\partial \bar{l}}{\partial \boldsymbol{\theta}'_0} \right) \end{aligned} \quad (2.8)$$

(see e.g., Ogasawara, 2015a) and the general formulas of c are available (see Konishi & Kitagawa, 2003 and Ogasawara, 2015a) but are not used in this paper.

Under correct model specification, $-\mathbf{\Lambda} = \mathbf{\Gamma} = \mathbf{I}$, where \mathbf{I} is the Fisher information matrix per observation. For $n^{-1}\text{TIC}_W$ addressed in the introductory section, the estimator $\widehat{b}_W^{(1)}$ or $\widehat{b}_W^{(2)}$ for b in (2.8) is used as follows:

$$\begin{aligned} \widehat{b}_W^{(1)} &\equiv 2 \operatorname{tr}(\widehat{\mathbf{\Lambda}}_W^{-1} \widehat{\mathbf{\Gamma}}_W) \\ &\equiv 2 \operatorname{tr} \left\{ \left(\left. \frac{\partial^2 \bar{l}}{\partial \boldsymbol{\theta} \partial \boldsymbol{\theta}'} \right|_{\boldsymbol{\theta}=\widehat{\boldsymbol{\theta}}_W} \right)^{-1} n^{-1} \sum_{j=1}^n \left. \frac{\partial l_j}{\partial \boldsymbol{\theta}} \right|_{\boldsymbol{\theta}=\widehat{\boldsymbol{\theta}}_W} \left. \frac{\partial l_j}{\partial \boldsymbol{\theta}'} \right|_{\boldsymbol{\theta}=\widehat{\boldsymbol{\theta}}_W} \right\} \\ &\equiv 2 \operatorname{tr} \left\{ \left(\frac{\partial^2 \bar{l}}{\partial \widehat{\boldsymbol{\theta}}_W \partial \widehat{\boldsymbol{\theta}}'_W} \right)^{-1} n^{-1} \sum_{j=1}^n \frac{\partial l_j}{\partial \widehat{\boldsymbol{\theta}}_W} \frac{\partial l_j}{\partial \widehat{\boldsymbol{\theta}}'_W} \right\}, \quad (2.9) \\ \widehat{b}_W^{(2)} &\equiv 2 \operatorname{tr}(-\widehat{\mathbf{I}}_W^{(\mathbf{\Lambda})-1} \widehat{\mathbf{I}}_W^{(\mathbf{\Gamma})}) \\ &\equiv 2 \operatorname{tr} \left(\left[\left\{ E_g \left(\frac{\partial^2 \bar{l}}{\partial \boldsymbol{\theta} \partial \boldsymbol{\theta}'} \right) \right\} \right]_{\boldsymbol{\theta}=\widehat{\boldsymbol{\theta}}_W}^{-1} E_g \left(\frac{\partial l_j}{\partial \boldsymbol{\theta}} \frac{\partial l_j}{\partial \boldsymbol{\theta}'} \right) \right)_{\boldsymbol{\theta}=\widehat{\boldsymbol{\theta}}_W}, \end{aligned}$$

where $\widehat{b}_W^{(2)}$ is generally simpler than $\widehat{b}_W^{(1)}$ though $\widehat{b}_W^{(2)}$ is not always available whereas $\widehat{b}_W^{(1)}$ is available even when $g(\mathbf{X} | \boldsymbol{\zeta}_0)$ is unknown.

Define

$$n^{-1}\text{TIC}_{W(k)}^{(j)} = -2\widehat{l}_W - n^{-1}k\widehat{b}_W^{(j)} \quad (j = 1, 2). \quad (2.10)$$

When $k = 1$ and $W = ML$, (2.10) becomes

$$n^{-1}\text{TIC}_{ML(1)}^{(j)} = n^{-1}\text{TIC}_{ML}^{(j)} = n^{-1}\text{TIC}^{(j)} \quad (j = 1, 2), \quad (2.11)$$

where $\text{TIC}^{(2)} = \text{TIC}$ is the original definition of TIC by Takeuchi (1976). Ogasawara (2015a) showed that under regularity conditions $-2\widehat{\bar{l}}_{\text{W}}$ is expanded as

$$\begin{aligned}
& -2\widehat{\bar{l}}_{\text{W}} = -2(\bar{l}_0^*)_{O(1)} - 2(\bar{l}_0 - \bar{l}_0^*)_{O_p(n^{-1/2})} + \left(\frac{\partial \bar{l}}{\partial \boldsymbol{\theta}'_0} \boldsymbol{\Lambda}^{-1} \frac{\partial \bar{l}}{\partial \boldsymbol{\theta}_0} \right)_{O_p(n^{-1})} \\
& + \left[-\text{vec}'(\mathbf{M}) \left(\boldsymbol{\Lambda}^{-1} \frac{\partial \bar{l}}{\partial \boldsymbol{\theta}_0} \right)^{\langle 2 \rangle} + \frac{1}{3} \text{vec}'\{\text{E}_g(\mathbf{J}_0^{(3)})\} \left(\boldsymbol{\Lambda}^{-1} \frac{\partial \bar{l}}{\partial \boldsymbol{\theta}_0} \right)^{\langle 3 \rangle} \right]_{O_p(n^{-3/2})} \\
& + O_p(n^{-2}) \\
& \equiv -2(\bar{l}_0^*)_{O(1)} + \sum_{j=1}^3 (\bar{l}_{\text{ML}}^{(j)})_{O_p(n^{-j/2})} + O_p(n^{-2}), \tag{2.12}
\end{aligned}$$

where $\bar{l}_0 = \bar{l}(\boldsymbol{\theta}_0 | \mathbf{X})$, $\mathbf{M} = \frac{\partial^2 \bar{l}}{\partial \boldsymbol{\theta}_0 \partial \boldsymbol{\theta}'_0} - \text{E}_g \left(\frac{\partial^2 \bar{l}}{\partial \boldsymbol{\theta}_0 \partial \boldsymbol{\theta}'_0} \right)$, $\mathbf{J}_0^{(3)} = \frac{\partial^3 \bar{l}}{\partial \boldsymbol{\theta}_0 (\partial \boldsymbol{\theta}'_0)^{\langle 2 \rangle}}$, $\text{vec}(\cdot)$ is the vectorizing operator stacking the columns of a matrix sequentially, $\text{vec}'(\cdot) = \{\text{vec}(\cdot)\}'$, $\mathbf{x}^{\langle k \rangle} = \mathbf{x} \otimes \cdots \otimes \mathbf{x}$ (k times of \mathbf{x}) is the k -fold Kronecker product of \mathbf{x} and $(\cdot)_{O_p(n^{-j/2})}$ indicates that (\cdot) is of order $O_p(n^{-j/2})$ for clarity. Equation (2.12) shows that $-2\widehat{\bar{l}}_{\text{W}}$ is equal to $-2\widehat{\bar{l}}_{\text{ML}}$ up to order $O_p(n^{-3/2})$. Similarly, it can be shown that $\widehat{b}_{\text{W}}^{(j)} = \widehat{b}_{\text{ML}}^{(j)} + O_p(n^{-1})$ ($j = 1, 2$) due to $\widehat{\boldsymbol{\theta}}_{\text{W}} = \widehat{\boldsymbol{\theta}}_{\text{ML}} + O_p(n^{-1})$ (see Ogasawara, 2014, Equation (2.4)). The vector $\widehat{\boldsymbol{\theta}}_{\text{ML}}$ is expanded as

$$\begin{aligned}
\widehat{\boldsymbol{\theta}}_{\text{ML}} - \boldsymbol{\theta}_0 &= \sum_{j=1}^3 (\boldsymbol{\Lambda}^{(j)})_{O(1)} (\mathbf{l}_0^{(j)})_{O_p(n^{-j/2})} + O_p(n^{-2}), \tag{2.13} \\
\boldsymbol{\Lambda}^{(1)} \mathbf{l}_0^{(1)} &= -\boldsymbol{\Lambda}^{-1} \frac{\partial \bar{l}}{\partial \boldsymbol{\theta}_0}, \\
\boldsymbol{\Lambda}^{(2)} \mathbf{l}_0^{(2)} &= \boldsymbol{\Lambda}^{-1} \mathbf{M} \boldsymbol{\Lambda}^{-1} \frac{\partial \bar{l}}{\partial \boldsymbol{\theta}_0} - \frac{1}{2} \boldsymbol{\Lambda}^{-1} \text{E}_g(\mathbf{J}_0^{(3)}) \left(\boldsymbol{\Lambda}^{-1} \frac{\partial \bar{l}}{\partial \boldsymbol{\theta}_0} \right)^{\langle 2 \rangle}, \\
\mathbf{l}_0^{(1)} &= \frac{\partial \bar{l}}{\partial \boldsymbol{\theta}_0}, \mathbf{l}_0^{(2)} = \left\{ \text{v}'(\mathbf{M}) \otimes \frac{\partial \bar{l}}{\partial \boldsymbol{\theta}'_0}, \left(\frac{\partial \bar{l}}{\partial \boldsymbol{\theta}'_0} \right)^{\langle 2 \rangle} \right\}', \boldsymbol{\Lambda}^{(1)} = -\boldsymbol{\Lambda}^{-1},
\end{aligned}$$

and $\boldsymbol{\Lambda}^{(2)}$ is implicitly defined by (2.13), where $\text{v}(\cdot)$ is the vectorizing operator taking the non-duplicated elements of a symmetric matrix (Ogasawara, 2010, Equation (2.4)).

Then, using the expansion of $\widehat{\boldsymbol{\theta}}_{\text{ML}} - \boldsymbol{\theta}_0$ up to order $O_p(n^{-1/2})$ in (2.13), we have

$$\begin{aligned} \widehat{b}_{\text{W}}^{(1)} &= 2 \operatorname{tr}(\boldsymbol{\Lambda}^{-1} \boldsymbol{\Gamma}) \\ &\quad - 2 \operatorname{tr} \left[\left[\boldsymbol{\Lambda}^{-1} \mathbf{M} \boldsymbol{\Lambda}^{-1} - \boldsymbol{\Lambda}^{-1} \mathbf{E}_g(\mathbf{J}_0^{(3)}) \left\{ \boldsymbol{\Lambda}^{-1} \otimes \left(\boldsymbol{\Lambda}^{-1} \frac{\partial \bar{l}}{\partial \boldsymbol{\theta}_0} \right) \right\} \right] \boldsymbol{\Gamma} \right. \\ &\quad \left. - \boldsymbol{\Lambda}^{-1} \left\{ \mathbf{M}_G - \sum_{j=1}^q \mathbf{E}_g(\mathbf{G}_{0(j)}^{(3)}) \left(\boldsymbol{\Lambda}^{-1} \frac{\partial \bar{l}}{\partial \boldsymbol{\theta}_0} \right)_j \right\} \right] \right]_{(A)O_p(n^{-1/2})} + O_p(n^{-1}), \end{aligned} \quad (2.14)$$

where

$$\begin{aligned} \mathbf{M}_G &= n^{-1} \sum_{j=1}^n \frac{\partial l_j}{\partial \boldsymbol{\theta}_0} \frac{\partial l_j}{\partial \boldsymbol{\theta}'_0} - n \mathbf{E}_g \left(\frac{\partial \bar{l}}{\partial \boldsymbol{\theta}_0} \frac{\partial \bar{l}}{\partial \boldsymbol{\theta}'_0} \right) \equiv \mathbf{G}_0 - \mathbf{E}_g(\mathbf{G}_0) = O_p(n^{-1/2}), \\ \mathbf{G}_{0(j)}^{(3)} &= \frac{\partial \mathbf{G}_0}{\partial (\boldsymbol{\theta}_0)_j} \quad (j = 1, \dots, q), \end{aligned} \quad (2.15)$$

$(\cdot)_j$ indicates the j -th element of a vector and $\left[\begin{smallmatrix} \cdot \\ (A) \end{smallmatrix} \right]$ is for ease of finding correspondence. On the other hand, $\widehat{b}_{\text{W}}^{(2)}$ for $n^{-1} \text{TIC}_{\text{W}}^{(2)}$ is given by omitting \mathbf{M} and \mathbf{M}_G in (2.14) i.e.,

$$\begin{aligned} \widehat{b}_{\text{W}}^{(2)} &= 2 \operatorname{tr}(\boldsymbol{\Lambda}^{-1} \boldsymbol{\Gamma}) + 2 \operatorname{tr} \left[\left[\boldsymbol{\Lambda}^{-1} \mathbf{E}_g(\mathbf{J}_0^{(3)}) \left\{ \boldsymbol{\Lambda}^{-1} \otimes \left(\boldsymbol{\Lambda}^{-1} \frac{\partial \bar{l}}{\partial \boldsymbol{\theta}_0} \right) \right\} \right] \boldsymbol{\Gamma} \right. \\ &\quad \left. - \boldsymbol{\Lambda}^{-1} \sum_{j=1}^q \mathbf{E}_g(\mathbf{G}_{0(j)}^{(3)}) \left(\boldsymbol{\Lambda}^{-1} \frac{\partial \bar{l}}{\partial \boldsymbol{\theta}_0} \right)_j \right] + O_p(n^{-1}). \end{aligned} \quad (2.16)$$

Define $\text{MSE}(\cdot)_{\rightarrow O(n^{-2})}$ as the MSE of (\cdot) up to order $O(n^{-2})$. Then,

$$\begin{aligned} &\text{MSE}(n^{-1} \text{TIC}_{\text{W}(k)}^{(j)})_{\rightarrow O(n^{-2})} \\ &= \text{avar}_g(-2\widehat{l}_{\text{W}})_{\rightarrow O(n^{-2})} \\ &\quad + n^{-2} \{ (1-k)^2 b^2 - 2kn \text{cov}_g(\widehat{b}_{\text{W}}^{(j)}, -2\widehat{l}_{\text{W}}) \} (j = 1, 2), \end{aligned} \quad (2.17)$$

where $\text{avar}_g(-\widehat{2\bar{l}}_W)_{\rightarrow O(n^{-2})}$ is the higher-order asymptotic variance of $-\widehat{2\bar{l}}_W$ up to order $O(n^{-2})$ with $\text{avar}_g(-\widehat{2\bar{l}}_W)_{\rightarrow O(n^{-2})} = \text{avar}_g(-\widehat{2\bar{l}}_{ML})_{\rightarrow O(n^{-2})}$ due to (2.12), and from (2.14)

$$\begin{aligned}
 & \text{nacov}_g(\widehat{b}_W^{(1)}, -\widehat{2\bar{l}}_W) = \text{nacov}_g(\widehat{b}_W^{(1)}, \bar{l}_W^{(1)}) = \text{nacov}_g(\widehat{b}_{ML}^{(1)}, -2\bar{l}_0) \\
 & = \text{nacov}_g \left[\underset{(A)}{-2 \text{tr}} \left[\underset{(B)}{\left[\mathbf{\Lambda}^{-1} \mathbf{M} \mathbf{\Lambda}^{-1} - \mathbf{\Lambda}^{-1} \mathbf{E}_g(\mathbf{J}_0^{(3)}) \left\{ \mathbf{\Lambda}^{-1} \otimes \left(\mathbf{\Lambda}^{-1} \frac{\partial \bar{l}}{\partial \boldsymbol{\theta}_0} \right) \right\} \right]} \right] \mathbf{\Gamma} \right. \\
 & \quad \left. - \mathbf{\Lambda}^{-1} \left\{ \mathbf{M}_G - \sum_{j=1}^q \mathbf{E}_g(\mathbf{G}_{0(j)}^{(3)}) \left(\mathbf{\Lambda}^{-1} \frac{\partial \bar{l}}{\partial \boldsymbol{\theta}_0} \right)_j \right\} \right] \underset{(B)}, -2\bar{l}_0 \right] \underset{(A)}{,} \\
 & = 4 \left[\underset{(A)}{\text{vec}'(\mathbf{\Lambda}^{-1} \mathbf{\Gamma} \mathbf{\Lambda}^{-1}) n \mathbf{E}_g\{\text{vec}(\mathbf{M}), \bar{l}_0\}} \right. \\
 & \quad - \text{vec}'\{\mathbf{E}_g(\mathbf{J}_0^{(3)})\} \left[\text{vec}(\mathbf{\Lambda}^{-1} \mathbf{\Gamma} \mathbf{\Lambda}^{-1}) \otimes \left\{ \mathbf{\Lambda}^{-1} n \mathbf{E}_g \left(\frac{\partial \bar{l}}{\partial \boldsymbol{\theta}_0} \bar{l}_0 \right) \right\} \right] \\
 & \quad - \text{tr}\{\mathbf{\Lambda}^{-1} n \mathbf{E}_g(\mathbf{M}_G \bar{l}_0)\} \\
 & \quad \left. + \text{tr} \left[\mathbf{\Lambda}^{-1} \sum_{j=1}^q \mathbf{E}_g(\mathbf{G}_{0(j)}^{(3)}) \left\{ \mathbf{\Lambda}^{-1} n \mathbf{E}_g \left(\frac{\partial \bar{l}}{\partial \boldsymbol{\theta}_0} \bar{l}_0 \right) \right\}_j \right] \right] \underset{(A)}. \tag{2.18}
 \end{aligned}$$

where $\text{acov}_g(\cdot)$ is the asymptotic covariance of order $O(n^{-1})$ for two variables in parentheses under possible model misspecification.

The result for $\widehat{b}_W^{(2)}$ corresponding to (2.18) is given by

$$\begin{aligned}
 & \text{nacov}_g(\widehat{b}_W^{(2)}, -\widehat{2\bar{l}}_W) = \text{nacov}_g(\widehat{b}_W^{(2)}, \bar{l}_W^{(1)}) = \text{nacov}_g(\widehat{b}_{ML}^{(2)}, -2\bar{l}_0) \\
 & = 4 \left[\underset{(A)}{- \text{vec}'\{\mathbf{E}_g(\mathbf{J}_0^{(3)})\}} \left[\text{vec}(\mathbf{\Lambda}^{-1} \mathbf{\Gamma} \mathbf{\Lambda}^{-1}) \otimes \left\{ \mathbf{\Lambda}^{-1} n \mathbf{E}_g \left(\frac{\partial \bar{l}}{\partial \boldsymbol{\theta}_0} \bar{l}_0 \right) \right\} \right] \right. \\
 & \quad \left. + \text{tr} \left[\mathbf{\Lambda}^{-1} \sum_{j=1}^q \mathbf{E}_g(\mathbf{G}_{0(j)}^{(3)}) \left\{ \mathbf{\Lambda}^{-1} n \mathbf{E}_g \left(\frac{\partial \bar{l}}{\partial \boldsymbol{\theta}_0} \bar{l}_0 \right) \right\}_j \right] \right] \underset{(A)}. \tag{2.19}
 \end{aligned}$$

Then, from (2.17) to (2.19), we have

Theorem 1. Under regularity conditions of (2.12) and (2.3) with possible model misspecification, the constant k minimizing the MSE up to order $O_p(n^{-2})$ for $n^{-1}TIC_{W(k)}^{(j)}$ denoted by $k_{\min}^{(Tj)}$ ($j = 1, 2$) defined by (2.10) with (2.9) using the WSE $\widehat{\boldsymbol{\theta}}_W$ is

$$k_{\min}^{(Tj)} = 1 + b^{-2}nacov_g(\widehat{b}_{ML}^{(j)}, -2\bar{l}_0), \quad (2.20)$$

which is equal to $k_{\min}^{(Tj)}$ by $\widehat{\boldsymbol{\theta}}_{ML}$. The minimized value is

$$\begin{aligned} MSE(n^{-1}TIC_{W(k_{\min}^{(Tj)})}^{(j)})_{\rightarrow O(n^{-2})} &= MSE(n^{-1}TIC_{ML(k_{\min}^{(Tj)})}^{(j)})_{\rightarrow O(n^{-2})} \\ &= avar_g(-2\widehat{l}_{ML})_{\rightarrow O(n^{-2})} + n^{-2}b^2\{1 - (k_{\min}^{(Tj)})^2\} \\ &= avar_g(-2\widehat{l}_{ML})_{\rightarrow O(n^{-2})} \\ &\quad - n^{-2}[2nacov_g(\widehat{b}_{ML}^{(j)}, -2\bar{l}_0) + b^{-2}\{nacov_g(\widehat{b}_{ML}^{(j)}, -2\bar{l}_0)\}^2] \quad (j = 1, 2). \end{aligned} \quad (2.21)$$

Theorem 1 shows that generally $MSE(n^{-1}TIC_{W(k_{\min}^{(Tj)})}^{(j)})_{\rightarrow O(n^{-2})}$ is smaller than $MSE(n^{-1}TIC_W^{(j)})_{\rightarrow O(n^{-2})}$ ($= MSE(n^{-1}TIC_{ML}^{(j)})_{\rightarrow O(n^{-2})} = avar_g(-2\widehat{l}_{ML})_{\rightarrow O(n^{-2})} - n^{-2}\{2nacov_g(\widehat{b}_{ML}^{(j)}, -2\bar{l}_0)\}_{O(1)}$) by $n^{-2}b^{-2}\{nacov_g(\widehat{b}_{ML}^{(j)}, -2\bar{l}_0)\}^2$ ($j = 1, 2$). Note that for the derivation of $k_{\min}^{(Tj)}$ ($j = 1, 2$), $avar_g(-2\widehat{l}_{ML})_{\rightarrow O(n^{-2})}$ was not used although $MSE(n^{-1}TIC_{W(k_{\min}^{(Tj)})}^{(j)})_{\rightarrow O(n^{-2})}$ depends on this value, which is somewhat complicated but can be obtained, if necessary, using Ogasawara (2015a). Note also that generally $k_{\min}^{(Tj)}$ ($j = 1, 2$) depend on $\boldsymbol{\theta}_0$ and is not directly available in practice. This issue will be addressed in a discussion section.

Under correct model specification, the bias adjustment of $n^{-1}AIC_W$ as $n^{-1}AIC_{W(k)} = -2\widehat{l}_W + n^{-1}2kq$ does not decrease the original value $MSE(n^{-1}AIC_W)_{\rightarrow O(n^{-2})}$ since by this adjustment the variance is unchanged whereas the absolute value of the bias up to order $O(n^{-1})$ when $k \neq 1$ is larger than that of $n^{-1}AIC_W$. Under possible model misspecification, the optimal k minimizing $MSE(n^{-1}AIC_{W(k)})_{\rightarrow O(n^{-2})}$ is $k_{\min}^{(A)} = -b/(2q)$, which gives $n^{-1}AIC_{W(k_{\min}^{(A)})} = -2\widehat{l}_W - n^{-1}b$. This bias-adjusted $n^{-1}AIC_W$ is equal to $n^{-1}TIC_W^{(j)}$ ($j = 1, 2$) when $\widehat{b}_W^{(j)}$ is replaced by b though b is not directly available in practice. Note, however, that $MSE(n^{-1}TIC_{W(k)}^{(j)})_{\rightarrow O(n^{-2})}$ and $MSE(n^{-1}AIC_{W(k)})_{\rightarrow O(n^{-2})}$ are generally different, which is shown as follows.

Corollary 1. *Under the same conditions as in Theorem 1,*

$$\begin{aligned}
 &MSE(n^{-1}TIC_{W(k)}^{(j)})_{\rightarrow O(n^{-2})} - MSE(n^{-1}AIC_{W(k)})_{\rightarrow O(n^{-2})} \\
 &= avar_g(-2\widehat{l}_{ML})_{\rightarrow O(n^{-2})} + n^{-2}\{(1-k)^2b^2 - 2knacov_g(\widehat{b}_{ML}^{(j)}, -2\bar{l}_0)\} \\
 &\quad - avar_g(-2\widehat{l}_{ML})_{\rightarrow O(n^{-2})} - n^{-2}(b+2kq)^2 \tag{2.22} \\
 &= -n^{-2}\{2k(b+2q)b + (-b^2+4q^2)k^2 + 2knacov_g(\widehat{b}_{ML}^{(j)}, -2\bar{l}_0)\} \quad (j = 1, 2).
 \end{aligned}$$

Note that the sign of (2.22) is indeterminable. When a model is true, consequently when $b = -2q$, if $nacov_f(\widehat{b}_{ML}^{(j)}, -2\bar{l}_0)$ is positive, where $nacov_f(\cdot)$ indicates the asymptotic covariance under the true model, a positive k gives smaller $MSE(n^{-1}TIC_{W(k)}^{(j)})_{\rightarrow O(n^{-2})}$ than $MSE(n^{-1}AIC_{W(k)})_{\rightarrow O(n^{-2})}$. An example of the last result is $k = 1$ which gives the usual $n^{-1}AIC_W$ and $n^{-1}TIC_W^{(j)}$ ($j = 1, 2$). That is, in this case with $nacov_f(\widehat{b}_{ML}^{(j)}, -2\bar{l}_0) > 0$, $MSE(n^{-1}TIC_W^{(j)})_{\rightarrow O(n^{-2})}$ ($j = 1, 2$) are smaller than $MSE(n^{-1}AIC_W)_{\rightarrow O(n^{-2})}$ by $n^{-2}\{2nacov_f(\widehat{b}_{ML}^{(j)}, -2\bar{l}_0)\}$. On the other hand, when $nacov_f(\widehat{b}_{ML}^{(j)}, -2\bar{l}_0) < 0$, an opposite result is obtained. When $nacov_f(\widehat{b}_{ML}^{(j)}, -2\bar{l}_0) = 0$, $MSE(n^{-1}TIC_W^{(j)})_{\rightarrow O(n^{-2})} = MSE(n^{-1}AIC_W)_{\rightarrow O(n^{-2})}$ ($j = 1, 2$).

3 Optimal Information Criteria by Shrinkage

3.1. $n^{-1}OIC_{W(k)}$. In the expression $D + a$ of (1.1), consider the case of $D = n(-2\widehat{l}_W)$ and $a = k2\widehat{l}_W$ with positive k . Divide $D + a$ by n and define this as a candidate optimal information criterion:

$$n^{-1}OIC_{W(k)} \equiv (1 - n^{-1}k)(-2\widehat{l}_W), \tag{3.1}$$

which is a shrinkage estimator of $-2E_g(\widehat{l}_W^*)$. Since

$$\begin{aligned}
 &MSE(n^{-1}OIC_{W(k)})_{\rightarrow O(n^{-2})} \\
 &= avar_g(-2\widehat{l}_W)_{\rightarrow O(n^{-2})} + n^{-2}\{(b+2k\bar{l}_0^*)^2 - 8knvar_g(\bar{l}_0)\}, \tag{3.2}
 \end{aligned}$$

we have

Theorem 2. *Under the same conditions as in Theorem 1, the constant k minimizing $MSE(n^{-1}OIC_{W(k)})_{\rightarrow O(n^{-2})}$ is*

$$k_{\min} = \frac{1}{2(\bar{l}_0^*)^2}\{-\bar{l}_0^*b + 2nvar_g(\bar{l}_0)\} \tag{3.3}$$

and the minimized value is

$$\begin{aligned}
 & \text{MSE}(n^{-1}\text{OIC}_{W(k_{\min})})_{\rightarrow O(n^{-2})} \\
 &= \text{avar}_g(-2\widehat{\bar{l}}_W)_{\rightarrow O(n^{-2})} + n^{-2}\{b^2 - 4(\bar{l}_0^*)^2 k_{\min}^2\} \\
 &= \text{avar}_g(-2\widehat{\bar{l}}_W)_{\rightarrow O(n^{-2})} - n^{-2}4 \left[-\frac{b}{\bar{l}_0^*} n\text{var}_g(\bar{l}_0) + \frac{1}{(\bar{l}_0^*)^2} \{n\text{var}_g(\bar{l}_0)\}^2 \right]. \quad (3.4)
 \end{aligned}$$

Note that $\text{MSE}(n^{-1}\text{OIC}_{W(k_{\min})})_{\rightarrow O(n^{-2})}$ is smaller than $\text{MSE}(-2\widehat{\bar{l}}_W)_{\rightarrow O(n^{-2})}$ by $4(\bar{l}_0^*)^2 k_{\min}^2 = (\bar{l}_0^*)^{-2} \{-\bar{l}_0^* b + 2n\text{var}_g(\bar{l}_0)\}^2$. Comparison of the results of $n^{-1}\text{OIC}_{W(k_{\min})}$, $n^{-1}\text{AIC}_W$ and $n^{-1}\text{TIC}_W^{(j)}$ ($j = 1, 2$) are given as follows.

$$\begin{aligned}
 & \text{MSE}(n^{-1}\text{AIC}_W)_{\rightarrow O(n^{-2})} - \text{MSE}(n^{-1}\text{OIC}_{W(k_{\min})})_{\rightarrow O(n^{-2})} \\
 &= n^{-2}(b + 2q)^2 - n^{-2}\{b^2 - 4(\bar{l}_0^*)^2 k_{\min}^2\} \\
 &= n^{-2} \left[(b + 2q)^2 + 4 \left\{ -\frac{b}{\bar{l}_0^*} n\text{var}_g(\bar{l}_0) + \frac{1}{(\bar{l}_0^*)^2} \{n\text{var}_g(\bar{l}_0)\}^2 \right\} \right]. \quad (3.5)
 \end{aligned}$$

The sign of $b = \text{tr}(\mathbf{\Lambda}^{-1}\mathbf{\Gamma})$ is usually negative while \bar{l}_0^* can also be negative. So, when b/\bar{l}_0^* is negative or the value of positive b/\bar{l}_0^* is relatively small, (3.5) is positive especially under model misspecification with $(b + 2q)^2 > 0$ (note that $b = -2q$ under correct model specification). Similarly, we have

$$\begin{aligned}
 & \text{MSE}(n^{-1}\text{TIC}_W^{(j)})_{\rightarrow O(n^{-2})} - \text{MSE}(n^{-1}\text{OIC}_{W(k_{\min})})_{\rightarrow O(n^{-2})} \\
 &= n^{-2}4n\text{acov}_g(\widehat{b}_{\text{ML}}^{(j)}, \bar{l}_0) - n^{-2}\{b^2 - 4(\bar{l}_0^*)^2 k_{\min}^2\} \\
 &= n^{-2}4 \left[n\text{acov}_g(\widehat{b}_{\text{ML}}^{(j)}, \bar{l}_0) - \frac{b}{\bar{l}_0^*} n\text{var}_g(\bar{l}_0) + \frac{1}{(\bar{l}_0^*)^2} \{n\text{var}_g(\bar{l}_0)\}^2 \right] \quad (j = 1, 2). \quad (3.6)
 \end{aligned}$$

3.2. $n^{-1}\text{OIC}_{W(k)}^{(A)}$. In Section 3.1, the bias adjustment of $-2\widehat{\bar{l}}_W$ was not explicitly considered as for $n^{-1}\text{AIC}_{W(k)}$ and $n^{-1}\text{TIC}_{W(k)}^{(j)}$ ($j = 1, 2$) in Section 2, though k_{\min} in Theorem 2 depends on b . In this section, a shrinkage estimator using $n^{-1}\text{AIC}_W$ is considered. Define

$$\begin{aligned}
 n^{-1}\text{OIC}_{W(k)}^{(A)} &\equiv (1 - n^{-1}k)n^{-1}\text{AIC}_W = (1 - n^{-1}k)(-2\widehat{\bar{l}}_W + n^{-1}2q) \\
 &= n^{-1}\text{AIC}_W + n^{-1}2k\widehat{\bar{l}}_W + O(n^{-2}). \quad (3.7)
 \end{aligned}$$

Let $l_{0j} = l_j |_{\theta=\theta_0} = f(\mathbf{x}_j^* | \theta_0)$ ($j = 1, \dots, n$) (see (2.1)) and $c_v(\cdot)$ as the coefficient of variation of a variable. Then, since

$$\begin{aligned} \text{MSE}(n^{-1}\text{OIC}_{W(k)}^{(A)})_{\rightarrow O(n-2)} &= \text{avar}_g(-2\widehat{l}_W)_{\rightarrow O(n-2)} + n^{-2}\{(b + 2q + 2k\bar{l}_0^*)^2 - 8kn\text{var}_g(\bar{l}_0)\}, \end{aligned} \tag{3.8}$$

we have

Theorem 3. *Under the same conditions as in Theorem 1, the constant k minimizing $\text{MSE}(n^{-1}\text{OIC}_{W(k)}^{(A)})_{\rightarrow O(n-2)}$ is*

$$k_{\min}^{(OA)} = -\frac{b + 2q}{2\bar{l}_0^*} + \frac{\text{var}_g(l_{0j})}{\{E_g(l_{0j})\}^2} = -\frac{b + 2q}{2\bar{l}_0^*} + c_v^2(l_{0j}) \tag{3.9}$$

and the minimized value is

$$\begin{aligned} \text{MSE}(n^{-1}\text{OIC}_{W(k_{\min}^{(OA)})}^{(A)})_{\rightarrow O(n-2)} &= \text{avar}_g(-2\widehat{l}_W)_{\rightarrow O(n-2)} + n^{-2}\{(b + 2q)^2 - 4(\bar{l}_0^*)^2(k_{\min}^{(OA)})^2\} \\ &= \text{avar}_g(-2\widehat{l}_W)_{\rightarrow O(n-2)} + n^{-2}4\{\bar{l}_0^*(b + 2q)c_v^2(l_{0j}) - (\bar{l}_0^*)^2c_v^4(l_{0j})\}. \end{aligned} \tag{3.10}$$

Theorem 3 with (3.8) shows that $\text{MSE}(n^{-1}\text{OIC}_{W(k_{\min}^{(OA)})}^{(A)})_{\rightarrow O(n-2)}$ is smaller than $\text{MSE}(n^{-1}\text{AIC}_W)_{\rightarrow O(n-2)}$ by

$$n^{-2}4(\bar{l}_0^*)^2(k_{\min}^{(OA)})^2 = n^{-2}\{-(b + 2q) + 2\bar{l}_0^*c_v^2(l_{0j})\}^2.$$

3.3. $n^{-1}\text{OIC}_{W(k)}^{(Tj)}$. In this subsection, shrinkage estimators using $n^{-1}\text{TIC}_W^{(j)}$ ($j = 1, 2$) are considered. Define

$$\begin{aligned} n^{-1}\text{OIC}_{W(k)}^{(Tj)} &= (1 - n^{-1}k)n^{-1}\text{TIC}_W^{(j)} = (1 - n^{-1}k)(-2\widehat{l}_W - n^{-1}\widehat{b}_W^{(j)}) \\ &= n^{-1}\text{TIC}_W^{(j)} + n^{-1}2k\widehat{l}_W + O_p(n^{-2}) \quad (j = 1, 2). \end{aligned} \tag{3.11}$$

Noting that

$$\begin{aligned} \text{MSE}(n^{-1}\text{OIC}_{W(k)}^{(Tj)})_{\rightarrow O(n-2)} &= \text{MSE}(n^{-1}\text{TIC}_W^{(j)})_{\rightarrow O(n-2)} + n^{-2}4\{(\bar{l}_0^*)^2k^2 - 2kn\text{var}_g(\bar{l}_0)\} \quad (j = 1, 2), \end{aligned} \tag{3.12}$$

we have

Theorem 4. *Under the same conditions as in Theorem 1, the constant k minimizing $MSE(n^{-1}OIC_{W(k)}^{(Tj)}) \rightarrow O(n^{-2})$ is*

$$k_{\min}^{(OTj)} = \frac{n \text{var}_g(\bar{l}_0)}{(\bar{l}_0^*)^2} = \frac{\text{var}_g(l_{0i})}{\{E_g(l_{0i})\}^2} = c_v^2(l_{0i}) \tag{3.13}$$

and the minimized value is

$$\begin{aligned} &MSE(n^{-1}OIC_{W(k_{\min}^{(OTj)})}^{(Tj)}) \rightarrow O(n^{-2}) \\ &= MSE(n^{-1}TIC_W^{(j)}) \rightarrow O(n^{-2}) - n^{-2}4(\bar{l}_0^*)^2 c_v^4(l_{0i}) (j = 1, 2). \end{aligned} \tag{3.14}$$

Note that $k_{\min}^{(OT1)} = k_{\min}^{(OT2)}$ in (3.13). Theorem 4 shows that $MSE(n^{-1}OIC_{W(k_{\min}^{(OTj)})}^{(Tj)}) \rightarrow O(n^{-2})$ is smaller than $MSE(n^{-1}TIC_W^{(j)}) \rightarrow O(n^{-2})$ by $n^{-2}4(\bar{l}_0^*)^2 c_v^4(l_{0i})$ ($j = 1, 2$).

Define $\hat{\theta}_u (= O_p(1))$ as an asymptotically unbiased estimator up to order $O(n^{-1})$ with the corresponding population value being θ_0 . Then, it is known that the multiplicative constant c^* for $\hat{\theta}_u$ minimizing the MSE of $c^*\hat{\theta}_u$ up to order $O(n^{-2})$ is given by $c_{\min}^* = 1 - n^{-1}\theta_0^{-2}n \text{avar}_g(\hat{\theta}_u)$ (Ogasawara, 2015b, the paragraph including Equation (2.6)). Since $n^{-1}TIC_W^{(j)}$ ($j = 1, 2$) are asymptotically unbiased up to order $O(n^{-1})$ and $k_{\min}^{(OTj)} = c_v^2(l_{0i})$, $n^{-1}OIC_{W(k_{\min}^{(OTj)})}^{(Tj)} = (1 - n^{-1}k_{\min}^{(OTj)})n^{-1}TIC_W^{(j)}$ ($j = 1, 2$) in (3.11) is seen as a special case of $c_{\min}^*\hat{\theta}_u$. Note also that

$$\begin{aligned} n^{-1}OIC_{W(k_{\min}^{(OTj)})}^{(Tj)} &= (1 - n^{-1}k_{\min}^{(OTj)})n^{-1}TIC_W^{(j)} = \frac{n^{-1}TIC_W^{(j)}}{1 + n^{-1}k_{\min}^{(OTj)}} + O_p(n^{-2}) \\ &= \frac{n^{-1}TIC_W^{(j)}}{1 + n^{-1}c_v^2(l_{0i})} + O_p(n^{-2}) \quad (j = 1, 2). \end{aligned} \tag{3.15}$$

Then, we find that the factor $1/\{1 + n^{-1}c_v^2(l_{0i})\}$ in (3.15) is the multiplicative constant for \bar{l}_0 exactly minimizing the MSE of \bar{l}_0 , where \bar{l}_0 is seen as a usual unbiased sample mean (Gruber, 1998, Section 1.6; Ogasawara, 2015b, Section 3.1).

It is of interest to see that $n^{-1}OIC_{W(k_{\min})}$ given in Section 3.1 is asymptotically equivalent to $n^{-1}OIC_{W(k_{\min}^{(OTj)})}^{(Tj)}$ ($j = 1, 2$) as follows.

Corollary 2. *Under the same conditions as in Theorem 1,*

$$n^{-1}OIC_{W(k_{\min})} = n^{-1}OIC_{W(k_{\min}^{(OTj)})}^{(Tj)} + O_p(n^{-3/2}) \quad (j = 1, 2). \tag{3.16}$$

PROOF.

$$\begin{aligned}
 n^{-1}\text{OIC}_{W(k_{\min})} &= (1 - n^{-1}k_{\min})(-\widehat{2\bar{l}}_W) \\
 &= \left[1 - \frac{n^{-1}}{2(\bar{l}_0^*)^2} \{-\bar{l}_0^*b + 2n\text{var}_g(\bar{l}_0)\} \right] (-\widehat{2\bar{l}}_W) \\
 &= -\widehat{2\bar{l}}_W - n^{-1}b - n^{-1} \frac{n\text{var}_g(\bar{l}_0)}{(\bar{l}_0^*)^2} (-\widehat{2\bar{l}}_W) + O_p(n^{-3/2}) \\
 &= \{1 - n^{-1}c_v^2(l_{0i})\}n^{-1}\text{TIC}_W^{(j)} + O_p(n^{-3/2}) \\
 &= n^{-1}\text{OIC}_{W(k_{\min}^{(\text{OT}j)})}^{(\text{T}j)} + O_p(n^{-3/2}) \quad (j = 1, 2). \tag{3.17}
 \end{aligned}$$

□

For comparison of the results of $n^{-1}\text{OIC}_{W(k_{\min}^{(\text{OA})})}^{(\text{A})}$ and $n^{-1}\text{OIC}_{W(k_{\min}^{(\text{OT}j)})}^{(\text{T}j)}$, we have from (3.10) and (3.14) with (2.17)

$$\begin{aligned}
 \text{MSE}(n^{-1}\text{OIC}_{W(k_{\min}^{(\text{OA})})}^{(\text{A})})_{\rightarrow O(n^{-2})} - \text{MSE}(n^{-1}\text{OIC}_{W(k_{\min}^{(\text{OT}j)})}^{(\text{T}j)})_{\rightarrow O(n^{-2})} \\
 = n^{-2}\{4\bar{l}_0^*(b + 2q)c_v^2(l_{0i}) - 4n\text{acov}_g(\widehat{b}_{\text{ML}}^{(j)}, \bar{l}_0)\} (j = 1, 2). \tag{3.18}
 \end{aligned}$$

The sign of (3.18) is indeterminable. Under correct model specification (3.18) becomes $-n^{-2}4n\text{acov}_g(\widehat{b}_{\text{ML}}^{(j)}, \bar{l}_0)$, which can be positive, negative or zero.

4 Examples

In this section, three simple examples are shown for illustration of computational aspects of the optimal information criteria. The first two examples are given under possible model misspecification using $\widehat{\theta}_{\text{ML}}$. The third example is shown under correct model specification when $\widehat{\theta}_W$ is used.

Example 1. The exponential distribution is used when the gamma distribution holds, where the shape parameter α is not necessarily equal to 1. That is,

$$\begin{aligned}
 f(x^* &= x \mid \lambda_0) = \lambda_0 \exp(-\lambda_0 x) (x > 0), \\
 g(x^* &= x \mid \lambda_1, \alpha) = x^{\alpha-1} \lambda_1 \exp(-\lambda_1 x) / \Gamma(\alpha) (x > 0), \\
 \widehat{\theta}_{\text{ML}} &= 1/\bar{x}^*, \theta_0 = \lambda_0, \zeta_0 = (\lambda_1, \alpha)', \tag{4.1}
 \end{aligned}$$

\bar{x}^* is the usual sample mean,

$$\begin{aligned}
 E_g(\bar{x}^*) &= \alpha/\lambda_1 = 1/\lambda_0, \lambda_0 = \lambda_1/\alpha, l_{0j} = -\lambda_0 x_j^* + \log \lambda_0, \\
 \bar{l}_0^* &= E_g(l_{0j}) = -\lambda_0(\alpha/\lambda_1) + \log \lambda_0 = \log(\lambda_1/\alpha) - 1,
 \end{aligned}$$

$$n\text{var}_g(\bar{l}_0) = \text{var}_g(l_{0j}) = \lambda_0^2 \text{var}_g(x_j^*) = \lambda_0^2(\alpha/\lambda_1^2) = (\lambda_1/\alpha)^2(\alpha/\lambda_1^2) = 1/\alpha,$$

$$c_v^2(l_{0j}) = \frac{1}{\alpha\{\log(\lambda_1/\alpha) - 1\}^2}, \quad \frac{\partial \bar{l}}{\partial \theta_0} = \frac{1}{\lambda_0} - \bar{x}^*,$$

$$\mathbf{\Lambda} = \lambda = \frac{\partial^2 \bar{l}}{\partial \theta_0^2} = -\frac{1}{\lambda_0^2} = -\frac{\alpha^2}{\lambda_1^2},$$

$$\mathbf{\Gamma} = \gamma = n\text{E}_g \left\{ \left(\frac{\partial \bar{l}}{\partial \theta_0} \right)^2 \right\} = \text{E}_g \left\{ \left(\frac{1}{\lambda_0} - x_j^* \right)^2 \right\} = \frac{\alpha}{\lambda_1^2},$$

$$b = 2 \text{tr}(\mathbf{\Lambda}^{-1} \mathbf{\Gamma}) = 2\lambda^{-1}\gamma = -\frac{2}{\alpha},$$

$$n^{-1}\text{AIC}_{\text{ML}} = -2\widehat{l}_{\text{ML}} + n^{-1}2, \quad n^{-1}\text{TIC}_{\text{ML}}^{(2)} = -2\widehat{l}_{\text{ML}} + n^{-1}2\widehat{\alpha}^{-1}.$$

Note that in this example $n^{-1}\text{TIC}_{\text{ML}}^{(1)}$ is not used since $n^{-1}\text{TIC}_{\text{ML}}^{(2)}$ is available.

(1) $n^{-1}\text{AIC}_{\text{ML}}$

The MSEs up to order $O(n^{-2})$ of $n^{-1}\text{AIC}_{\text{ML}}$ and $n^{-1}\text{OIC}_{\text{W}(k_{\min}^{\text{OA}})}^{(\text{A})}$ are compared. From (3.9) and (3.10),

$$\begin{aligned} k_{\min}^{(\text{OA})} &= -\frac{b+2q}{2\bar{l}_0^*} + c_v^2(l_{0j}) = -\frac{-(2/\alpha)+2}{2\{\log(\lambda_1/\alpha)-1\}} + c_v^2(l_{0j}) \\ &= \frac{\alpha^{-1}-1}{\log(\lambda_1/\alpha)-1} + \frac{1}{\alpha\{\log(\lambda_1/\alpha)-1\}^2} = \frac{(1-\alpha)\log(\lambda_1/\alpha)+\alpha}{\alpha\{\log(\lambda_1/\alpha)-1\}^2}, \end{aligned}$$

$$\begin{aligned} \text{MSE}(n^{-1}\text{AIC}_{\text{ML}})_{\rightarrow O(n^{-2})} &= \text{avar}_g(-2\widehat{l}_{\text{ML}})_{\rightarrow O(n^{-2})} + n^{-2}(b+2q)^2 \\ &= \text{avar}_g(-2\widehat{l}_{\text{ML}})_{\rightarrow O(n^{-2})} + n^{-2}4(1-\alpha^{-1})^2, \end{aligned}$$

$$\begin{aligned} \text{MSE}(n^{-1}\text{OIC}_{\text{ML}(k_{\min}^{\text{OA}})}^{(\text{A})})_{\rightarrow O(n^{-2})} &= \text{avar}_g(-2\widehat{l}_{\text{ML}})_{\rightarrow O(n^{-2})} + n^{-2}4\{\bar{l}_0^*(b+2q)c_v^2(l_{0j}) - (\bar{l}_0^*)^2 c_v^4(l_{0j})\} \\ &= \text{avar}_g(-2\widehat{l}_{\text{ML}})_{\rightarrow O(n^{-2})} \\ &\quad + n^{-2}4 \left[2(1-\alpha^{-1}) \frac{1}{\alpha\{\log(\lambda_1/\alpha)-1\}} - \frac{1}{\alpha^2\{\log(\lambda_1/\alpha)-1\}^2} \right] \\ &= \text{avar}_g(-2\widehat{l}_{\text{ML}})_{\rightarrow O(n^{-2})} + n^{-2}4 \frac{2(\alpha-1)\log(\lambda_1/\alpha)-2\alpha+1}{\alpha^2\{\log(\lambda_1/\alpha)-1\}^2}, \end{aligned}$$

(4.2)

where $\text{avar}_g(-2\widehat{l}_{\text{ML}}) \rightarrow O(n^{-2}) = n^{-1}4\alpha^{-1} + n^{-2}2\alpha^{-2}$ (see Ogasawara, 2015a). That is, $\text{MSE}(n^{-1}\text{OIC}_{\text{ML}(k_{\min}^{\text{(OA)})}}^{(\text{A})}) \rightarrow O(n^{-2})$ is smaller than $\text{MSE}(n^{-1}\text{AIC}_{\text{ML}}) \rightarrow O(n^{-2})$ by

$$\begin{aligned} n^{-2}\{-(b+2q) + 2\bar{l}_0^*c_v^2(l_{0j})\}^2 &= n^{-2}\left[-2(1-\alpha^{-1}) + \frac{2}{\alpha\{\log(\lambda_1/\alpha)-1\}}\right]^2 \\ &= n^{-2}4\left[\frac{(\alpha-1)\log(\lambda_1/\alpha)-\alpha}{\alpha\{\log(\lambda_1/\alpha)-1\}}\right]^2. \end{aligned}$$

(2) $n^{-1}\text{OIC}_{\text{ML}(k)}$ and $n^{-1}\text{TIC}_{\text{ML}}^{(2)}$

The constant k minimizing $\text{MSE}(n^{-1}\text{OIC}_{\text{W}(k)}) \rightarrow O(n^{-2})$ with $n^{-1}\text{OIC}_{\text{W}(k)} = (1-n^{-1}k)(-2\widehat{l}_{\text{W}})$ (see (3.1)) is given from (3.3) as

$$\begin{aligned} k_{\min} &= \frac{1}{2(\bar{l}_0^*)^2}\{-\bar{l}_0^*b + 2n\text{var}_g(\bar{l}_0)\} = -\frac{1}{2\bar{l}_0^*}(2\lambda^{-1}\gamma) + c_v^2(l_{0i}) \\ &= \frac{1}{\alpha\{\log(\lambda_1/\alpha)-1\}} + \frac{1}{\alpha\{\log(\lambda_1/\alpha)-1\}^2} \\ &= \frac{\log(\lambda_1/\alpha)}{\alpha\{\log(\lambda_1/\alpha)-1\}^2}. \end{aligned} \quad (4.3)$$

The constant k minimizing $\text{MSE}(n^{-1}\text{OIC}_{\text{ML}(k)}^{(\text{T2})}) \rightarrow O(n^{-2})$ with $n^{-1}\text{OIC}_{\text{ML}(k)}^{(\text{T2})} = (1-n^{-2}k)n^{-1}\text{TIC}_{\text{ML}}^{(2)}$ (see (3.11) and (3.13)) is

$$k_{\min}^{(\text{OT2})} = c_v^2(l_{0i}) = \frac{1}{\alpha\{\log(\lambda_1/\alpha)-1\}^2}. \quad (4.4)$$

From (4.3) and (4.4),

$$k_{\min} = -\frac{b}{2\bar{l}_0^*} + k_{\min}^{(\text{OT2})} = \frac{1}{\alpha\{\log(\lambda_1/\alpha)-1\}} + k_{\min}^{(\text{OT2})}. \quad (4.5)$$

From (3.14), $\text{MSE}(n^{-1}\text{OIC}_{\text{ML}(k_{\min}^{(\text{OT2})})}^{(\text{T2})}) \rightarrow O(n^{-2})$ is smaller than $\text{MSE}(n^{-1}\text{TIC}_{\text{ML}}^{(2)}) \rightarrow O(n^{-2})$ by

$$n^{-2}4(\bar{l}_0^*)^2c_v^4(l_{0i}) = n^{-2}\frac{4\{\text{var}_g(l_{0i})\}^2}{(\bar{l}_0^*)^2} = n^{-2}\frac{4}{\alpha^2\{\log(\lambda_1/\alpha)-1\}^2}. \quad (4.6)$$

From (3.18) and $4nacov_g(\widehat{b}_{ML}^{(2)}, \bar{l}_0) = 0$ in this example (see Ogasawara, 2015a),

$$\begin{aligned} & \text{MSE}(n^{-1}\text{OIC}_{ML(k_{\min}^{(OA)})}^{(A)})_{\rightarrow O(n^{-2})} - \text{MSE}(n^{-1}\text{OIC}_{ML(k_{\min}^{(OT2)})}^{(T2)})_{\rightarrow O(n^{-2})} \\ &= n^{-2}4\bar{l}_0^*(b + 2q)c_v^2(l_{0i}) = n^{-2}4(2 - 2\alpha^{-1})\frac{1}{\alpha\{\log(\lambda_1/\alpha) - 1\}} \\ &= n^{-2}\frac{8(\alpha - 1)}{\alpha^2\{\log(\lambda_1/\alpha) - 1\}}, \end{aligned} \tag{4.7}$$

whose sign is indeterminable.

From (3.5).

$$\begin{aligned} & \text{MSE}(n^{-1}\text{AIC}_{ML})_{\rightarrow O(n^{-2})} - \text{MSE}(n^{-1}\text{OIC}_{ML(k_{\min})})_{\rightarrow O(n^{-2})} \\ &= n^{-2}\left[(b + 2q)^2 + 4\left\{ -\frac{b}{\bar{l}_0^*}n\text{var}_g(\bar{l}_0) + \frac{1}{(\bar{l}_0^*)^2}\{n\text{var}_g(\bar{l}_0)\}^2 \right\} \right] \\ &= n^{-2}\left[4(1 - \alpha^{-1})^2 + 4\left\{ -\frac{1}{\log(\lambda_1/\alpha) - 1}\left(-\frac{2}{\alpha}\right)\frac{1}{\alpha} + \frac{1}{\alpha^2\{\log(\lambda_1/\alpha) - 1\}^2} \right\} \right] \\ &= n^{-2}4\left[(1 - \alpha^{-1})^2 + \frac{1}{\alpha^2}\left\{ \frac{2}{\log(\lambda_1/\alpha) - 1} + \frac{1}{\{\log(\lambda_1/\alpha) - 1\}^2} \right\} \right]. \end{aligned} \tag{4.8}$$

When λ_1/α is sufficiently large, (4.8) is positive.

From (3.6) and $nacov_g(\widehat{b}_{ML}^{(2)}, \bar{l}_0) = 0$,

$$\begin{aligned} & \text{MSE}(n^{-1}\text{TIC}_{ML}^{(2)})_{\rightarrow O(n^{-2})} - \text{MSE}(n^{-1}\text{OIC}_{ML(k_{\min})})_{\rightarrow O(n^{-2})} \\ &= n^{-2}4\left\{ -\frac{b}{\bar{l}_0^*}n\text{var}_g(\bar{l}_0) + \frac{1}{(\bar{l}_0^*)^2}\{n\text{var}_g(\bar{l}_0)\}^2 \right\} \\ &= n^{-2}\frac{4}{\alpha^2}\left\{ \frac{2}{\log(\lambda_1/\alpha) - 1} + \frac{1}{\{\log(\lambda_1/\alpha) - 1\}^2} \right\}. \end{aligned} \tag{4.9}$$

When $\log(\lambda_1/\alpha) - 1 > 0$, (4.9) is positive.

Example 2. In this example, the normal distribution with a known variance and the MLE of the mean is used although the true distribution is possibly non-normal with the known variance. That is,

$$\begin{aligned} f(x^* = x \mid \mu_0, \sigma^2) &= \frac{1}{\sqrt{2\pi}\sigma} \exp\left\{ -\frac{(x - \mu_0)^2}{2\sigma^2} \right\}, \\ \widehat{\theta}_{ML} = \bar{x}^*, \bar{l}_0^* &= -\frac{1}{2}\log(2\pi\sigma^2) - \frac{1}{2}, \text{var}_g(z^2) = \kappa_4 + 2, \end{aligned} \tag{4.10}$$

where $z = (x^* - \mu_0)/\sigma$, $\kappa_j = \kappa_j(z)$ is the j -th cumulant of z under possible non-normality,

$$\begin{aligned}
 c_v^2(l_{0i}) &= \frac{(\kappa_4 + 2)/4}{\{-\frac{1}{2}\log(2\pi\sigma^2) - \frac{1}{2}\}^2} = \frac{\kappa_4 + 2}{\{\log(2\pi\sigma^2) + 1\}^2}, \\
 \Lambda &= \lambda = -1/\sigma^2, \quad \gamma = 1/\sigma^2, \quad b = 2\lambda^{-1}\gamma = -2, \quad b + 2q = 0, \\
 n^{-1}\text{AIC}_{\text{ML}} &= -2\widehat{\bar{l}}_{\text{ML}} + n^{-1}2, \\
 n^{-1}\text{TIC}_{\text{ML}}^{(2)} &= -2\widehat{\bar{l}}_{\text{ML}} + n^{-1}2(-\lambda^{-1}\gamma) = -2\widehat{\bar{l}}_{\text{ML}} + n^{-1}2 = n^{-1}\text{AIC}_{\text{ML}}.
 \end{aligned}
 \tag{4.11}$$

As before $n^{-1}\text{TIC}_{\text{ML}}^{(1)}$ is not used. Note that in this example, $n^{-1}\text{AIC}_{\text{ML}} = n^{-1}\text{TIC}_{\text{ML}}^{(2)}$ even under non-normality.

From (3.9) and (3.10),

$$\begin{aligned}
 k_{\min}^{(\text{OA})} &= k_{\min}^{(\text{OT2})} = -\frac{b + 2q}{2\bar{l}_0^*} + c_v^2(l_{0i}) = c_v^2(l_{0i}), \\
 \text{MSE}(n^{-1}\text{AIC}_{\text{ML}}) &= \text{MSE}(n^{-1}\text{TIC}_{\text{ML}}^{(2)}) \\
 &= \text{avar}_g(-2\widehat{\bar{l}}_{\text{ML}})_{\rightarrow O(n^{-2})} + n^{-2}(b + 2q)^2 = \text{avar}_g(-2\widehat{\bar{l}}_{\text{ML}})_{\rightarrow O(n^{-2})}, \\
 \text{MSE}(n^{-1}\text{OIC}_{\text{ML}(k_{\min}^{(\text{OA})})}^{(\text{A})})_{\rightarrow O(n^{-2})} &= \text{MSE}(n^{-1}\text{OIC}_{\text{ML}(k_{\min}^{(\text{OT2})})}^{(\text{T2})})_{\rightarrow O(n^{-2})} \\
 &= \text{avar}_g(-2\widehat{\bar{l}}_{\text{ML}})_{\rightarrow O(n^{-2})} - n^{-2}4(\bar{l}_0^*)^2 c_v^4(l_{0i}) \\
 &= \text{avar}_g(-2\widehat{\bar{l}}_{\text{ML}})_{\rightarrow O(n^{-2})} - n^{-2} \frac{(\kappa_4 + 2)^2}{\{\log(2\pi\sigma^2) + 1\}^2}.
 \end{aligned}
 \tag{4.12}$$

That is, $\text{MSE}(n^{-1}\text{OIC}_{\text{ML}(k_{\min}^{(\text{OA})})}^{(\text{A})})_{\rightarrow O(n^{-2})} (= \text{MSE}(n^{-1}\text{OIC}_{\text{ML}(k_{\min}^{(\text{OT2})})}^{(\text{T2})})_{\rightarrow O(n^{-2})})$ is smaller than $\text{MSE}(n^{-1}\text{AIC}_{\text{ML}}) (= \text{MSE}(n^{-1}\text{TIC}_{\text{ML}}^{(2)}))$ by $n^{-2} \frac{(\kappa_4 + 2)^2}{\{\log(2\pi\sigma^2) + 1\}^2}$.

From (3.3) to (3.6),

$$\begin{aligned}
 k_{\min} &= \frac{1}{2(\bar{l}_0^*)^2} \{-\bar{l}_0^*b + 2n\text{var}_g(\bar{l}_0)\} = \frac{1}{\bar{l}_0^*} + c_v^2(l_{0i}) \\
 &= \frac{1}{\bar{l}_0^*} + k_{\min}^{(\text{OA})} = \frac{1}{\bar{l}_0^*} + k_{\min}^{(\text{OT2})}, \\
 \text{MSE}(n^{-1}\text{AIC}_{\text{ML}})_{\rightarrow O(n^{-2})} &- \text{MSE}(n^{-1}\text{OIC}_{\text{ML}(k_{\min})})_{\rightarrow O(n^{-2})} \\
 &= \text{MSE}(n^{-1}\text{TIC}_{\text{ML}}^{(2)})_{\rightarrow O(n^{-2})} - \text{MSE}(n^{-1}\text{OIC}_{\text{ML}(k_{\min})})_{\rightarrow O(n^{-2})} \\
 &= n^{-2} \left[(b + 2q)^2 + 4 \left\{ -\frac{b}{\bar{l}_0^*} n\text{var}_g(\bar{l}_0) + \frac{1}{(\bar{l}_0^*)^2} \{n\text{var}_g(\bar{l}_0)\}^2 \right\} \right]
 \end{aligned}$$

$$\begin{aligned}
&= n^{-2} 4 \left[\left\{ \frac{1}{2} \log(2\pi\sigma^2) + \frac{1}{2} \right\}^{-1} (-2)^{\frac{\kappa_4 + 2}{4}} \right. \\
&\quad \left. + \left\{ \frac{1}{2} \log(2\pi\sigma^2) + \frac{1}{2} \right\}^{-2} \frac{1}{16} (\kappa_4 + 2)^2 \right] \\
&= n^{-2} [-4\{\log(2\pi\sigma^2) + 1\}^{-1}(\kappa_4 + 2) + \{\log(2\pi\sigma^2) + 1\}^{-2}(\kappa_4 + 2)^2]. \tag{4.13}
\end{aligned}$$

Since $\kappa_4 + 2 = 4n\text{var}_g(\bar{l}_0) > 0$, when $\bar{l}_0^* = -\{\log(2\pi\sigma^2) + 1\}/2 > 0$, (4.13) is positive.

Example 3. In this example, it is assumed that the Bernoulli distribution holds, where the WSE of the logit with $a^*/2$ pseudocounts for each of two categories is used:

$$\begin{aligned}
\Pr(x^* = x | \theta) &= \pi^x (1 - \pi)^{1-x} \quad (x = 0, 1), \quad \pi = \frac{1}{1 + \exp(-\theta)}, \\
l_{0i} &= x_i^* \log \frac{\pi_0}{1 - \pi_0} + \log(1 - \pi_0) = x_i^* \theta_0 + \log(1 - \pi_0), \quad \pi_0 = \frac{1}{1 + \exp(-\theta_0)}, \\
\hat{\theta}_{\text{ML}} &= \log \frac{\bar{x}^*}{1 - \bar{x}^*}, \quad \hat{\theta}_{\text{W}} = \log \frac{\bar{x}^* + n^{-1}0.5a^*}{1 - \bar{x}^* + n^{-1}0.5a^*}, \\
\mathbf{\Lambda} = \lambda &= -\pi_0(1 - \pi_0) \equiv -\bar{i}_0, \quad \gamma = \bar{i}_0, \quad b = 2\lambda^{-1}\gamma = -2, \quad b + 2q = 0, \\
\bar{l}_0^* &= \pi_0\theta_0 + \log(1 - \pi_0), \quad \text{var}_f(l_{0i}) = \theta_0^2 \bar{i}_0, \\
c_{\text{V}}^2(l_{0i}) &= \frac{\theta_0^2 \bar{i}_0}{\{\pi_0\theta_0 + \log(1 - \pi_0)\}^2}, \quad n^{-1}\text{AIC}_{\text{W}} = n^{-1}\text{TIC}_{\text{W}}^{(2)} = -2\hat{l}_{\text{W}} + n^{-1}2. \tag{4.14}
\end{aligned}$$

From (3.9) and (3.10),

$$\begin{aligned}
k_{\text{min}}^{(\text{OA})} &= k_{\text{min}}^{(\text{OT2})} = -\frac{b + 2q}{2\bar{l}_0^*} + c_{\text{V}}^2(l_{0i}) = c_{\text{V}}^2(l_{0i}), \\
\text{MSE}(n^{-1}\text{AIC}_{\text{W}}) &= \text{MSE}(n^{-1}\text{TIC}_{\text{W}}^{(2)}) \\
&= \text{avar}_f(-2\hat{l}_{\text{ML}})_{\rightarrow O(n-2)} + n^{-2}(b + 2q)^2 = \text{avar}_f(-2\hat{l}_{\text{ML}})_{\rightarrow O(n-2)}, \\
\text{MSE}(n^{-1}\text{OIC}_{\text{W}(k_{\text{min}}^{(\text{OA})})}^{(\text{A})})_{\rightarrow O(n-2)} &= \text{MSE}(n^{-1}\text{OIC}_{\text{W}(k_{\text{min}}^{(\text{OT2})})}^{(\text{T2})})_{\rightarrow O(n-2)} \\
&= \text{avar}_f(-2\hat{l}_{\text{ML}})_{\rightarrow O(n-2)} + n^{-2}\{(b + 2q)^2 - 4(\bar{l}_0^*)^2(k_{\text{min}}^{(\text{OA})})^2\} \\
&= \text{avar}_f(-2\hat{l}_{\text{ML}})_{\rightarrow O(n-2)} - n^{-2}4(\bar{l}_0^*)^2 c_{\text{V}}^4(l_{0i}) \\
&= \text{avar}_f(-2\hat{l}_{\text{ML}})_{\rightarrow O(n-2)} - n^{-2} \frac{4\theta_0^4 \bar{i}_0^2}{\{\pi_0\theta_0 + \log(1 - \pi_0)\}^2}, \tag{4.15}
\end{aligned}$$

which shows that $\text{MSE}(n^{-1}\text{OIC}_{\text{W}(k_{\min}^{\text{OA}})}^{(\text{A})})_{\rightarrow O(n-2)}$ ($=\text{MSE}(n^{-1}\text{OIC}_{\text{W}(k_{\min}^{\text{OT2}})}^{(\text{T2})})_{\rightarrow O(n-2)}$) is smaller than $\text{MSE}(n^{-1}\text{AIC}_{\text{W}})$ ($=\text{MSE}(n^{-1}\text{TIC}_{\text{W}}^{(2)})$) by $n^{-2} \frac{4\theta_0^4 \bar{l}_0^{-2}}{\{\pi_0\theta_0 + \log(1 - \pi_0)\}^2}$.

From (3.3) to (3.6),

$$\begin{aligned}
 k_{\min} &= \frac{1}{2(\bar{l}_0^*)^2} \{-\bar{l}_0^* b + 2n\text{var}_f(\bar{l}_0)\} = \frac{1}{\bar{l}_0^*} + c_v^2(l_{0i}) \\
 &= \frac{1}{\bar{l}_0^*} + k_{\min}^{(\text{OA})} = \frac{1}{\bar{l}_0^*} + k_{\min}^{(\text{OT2})} \\
 &= \frac{1}{\pi_0\theta_0 + \log(1 - \pi_0)} + \frac{\theta_0^2 \bar{l}_0}{\{\pi_0\theta_0 + \log(1 - \pi_0)\}^2}, \\
 \text{MSE}(n^{-1}\text{AIC}_{\text{W}})_{\rightarrow O(n-2)} - \text{MSE}(n^{-1}\text{OIC}_{\text{W}(k_{\min})})_{\rightarrow O(n-2)} \\
 &= \text{MSE}(n^{-1}\text{TIC}_{\text{W}}^{(2)})_{\rightarrow O(n-2)} - \text{MSE}(n^{-1}\text{OIC}_{\text{W}(k_{\min})})_{\rightarrow O(n-2)} \\
 &= n^{-2} \left[(b + 2q)^2 + 4 \left\{ -\frac{b}{\bar{l}_0^*} n\text{var}_f(\bar{l}_0) + \frac{1}{(\bar{l}_0^*)^2} \{n\text{var}_f(\bar{l}_0)\}^2 \right\} \right] \\
 &= n^{-2} 4 \left[-\frac{1}{\pi_0\theta_0 + \log(1 - \pi_0)} (-2)\theta_0^2 \bar{l}_0 + \frac{\theta_0^4 \bar{l}_0^2}{\{\pi_0\theta_0 + \log(1 - \pi_0)\}^2} \right] \\
 &= n^{-2} \left[\frac{8\theta_0^2 \bar{l}_0}{\pi_0\theta_0 + \log(1 - \pi_0)} + \frac{4\theta_0^4 \bar{l}_0^2}{\{\pi_0\theta_0 + \log(1 - \pi_0)\}^2} \right]. \tag{4.16}
 \end{aligned}$$

Note that the above asymptotic results by $\hat{\theta}_{\text{W}}$ are the same as those by $\hat{\theta}_{\text{ML}}$ with $\text{avar}_f(-2\hat{l}_{\text{W}})_{\rightarrow O(n-2)} = \text{avar}_f(-2\hat{l}_{\text{ML}})_{\rightarrow O(n-2)}$ (see (2.12)).

5 A Simulation for Model Selection in Normal Linear Regression

Since information criteria are used primarily for selecting appropriate models, a simulation for model selection is carried out in the case of usual normal-theory linear regression with non-stochastic regressors or covariates under normality and non-normality. Let \mathbf{y} be the $n \times 1$ random vector with

$$\mathbf{y} = \mathbf{X}\beta_0 + \mathbf{e} \text{ and } \mathbf{e} \sim N(\mathbf{0}, \sigma_0^2 \mathbf{I}_{(n)}), \tag{5.1}$$

where n is the sample size, \mathbf{X} is a $n \times p$ design matrix, β_0 is a $p \times 1$ vector of population regression coefficients, \mathbf{e} is a $n \times 1$ random vector of errors, σ_0^2 is the common variance for the errors and $\mathbf{I}_{(n)}$ is the $n \times n$ identity

matrix. Using the likelihood of $\boldsymbol{\theta} = (\boldsymbol{\beta}', \sigma^2)'$ corresponding to the population counterpart $\boldsymbol{\theta}_0 = (\boldsymbol{\beta}'_0, \sigma_0^2)'$ i.e.,

$$L(\boldsymbol{\theta} | \mathbf{y}, \mathbf{X}) = (2\pi\sigma^2)^{-n/2} \exp\{-\frac{1}{2\sigma^2}(\mathbf{y} - \mathbf{X}\boldsymbol{\beta})'(\mathbf{y} - \mathbf{X}\boldsymbol{\beta})\}, \quad (5.2)$$

we have

$$\begin{aligned} \frac{\partial \bar{l}}{\partial \boldsymbol{\theta}_0} &= \left\{ \frac{n^{-1}}{\sigma_0^2}(\mathbf{y} - \mathbf{X}\boldsymbol{\beta}_0)' \mathbf{X}, -\frac{1}{2\sigma_0^2} + \frac{n^{-1}}{2\sigma_0^4}(\mathbf{y} - \mathbf{X}\boldsymbol{\beta}_0)'(\mathbf{y} - \mathbf{X}\boldsymbol{\beta}_0) \right\}' \\ &= \left(\frac{n^{-1}}{\sigma_0^2} \mathbf{e}' \mathbf{X}, -\frac{1}{2\sigma_0^2} + \frac{n^{-1}}{2\sigma_0^4} \mathbf{e}' \mathbf{e} \right)', \end{aligned} \quad (5.3)$$

and

$$\begin{aligned} \widehat{\boldsymbol{\theta}}_{\text{ML}} &= (\widehat{\boldsymbol{\beta}}'_{\text{ML}}, \widehat{\sigma}_{\text{ML}}^2)' = \{\mathbf{y}' \mathbf{X} (\mathbf{X}' \mathbf{X})^{-1}, n^{-1} \mathbf{y}' (\mathbf{I}_{(n)} - \mathbf{P}) \mathbf{y}\}' \\ &= \{\boldsymbol{\beta}'_0 + \mathbf{e}' \mathbf{X} (\mathbf{X}' \mathbf{X})^{-1}, n^{-1} \mathbf{e}' (\mathbf{I}_{(n)} - \mathbf{P}) \mathbf{e}\}', \end{aligned} \quad (5.4)$$

where $\mathbf{P} = \mathbf{X}(\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'$ is a projection matrix.

In this section, $n^{-1}\text{AIC}$, $n^{-1}\text{CAIC}$, $n^{-1}\text{TIC}^{(1)}$ and $n^{-1}\text{TIC}_{\text{ML}(k)}^{(1)}$ with various k 's including $k_{\min}^{(\text{T1})}$ are considered using $\widehat{\boldsymbol{\theta}}_{\text{ML}}$ of (5.4), where CAIC is a corrected AIC (Sugiura, 1978) to have an exactly unbiased estimator of $E_f(-2\bar{l}_{\text{ML}}^*)$ under correct model specification or under over specification with true and redundant regressors. Note that in the case of simulation in this section $n^{-1}\text{OIC}_{\text{ML}(k)}^{(\text{A})}$ and $n^{-1}\text{OIC}_{\text{ML}(k)}^{(\text{T}j)}$ ($j = 1, 2$) give the same model selection as those by $n^{-1}\text{AIC}_{\text{ML}}$ and $n^{-1}\text{TIC}_{\text{ML}}^{(j)}$ ($j = 1, 2$), respectively due to proportionality.

Since $-2\bar{l}_{\text{ML}} = \log(2\pi\widehat{\sigma}_{\text{ML}}^2) + 1$, we obtain

$$n^{-1}\text{AIC} = \log(2\pi\widehat{\sigma}_{\text{ML}}^2) + 1 + n^{-1}2(p+1). \quad (5.5)$$

On the other hand $n^{-1}\text{CAIC}$ is defined by

$$\begin{aligned} n^{-1}\text{CAIC} &= \log(2\pi\widehat{\sigma}_{\text{ML}}^2) + \frac{n+p}{n-p-2} = \log(2\pi\widehat{\sigma}_{\text{ML}}^2) + 1 + \frac{2(p+1)}{n-p-2} \\ &> \log(2\pi\widehat{\sigma}_{\text{ML}}^2) + 1 + n^{-1}2(p+1) = n^{-1}\text{AIC}. \end{aligned} \quad (5.6)$$

For the TIC,

$$n^{-1}\text{TIC}^{(1)} = n^{-1}\text{TIC}_{\text{ML}}^{(1)} = -2\bar{l}_{\text{ML}} - n^{-1}2 \text{tr}(\widehat{\boldsymbol{\Lambda}}_{\text{ML}}^{-1} \widehat{\boldsymbol{\Gamma}}_{\text{ML}}) \quad (5.7)$$

is used under normality and non-normality. It is assumed that even under non-normality, $E_g(\mathbf{e}) = \mathbf{0}$ and $E_g(\mathbf{e}\mathbf{e}') = \sigma_0^2 \mathbf{I}_{(n)}$ hold. In (5.7),

$$\begin{aligned} \widehat{\Lambda}_{\text{ML}} &= \frac{\partial^2 \bar{l}}{\partial \widehat{\boldsymbol{\theta}}_{\text{ML}} \partial \widehat{\boldsymbol{\theta}}'_{\text{ML}}} \\ &= \begin{bmatrix} -n^{-1} \mathbf{X}'\mathbf{X}/\widehat{\sigma}_{\text{ML}}^2 & -n^{-1} \mathbf{X}'(\mathbf{y} - \mathbf{X}\widehat{\boldsymbol{\beta}}_{\text{ML}})/\widehat{\sigma}_{\text{ML}}^4 \\ -n^{-1}(\mathbf{y} - \mathbf{X}\widehat{\boldsymbol{\beta}}_{\text{ML}})'\mathbf{X}/\widehat{\sigma}_{\text{ML}}^4 & \{1/(2\widehat{\sigma}_{\text{ML}}^4)\} \\ & -n^{-1}(\mathbf{y} - \mathbf{X}\widehat{\boldsymbol{\beta}}_{\text{ML}})'(\mathbf{y} - \mathbf{X}\widehat{\boldsymbol{\beta}}_{\text{ML}})/\widehat{\sigma}_{\text{ML}}^6 \end{bmatrix} \\ &= \begin{bmatrix} -n^{-1} \mathbf{X}'\mathbf{X}/\widehat{\sigma}_{\text{ML}}^2 & \mathbf{0} \\ \mathbf{0}' & -1/(2\widehat{\sigma}_{\text{ML}}^4) \end{bmatrix} \end{aligned} \tag{5.8}$$

and

$$\begin{aligned} \widehat{\Gamma}_{\text{ML}} &= n^{-1} \sum_{j=1}^n \frac{\partial l_j}{\partial \widehat{\boldsymbol{\theta}}_{\text{ML}}} \frac{\partial l_j}{\partial \widehat{\boldsymbol{\theta}}'_{\text{ML}}} \\ &= n^{-1} \sum_{j=1}^n \begin{bmatrix} (y_j - \widehat{y}_j)\mathbf{x}_j/\widehat{\sigma}_{\text{ML}}^2 \\ -\frac{1}{2\widehat{\sigma}_{\text{ML}}^2} + \frac{(y_j - \widehat{y}_j)^2}{2\widehat{\sigma}_{\text{ML}}^4} \end{bmatrix} \begin{bmatrix} \frac{y_j - \widehat{y}_j}{\widehat{\sigma}_{\text{ML}}^2} \mathbf{x}'_j, -\frac{1}{2\widehat{\sigma}_{\text{ML}}^2} + \frac{(y_j - \widehat{y}_j)^2}{2\widehat{\sigma}_{\text{ML}}^4} \end{bmatrix} \\ &= n^{-1} \sum_{j=1}^n \begin{bmatrix} \frac{[\{(\mathbf{I}_{(n)} - \mathbf{P})\mathbf{e}\}_j]^2 \mathbf{x}_j \mathbf{x}'_j/\widehat{\sigma}_{\text{ML}}^4}{[\{(\mathbf{I}_{(n)} - \mathbf{P})\mathbf{e}\}_j]^3 \mathbf{x}'_j} & \frac{[\{(\mathbf{I}_{(n)} - \mathbf{P})\mathbf{e}\}_j]^3 \mathbf{x}_j/(2\widehat{\sigma}_{\text{ML}}^6)}{-\frac{1}{4\widehat{\sigma}_{\text{ML}}^4} + \frac{[\{(\mathbf{I}_{(n)} - \mathbf{P})\mathbf{e}\}_j]^4}{4\widehat{\sigma}_{\text{ML}}^8}} \end{bmatrix}, \end{aligned} \tag{5.9}$$

where $\widehat{y}_j = (\widehat{\mathbf{y}})_j = (\mathbf{P}\mathbf{y})_j$ and \mathbf{x}'_j is the j -th row of \mathbf{X} ($j = 1, \dots, n$). From (5.7) to (5.9),

$$\begin{aligned} n^{-1} \text{TIC}^{(1)} &= -2\widehat{l}_{\text{ML}} - n^{-1} \widehat{b}_{\text{ML}}^{(1)} \\ &= -2\widehat{l}_{\text{ML}} - n^{-1} \underset{(A)}{\left[-2 \text{tr} \left[\underset{(B)}{(\mathbf{X}'\mathbf{X})^{-1}} \sum_{j=1}^n \left\{ [(\mathbf{I}_{(n)} - \mathbf{P})\mathbf{e}]_j \right\}^2 \mathbf{x}_j \mathbf{x}'_j / \widehat{\sigma}_{\text{ML}}^2 \right] \right]} \underset{(B)}{\left[\right]} \\ &\quad + 1 - n^{-1} \sum_{j=1}^n \left\{ [(\mathbf{I}_{(n)} - \mathbf{P})\mathbf{e}]_j \right\}^4 / \widehat{\sigma}_{\text{ML}}^4 \underset{(A)}{\left[\right]} \end{aligned} \tag{5.10}$$

follows.

For $k_{\min}^{(\text{T1})}$, we derive $\text{nacov}_g(\widehat{b}_{\text{ML}}^{(1)}, -2\widehat{l}_{\text{ML}})$ under possible non-normality. Note that

$$\text{nacov}_g(\widehat{b}_{\text{ML}}^{(1)}, -2\widehat{l}_{\text{ML}}) = \text{nacov}_g(\widehat{b}_{\text{ML}}^{(1)}, -2\bar{l}_0) = \text{nacov}_g(\widehat{b}_{\text{ML}}^{(1)}, n^{-1} \mathbf{e}'\mathbf{e}/\sigma_0^2) \tag{5.11}$$

On the other hand,

$$\begin{aligned} \text{nacov}_g(\widehat{b}_{\text{ML}}^{(1)}, -2\widehat{l}_{\text{ML}}) &= \text{nacov}_g\{\widehat{b}_{\text{ML}}^{(1)}, \log(2\pi\widehat{\sigma}_{\text{ML}}^2) + 1\} \\ &= \text{nacov}_g(\widehat{b}_{\text{ML}}^{(1)}, \log\widehat{\sigma}_{\text{ML}}^2) = \text{nacov}_g(\widehat{b}_{\text{ML}}^{(1)}, \widehat{\sigma}_{\text{ML}}^2/\sigma_0^2). \end{aligned} \quad (5.12)$$

Note also that $\{(\mathbf{I}_{(n)} - \mathbf{P})\mathbf{e}\}_j = e_j - (\mathbf{P}\mathbf{e})_j$, where $e_j = O_p(1)$ and $(\mathbf{P}\mathbf{e})_j = O_p(n^{-1/2})$ under some regularity conditions since $(\mathbf{P}\mathbf{e})_j$ is a weighted mean of e_1, \dots, e_n . Define $\bar{z} = e_1/\sigma_0$ and $\kappa_j = \kappa_j(\bar{z})$ is the j -th cumulant of \bar{z} . Then,

$$\begin{aligned} &\text{nacov}_g(\widehat{b}_{\text{ML}}^{(1)}, -2\widehat{l}_{\text{ML}}) \\ &= \text{nacov}_g \left[\underset{(A)}{-2 \operatorname{tr}[(\mathbf{X}'\mathbf{X})^{-1} \sum_{j=1}^n e_j^2 \mathbf{x}_j \mathbf{x}_j' / \widehat{\sigma}_{\text{ML}}^2]} - n^{-1} \sum_{j=1}^n e_j^4 / \widehat{\sigma}_{\text{ML}}^4, \underset{(A)}{n^{-1} \mathbf{e}'\mathbf{e} / \sigma_0^2} \right] \\ &= -n2 \operatorname{tr} \left[(\mathbf{X}'\mathbf{X})^{-1} n^{-1} \sum_{j=1}^n \left\{ \operatorname{cov}_g \left(\frac{e_j^2}{\sigma_0^2}, \frac{\mathbf{e}'\mathbf{e}}{\sigma_0^2} \right) - \frac{\mathbf{E}_g(e_j^2)}{\sigma_0^6} \operatorname{navar}_g(\widehat{\sigma}_{\text{ML}}^2) \right\} \mathbf{x}_j \mathbf{x}_j' \right] \\ &\quad - nn^{-2} \sum_{j=1}^n \left\{ \operatorname{cov}_g \left(\frac{e_j^4}{\sigma_0^4}, \frac{\mathbf{e}'\mathbf{e}}{\sigma_0^2} \right) - \frac{2\mathbf{E}_g(e_j^4)}{\sigma_0^8} \operatorname{navar}_g(\widehat{\sigma}_{\text{ML}}^2) \right\} \\ &= -2 \operatorname{tr}[\mathbf{I}_{(p)}\{(\kappa_4 + 2) - (\kappa_4 + 2)\}] - \{\mathbf{E}_g(\bar{z}^6) - (\kappa_4 + 3) - 2(\kappa_4 + 3)(\kappa_4 + 2)\} \\ &= -(\kappa_6 + 15\kappa_4 + 10\kappa_3^2 + 15) + (2\kappa_4 + 5)(\kappa_4 + 3) \\ &= -(\kappa_6 + 4\kappa_4 + 10\kappa_3^2) + 2\kappa_4^2. \end{aligned} \quad (5.13)$$

From (2.20) and (5.13), we find that $k_{\min}^{(\text{T1})} = 1$ and consequently, $n^{-1} \text{TIC}_{\text{ML}(k_{\min}^{(\text{T1})})}^{(1)} = n^{-1} \text{TIC}^{(1)}$ under normality.

For $b = 2 \operatorname{tr}(\mathbf{\Lambda}^{-1}\mathbf{\Gamma})$, we have

$$\mathbf{\Lambda} = \begin{bmatrix} -n^{-1} \mathbf{X}'\mathbf{X} / \sigma_0^2 & \mathbf{0} \\ \mathbf{0}' & -1 / (2\sigma_0^4) \end{bmatrix} = -\mathbf{I}_0 \quad (5.14)$$

even under non-normality and

$$\begin{aligned} \mathbf{\Gamma} &= \mathbf{E}_g \left(n^{-1} \sum_{j=1}^n \frac{\partial l_j}{\partial \boldsymbol{\theta}_0} \frac{\partial l_j}{\partial \boldsymbol{\theta}_0'} \right) \\ &= n^{-1} \mathbf{E}_g \left[\sum_{j=1}^n \begin{bmatrix} e_j \mathbf{x}_j / \sigma_0^2 \\ -\frac{1}{2\sigma_0^2} + \frac{e_j^2}{2\sigma_0^4} \end{bmatrix} \begin{bmatrix} \frac{e_j}{\sigma_0^2} \mathbf{x}_j', -\frac{1}{2\sigma_0^2} + \frac{e_j^2}{2\sigma_0^4} \end{bmatrix} \right] \\ &= \begin{bmatrix} n^{-1} \mathbf{X}'\mathbf{X} / \sigma_0^2 & n^{-1} \sum_{j=1}^n \kappa_3 \mathbf{x}_j / (2\sigma_0^3) \\ n^{-1} \sum_{j=1}^n \kappa_3 \mathbf{x}_j' / (2\sigma_0^3) & (\kappa_4 + 2) / (4\sigma_0^4) \end{bmatrix}, \end{aligned} \quad (5.15)$$

which gives

$$b = 2 \operatorname{tr}(\mathbf{\Lambda}^{-1}\mathbf{\Gamma}) = -2 \operatorname{tr} \begin{pmatrix} \mathbf{I}^{(p)} & \mathbf{0} \\ \mathbf{0}' & (\kappa_4 + 2)/2 \end{pmatrix} = -(2p + 2 + \kappa_4). \quad (5.16)$$

In the simulation, two sets of population values are used: [1] $p = 2$, $\beta_0 = (2, 1)'$, $\sigma_0^2 = 1$, $n = 20$ and 50 , [2] $p = 5$, $\beta_0 = (2, 1, 1, 1, 1)'$, $\sigma_0^2 = 1$, $n = 50$. The covariates are generated by the independent standard normal distributions except the first unit covariate corresponding to an intercept which is always included. In the first case, candidate models with $p = 1$ (only an intercept), $2, 3, 4$ are used while in the second case, those with $p = 1, 2, \dots, 8$ are used. For non-normal distributions of e_j , the chi-square distributions with 3 and 1 degrees of freedom followed by standardization to have mean zero and variance σ_0^2 are used. Their skewnesses are $2\sqrt{2}/\sqrt{3}$ and $2\sqrt{2}$, the excess kurtoses are 4 and 12 , and the k_6 's are $160/3$ and 480 for the two non-normal distributions, respectively. For $n^{-1}\operatorname{TIC}_{\text{ML}(k)}^{(1)}$, $k = k_{\min}^{(\text{T1})}, 2k_{\min}^{(\text{T1})}, 4k_{\min}^{(\text{T1})}, 5k_{\min}^{(\text{T1})}$ are used after some preliminary investigation.

Tables 1 and 2 show the proportions of correct model selection by the 7 information criteria when the 4 and 8 candidate models are available, respectively. In Table 1, with $n = 20$ $n^{-1}\text{CAIC}$ shows the best results. The added proportions of correct model selection over $n^{-1}\text{AIC}$ are substantial even under non-normality. On the other hand, usual or non-adjusted $n^{-1}\operatorname{TIC}^{(1)}$ gives poor results. When $k_{\min}^{(\text{T1})}$ is used for $n^{-1}\operatorname{TIC}_{\text{ML}(k)}^{(1)}$, the proportions of correct model selection have increased under non-normality. However, the increased proportions are still smaller than those by those by $n^{-1}\text{CAIC}$. When $k = 2k_{\min}^{(\text{T1})}$ or $4k_{\min}^{(\text{T1})}$ is used, the correct proportions have surprisingly increased especially when $n = 50$. In Table 2, under normality the correct proportion by $n^{-1}\operatorname{TIC}_{\text{ML}(4k_{\min}^{(\text{T1})})}^{(1)}$ is as large as 99.4 %.

It is known that AIC tends to select overspecified models or models with redundant regressor(s) (Hurvich and Tsai, 1989; Fujikoshi and Satoh, 1997). Tables 1 and 2 repeat this tendency. We find that in the tables, $n^{-1}\text{CAIC}$ also has this tendency. On the other hand, when k is increased from $k = k_{\min}^{(\text{T1})}$ to $k = 5k_{\min}^{(\text{T1})}$, it is clearly shown that the opposite tendency of selecting underspecified models gradually appears. The information criteria with the largest correct proportions in the tables seem to select under and overspecified models in a balanced way. Although $k_{\min}^{(\text{T1})}$ is not directly available in practice since $k_{\min}^{(\text{T1})}$ depends on θ_0 , the results in the tables are

encouraging in that when information about $k_{\min}^{(T1)}$ and optimal c for $ck_{\min}^{(T1)}$ accumulates using e.g., simulations, it is expected that the proportions of correct model selection can be increased as in Tables 1 and 2.

6 A Real Example for Model Selection

A real example is shown for model selection in linear regression using the house price data analyzed by Gilmour (1996, Table 1) in the context of model selection by Mallows' C_P and its modification, where the dependent variable is the house price and nine regressors are available with $n = 24$: x_1 taxes (\$/1000), x_2 number of baths, x_3 lot size (ft²/1000), x_4 living space (ft²/1000), x_5 number of garage stalls, x_6 number of rooms, x_7 number of bed rooms, x_8 age (years) and x_9 number of fireplaces. An intercept is always included in candidate models. That is, the number of possible models considered is $2^9 = 512$, where the model with an intercept only is also considered as a candidate.

Six information criteria $n^{-1}AIC$, $n^{-1}CAIC$, $n^{-1}TIC^{(1)}$ and $n^{-1}TIC_{ML}^{(1)(k)}$ with $k = k_{\min}^{(T1)}, 2k_{\min}^{(T1)}, 4k_{\min}^{(T1)}$ are used for model selection. For $k_{\min}^{(T1)}$, two conditions for the distributions of e_i are illustrated: the normal and chi-square (3 degrees of freedom) distributions, where the latter distribution is employed as a substantially non-normal distribution. For σ_0^2 , the unbiased estimator $n\hat{\sigma}_{ML}^2/(n - p)$ in the full model with $p = 10$ is used with the assumption that the model includes the true set of regressors as in Mallows' C_P . The value $-2E_f(\hat{l}_{ML}^*)$ under normality is given by

$$-2E_f(\hat{l}_{ML}^*) = \log \frac{4\pi\sigma_0^2}{n} + \psi\left(\frac{n - p}{2}\right) + \frac{n + p}{n - p - 2} \tag{6.1}$$

(see Davies et al., 2006, Equation (3.5)), where $\psi(\cdot)$ is the digamma function, and the last term $(n + p)/(n - p - 2)$ is equal to the correction term plus 1 in $n^{-1}CAIC$ (see (5.6)).

Figures 1 and 2 show the values of the six information criteria under normality and non-normality, respectively, where the horizontal axis is for the number of regressors including an intercept. A fine horizontal line indicates the value 5.168 of (6.1), when $p = 2$, which is also used under non-normality for comparison. Note that Gilmour (1996, p.54) concluded that the model of $p = 2$ with regressor x_1 and an intercept is the most appropriate model. In Fig. 1, $n^{-1}TIC^{(1)} = n^{-1}TIC_{ML}^{(1)(k_{\min}^{(T1)})}$ since $k_{\min}^{(T1)} = 1$ under normality.

The models selected by the minimum values of the six information criteria in Fig. 1 are $(1, x_1, x_2, x_5, x_7), (1, x_1, x_2), (1, x_1, x_2, x_4, x_5, x_7),$

Table 1: 1,000 times the proportions of model selection when the true number of regressors is 2 with 1,000 replications

	$n = 20$				$n = 50$			
	$p = 1$	2*	3	4	$p = 1$	2*	3	4
Normal	$(k_{\min}^{(T1)} = 1 \quad 1 \quad 1 \quad 1)$							
n^{-1} AIC	11	685	177	127	0	779	134	87
n^{-1} CAIC	25	<u>824</u>	114	37	0	833	116	51
n^{-1} TIC ⁽¹⁾	7	618	210	165	0	736	152	112
n^{-1} TIC _{ML(k)} ⁽¹⁾								
$k = k_{\min}^{(T1)}$	7	618	210	165	0	736	152	112
$k = 2k_{\min}^{(T1)}$	39	806	108	47	0	928	53	19
$k = 4k_{\min}^{(T1)}$	169	808	19	4	5	981	14	0
$k = 5k_{\min}^{(T1)}$	274	716	10	0	11	<u>982</u>	7	0
Chi-square ($df = 3$)	$(k_{\min}^{(T1)} = 0 \quad .36 \quad .56 \quad .67)$							
n^{-1} AIC	11	720	151	118	0	761	148	91
n^{-1} CAIC	28	<u>838</u>	95	39	0	818	127	55
n^{-1} TIC ⁽¹⁾	28	616	184	172	5	703	176	116
n^{-1} TIC _{ML(k)} ⁽¹⁾								
$k = k_{\min}^{(T1)}$	27	727	137	109	1	820	115	64
$k = 2k_{\min}^{(T1)}$	94	834	49	23	9	<u>959</u>	24	8
$k = 4k_{\min}^{(T1)}$	312	681	7	0	64	936	0	0
$k = 5k_{\min}^{(T1)}$	435	561	4	0	106	894	0	0
Chi-square ($df = 1$)	$(k_{\min}^{(T1)} = -.25 \quad .01 \quad .20 \quad .34)$							
n^{-1} AIC	27	715	173	85	0	743	149	108
n^{-1} CAIC	39	824	109	28	0	801	126	73
n^{-1} TIC ⁽¹⁾	66	572	201	161	12	656	186	146
n^{-1} TIC _{ML(k)} ⁽¹⁾								
$k = k_{\min}^{(T1)}$	18	672	194	116	0	778	143	79
$k = 2k_{\min}^{(T1)}$	57	812	98	33	5	927	52	16
$k = 4k_{\min}^{(T1)}$	133	<u>841</u>	22	4	24	<u>971</u>	5	0
$k = 5k_{\min}^{(T1)}$	178	806	15	1	40	958	2	0

Note. 2* = the true value. An underscore indicates the largest proportion of correct model selection under each condition

Table 2: 1,000 times the proportions of model selection when the true number of regressors is 5 with 1,000 replications

$n = 50$	$p = 1$	2	3	4	5*	6	7	8	
Normal	$(k_{\min}^{(T1)} = 1$	1	1	1	1	1	1	1)	
n^{-1} AIC		0	0	0	0	736	120	78	66
n^{-1} CAIC		0	0	0	0	849	77	47	27
n^{-1} TIC ⁽¹⁾		0	0	0	0	707	128	85	80
n^{-1} TIC _{ML(k)} ⁽¹⁾									
$k = k_{\min}^{(T1)}$		0	0	0	0	707	128	85	80
$k = 2k_{\min}^{(T1)}$		0	0	0	0	925	41	25	9
$k = 4k_{\min}^{(T1)}$		0	0	0	0	<u>994</u>	5	1	0
$k = 5k_{\min}^{(T1)}$		0	0	1	4	993	2	0	0
Chi-square ($df = 3$)	$(k_{\min}^{(T1)} = 0$.36	.56	.67	.75	.80	.84	.87)	
n^{-1} AIC		0	0	0	0	706	139	74	81
n^{-1} CAIC		0	0	0	0	815	106	46	33
n^{-1} TIC ⁽¹⁾		0	0	0	0	642	149	104	105
n^{-1} TIC _{ML(k)} ⁽¹⁾									
$k = k_{\min}^{(T1)}$		0	0	0	0	729	134	70	67
$k = 2k_{\min}^{(T1)}$		1	1	2	5	922	47	16	6
$k = 4k_{\min}^{(T1)}$		6	3	6	27	<u>954</u>	4	0	0
$k = 5k_{\min}^{(T1)}$		13	4	23	48	911	1	0	0
Chi-square ($df = 1$)	$(k_{\min}^{(T1)} = -.25$.01	.20	.34	.44	.53	.59	.64)	
n^{-1} AIC		0	0	0	0	703	147	85	65
n^{-1} CAIC		0	0	0	0	819	110	43	28
n^{-1} TIC ⁽¹⁾		0	0	3	6	602	175	111	103
n^{-1} TIC _{ML(k)} ⁽¹⁾									
$k = k_{\min}^{(T1)}$		0	0	1	0	739	134	66	60
$k = 2k_{\min}^{(T1)}$		1	1	6	8	928	42	6	8
$k = 4k_{\min}^{(T1)}$		4	4	16	39	<u>934</u>	3	0	0
$k = 5k_{\min}^{(T1)}$		5	5	42	48	899	1	0	0

Note. 5* = the true value. An underscore indicates the largest proportion of correct model selection under each condition

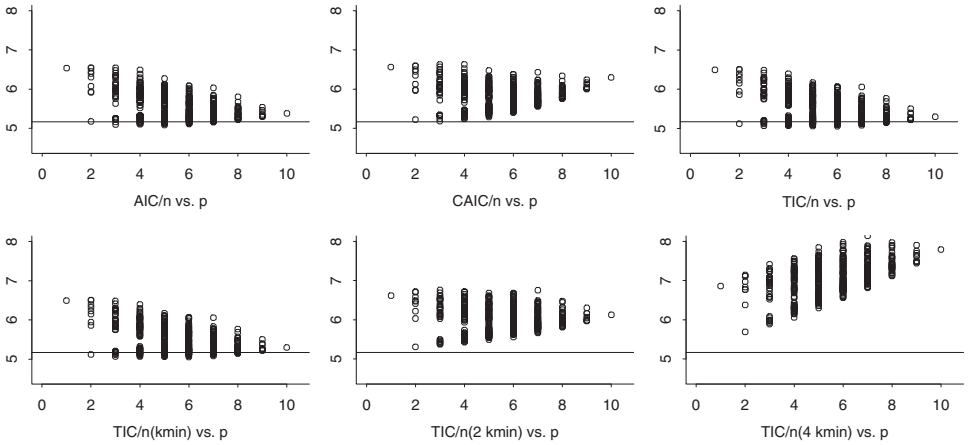


Figure 1: Information criteria for the house price data under normality

$(1, x_1, x_2, x_4, x_5, x_7), (1, x_1)$ and $(1, x_1)$ for $n^{-1}AIC$, $n^{-1}CAIC$, $n^{-1}TIC^{(1)}$, $n^{-1}TIC^{(1)}_{ML(k_{min}^{(T1)})}$, $n^{-1}TIC^{(1)}_{ML(2k_{min}^{(T1)})}$ and $n^{-1}TIC^{(1)}_{ML(4k_{min}^{(T1)})}$, respectively, where e.g., $(1, x_1)$ indicates the model with regressors constant 1 for an intercept and x_1 . The corresponding results in Fig. 2 are $(1, x_1, x_2, x_5, x_7), (1, x_1, x_2), (1, x_1, x_2, x_4, x_5, x_7), (1, x_1, x_2), (1, x_1)$ and $(1, x_1)$, where the results are unchanged except that $(1, x_1, x_2)$ by $n^{-1}TIC^{(1)}_{ML(k_{min}^{(T1)})}$ was $(1, x_1, x_2, x_4, x_5, x_7)$ in Fig. 1.

The appropriate model with $(1, x_1)$ is selected only by $n^{-1}TIC^{(1)}_{ML(2k_{min}^{(T1)})}$ and $n^{-1}TIC^{(1)}_{ML(4k_{min}^{(T1)})}$ in Figs. 1 and 2, which is consistent with the results in Section 5. Selection of the appropriate model was made possible by penalizing relatively complicated models. It is seen from the figures that $n^{-1}CAIC$ also has this tendency in comparison with $n^{-1}AIC$ though the correction is not sufficient.

7 Discussion

As mentioned earlier, optimal k values minimizing the associated MSEs generally depend on θ_0 , which is unknown in practice. The problem of k depending on θ_0 is similar to that choosing optimal ridge parameters depending on the population parameters (Hoerl & Kennard, 1970). Since θ_0 can be estimated, the optimal k can also be estimated using the same data as those for $\hat{\theta}_W$. However, if we use the estimated k , the optimality does not generally hold.

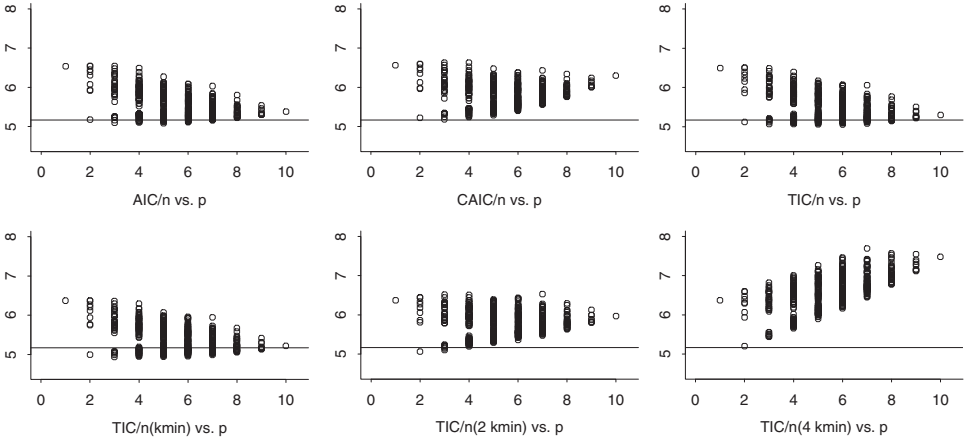


Figure 2: Information criteria for the house price data under chi-square (3 df)

Let an information criterion

$$IC_{W(k)} \equiv -2\widehat{l}_W + n^{-1}k\widehat{a} \tag{7.1}$$

and a similar one

$$IC_{W(\widehat{k})} \equiv -2\widehat{l}_W + n^{-1}\widehat{k}\widehat{a}, \tag{7.2}$$

where \widehat{k} is an estimated k using the same data. Then,

$$\begin{aligned} & \text{MSE}(IC_{W(\widehat{k})}) \rightarrow O(n^{-2}) \\ & = \text{MSE}(IC_{W(k)}) \rightarrow O(n^{-2}) + n^{-2}\{2anacov_g(-2\widehat{l}_W, \widehat{k})\}_{O(1)}. \end{aligned} \tag{7.3}$$

That is, except under the condition of $nacov_g(-2\widehat{l}_W, \widehat{k}) = 0$, the MSE up to order $O(n^{-2})$ changes although the MSE may decrease when $anacov_g(-2\widehat{l}_W, \widehat{k})$ is negative.

While it is difficult to overcome the above problem completely, several methods can be used. First, in a limited number of cases, optimal k 's do not depend on (the whole of) θ_0 . Note that in Example 2, the optimal k 's are not functions of $\theta_0 (= \mu_0)$. In Sections 5 and 6 for linear regression, $k_{\min}^{(T1)}$ does not depend on β_0 , but depends on the distributions of errors. Secondly, we often have some prior information about the values of θ_0 , which can give e.g., a lower bound of $k_{\min}^{(OA)}$. Since the MSEs of the information criteria up

to order $O(n^{-2})$ shown earlier with k are quadratic functions of k , we find that using the lower/upper bound decreases the original MSE. For instance, in Example 1 when $\alpha = 1$ or under correct model specification, (4.4) gives $k_{\min}^{(\text{OT}2)} = 1/(\log \lambda_1 - 1)^2 = 1/(\log \lambda_0 - 1)^2$. When it is known that an upper bound of λ_0 is $\lambda_0^*(> e)$, $1/(\log \lambda_0^* - 1)^2$ can be used as a lower bound for $k_{\min}^{(\text{OT}2)}$. Similarly, in Example 2, $k_{\min}^{(\text{OA})} = c_v^2(l_{0i}) = (\kappa_4 + 2)/\{\log(2\pi\sigma^2) + 1\}^2$ is a linear function of the excess kurtosis of z under arbitrary distributions. When a lower bound of $\kappa_4 + 2 (= \text{var}_g(z^2))$ is available, this gives a lower bound of $k_{\min}^{(\text{OA})}$.

Thirdly, as in cross validation, if another independent data set of size n^* is available for estimation of k , the estimated k , say \widehat{k}^* , can be used. Note that \widehat{k}^* does not change the optimality of using k since the expectation of $n^{-1}\widehat{k}^*\widehat{a}$ is the same as that of $n^{-1}k\widehat{a}$ up to order $O(n^{-1})$ under usual conditions with the term $\text{nacov}_g(-2\widehat{l}_W, \widehat{k}^*) = 0$ in (7.3). Note that the variance of $n^{-1}\widehat{k}^*\widehat{a}$ is of order $o(n^{-2})$, which can be asymptotically neglected when we consider the MSE up to order $O(n^{-2})$ for the associated information criterion of order $O_p(n^{-1})$.

Acknowledgments. This work was partially supported by a Grant-in-Aid for Scientific Research from the Japanese Ministry of Education, Culture, Sports, Science and Technology (JSPS KAKENHI, Grant No.26330031).

References

- AKAIKE, H. (1973). *Information theory and an extension of the maximum likelihood principle*. Académiai Kiado, Budapest. Proceedings of the 2nd international symposium on information theory, B. N. PETROV and F. CSÁKI (eds.), p. 267–281.
- CAVANAUGH, J.E. and SHUMWAY, R.H. (1997). A bootstrap variant of AIC for state-space model selection. *Statistica Sinica* **7**, 473–496.
- DAVIES, S.L., NEATH, A.A. and CAVANAUGH, J.E. (2006). Estimation optimality of corrected AIC and modified Cp in linear regression. *Int. Stat. Rev.* **74**, 161–168.
- FUJIKOSHI, Y. and SATOH, K. (1997). Modified AIC and C_p in multivariate linear regression. *Biometrika* **84**, 707–716.
- GILMOUR, S.G. (1996). The interpretation of Mallows's C_p -statistic. *J. R. Stat. Soc. D* **45**, 45–56.
- GRUBER, M.H.J. (1998). *Improving efficiency by shrinkage: The James-Stein and ridge regression estimators*. Marcel Dekker, New York.
- HOERL, A.E. and KENNARD, R.W. (1970). Ridge regression: Biased estimation for nonorthogonal problems. *Technometrics* **12**, 55–67.
- HURVICH, C.M. and TSAI, C.-L. (1989). Regression and time series model selection in small samples. *Biometrika* **76**, 297–307.
- ISHIGURO, M., SAKAMOTO, Y. and KITAGAWA, G. (1997). Bootstrapping log-likelihood and EIC, An extension of AIC. *Ann. Inst. Stat. Math.* **49**, 411–434.

- KISHINO, H. and HASEGAWA, M. (1989). Evaluation of the maximum likelihood estimate of the evolutionary tree topologies from DNA sequence data, and the branching order in hominoidea. *J. Mol. Evol.* **29**, 170–179.
- KONISHI, S. and KITAGAWA, G. (1996). Generalized information criteria in model selection. *Biometrika* **83**, 875–890.
- KONISHI, S. and KITAGAWA, G. (2003). Asymptotic theory for information criteria in model selection – functional approach. *Journal of Statistical Planning and Inference* **114**, 45–61.
- KONISHI, S. and KITAGAWA, G. (2008). *Information criteria and statistical modeling*. Springer, New-York.
- KULLBACK, S. and LEIBLER, R.A. (1951). On information and sufficiency. *Ann. Stat.* **22**, 79–86.
- MALLOWS, C.L. (1973). Some comments on C_P . *Technometrics* **15**, 661–675.
- OGASAWARA, H. (2010). Asymptotic expansions for the pivots using log-likelihood derivatives with an application in item response theory. *J. Multivar. Anal.* **101**, 2149–2167.
- OGASAWARA, H. (2013). Asymptotic cumulants of the estimator of the canonical parameter in the exponential family. *Journal of Statistical Planning and Inference* **143**, 2142–2150.
- OGASAWARA, H. (2014). Optimization of the Gaussian and Jeffreys power priors with emphasis on the canonical parameters in the exponential family. *Behaviormetrika* **41**, 195–223.
- OGASAWARA, H. (2015a) Asymptotic cumulants of some information criteria (2nd version). Discussion Paper, Center for Business Creation, Otaru University of Commerce, No. 174. <http://barrel.ih.otaru-uc.ac.jp/>.
- OGASAWARA, H. (2015b). Bias adjustment minimizing the asymptotic mean square error. *Communications in Statistics – Theory and Methods* **44**, 3503–3522.
- SCHWARZ, G. (1978). Estimating the dimension of a model. *Ann. Stat.* **6**, 461–464.
- SHIBATA, R. (1989). *Statistical aspects of model selection*. Springer, Berlin. From data to model, J. C. WILLEMS (ed.), p. 215–240.
- SHIBATA, R. (1997). Bootstrap estimate of Kullback-Leibler information for model selection. *Stat. Sin.* **7**, 375–394.
- SHIMODAIRA, H. and HASEGAWA, M. (1999). Multiple comparisons of log-likelihoods with applications to phylogenetic inference. *Mol. Biol. Evol.* **16**, 1114–1116.
- SUGIURA, N. (1978). Further analysis of the data by Akaike’s information criterion and the finite corrections. *Communications in Statistics – Theory and Methods* **7**, 13–26.
- TAKEUCHI, K. (1976). Distributions of information statistics and criteria of the goodness of models. *Math. Sci.* **153**, 12–18. (in Japanese).

HARUHIKO OGASAWARA
 DEPARTMENT OF INFORMATION AND
 MANAGEMENT SCIENCE
 OTARU UNIVERSITY OF COMMERCE
 3-5-21, OTARU 047-8501, JAPAN
 E-mail: hogasa@res.otaru-uc.ac.jp