# Bias correction of cross-validation criterion based on Kullback–Leibler information under a general condition

Hirokazu Yanagihara[a],[*], Tetsuji Tonda[b], Chieko Matsumoto[c]

[a]*Department of Social Systems and Management, Graduate School of Systems and Information Engineering, University of Tsukuba, 1-1-1 Tennodai, Tsukuba, Ibaraki 305-8573, Japan*

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

[c]*Department of Mathematics, Graduate School of Science, Hiroshima University, 1-3-1 Kagamiyama, Higashi-Hiroshima, Hiroshima 739-8526, Japan*

## Abstract

This paper deals with the bias correction of the cross-validation (CV) criterion to estimate the predictive Kullback–Leibler information. A bias-corrected CV criterion is proposed by replacing the ordinary maximum likelihood estimator with the maximizer of the adjusted log-likelihood function. The adjustment is just slight and simple, but the improvement of the bias is remarkable. The bias of the ordinary CV criterion is $O(n^{-1})$, but that of the bias-corrected CV criterion is $O(n^{-2})$. We verify that our criterion has smaller bias than the AIC, TIC, EIC and the ordinary CV criterion by numerical experiments.
© 2006 Elsevier Inc. All rights reserved.

[*] Corresponding author. Fax: +81 29 853 6451.
 *E-mail address:* yanagi@sk.tsukuba.ac.jp (H. Yanagihara).

## 1. Introduction

Let $\varphi(\mathbf{y})$ be the probability density function and $\mathbf{y}_i$ ($i = 1, \ldots, n$) be a $p \times 1$ observation vector drawn from $\varphi$, where $n$ is the sample size. The true model is expressed as

$$M^* : \quad \mathbf{y}_i \sim i.i.d. \; \varphi(\mathbf{y}_i), \quad (i = 1, \ldots, n). \tag{1}$$

We consider a family of parametric models $\mathcal{F} = \{f(\mathbf{y}|\boldsymbol{\theta}); \; \boldsymbol{\theta} \in \boldsymbol{\Theta} \subset \mathbb{R}^q\}$, where $\boldsymbol{\theta} = (\theta_1, \ldots, \theta_q)'$ is the $q$-dimensional vector of unknown parameters. A candidate model is expressed as

$$M : \quad \mathbf{y}_i \sim i.i.d. \; f(\mathbf{y}_i|\boldsymbol{\theta}), \quad (i = 1, \ldots, n). \tag{2}$$

Akaike's information criterion (AIC) proposed by Akaike [1,2] has been used universally for choosing the best model from the candidate models. It is well known that the AIC is an estimator of the risk based on the predictive Kullback–Leibler (K–L) information (see [8]), which measures the discrepancy between the true model $M^*$ and the candidate model $M$. However, AIC has a constant bias for the risk when the candidate model is misspecified (see e.g., [11,12]). Takeuchi [12] revaluated the bias correction term of AIC under the situation that $\mathcal{F}$ does not contain $\varphi(\mathbf{y})$, and proposed the Takeuchi's information criterion (TIC) by replacing the AIC's bias correction term with the revaluated term. The TIC is an asymptotically unbiased estimator for the risk if $\mathbf{y}_i$'s are *i.i.d.* However, Fujikoshi et al. [5] pointed out that TIC in normal regression models hardly corrects the bias in actual use, because its bias correction term mainly consists of an estimator of the fourth cumulant of the true distribution. Such an estimator tends to underestimate too much, even if the sample size $n$ is moderate (see [15]).

Like TIC, the cross-validation (CV) criterion proposed by Stone [10] is known as an asymptotically unbiased estimator for the risk (see [11]), although there are no estimators of higher-order cumulants in the CV criterion. Therefore, unlike TIC, the CV criterion can correct the bias efficiently. Using the better property of the CV criterion, Yanagihara [14,15] proposed new criteria which are partially constructed by the cross-validation method, and which are slightly influenced by the difference between $\varphi(\mathbf{y})$ and $f(\mathbf{y}|\boldsymbol{\theta})$. However, a bias for the risk exists also in the CV criterion. Fujikoshi et al. [4] corrected the biases of the CV criteria in normal multivariate regression and GMANOVA models. The purpose of our paper is to reduce the bias in the CV criterion under a general condition without adding any correction terms. We replaced the maximum likelihood estimator (MLE) of $\boldsymbol{\theta}$ with the maximizer of the adjusted log-likelihood function, and thus propose a bias-corrected CV (corrected CV, CCV) criterion. The adjustment is merely to change a weight of the weighted log-likelihood function, but the improvement of the bias is remarkable. The bias of the ordinary CV criterion is $O(n^{-1})$, but that of the CCV criterion is $O(n^{-2})$.

This paper is organized in the following way. In Section 2, we describe the risk based on the K–L information and usual information criteria. In Section 3, we state the derivation of CCV criterion and its asymptotic property. In Section 4, we verify that CCV criterion has smaller bias than other criteria, namely, the AIC, TIC, (the extended information criterion (EIC) [6]) and CV criterion by numerical experiments.

## 2. Risk and usual information criteria

Let $L(\boldsymbol{\theta}|\mathbf{Y}, \boldsymbol{d})$ be a weighted log-likelihood function on $f(\mathbf{y}_i|\boldsymbol{\theta})$ given by

$$L(\boldsymbol{\theta}|\mathbf{Y}, \boldsymbol{d}) = \sum_{i=1}^{n} d_i \log f(\mathbf{y}_i|\boldsymbol{\theta}), \tag{3}$$

where $Y = (y_1, \ldots, y_n)'$ and $d = (d_1, \ldots, d_n)'$. Let $\mathbf{1}_n$ be the $n \times 1$ vector, all of which elements are 1. For simplicity, let us express $L(\theta|Y) = L(\theta|Y, \mathbf{1}_n)$. The MLE of $\theta$ is obtained by maximizing the ordinary log-likelihood function $L(\theta|Y)$, i.e.,

$$\hat{\theta} = \arg\max_{\theta} L(\theta|Y). \tag{4}$$

Let $u_i$ be a $p \times 1$ future observation vector and $U = (u_1, \ldots, u_n)'$. We assume that $U$ is independent of $Y$ and $u_i$'s are independently and identically distributed according to $\varphi$. Note that the distribution of $u_i$ is the same as that of $y_i$. Then, the risk based on the predictive K–L information, which measures the discrepancy between the true model $M^*$ and the candidate model $M$ is defined by

$$R_{\mathrm{KL}} = E_{\mathbf{y}}^* E_{\mathbf{u}}^* \left[ -2L(\hat{\theta}|U) \right], \tag{5}$$

where $E^*$ means an expectation under the true model $M^*$.

The AIC proposed by Akaike [1,2] is a simple estimator of the risk $R_{\mathrm{KL}}$, and is given by

$$\mathrm{AIC} = -2L(\hat{\theta}|Y) + 2q. \tag{6}$$

However, if the candidate model is misspecified, AIC has a constant bias, i.e.,

$$B_{\mathrm{AIC}} = R_{\mathrm{KL}} - E_{\mathbf{y}}^*[\mathrm{AIC}] = O(1) \tag{7}$$

(see e.g., [11,12]). This is mainly because Akaike [1,2] derived AIC under the condition that the candidate model is correctly specified. Takeuchi [12] reevaluated the AIC's bias correction term, $2q$, under the situation that $\mathcal{F}$ does not contain $\varphi(y)$, and proposed TIC as follows: Let

$$g(y_i|\hat{\theta}) = \frac{\partial}{\partial\theta} \log f(y_i|\theta) \Big|_{\theta=\hat{\theta}}, \quad H(y_i|\hat{\theta}) = \frac{\partial^2}{\partial\theta\,\partial\theta'} \log f(y_i|\theta) \Big|_{\theta=\hat{\theta}}. \tag{8}$$

Under the candidate model $M$, TIC is given by

$$\mathrm{TIC} = -2L(\hat{\theta}|Y) + 2\operatorname{tr}(\hat{J}(\hat{\theta})^{-1}\hat{I}(\hat{\theta})), \tag{9}$$

where

$$\hat{J}(\hat{\theta}) = -\frac{1}{n} \sum_{i=1}^{n} H(y_i|\hat{\theta}), \quad \hat{I}(\hat{\theta}) = \frac{1}{n} \sum_{i=1}^{n} g(y_i|\hat{\theta})g(y_i|\hat{\theta})'. \tag{10}$$

Takeuchi [12] showed that TIC is an asymptotically unbiased estimator for the risk $R_{\mathrm{KL}}$ if $y_i$'s are *i.i.d.*, i.e.,

$$B_{\mathrm{TIC}} = R_{\mathrm{KL}} - E_{\mathbf{y}}^*[\mathrm{TIC}] = O(n^{-1}). \tag{11}$$

On the other hand, Stone [10] proposed the CV criterion in the following way. Let $\hat{\theta}_{[-i]}$ be the MLE of $\theta$ evaluated from $i$th jackknife sample. Note that $\hat{\theta}_{[-i]}$ is the maximizer of $\sum_{j\neq i}^{n} \log f(y_j|\theta)$, i.e.,

$$\hat{\theta}_{[-i]} = \arg\max_{\theta} L(\theta|Y, \mathbf{1}_n - e_i), \tag{12}$$

where $e_i$ is the $n \times 1$ vector whose $i$th element is 1 and the other elements are 0. The CV criterion is given by

$$\mathrm{CV} = -2 \sum_{i=1}^{n} \log f(y_i | \hat{\theta}_{[-i]}). \tag{13}$$

Stone [11] pointed out that the CV criterion is an asymptotically unbiased estimator for the risk. From the result in Stone [11], we can see that the TIC and CV criteria are asymptotically equivalent, i.e., $\mathrm{CV} = \mathrm{TIC} + O_p(n^{-1})$. Therefore, from Eq. (11), the bias of the CV criterion is given by

$$B_{\mathrm{CV}} = R_{\mathrm{KL}} - E_y^*[\mathrm{CV}] = O(n^{-1}). \tag{14}$$

We can see that the order of bias of CV is the same as that of TIC. However, Yanagihara [15] showed that the CV criterion in normal regression models has smaller bias than TIC by investigating the asymptotic expansions of biases for the risk. In order to get TIC, we must estimate higher-order cumulants. However, the ordinary estimators of higher-order cumulants tend to underestimate too much, even if the sample size $n$ is moderate. Consequently, TIC tends to have large bias. Needless to say, we can obtain the CV criterion without estimating higher-order cumulants.

## 3. Bias correction of the CV criterion

### 3.1. Asymptotic expansion of the bias of the CV criterion

In this section, we propose the bias-corrected CV (CCV) criterion by replacing $\hat{\theta}_{[-i]}$ in the CV criterion with the maximizer of another weighted log-likelihood function. First, in order to correct the bias of the CV criterion, we derive an asymptotic expansion of its bias up to the order $n^{-1}$. Let $\theta_0$ be the $q \times 1$ vector such that

$$E_y^* \left[ g(y|\theta_0) \right] = \mathbf{0}_q, \tag{15}$$

where $\mathbf{0}_q$ is the $q \times 1$ vector, all of which elements are 0. Note that $\hat{\theta}$ converges to $\theta_0$ almost surely as $n$ goes to infinity (see [13]). Then, we obtain an asymptotic expansion of the bias of the CV criterion up to the order $n^{-1}$ in the following theorem.

**Theorem 1.** *Under the regularity conditions, the bias of the CV criterion is expanded as*

$$B_{\mathrm{CV}} = R_{\mathrm{KL}} - E_y^*[\mathrm{CV}] = -\frac{1}{n} \operatorname{tr}(J(\theta_0)^{-1} I(\theta_0)) + O(n^{-2}), \tag{16}$$

*where*

$$J(\theta_0) = -E_y^* \left[ H(y|\theta_0) \right], \quad I(\theta_0) = E_y^* \left[ g(y|\theta_0) g(y|\theta_0)' \right].$$

**Proof.** Let

$$\hat{K}(\hat{\theta}) = -\frac{1}{n} \sum_{i=1}^{n} \left( \frac{\partial}{\partial \theta'} \otimes \frac{\partial^2}{\partial \theta \, \partial \theta'} \right) \log f(y_i | \theta) \Big|_{\theta = \hat{\theta}}.$$

Using the Taylor expansion, we obtain the perturbation expansion of $\hat{\theta}_{[-i]}$ as

$$\hat{\theta}_{[-i]} = \hat{\theta} - \frac{1}{n} z_{1,i} - \frac{1}{n^2} z_{2,i} + O_p(n^{-3}), \tag{17}$$

where

$$z_{1,i} = \hat{J}(\hat{\theta})^{-1} g(y_i|\hat{\theta}), \quad z_{2,i} = \hat{J}(\hat{\theta})^{-1} \left\{ \tfrac{1}{2}\hat{K}(\hat{\theta})\text{vec}(z_{1,i}z'_{1,i}) - H(y_i|\hat{\theta})z_{1,i} \right\}.$$

Here, $g(y_i|\hat{\theta})$ and $H(y_i|\hat{\theta})$ are given by (8), and $\hat{J}(\hat{\theta})$ is given by (10). Since the distributions of $y_i$ and $u_i$ are the same, the commutative equation $E_{\boldsymbol{y}}^*[\log f(y_i|\hat{\theta}_{[-i]})] = E_{\boldsymbol{y}}^* E_{\boldsymbol{u}}^*[\log f(u_i|\hat{\theta}_{[-i]})]$ holds. Therefore, using the Taylor expansion and Eq. (17), $E_{\boldsymbol{y}}^*[\text{CV}]$ is expanded as

$$E_{\boldsymbol{y}}^*[\text{CV}] = R_{\text{KL}} + R_1 + \frac{1}{n}R_2 + O(n^{-2}), \tag{18}$$

where

$$R_1 = \frac{2}{n}\sum_{i=1}^{n} E_{\boldsymbol{y}}^* E_{\boldsymbol{u}}^* \left[ g(u_i|\hat{\theta})'z_{1,i} \right],$$

$$R_2 = \frac{1}{n}\sum_{i=1}^{n} E_{\boldsymbol{y}}^* E_{\boldsymbol{u}}^* \left[ 2g(u_i|\hat{\theta})'z_{2,i} - z'_{1,i}H(u_i|\hat{\theta})z_{1,i} \right].$$

Note that $\sum_{i=1}^{n} z_{1,i} = \boldsymbol{0}_q$ because $\hat{\theta}$ is the MLE, i.e., $\sum_{i=1}^{n} g(y_i|\hat{\theta}) = \boldsymbol{0}_q$. Therefore, taking a conditional expectation of $g(u_i|\hat{\theta})$ for $Y$ as $\boldsymbol{\eta} = E_{\boldsymbol{u}}^*[g(u_i|\hat{\theta})|Y]$, we obtain

$$R_1 = \frac{2}{n}\sum_{i=1}^{n} E_{\boldsymbol{y}}^* \left[ E_{\boldsymbol{u}}^* \left[ g(u_i|\hat{\theta})'z_{1,i} \,\Big|\, Y \right] \right] = \frac{2}{n}\sum_{i=1}^{n} E_{\boldsymbol{y}}^* \left[ \boldsymbol{\eta}'z_{1,i} \right] = 0. \tag{19}$$

On the other hand, by using the equation $\hat{\theta} \overset{a.s.}{\to} \theta_0$ ($n \to \infty$), $R_2$ is expanded as

$$R_2 = \frac{1}{n}\sum_{i=1}^{n} E_{\boldsymbol{y}}^* E_{\boldsymbol{u}}^* \left[ 2g(u_i|\theta_0)'z_{2,i} - z'_{1,i}H(u_i|\theta_0)z_{1,i} \right] + O(n^{-1}).$$

From Eq. (15), the first term on the right side of the above equation disappears. Moreover, using equation $\hat{J}(\hat{\theta}) \overset{a.s.}{\to} J(\theta_0)$ and $\hat{I}(\hat{\theta}) \overset{a.s.}{\to} I(\theta_0)$ ($n \to \infty$), we derive the following equation:

$$\frac{1}{n}\sum_{i=1}^{n} E_{\boldsymbol{y}}^* E_{\boldsymbol{u}}^* \left[ z'_{1,i}H(u_i|\theta_0)z_{1,i} \right] = -\text{tr}(J(\theta_0)^{-1}I(\theta_0)) + O(n^{-1}).$$

Therefore, $R_2$ is expanded as

$$R_2 = \text{tr}(J(\theta_0)^{-1}I(\theta_0)) + O(n^{-1}). \tag{20}$$

Substituting Eqs. (19) and (20) into (18) yields

$$E_{\boldsymbol{y}}^*[\text{CV}] = R_{\text{KL}} + \frac{1}{n}\text{tr}(J(\theta_0)^{-1}I(\theta_0)) + O(n^{-2}). \tag{21}$$

Consequently, the result (16) in Theorem 1 is obtained. $\quad\square$

Because $-\log f(y|\theta)$ is strictly convex with respect to $\theta$, the matrix $H(y|\theta)$ is positive definite (see e.g., [9, p. 49]). It makes the inequality $\text{tr}(J(\theta_0)^{-1}I(\theta_0)) > 0$. Consequently, the CV criterion tends to overestimate for the risk $R_{\text{KL}}$.

### 3.2. Corrected CV criterion

Next, we propose a new criterion, CCV criterion, which always corrects the bias for the risk $R_{KL}$ to $O(n^{-2})$. Theoretically, we can correct the bias in the CV criterion by subtracting the term $n^{-1} \operatorname{tr}(\hat{J}(\hat{\theta})^{-1}\hat{I}(\hat{\theta}))$ from the CV criterion. However, we can easily forecast that the bias is not fully corrected by such a plug-in estimator because $\operatorname{tr}(\hat{J}(\hat{\theta})^{-1}\hat{I}(\hat{\theta}))$ must have a large bias for $\operatorname{tr}(J(\theta_0)^{-1}I(\theta_0))$, even if the sample size $n$ is moderate. The reason for this is the same as the reason that TIC does not reduce the bias enough in actual use. Therefore, we need to pre-pare other methods to correct the bias without estimating $\operatorname{tr}(J(\theta_0)^{-1}I(\theta_0))$. From Eq. (12), we notice that $\hat{\theta}_{[-i]}$ removes the influence of $y_i$ perfectly. However, we consider that the effect of $y_i$ should not be removed completely because $R_{KL}$ in (5) is not the predictive K–L infor-mation measuring the discrepancy between $\varphi(u)$ and $f(u|\hat{\theta}_{[-i]})$, but, rather, $\varphi(u)$ and $f(u|\hat{\theta})$. Thus, we use the estimator obtained by maximizing another weighted log-likelihood function, in which the influence of $y_i$ remains for a while. Consequently, we propose the following CCV criterion.

**Definition.** Let $\tilde{\theta}_i$ be the estimator of $\theta$ by maximizing the weighted log-likelihood function as

$$\tilde{\theta}_i = \arg\max_{\theta} L(\theta|Y, \mathbf{1}_n - c_n e_i), \tag{22}$$

where $c_n$ is any constant which can be expanded as $c_n = 1 - 1/(2n) + O(n^{-2})$. Then, we propose the bias-corrected CV (CCV) criterion as

$$\mathrm{CCV} = -2 \sum_{i=1}^{n} \log f(y_i|\tilde{\theta}_i). \tag{23}$$

We can see that any estimators of higher-order cumulants are not necessary for obtaining CCV criterion. However, CCV criterion always corrects the bias to $O(n^{-2})$, even though there is no term based on $\operatorname{tr}(\hat{J}(\hat{\theta})^{-1}\hat{I}(\hat{\theta}))$ in the formula (23). The order of bias of the CCV criterion is obtained in the following theorem.

**Theorem 2.** *Under the regularity conditions, the order of the bias of the CCV criterion is given by*

$$B_{CCV} = R_{KL} - E_y^*[\mathrm{CCV}] = O(n^{-2}). \tag{24}$$

**Proof.** From the definition of $\tilde{\theta}_i$ and the Taylor expansion, we can expand $\tilde{\theta}_i$ as

$$\tilde{\theta}_i = \hat{\theta}_{[-i]} - \frac{1}{2n^2} \left\{ \frac{1}{n-1} \sum_{j \neq i}^{n} H(y_j|\hat{\theta}_{[-i]}) \right\}^{-1} g(y_i|\hat{\theta}_{[-i]}) + O_p(n^{-3}).$$

Thus, the perturbation expansion of CCV is given by

$$\mathrm{CCV} = \mathrm{CV} + \frac{1}{n^2} \sum_{i=1}^{n} g(y_i|\hat{\theta}_{[-i]})' \left\{ \frac{1}{n-1} \sum_{j \neq i}^{n} H(y_j|\hat{\theta}_{[-i]}) \right\}^{-1} g(y_i|\hat{\theta}_{[-i]}) + O_p(n^{-2}).$$

Therefore, we calculate $E_{\boldsymbol{y}}^*[\text{CCV}]$ as

$$E_{\boldsymbol{y}}^*[\text{CCV}] = E_{\boldsymbol{y}}^*[\text{CV}] - \frac{1}{n}\,\text{tr}(\boldsymbol{J}(\boldsymbol{\theta}_0)^{-1}\boldsymbol{I}(\boldsymbol{\theta}_0)) + O(n^{-2}).$$

Substituting Eq. (21) into the above equation yields Eq. (24) in Theorem 2.  □

## 4. Numerical examination

First, we prepare another bias-corrected criterion constructed by the bootstrap method, which was named EIC by Ishiguro et al. [6]. Let $\boldsymbol{Y}_b^\star$ $(b = 1, \ldots, B)$ be the $b$th bootstrap sample by resampling, and $\hat{\boldsymbol{\theta}}_b^\star$ be the MLE of $\boldsymbol{\theta}$ evaluated from $\boldsymbol{Y}_b^\star$, i.e.,

$$\hat{\boldsymbol{\theta}}_b^\star = \arg\max_{\boldsymbol{\theta}} L(\boldsymbol{\theta}|\boldsymbol{Y}_b^\star). \tag{25}$$

Then, the EIC is given by

$$\text{EIC} = -2L(\hat{\boldsymbol{\theta}}|\boldsymbol{Y}) - \frac{2}{B}\sum_{b=1}^{B}\left\{L(\hat{\boldsymbol{\theta}}_b^\star|\boldsymbol{Y}) - L(\hat{\boldsymbol{\theta}}_b^\star|\boldsymbol{Y}_b^\star)\right\}, \tag{26}$$

(see e.g., [7]). Through the simulation, we compare the biases and frequencies of the selected model in our proposed CCV criterion, and also the AIC, TIC, EIC and CV criterion. In this paper, we deal with the selection of the best model from the candidate models having the elliptical distribution, i.e.,

$$f(\boldsymbol{y}_i|\boldsymbol{\theta}) = a_p|\boldsymbol{\Lambda}|^{-1/2}g((\boldsymbol{y}_i - \boldsymbol{\mu})'\boldsymbol{\Lambda}^{-1}(\boldsymbol{y}_i - \boldsymbol{\mu})), \quad (i = 1, \ldots, 20), \tag{27}$$

where $g(r)$ is a known non-negative function and $a_p$ is a positive constant depending on the dimension $p$ (see e.g., [3]). We choose the best $g(r)$ and $a_p$ from the candidate models by minimizing the information criteria. The candidate models considered are as follows:

*Model* 1: *Multivariate normal distribution*,

$$a_p = (2\pi)^{-p/2}, \quad g(r) = e^{-r/2}.$$

*Model* 2: *Multivariate logistic distribution*,

$$a_p = (2\pi)^{-p/2}\left\{\sum_{j=1}^{\infty}(-1)^{j-1}j^{1-p/2}\right\}^{-1}, \quad g(r) = \frac{e^{-r/2}}{\{1 + e^{-r/2}\}^2}.$$

*Model* 3: *Multivariate Cauchy distribution*,

$$a_p = \frac{\Gamma((p+1)/2)}{\pi^{(p+1)/2}}, \quad g(r) = (1 + r)^{-(p+1)/2},$$

where $\Gamma(x)$ is the gamma function.

Choosing the best model is equivalent to determining the best weight function in the $M$-estimation. Therefore, we will judge whether or not the robust estimation should be performed through minimizing the information criterion, since the normal distribution is included in the candidate

Table 1
Biases and frequencies of the selected model according to the criteria

| Distribution | Criterion | | $p = 2$ | | | $p = 6$ | | |
|---|---|---|---|---|---|---|---|---|
| | | | Model 1 | Model 2 | Model 3 | Model 1 | Model 2 | Model 3 |
| Normal | | Risk | 120.50[a] | 123.45 | 132.31 | 399.91[a] | 405.21 | 404.71 |
| | AIC | Bias | 2.30 | 5.54 | 1.77 | 37.43 | 43.81 | 15.90 |
| | | Frequency | (40.7) | (58.5) | (0.8) | (14.4) | (85.5) | (0.1) |
| | TIC | Bias | 3.07 | 5.34 | −0.60 | 42.04 | 49.58 | −11.84 |
| | | Frequency | (55.8) | (43.2) | (1.0) | (13.3) | (86.7) | (0.0) |
| | EIC | Bias | **−0.06** | 0.43 | −0.74 | 1.92 | 2.43 | **−0.24** |
| | | Frequency | (65.7) | (28.7) | (5.6) | (55.0) | (21.0) | (24.0) |
| | CV | Bias | −0.73 | −0.93 | −0.40 | −4.35 | −4.71 | −2.37 |
| | | Frequency | (71.4) | (21.6) | (7.1) | (57.4) | (4.0) | (38.6) |
| | CCV | Bias | −0.27 | **−0.30** | **−0.11** | **0.49** | **0.66** | −0.79 |
| | | Frequency | (71.0) | (22.5) | (6.5) | (62.8) | (5.5) | (31.7) |
| Laplace | | Risk | 125.27 | 135.24 | 121.32[a] | 422.27 | 434.96 | 387.84[a] |
| | AIC | Bias | 9.93 | 17.42 | 1.47 | 67.59 | 79.65 | 17.69 |
| | | Frequency | (62.8) | (14.0) | (23.3) | (57.2) | (36.4) | (6.4) |
| | TIC | Bias | 8.52 | 13.66 | −0.59 | 68.07 | 79.79 | −10.49 |
| | | Frequency | (65.4) | (9.0) | (25.6) | (60.7) | (38.9) | (0.4) |
| | EIC | Bias | 2.49 | 4.35 | −0.80 | 11.68 | 14.22 | **0.21** |
| | | Frequency | (48.1) | (7.0) | (44.9) | (14.1) | (3.6) | (82.3) |
| | CV | Bias | **−0.31** | **−0.57** | −0.39 | −6.27 | −7.12 | −1.01 |
| | | Frequency | (46.9) | (2.9) | (50.2) | (17.0) | (0.2) | (82.8) |
| | CCV | Bias | 0.71 | 0.92 | **−0.11** | 2.74 | 3.09 | 0.58 |
| | | Frequency | (47.9) | (3.1) | (49.0) | (19.9) | (0.3) | (79.8) |
| Chi-square | | Risk | 126.92 | 135.41 | 122.52[a] | 426.64 | 437.16 | 387.63[a] |
| | AIC | Bias | 12.06 | 18.76 | 2.82 | 71.44 | 82.00 | 17.15 |
| | | Frequency | (39.3) | (36.5) | (24.2) | (35.2) | (56.5) | (8.3) |
| | TIC | Bias | 10.61 | 15.74 | -0.31 | 71.96 | 82.90 | −13.39 |
| | | Frequency | (47.7) | (27.8) | (24.5) | (38.7) | (61.1) | (0.2) |
| | EIC | Bias | 2.94 | 4.76 | −0.51 | 9.74 | 12.46 | −2.21 |
| | | Frequency | (39.5) | (18.0) | (42.5) | (12.7) | (6.3) | (81.0) |
| | CV | Bias | −0.76 | −1.36 | **−0.11** | −11.56 | −13.06 | −3.72 |
| | | Frequency | (38.2) | (13.8) | (48.0) | (16.3) | (0.8) | (82.9) |
| | CCV | Bias | **0.65** | **0.67** | 0.20 | **−0.21** | **−0.34** | **−2.09** |
| | | Frequency | (38.8) | (14.4) | (46.8) | (20.1) | (0.9) | (79.0) |
| Log-normal | | Risk | 128.07 | 137.34 | 121.90[a] | 427.94 | 439.32 | 384.36[a] |
| | AIC | Bias | 13.85 | 21.38 | 2.71 | 75.71 | 87.04 | 17.12 |
| | | Frequency | (42.9) | (33.3) | (23.8) | (39.4) | (52.5) | (8.1) |
| | TIC | Bias | 12.31 | 18.17 | −0.03 | 76.01 | 87.57 | −12.89 |
| | | Frequency | (49.7) | (24.9) | (25.4) | (43.7) | (55.4) | (0.9) |
| | EIC | Bias | 4.68 | 7.00 | −0.19 | 14.07 | 16.86 | **−1.42** |
| | | Frequency | (41.5) | (16.2) | (42.3) | (13.5) | (4.6) | (81.9) |
| | CV | Bias | **−0.24** | **−0.83** | **0.17** | −12.92 | −14.57 | −3.28 |
| | | Frequency | (41.1) | (11.3) | (47.6) | (15.7) | (0.7) | (83.6) |
| | CCV | Bias | 1.51 | 1.70 | 0.47 | **−0.31** | **−0.41** | −1.65 |
| | | Frequency | (41.6) | (11.9) | (46.5) | (19.9) | (1.2) | (78.9) |

[a] Denotes the smallest risk in all the candidate models, and the smallest bias in all the criteria is in bold.

Table 2
Biases of CV, CCV, CCV′ and CCV″ criteria

| Distribution | Criterion | $p = 2$ | | | $p = 6$ | | | Average |
|---|---|---|---|---|---|---|---|---|
| | | Model 1 | Model 2 | Model 3 | Model 1 | Model 2 | Model 3 | |
| Normal | CV | −0.54 | −0.73 | −0.48 | −6.39 | −7.01 | −2.71 | −2.98 |
| | CCV | **−0.08** | **−0.09** | −0.19 | **−1.49** | **−1.57** | −1.14 | **−0.76** |
| | CCV′ | −0.31 | −0.47 | **−0.17** | −5.16 | −5.81 | **−0.67** | −2.10 |
| | CCV″ | −0.21 | −0.31 | −0.18 | −3.99 | −4.43 | −0.90 | −1.67 |
| Laplace | CV | −1.20 | −1.79 | −0.31 | −8.61 | −9.62 | −2.75 | −4.04 |
| | CCV | **−0.17** | **−0.29** | −0.02 | **0.16** | **0.31** | −1.16 | **−0.20** |
| | CCV′ | −0.91 | −1.45 | **−0.01** | −7.27 | −8.27 | **−0.69** | −3.10 |
| | CCV″ | −0.68 | −1.09 | **−0.01** | −5.43 | −6.13 | −0.93 | −2.38 |
| Chi-square | CV | −0.57 | −1.13 | **−0.03** | −10.78 | −12.07 | −2.86 | −4.57 |
| | CCV | 0.79 | 0.81 | 0.28 | **0.36** | **0.41** | −1.23 | **0.24** |
| | CCV′ | −0.29 | −0.80 | 0.30 | −9.44 | −10.74 | **−0.74** | −3.62 |
| | CCV″ | **−0.07** | **−0.38** | 0.30 | −7.35 | −8.34 | −0.99 | −2.80 |
| Log-normal | CV | −1.14 | −2.01 | −0.35 | −12.91 | −14.78 | −2.95 | −5.69 |
| | CCV | 0.65 | **0.58** | −0.05 | **−0.35** | **−0.67** | −1.33 | **−0.20** |
| | CCV′ | −0.85 | −1.68 | **−0.03** | −11.57 | −13.45 | **−0.86** | −4.74 |
| | CCV″ | **−0.53** | −1.20 | −0.04 | −9.36 | −10.90 | −1.10 | −3.85 |

The smallest bias in all the criteria is in bold.

models. Let $m$ and $S$ be the $p \times 1$ vector and $p \times p$ matrix obtained by maximizing the weighted log-likelihood function $L(\theta|Y, d)$ in (3), i.e.,

$$m = \frac{1}{\text{tr}(WD)} Y'DW1_n, \quad S = -\frac{2}{\text{tr}(D)} (Y - 1_n m')'DW(Y - 1_n m'), \quad (28)$$

where $D = \text{diag}(d_1, \ldots, d_n)$ and $W = \text{diag}(w(r_1), \ldots, w(r_n))$. Here, $w(r) = \{dg(r)/dr\}/g(r)$ and $r_i = (y_i - m)'S^{-1}(y_i - m)$. We can obtain $\hat{\theta}$, $\hat{\theta}_{[-i]}$, $\tilde{\theta}_i$ and $\hat{\theta}_b^\star$ from formula (28). On the other hand, we prepare the following four distributions for the true model.

- *Normal distribution*: Each of the $p$ variables is generated independently from $N(0, 1)$ ($\kappa_{3,3}^{(1)} = \kappa_{3,3}^{(2)} = 0$ and $\kappa_4^{(1)} = 0$),
- *Laplace distribution*: Each of the $p$ variables is generated independently from the Laplace distribution $L(0, 1)$ divided by the standard deviation $\sqrt{2}$ ($\kappa_{3,3}^{(1)} = \kappa_{3,3}^{(2)} = 0$ and $\kappa_4^{(1)} = 2p$),
- *Chi-square distribution*: Each of the $p$ variables is generated independently from the $\chi^2$ distribution with 3 degrees of freedom standardized by the mean 3 and standard deviation $\sqrt{6}$ ($\kappa_{3,3}^{(1)} = \kappa_{3,3}^{(2)} \approx 1.63 \times p$ and $\kappa_4^{(1)} = 4p$),
- *Log-normal distribution*: Each of the $p$ variables is generated independently from a log-normal distribution $LN(0, 1/4)$ standardized by the mean $e^{1/8}$ and standard deviation $e^{1/8}\sqrt{e^{1/4} - 1}$ ($\kappa_{3,3}^{(1)} = \kappa_{3,3}^{(2)} \approx 1.71 \times p$ and $\kappa_4^{(1)} \approx 8.90 \times p$).

Table 1 lists the average risk, the biases of the CCV criterion along with the AIC, TIC, EIC and CV criterion, and the frequencies of the model selected by the criteria in the cases of $p = 2$ and 6. These average values were obtained after 10,000 iterations, and the EIC was obtained by

resampling 100 times. Moreover, we use $c_n = \sqrt{n/(n+1)}$ as the constant for adjusting $\tilde{\boldsymbol{\theta}}_i$. From the table, we can see that the biases of AIC were large in all the cases. TIC hardly corrected the bias in Models 1 and 2. On the other hand, the biases of the CV criterion were smaller than the biases of AIC and TIC. Especially, the biases of the CV criterion were smaller than the biases of EIC in most cases. Moreover, when we use the CV criterion for model selection, the frequencies of the model with the smallest risk selected was the highest in all the criteria. However, the bias of the CV criterion became large when the dimension $p$ increased. We can see that CCV corrected the bias efficiently.

Next, we compared several methods for correcting the bias in the CV criterion. We prepared the following two different bias-corrected CV criteria from the CCV criterion, which were obtained by adding some bias correction terms

$$\mathrm{CCV}' = \mathrm{CV} - \frac{1}{n}\,\mathrm{tr}(\hat{\boldsymbol{J}}(\hat{\boldsymbol{\theta}})^{-1}\hat{\boldsymbol{I}}(\hat{\boldsymbol{\theta}})), \quad \mathrm{CCV}'' = \left(1 - \frac{1}{2n}\right)\mathrm{CV} - \frac{1}{n}\,L(\hat{\boldsymbol{\theta}}|\boldsymbol{Y}).$$

Note that the CCV' and CCV'' criteria correct the biases to $O(n^{-2})$ as well as the CCV criterion. Table 2 shows the biases of the CV, CCV, CCV' and CCV'' criteria. From the table, we can see that CCV' and CCV'' did not reduce the bias fully when the bias is large. Therefore, the methods for reducing the bias by adding correction terms should not be used for bias correction. We have studied several other models and have obtained similar results.

## Acknowledgments

## References

[1] 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.

[2] H. Akaike, A new look at the statistical model identification, IEEE Trans. Automat. Control AC-19 (1974) 716–723.

[3] K.T. Fang, S. Kotz, K.W. Ng, Symmetric Multivariate and Related Distributions, Chapman & Hall, CRC, London, 1990.

[4] Y. Fujikoshi, T. Noguchi, M. Ohtaki, H. Yanagihara, Corrected versions of cross-validation criteria for selecting multivariate regression and growth curve models, Ann. Inst. Statist. Math. 55 (2003) 537–553.

[5] Y. Fujikoshi, H. Yanagihara, H. Wakaki, Bias corrections of some criteria for selection multivariate linear regression models in a general case, Amer. J. Math. Management Sci. 25 (2005) 221–258.

[6] M. Ishiguro, Y. Sakamoto, G. Kitagawa, Bootstrapping log likelihood and EIC, an extension of AIC, Ann. Inst. Statist. Math. 49 (1997) 411–434.

[7] S. Konishi, Statistical model evaluation and information criteria, in: S. Ghosh (Ed.), Multivariate Analysis, Design of Experiments, and Survey Sampling, Marcel Dekker, New York, 1999, pp. 369–399.

[8] S. Kullback, R.A. Leibler, On information and sufficiency, Ann. Math. Statist. 22 (1951) 79–86.

[9] E.L. Lehmann, G. Casella, Theory of Point Estimation, second ed., Springer, New York, 1998.

[10] M. Stone, Cross-validatory choice and assessment of statistical predictions, J. Roy. Statist. Soc. Ser. B 36 (1974) 111–147.

[11] M. Stone, An asymptotic equivalence of choice of model by cross-validation and Akaike's criterion, J. Roy. Statist. Soc. Ser. B 39 (1977) 44–47.

[12] K. Takeuchi, Distribution of information statistics and criteria for adequacy of models, Math. Sci. 153 (1976) 12–18 (in Japanese).

[13] H. White, Maximum likelihood estimation of misspecified models, Econometrica 50 (1982) 1–25.

[14] H. Yanagihara, Selection of covariance structure models in nonnormal data by using information criterion: an application to data from the survey of the Japanese national character, Proc. Inst. Statist. Math. 53 (2005) 133–157 (in Japanese).

[15] H. Yanagihara, Corrected version of *AIC* for selecting multivariate normal linear regression models in a general nonnormal case, J. Multivariate Anal. 97 (2006) 1070–1089.

[16] H. Yanagihara, A family of estimators for multivariate kurtosis in a nonnormal linear regression model, J. Multivariate Anal., in press.