

# A Class of Model Selection Criteria Based on Cross-Validation Method

(Last Modified: January 13, 2006)

Hirokazu YANAGIHARA<sup>1</sup>, Ke-Hai YUAN<sup>2</sup>, Hironori FUJISAWA<sup>3</sup>

AND

Kentaro HAYASHI<sup>4</sup>

<sup>1</sup>*Department of Mathematics, Graduate School of Science, Hiroshima University  
1-3-1 Kagamiyama, Higashi-Hiroshima, Hiroshima 739-8626, Japan*

<sup>2</sup>*Department of Psychology, University of Notre Dame  
Notre Dame, Indiana 46556, USA*

<sup>3</sup>*Department of Mathematical Analysis and Statistical Inference  
The Institute of Statistical Mathematics  
4-6-7 Minami-Azabu, Minato-ku, Tokyo 106-8569, Japan*

<sup>4</sup>*Department of Psychology, University of Hawaii at Manoa  
Honolulu, Hawaii 96822, USA*

## Abstract

In this paper, we define a class of cross-validators model selection criteria as an estimator of the predictive risk based on a discrepancy between a candidate model and the true model. For a vector of unknown parameters,  $n$  estimators are required for the definition of the class, where  $n$  is the sample size. The  $i$ th estimator ( $i = 1, \dots, n$ ) is obtained by minimizing a weighted discrepancy function in which the  $i$ th observation has a weight of  $1 - \lambda$  and others have weight of 1. Cross-validators model selection criteria in the class are specified by the individual  $\lambda$ . The sample discrepancy function and the ordinary cross-validation (CV) criterion are special cases of the class. One may choose  $\lambda$  to minimize the biases. The

---

<sup>1</sup>Corresponding author. E-mail: yanagi@math.sci.hiroshima-u.ac.jp

optimal  $\lambda$  makes the bias-corrected CV (CCV) criterion a second-order unbiased estimator for the risk, while the ordinary CV criterion is a first-order unbiased estimator of the risk.

*AMS 2000 subject classifications.* Primary 62H25; Secondary 62F07.

*Key words:* Asymptotic expansion, Bias correction, Cross-validation criterion, Model misspecification, Model selection, Predictive discrepancy, Sample discrepancy function, Structural equation model.

## 1. Introduction

Let  $\mathbf{y}_1, \dots, \mathbf{y}_n$  be a random sample from a  $p$ -dimensional population  $\mathbf{y}$  whose probability density function  $\varphi(\mathbf{y})$  is unknown. Nevertheless, the true model can be expressed as

$$M^* : \mathbf{y}_1, \dots, \mathbf{y}_n \sim i.i.d. \varphi(\mathbf{y}). \quad (1)$$

Consider a family of parametric models  $\mathcal{F} = \{f(\mathbf{y}|\boldsymbol{\theta}); \boldsymbol{\theta} \in \Theta \subset \mathbb{R}^q\}$ , where  $\boldsymbol{\theta} = (\theta_1, \dots, \theta_q)'$  is a  $q$ -dimensional vector of unknown parameters. This implies that a candidate model is given by

$$M : \mathbf{y}_1, \dots, \mathbf{y}_n \sim i.i.d. f(\mathbf{y}|\boldsymbol{\theta}). \quad (2)$$

Most model selection criteria (or information criteria) for determining the best model among all candidate models are estimators of the predictive risk based on the discrepancy between the candidate model and the true model. For example, Akaike's information criterion (AIC; Akaike, 1973, 1974), Takeuchi's (1976) bias-corrected information criterion (TIC), and the extended information criterion (EIC; Ishiguro, Sakamoto, & Kitagawa, 1997) are estimators of the predictive Kullback-Leibler (K-L) discrepancy (Kullback & Leibler, 1951). On the other hand, the cross-validation (CV) criterion (Stone, 1974, 1977) serves as an estimator of the predictive risk based on an arbitrary discrepancy, e.g., K-L discrepancy,  $L_2$  distance, and density power divergence (Basu *et al.*, 1998).

In this paper, we propose a class of model selection criteria by the cross-validatory method. For a given discrepancy,  $n$  estimators of  $\boldsymbol{\theta}$  are required to define the class. The  $i$ th estimator ( $i = 1, \dots, n$ ) is obtained by minimizing a weighted discrepancy function in which the  $i$ th observation has a weight of  $1 - \lambda$  ( $0 \leq \lambda \leq 1$ ) and others have weights of 1. Each  $\lambda$  represents a cross-validatory model selection criteria. The sample discrepancy function and the ordinary CV criterion correspond to  $\lambda = 0$  and 1, respectively. From the viewpoint of second order asymptotics for biases, the optimal  $\lambda$  can be expanded as  $\lambda = 1 - 1/(2n) + O(n^{-2})$ .

The optimal  $\lambda$  yield a bias-corrected CV (CCV) criterion that corrects the bias to  $O(n^{-2})$  while the bias of the ordinary CV criterion is  $O(n^{-1})$ . The CCV criterion extends the result of Yanagihara, Tonda, and Matsumoto (2006), which consists of the K-L discrepancy.

In Section 2, we define the class of cross-validatory model selection criteria and study its mathematical properties. In Section 3, we describe other model selection criteria and their properties. In Section 4, the developed criteria will be applied to selecting structural equation models (SEM) under the normal distribution assumption. In Section 5, via the Monte Carlo method, we check the mathematical properties of the developed model selection criteria and compare CV and CCV criteria with other criteria such as AIC, TIC, and EIC. Conclusions and discussion are given in Section 6. Technical details are provided in an appendix.

## 2. A Class of Cross-Validatory Model Selection Criteria

Suppose that  $\psi(\mathbf{y}|\boldsymbol{\theta})$  is a discrepancy function for the candidate model  $M$  in (2), which is typically a function of  $f(\mathbf{y}|\boldsymbol{\theta})$ . Let  $\Psi(\boldsymbol{\theta}|\mathbf{Y}, \mathbf{w})$  be a weighted discrepancy function defined by

$$\Psi(\boldsymbol{\theta}|\mathbf{Y}, \mathbf{w}) = \sum_{i=1}^n w_i \psi(\mathbf{y}_i|\boldsymbol{\theta}), \quad (3)$$

where  $\mathbf{Y} = (\mathbf{y}_1, \dots, \mathbf{y}_n)'$  and  $\mathbf{w} = (w_1, \dots, w_n)'$ . Then an estimator of  $\boldsymbol{\theta}$  is obtained by minimizing the discrepancy function  $\Psi(\boldsymbol{\theta}|\mathbf{Y}, \mathbf{w})$  in (3), i.e.,

$$\hat{\boldsymbol{\theta}}(\mathbf{w}) = \arg \min_{\boldsymbol{\theta}} \Psi(\boldsymbol{\theta}|\mathbf{Y}, \mathbf{w}). \quad (4)$$

It is easy to see that  $\hat{\boldsymbol{\theta}}(\mathbf{w})$  is the maximum likelihood estimator (MLE) of  $\boldsymbol{\theta}$  when  $\psi(\mathbf{y}|\boldsymbol{\theta}) = -\log f(\mathbf{y}|\boldsymbol{\theta})$  and  $\mathbf{w} = \mathbf{1}_n = (1, \dots, 1)'$ . When

$$\psi(\mathbf{y}|\boldsymbol{\theta}) = -\frac{1}{\beta} f(\mathbf{y}|\boldsymbol{\theta})^\beta + \frac{1}{1+\beta} \int \{f(\mathbf{x}|\boldsymbol{\theta})\}^{1+\beta} d\mathbf{x},$$

$\Psi(\boldsymbol{\theta}|\mathbf{Y}, \mathbf{w})$  is the density power divergence (Basu *et al.*, 1998). An application of the power density divergence can be found in Fujisawa and Eguchi (2006). For simplicity, we use  $\Psi(\boldsymbol{\theta}|\mathbf{Y}) = \Psi(\boldsymbol{\theta}|\mathbf{Y}, \mathbf{1}_n)$  and  $\hat{\boldsymbol{\theta}} = \hat{\boldsymbol{\theta}}(\mathbf{1}_n)$ . Furthermore, we write  $\hat{\boldsymbol{\theta}}_{[-i]} = \hat{\boldsymbol{\theta}}(\mathbf{1}_n - \mathbf{e}_i)$ , where  $\mathbf{e}_i$  is a  $n \times 1$  vector whose  $i$ th element is 1, and the other elements are 0. Notice that  $\hat{\boldsymbol{\theta}}_{[-i]}$  becomes the jackknife estimator evaluated from the  $i$ th jackknife sample, which is obtained from  $\mathbf{Y}$  by deleting  $\mathbf{y}_i$ . Let  $\mathbf{u}_1, \dots, \mathbf{u}_n$  be  $p \times 1$  independent random vectors from  $\mathbf{u} \sim \varphi(\mathbf{u})$  with  $\mathbf{U} = (\mathbf{u}_1, \dots, \mathbf{u}_n)'$ , which is also independent from  $\mathbf{Y}$ . Notice that  $\hat{\boldsymbol{\theta}}$  is a function of

**Y**. We define a risk based on the predictive discrepancy  $\Psi(\boldsymbol{\theta}|\cdot)$  as

$$R_{\text{PD}} = E_{\mathbf{y}}^* E_{\mathbf{u}}^* \left[ \Psi(\hat{\boldsymbol{\theta}}|\mathbf{U}) \right] = n E_{\mathbf{y}}^* E_{\mathbf{u}}^* \left[ \psi(\mathbf{u}|\hat{\boldsymbol{\theta}}) \right], \quad (5)$$

where  $E_{\mathbf{y}}^*$  and  $E_{\mathbf{u}}^*$  are expectations under the true model  $M^*$  in (1) with respect to  $\mathbf{y}$  and  $\mathbf{u}$ , respectively. In model selection based on  $\psi(\mathbf{y}|\boldsymbol{\theta})$ , we regard the model having the smallest  $R_{\text{PD}}$  as the best model, which is typically different from the true model. In many contexts of statistical modeling, the aim is to determine the best model. Obtaining an unbiased estimator of  $R_{\text{PD}}$  will allow us to correctly evaluate the discrepancy between data and model, which will further facilitate the selection of the best model. The simplest estimator of  $R_{\text{PD}}$  is the sample discrepancy function  $\Psi(\hat{\boldsymbol{\theta}}|\mathbf{Y})$ . The CV criterion proposed by Stone (1974, 1977),

$$\text{CV} = \sum_{i=1}^n \psi(\mathbf{y}_i|\hat{\boldsymbol{\theta}}_{[-i]}), \quad (6)$$

is also an estimator of  $R_{\text{PD}}$ . Let

$$\mathbf{g}(\mathbf{y}|\boldsymbol{\theta}) = \left. \frac{\partial}{\partial \boldsymbol{\theta}} \psi(\mathbf{y}|\boldsymbol{\theta}) \right|_{\boldsymbol{\theta}=\boldsymbol{\theta}}, \quad \mathbf{H}(\mathbf{y}|\boldsymbol{\theta}) = \left. \frac{\partial^2}{\partial \boldsymbol{\theta} \partial \boldsymbol{\theta}'} \psi(\mathbf{y}|\boldsymbol{\theta}) \right|_{\boldsymbol{\theta}=\boldsymbol{\theta}}, \quad (7)$$

and

$$\mathbf{r}(\boldsymbol{\theta}) = E_{\mathbf{y}}^*[\mathbf{g}(\mathbf{y}|\boldsymbol{\theta})], \quad \mathbf{I}(\boldsymbol{\theta}) = E_{\mathbf{y}}^*[\mathbf{g}(\mathbf{y}|\boldsymbol{\theta})\mathbf{g}(\mathbf{y}|\boldsymbol{\theta})'], \quad \mathbf{J}(\boldsymbol{\theta}) = E_{\mathbf{y}}^*[\mathbf{H}(\mathbf{y}|\boldsymbol{\theta})]. \quad (8)$$

Suppose that  $\boldsymbol{\theta}_0$  is a  $q \times 1$  vector such that  $\hat{\boldsymbol{\theta}} \xrightarrow{a.s.} \boldsymbol{\theta}_0$  as  $n \rightarrow \infty$ . Under proper conditions, as specified in White (1982),  $\boldsymbol{\theta}_0$  satisfies

$$\mathbf{r}(\boldsymbol{\theta}_0) = \mathbf{0}_q, \quad (9)$$

where  $\mathbf{0}_q$  is a vector of  $q$  zeros. Notice that  $\mathbf{I}(\boldsymbol{\theta}_0)$  is called the Fisher's information matrix when  $\psi(\mathbf{y}|\boldsymbol{\theta}) = -\log f(\mathbf{y}|\boldsymbol{\theta})$ . Because  $\boldsymbol{\theta}_0$  is a local minimum of  $E_{\mathbf{y}}^*[\psi(\mathbf{y}|\boldsymbol{\theta})]$ , equation (9) leads to a natural assumption that  $\mathbf{J}(\boldsymbol{\theta}_0)$  is positive definite.

Let  $\hat{\boldsymbol{\theta}}_i(\lambda)$  ( $0 \leq \lambda \leq 1$ ) be the estimator of  $\boldsymbol{\theta}$ , which is obtained by minimizing the weighted discrepancy function  $\Psi(\boldsymbol{\theta}|\mathbf{Y}, \mathbf{1}_n - \lambda \mathbf{e}_i)$ , i.e.,  $\hat{\boldsymbol{\theta}}_i(\lambda) = \hat{\boldsymbol{\theta}}(\mathbf{1}_n - \lambda \mathbf{e}_i)$ . Notice that, with weight  $\mathbf{1}_n - \lambda \mathbf{e}_i$ , the effect of the  $i$ th observation  $\mathbf{y}_i$  on  $\hat{\boldsymbol{\theta}}_i(\lambda)$  decreases as  $\lambda$  increases. The estimator  $\hat{\boldsymbol{\theta}}_i(\lambda)$  includes the ordinary estimator and the  $i$ th jackknife estimator as special cases, i.e.,  $\hat{\boldsymbol{\theta}}_i(0) = \hat{\boldsymbol{\theta}}$  and  $\hat{\boldsymbol{\theta}}_i(1) = \hat{\boldsymbol{\theta}}_{[-i]}$ . Replacing  $\hat{\boldsymbol{\theta}}$  by  $\hat{\boldsymbol{\theta}}_i(\lambda)$ , we define the following cross-validatory model selection criterion:

$$\text{CV}(\lambda) = \sum_{i=1}^n \psi(\mathbf{y}_i|\hat{\boldsymbol{\theta}}_i(\lambda)), \quad (0 \leq \lambda \leq 1). \quad (10)$$

Let

$$\mathcal{G}_\lambda = \{\mathbf{E}_\mathbf{y}^*[\text{CV}(\lambda)] \mid 0 \leq \lambda \leq 1\},$$

and

$$R_1 = \sum_{i=1}^n \mathbf{E}_\mathbf{y}^* \left[ \mathbf{r}(\hat{\boldsymbol{\theta}})' (\hat{\boldsymbol{\theta}}_{[-i]} - \hat{\boldsymbol{\theta}}) \right], \quad R_2 = \sum_{i=1}^n \mathbf{E}_\mathbf{y}^* \left[ (\hat{\boldsymbol{\theta}}_{[-i]} - \hat{\boldsymbol{\theta}})' \mathbf{J}(\bar{\boldsymbol{\theta}}_i(\delta_i)) (\hat{\boldsymbol{\theta}}_{[-i]} - \hat{\boldsymbol{\theta}}) \right], \quad (11)$$

where

$$\bar{\boldsymbol{\theta}}_i(\delta_i) = \hat{\boldsymbol{\theta}} + \delta_i(\hat{\boldsymbol{\theta}}_{[-i]} - \hat{\boldsymbol{\theta}}), \quad (i = 1, \dots, n), \quad (12)$$

with  $\delta_i \in (0, 1)$ . The following theorem characterizes the properties of  $\text{CV}(\lambda)$  (the proof is given in Appendix A.1).

**THEOREM 1.** *The model selection criterion  $\text{CV}(\lambda)$  has the following properties:*

1.  $\text{CV}(0) = \Psi(\hat{\boldsymbol{\theta}}|\mathbf{Y})$  and  $\text{CV}(1) = \text{CV}$ .
2.  $\text{CV}(\lambda)$  is an increasing function of  $\lambda$ .
3.  $n\mathbf{E}_\mathbf{y}^*[\psi(\mathbf{y}|\boldsymbol{\theta}_0)] \in \mathcal{G}_\lambda$  when  $\boldsymbol{\theta}_0$  is a global minimum of  $\mathbf{E}_\mathbf{y}^*[\psi(\mathbf{y}|\boldsymbol{\theta})]$ .
4.  $R_{\text{PD}} \in \mathcal{G}_\lambda$  when  $R_1 + R_2/2 \geq 0$ .

Appendix A.2 provides the detail leading to  $R_1 = O(n^{-2})$  and  $R_2 = \gamma_1 + O(n^{-2})$ , where  $\gamma_1$  is given by

$$\gamma_1 = \text{tr}\{\mathbf{I}(\boldsymbol{\theta}_0)\mathbf{J}(\boldsymbol{\theta}_0)^{-1}\}. \quad (13)$$

Because  $\mathbf{J}(\boldsymbol{\theta}_0)$  is positive definite,  $\gamma_1$  is positive. Thus,  $R_1 + R_2/2 \geq 0$  asymptotically holds. Consequently,  $R_{\text{PD}} \in \mathcal{G}_\lambda$  when  $n$  is adequate. When  $\psi(\mathbf{y}|\boldsymbol{\theta})$  is a strictly convex function of  $\boldsymbol{\theta}$ ,  $\mathbf{H}(\mathbf{y}|\boldsymbol{\theta})$  is positive definite for any  $\boldsymbol{\theta}$  and  $\mathbf{y}$  (see e.g., Lehmann & Casella, 1998, p. 49). Then,  $\mathbf{J}(\boldsymbol{\theta})$  will be positive definite for any  $\boldsymbol{\theta}$ . This directly implies that  $R_2 > 0$  when  $\psi(\mathbf{y}|\boldsymbol{\theta})$  is a strictly convex function of  $\boldsymbol{\theta}$ . Thus,  $R_{\text{PD}} \in \mathcal{G}_\lambda$  when  $R_1 \geq 0$ . Although the order of  $R_1$  is  $O(n^{-2})$ ,  $R_1 \geq 0$  holds under special cases, as in the following example.

**EXAMPLE 1.** Suppose that the candidate model  $M$  and the true model  $M^*$  are given by

$$\begin{aligned} M : \quad & \mathbf{y}_1, \dots, \mathbf{y}_n \sim i.i.d. \text{N}_p(\boldsymbol{\mu}, \boldsymbol{\Sigma}), \\ M^* : \quad & \mathbf{y}_1, \dots, \mathbf{y}_n \sim i.i.d. \text{E}[\mathbf{y}] = \boldsymbol{\mu}^* \text{ and } \text{Cov}[\mathbf{y}] = \boldsymbol{\Sigma}^*. \end{aligned}$$

If the K-L discrepancy is used to define  $\text{CV}(\lambda)$ , Appendix A.3 shows that  $R_1 > 0$  always holds. Thus,  $R_{\text{PD}} \in \mathcal{G}_\lambda$ .

An important issue is how to choose  $\lambda$ . It follows from Theorem 1 that, when  $R_1 + R_2/2 \geq 0$ , a  $\lambda_0$  exists such that  $\text{E}_{\mathbf{y}}^*[\text{CV}(\lambda_0)] = R_{\text{PD}}$ . However, since  $\lambda_0$  depends on the unknown distribution  $\varphi(\mathbf{y})$ , it is very difficult to find the exact  $\lambda_0$ . Even if we can obtain  $\lambda_0$  somehow, it may be difficult to put it to practice. This is because the optimal  $\lambda_0$  may depend on cumulants of  $\varphi(\mathbf{y})$ . It is difficult to obtain good estimates of higher-order cumulants even when  $n$  is relatively large (see Yanagihara (2007) for the case of kurtosis). Thus, an estimator of  $\lambda$  that does not involve higher-order cumulants is preferable. Let  $\beta_{abcd} = \text{E}_{\mathbf{y}}^*[\partial^4 \psi(\mathbf{y}|\boldsymbol{\theta}) / (\partial \theta_a \partial \theta_b \partial \theta_c \partial \theta_d)]|_{\boldsymbol{\theta}=\boldsymbol{\theta}_0}$  ( $1 \leq a, b, c, d \leq q$ ). The following theorem characterizes the bias of  $\text{CV}(\lambda)$  (the proof is given in Appendix A.4).

**THEOREM 2.** *Under the condition  $|\beta_{abcd}| < \infty$  ( $1 \leq a, b, c, d \leq q$ ), the bias of  $\text{CV}(\lambda)$  is characterized as*

$$R_{\text{PD}} - \text{E}_{\mathbf{y}}^*[\text{CV}(\lambda)] = \begin{cases} (1 - \lambda)\gamma_1 + O(n^{-1}) & (\lambda \text{ is independent of } n) \\ \{1 - \lambda - 1/(2n)\}\gamma_1 + O(n^{-2}) & (\lambda = 1 + O(n^{-1})) \end{cases}, \quad (14)$$

where  $\gamma_1$  is given by (13).

The moment condition in Theorem 2 (also Theorems A.1 and A.2) may be weakened as in Hall (1987). If  $\lambda = 1 - 1/(2n)$ , then the  $O(n^{-1})$  term in the bias of  $\text{CV}(\lambda)$  in Theorem 2 vanishes. Thus, using second-order asymptotics, the optimal value of  $\lambda$  is  $\lambda = 1 - 1/(2n) + O(n^{-2})$ . Based on this, we propose a bias-corrected CV (CCV) criterion as in the following theorem.

**THEOREM 3.** *Let  $a_n \in (0, 1)$  that can be expanded as  $a_n = 1 - 1/(2n) + O(n^{-2})$ , and*

$$\text{CCV} = \text{CV}(a_n) = \sum_{i=1}^n \psi(\mathbf{y}_i | \hat{\boldsymbol{\theta}}_i(a_n)). \quad (15)$$

*Then the bias of the CCV criterion is  $O(n^{-2})$ , while the bias in the ordinary CV criterion is  $O(n^{-1})$ . Because  $a_n < 1$ , the CCV criterion is always smaller than the ordinary CV criterion.*

Notice that the CCV in (15) coincides with the CCV criterion in Yanagihara, Tonda, and Matsumoto (2006) when  $\psi(\mathbf{y}|\boldsymbol{\theta}) = -2 \log f(\mathbf{y}|\boldsymbol{\theta})$ .

Since our assumption is that  $\mathbf{y}_1, \dots, \mathbf{y}_n$  are *i.i.d.*, it may seem that Theorem 3 does not apply to selecting explanatory variables in regression models, which are widely used in data analysis. Let  $\mathbf{y} = (\mathbf{z}', \mathbf{x}')'$ , where  $\mathbf{z}$  is the vector of response variables and  $\mathbf{x}$  is the vector of explanatory variables. Then our result immediately applies to the regression model. In order to calculate  $\text{CV}(\lambda)$ , it is often necessary to obtain each  $\hat{\boldsymbol{\theta}}_i(\lambda)$ . However,  $\text{CV}(\lambda)$  in the linear regression model under the normal distribution assumption can be derived using  $\hat{\boldsymbol{\theta}}$  alone, as in the following example.

EXAMPLE 2. Let  $\mathbf{z}$  and  $\mathbf{x}$  be  $m \times 1$  and  $k \times 1$  vectors and  $\mathbf{y} = (\mathbf{z}', \mathbf{x}')'$ . Suppose that the candidate model  $M$  and the true model  $M^*$  are given by

$$\begin{aligned} M : \quad \mathbf{z}_i | \mathbf{x}_i &\sim N_m(\boldsymbol{\Xi}' \tilde{\mathbf{x}}_i, \boldsymbol{\Gamma}), \\ M^* : \quad \mathbf{y}_1, \dots, \mathbf{y}_n &\sim i.i.d. \text{ E}[\mathbf{y}] = \boldsymbol{\mu}^* \text{ and } \text{Cov}[\mathbf{y}] = \boldsymbol{\Sigma}^*, \end{aligned}$$

where  $\tilde{\mathbf{x}}_i = (1, \mathbf{x}'_i)'$ . Notice that the MLEs of  $\boldsymbol{\Xi}$  and  $\boldsymbol{\Gamma}$  are  $\hat{\boldsymbol{\Xi}} = (\tilde{\mathbf{X}}' \tilde{\mathbf{X}})^{-1} \tilde{\mathbf{X}}' \mathbf{Z}$  and  $\hat{\boldsymbol{\Gamma}} = \mathbf{Z}' \{ \mathbf{I}_n - \tilde{\mathbf{X}} (\tilde{\mathbf{X}}' \tilde{\mathbf{X}})^{-1} \tilde{\mathbf{X}}' \} \mathbf{Z} / n$ , where  $\mathbf{Z} = (\mathbf{z}_1, \dots, \mathbf{z}_n)'$  and  $\tilde{\mathbf{X}} = (\tilde{\mathbf{x}}_1, \dots, \tilde{\mathbf{x}}_n)'$ . Then,  $\text{CV}(\lambda)$  in the case of  $\psi(\mathbf{y} | \boldsymbol{\theta}) = -2 \log f(\mathbf{y} | \boldsymbol{\theta})$  is given by

$$\begin{aligned} \text{CV}(\lambda) = \quad & n \log |\hat{\boldsymbol{\Gamma}}| + nm \log \left( \frac{2n\pi}{n-\lambda} \right) + \sum_{i=1}^n \log \left\{ 1 - \frac{\lambda \hat{r}_i^2}{n(1-\lambda c_i)} \right\} \\ & + \left( 1 - \frac{\lambda}{n} \right) \sum_{i=1}^n \frac{\hat{r}_i^2}{(1-\lambda c_i)^2} \left\{ 1 - \frac{\lambda \hat{r}_i^2}{n(1-\lambda c_i)} \right\}^{-1}, \end{aligned}$$

where  $c_i = \tilde{\mathbf{x}}_i' (\tilde{\mathbf{X}}' \tilde{\mathbf{X}})^{-1} \tilde{\mathbf{x}}_i$  and  $\hat{r}_i^2 = (\mathbf{z}_i - \hat{\boldsymbol{\Xi}}' \tilde{\mathbf{x}}_i)' \hat{\boldsymbol{\Gamma}}^{-1} (\mathbf{z}_i - \hat{\boldsymbol{\Xi}}' \tilde{\mathbf{x}}_i)$ .

Yanagihara, Kamo, and Tonda (2006) proposed a second-order bias-corrected AIC, called  $\text{CAIC}_J$ , in multivariate linear models. The order of the bias of  $\text{CAIC}_J$  is the same as that of CCV. However,  $\text{CAIC}_J$  was obtained under the assumption that the explanatory variables  $\mathbf{x}$  are nonstochastic, while the condition here is that both the explanatory variables  $\mathbf{x}$  and the response variables  $\mathbf{z}$  are stochastic.

For linear regression models, the well-known CV criterion is defined by the predicted residual sum of squares. Our general formula also applies to this case, and the  $\text{CV}(\lambda)$  is given by the following example.

EXAMPLE 3. Let  $\mathbf{x}$  be a  $k \times 1$  vector and  $\mathbf{y} = (z, \mathbf{x}')'$ . Suppose that the candidate model  $M$  and the true model  $M^*$  are

$$M : \quad \text{E}[z_i | \mathbf{x}_i] = \boldsymbol{\beta}' \tilde{\mathbf{x}}_i,$$

$$M^* : \mathbf{y}_1, \dots, \mathbf{y}_n \sim i.i.d. \text{ E}[\mathbf{y}] = \boldsymbol{\mu}^* \text{ and Cov}[\mathbf{y}] = \boldsymbol{\Sigma}^*,$$

where  $\tilde{\mathbf{x}}_i = (1, \mathbf{x}'_i)'$ . Notice that the least square estimator of  $\boldsymbol{\beta}$  is given by  $\hat{\boldsymbol{\beta}} = (\tilde{\mathbf{X}}' \tilde{\mathbf{X}})^{-1} \tilde{\mathbf{X}}' \mathbf{z}$ , where  $\mathbf{z} = (z_1, \dots, z_n)'$  and  $\tilde{\mathbf{X}} = (\tilde{\mathbf{x}}_1, \dots, \tilde{\mathbf{x}}_n)'$ . Thus,  $\text{CV}(\lambda)$  in the case of the predicted residual sum of squares is given by

$$\text{CV}(\lambda) = \sum_{i=1}^n \left\{ \frac{z_i - \hat{\boldsymbol{\beta}}' \tilde{\mathbf{x}}_i}{1 - \lambda \tilde{\mathbf{x}}_i' (\tilde{\mathbf{X}}' \tilde{\mathbf{X}})^{-1} \tilde{\mathbf{x}}_i} \right\}^2.$$

### 3. Other Model Selection Criteria

In this section, we discuss other criteria for selecting the best model among all the candidate models using the general discrepancy function  $\psi(\mathbf{y}|\boldsymbol{\theta})$ . The AIC-type criterion can be defined by adding the number of parameters to the sample discrepancy function as

$$\text{AIC} = \Psi(\hat{\boldsymbol{\theta}}|\mathbf{Y}) + q. \quad (16)$$

However, unless  $\psi(\mathbf{y}|\boldsymbol{\theta}) = -\log f(\mathbf{y}|\boldsymbol{\theta})$  and  $\mathcal{F}$  contains  $\varphi(\mathbf{y})$ , (16) has a constant bias in estimating  $R_{\text{PD}}$ . The TIC-type criterion corrects the bias of the AIC-type criterion, reducing the bias to  $O(n^{-1})$ . The TIC-type criterion is given by

$$\text{TIC} = \Psi(\hat{\boldsymbol{\theta}}|\mathbf{Y}) + \text{tr}\{\hat{\mathbf{I}}(\hat{\boldsymbol{\theta}})\hat{\mathbf{J}}(\hat{\boldsymbol{\theta}})^{-1}\}, \quad (17)$$

where

$$\hat{\mathbf{I}}(\hat{\boldsymbol{\theta}}) = \frac{1}{n} \sum_{i=1}^n \mathbf{g}(\mathbf{y}_i|\hat{\boldsymbol{\theta}})\mathbf{g}(\mathbf{y}_i|\hat{\boldsymbol{\theta}})', \quad \hat{\mathbf{J}}(\hat{\boldsymbol{\theta}}) = \frac{1}{n} \sum_{i=1}^n \mathbf{H}(\mathbf{y}_i|\hat{\boldsymbol{\theta}}), \quad (18)$$

with  $\mathbf{g}(\cdot|\cdot)$  and  $\mathbf{H}(\cdot|\cdot)$  being given by (7). Although the order of the bias in TIC is the same as that in the CV criterion, the bias of TIC tends to be larger than that of CV because  $\text{tr}\{\hat{\mathbf{I}}(\hat{\boldsymbol{\theta}})\hat{\mathbf{J}}(\hat{\boldsymbol{\theta}})^{-1}\}$  may contain a large bias. Actually, Theorems A.1 and A.2 in Appendix A.2 show that the  $n^{-1}$  term of the bias in TIC contains more terms of higher-order moments than that of the CV criterion.

The bootstrap method can also correct the bias of the AIC-type criterion. The resulting criterion is called the EIC-type criterion. Let  $\mathbf{y}_{b,1}^*, \dots, \mathbf{y}_{b,n}^*$  be the  $b$ th bootstrap resample from  $\mathbf{Y}$  ( $b = 1, \dots, B$ ) and  $\hat{\boldsymbol{\theta}}_b^*$  be the estimator of  $\boldsymbol{\theta}$ , where

$$\hat{\boldsymbol{\theta}}_b^* = \arg \min_{\boldsymbol{\theta}} \sum_{i=1}^n \psi(\mathbf{y}_{b,i}^*|\boldsymbol{\theta}).$$



Replacing the log-likelihood function by the discrepancy function  $\psi(\mathbf{y}|\boldsymbol{\theta})$  in the formula of Konishi (1999), the EIC-type criterion can be defined by

$$\text{EIC} = \Psi(\hat{\boldsymbol{\theta}}|\mathbf{Y}) + \frac{1}{B} \sum_{b=1}^B \left\{ \sum_{i=1}^n \psi(\mathbf{y}_i|\hat{\boldsymbol{\theta}}_b^*) - \sum_{i=1}^n \psi(\mathbf{y}_{b,i}^*|\hat{\boldsymbol{\theta}}_b^*) \right\}. \quad (19)$$

By using random vectors distributed according to the multinomial distribution, we can rewrite the definition of EIC in (19). Let  $\mathbf{d}_b = (d_{b1}, \dots, d_{bn})'$  ( $b = 1, \dots, B$ ) be random samples of size  $n$  from the multinomial distribution  $\text{Multi}_n(n; 1/n, \dots, 1/n)$ . Then, the EIC in (19) is equivalent to the following formula (the derivation is given in Appendix A.5):

$$\text{EIC} = \Psi(\hat{\boldsymbol{\theta}}|\mathbf{Y}) + \frac{1}{B} \sum_{b=1}^B \Psi(\hat{\boldsymbol{\theta}}(\mathbf{d}_b)|\mathbf{Y}, \mathbf{1}_n - \mathbf{d}_b). \quad (20)$$

where  $\hat{\boldsymbol{\theta}}(\cdot)$  is given by (4). Because the bias of EIC is  $O(n^{-1})$ , the order of bias in EIC is the same as those in TIC and the ordinary CV criterion. However, since EIC does not contain the term  $\text{tr}\{\hat{\mathbf{I}}(\hat{\boldsymbol{\theta}})\hat{\mathbf{J}}(\hat{\boldsymbol{\theta}})^{-1}\}$ , the bias of EIC tends to be smaller than that of TIC. On the other hand, EIC involves more computation than the CV criterion. Furthermore, EIC may behave poorly when the sample size is small and the number of parameters is large. Caution is needed when using EIC with small samples.

## 4. Application to Selecting Structural Equation Models Under the Normal Distribution Assumption

SEM is a multivariate statistical technique designed to model the covariance matrix by a structure with relatively few parameters (see e.g., Lee & Kontoghiorghes, 2007; Yuan & Bentler, 2007). The normal distribution assumption is typically used in the practice of SEM and is the default option of all statistical software (AMOS, EQS, LISREL, Mplus, SAS Calis). We will obtain the analytical expression of  $\text{CV}(\lambda)$  when the candidate model is from the normal family while the true model is unknown.

Let the candidate model  $M$  and the true model  $M^*$  be

$$\begin{aligned} M &: \mathbf{y}_1, \dots, \mathbf{y}_n \sim i.i.d. N_p(\boldsymbol{\mu}, \boldsymbol{\Sigma}(\boldsymbol{\xi})), \\ M^* &: \mathbf{y}_1, \dots, \mathbf{y}_n \sim i.i.d. E[\mathbf{y}] = \boldsymbol{\mu}^* \text{ and } \text{Cov}[\mathbf{y}] = \boldsymbol{\Sigma}^*, \end{aligned} \quad (21)$$

where  $\boldsymbol{\mu} = (\mu_1, \dots, \mu_p)'$  and  $\boldsymbol{\xi} = (\xi_1, \dots, \xi_q)'$  are  $p \times 1$  and  $q \times 1$  unknown vectors of parameters, respectively, and the true distribution of  $\mathbf{y}$  is unknown. Consider the K-L

discrepancy with

$$\psi(\mathbf{y}|\boldsymbol{\theta}) = -2 \log f(\mathbf{y}|\boldsymbol{\theta}) = p \log(2\pi) + \log |\boldsymbol{\Sigma}(\boldsymbol{\xi})| + (\mathbf{y} - \boldsymbol{\mu})' \boldsymbol{\Sigma}(\boldsymbol{\xi})^{-1} (\mathbf{y} - \boldsymbol{\mu}), \quad (22)$$

where  $\boldsymbol{\theta} = (\boldsymbol{\mu}', \boldsymbol{\xi}')'$ . Let

$$\bar{\mathbf{y}} = \frac{1}{n} \sum_{i=1}^n \mathbf{y}_i, \quad \mathbf{S} = \frac{1}{n} \sum_{i=1}^n (\mathbf{y}_i - \bar{\mathbf{y}})(\mathbf{y}_i - \bar{\mathbf{y}})', \quad (23)$$

and

$$F(\boldsymbol{\xi}|\mathbf{A}) = \log |\boldsymbol{\Sigma}(\boldsymbol{\xi})| + \text{tr} \{ \mathbf{A} \boldsymbol{\Sigma}(\boldsymbol{\xi})^{-1} \}. \quad (24)$$

Then, the  $\text{CV}(\lambda)$  defined in (10) is given by the following corollary (the proof is given in Appendix A.6).

**COROLLARY 1.** *The  $\text{CV}(\lambda)$  under the candidate model  $M$  in (21) is given by*

$$\begin{aligned} \text{CV}(\lambda) &= np \log(2\pi) \\ &+ \sum_{i=1}^n \left\{ \log |\boldsymbol{\Sigma}(\hat{\boldsymbol{\theta}}_i(\lambda))| + \left( \frac{n}{n-\lambda} \right)^2 (\mathbf{y}_i - \bar{\mathbf{y}})' \boldsymbol{\Sigma}(\hat{\boldsymbol{\xi}}_i(\lambda))^{-1} (\mathbf{y}_i - \bar{\mathbf{y}}) \right\}, \end{aligned} \quad (25)$$

where  $\hat{\boldsymbol{\xi}}_i(\lambda)$  is the estimator of  $\boldsymbol{\xi}$  defined by

$$\hat{\boldsymbol{\xi}}_i(\lambda) = \arg \min_{\boldsymbol{\xi}} F(\boldsymbol{\xi}|\mathbf{S}_i(\lambda)), \quad (26)$$

with

$$\mathbf{S}_i(\lambda) = \frac{n}{n-\lambda} \left\{ \mathbf{S} - \frac{\lambda}{n-\lambda} (\mathbf{y}_i - \bar{\mathbf{y}})(\mathbf{y}_i - \bar{\mathbf{y}})' \right\}. \quad (27)$$

The following corollary provides the analytical expression for other model selection criteria.

**COROLLARY 2.** *Let  $\hat{\boldsymbol{\xi}}$  be the estimator of  $\boldsymbol{\xi}$  such that*

$$\hat{\boldsymbol{\xi}} = \arg \min_{\boldsymbol{\xi}} F(\boldsymbol{\xi}|\mathbf{S}). \quad (28)$$

*Then, AIC, TIC, and EIC under the candidate model  $M$  in (21) are given by*

$$\text{AIC} = nF(\hat{\boldsymbol{\xi}}|\mathbf{S}) + np \log(2\pi) + 2(p+q);$$

$$\begin{aligned} \text{TIC} &= \text{AIC} - 2(p+q) + 2F(\hat{\boldsymbol{\xi}}|\mathbf{S}) \\ &+ \text{tr} \left\{ \hat{\boldsymbol{\Omega}}(\boldsymbol{\Sigma}(\hat{\boldsymbol{\xi}})^{-1} \otimes \boldsymbol{\Sigma}(\hat{\boldsymbol{\xi}})^{-1}) \mathbf{Q}(\hat{\boldsymbol{\xi}}|\mathbf{S}) (\boldsymbol{\Sigma}(\hat{\boldsymbol{\xi}})^{-1} \otimes \boldsymbol{\Sigma}(\hat{\boldsymbol{\xi}})^{-1}) \right\} \\ &- \text{vec}(\mathbf{S})' (\boldsymbol{\Sigma}(\hat{\boldsymbol{\xi}})^{-1} \otimes \boldsymbol{\Sigma}(\hat{\boldsymbol{\xi}})^{-1}) \mathbf{Q}(\hat{\boldsymbol{\xi}}|\mathbf{S}) (\boldsymbol{\Sigma}(\hat{\boldsymbol{\xi}})^{-1} \otimes \boldsymbol{\Sigma}(\hat{\boldsymbol{\xi}})^{-1}) \text{vec}(\mathbf{S}), \end{aligned}$$

where

$$\begin{aligned}\hat{\Omega} &= \frac{1}{n} \sum_{i=1}^n \text{vec}((\mathbf{y}_i - \bar{\mathbf{y}})(\mathbf{y}_i - \bar{\mathbf{y}})') \text{vec}((\mathbf{y}_i - \bar{\mathbf{y}})(\mathbf{y}_i - \bar{\mathbf{y}})')', \\ \mathbf{Q}(\hat{\boldsymbol{\xi}}|\mathbf{S}) &= \left\{ \frac{\partial}{\partial \boldsymbol{\xi}'} \text{vec}(\boldsymbol{\Sigma}(\boldsymbol{\xi})) \right\} \left\{ \frac{\partial^2}{\partial \boldsymbol{\xi}' \partial \boldsymbol{\xi}'} F(\boldsymbol{\xi}|\mathbf{S}) \right\}^{-1} \left\{ \frac{\partial}{\partial \boldsymbol{\xi}} \text{vec}(\boldsymbol{\Sigma}(\boldsymbol{\xi}))' \right\} \Bigg|_{\boldsymbol{\xi}=\hat{\boldsymbol{\xi}}}; \\ \text{EIC} &= \text{AIC} - 2(p + q) + \frac{n}{B} \sum_{b=1}^B \text{tr} \left\{ \mathbf{V}(\mathbf{d}_b) \boldsymbol{\Sigma}(\hat{\boldsymbol{\xi}}(\mathbf{d}_b))^{-1} \right\},\end{aligned}\quad (29)$$

where

$$\hat{\boldsymbol{\xi}}(\mathbf{d}_b) = \arg \min_{\boldsymbol{\xi}} F(\boldsymbol{\xi}|\mathbf{S}(\mathbf{d}_b)), \quad (30)$$

$$\mathbf{V}(\mathbf{d}_b) = \frac{1}{n} \mathbf{Y}' \left\{ \mathbf{I}_p - \text{diag}(\mathbf{d}_b) + \frac{1}{n} (2\mathbf{d}_b \mathbf{d}_b' - \mathbf{1}_n \mathbf{d}_b' - \mathbf{d}_b \mathbf{1}_n') \right\} \mathbf{Y}, \quad (31)$$

with

$$\mathbf{S}(\mathbf{d}_b) = \frac{1}{n} \mathbf{Y}' \left\{ \text{diag}(\mathbf{d}_b) - \frac{1}{n} \mathbf{d}_b \mathbf{d}_b' \right\} \mathbf{Y}. \quad (32)$$

The use of AIC for selecting the number of factors in the explanatory factor model was discussed by Akaike (1987). TIC for selecting SEM models under the normal distribution assumption was obtained by Yanagihara (2005). The details leading to the expression for EIC are provided in Appendix A.7.

## 5. Numerical Examinations

In this section, we verify the mathematical properties of model selection criteria using a Monte Carlo method. In particular, we compare CV and CCV criteria with AIC, TIC, and EIC. Bayesian information criterion (BIC; Schwarz, 1978) and the consistent Akaike's information criterion (CAIC; Bozdogan, 1987) are also frequently used for model selection, but their expectations do not convergence to  $R_{\text{PD}}$ . Thus, our study will not include BIC and CAIC.

In designing the Monte Carlo, we let the candidate distribution be multivariate normal as in the previous section, while the true distribution varies. Let  $\mathbf{y}$  be the  $6 \times 1$  vector defined

by  $\mathbf{y} = \Sigma^{*1/2}\boldsymbol{\varepsilon}$ , where

$$\Sigma^* = \begin{pmatrix} 2 & 1 & 0 & 0 & 0 & 0 \\ 1 & 2 & 1 & 0 & 0 & 0 \\ 0 & 1 & 2 & 0 & 0 & 0 \\ 0 & 0 & 1 & 2 & 1 & 0 \\ 0 & 0 & 0 & 1 & 2 & 1 \\ 0 & 0 & 0 & 0 & 1 & 6 \end{pmatrix}.$$

We use Mardia's (1970) multivariate skewnesses  $\kappa_{3,3}^{(1)}$  and  $\kappa_{3,3}^{(2)}$  and kurtosis  $\kappa_4^{(1)}$  to measure the departure of the candidate distribution from the true distribution. These are given by

$$\kappa_{3,3}^{(1)} = E[(\boldsymbol{\varepsilon}'_1 \boldsymbol{\varepsilon}_2)^3], \quad \kappa_{3,3}^{(2)} = E[(\boldsymbol{\varepsilon}'_1 \boldsymbol{\varepsilon}_1)(\boldsymbol{\varepsilon}'_1 \boldsymbol{\varepsilon}_2)(\boldsymbol{\varepsilon}'_2 \boldsymbol{\varepsilon}_2)], \quad \kappa_4^{(1)} = E[(\boldsymbol{\varepsilon}'_1 \boldsymbol{\varepsilon}_1)^2] - 48,$$

where  $\boldsymbol{\varepsilon}_1$  and  $\boldsymbol{\varepsilon}_2$  are independent random vectors having the same distribution of  $\boldsymbol{\varepsilon}$ .

Six populations or true models are created when the elements  $\varepsilon_j$  of  $\boldsymbol{\varepsilon} = (\varepsilon_1, \dots, \varepsilon_6)'$  are independently and identically distributed standardized variables from each of the following six distributions:

1. *Normal Distribution*:  $\varepsilon_j \sim N(0, 1)$ , ( $\kappa_{3,3}^{(1)} = \kappa_{3,3}^{(2)} = 0$  and  $\kappa_4^{(1)} = 0$ ).
2. *Laplace Distribution*:  $\varepsilon_j$  is generated from the Laplace distribution with mean 0 and standard deviation 1 ( $\kappa_{3,3}^{(1)} = \kappa_{3,3}^{(2)} = 0$  and  $\kappa_4^{(1)} = 18$ ).
3. *Uniform Distribution*:  $\varepsilon_j$  is generated from the uniform distribution on  $(-1, 1)$ , divided by the standard deviation  $1/\sqrt{3}$  ( $\kappa_{3,3}^{(1)} = \kappa_{3,3}^{(2)} = 0$  and  $\kappa_4^{(1)} = -7.2$ ).
4. *Skew-Laplace Distribution*:  $\varepsilon_j$  is generated from the skew-Laplace distribution with location parameter 0, dispersion parameter 1 and skew parameter 1, standardized by mean  $3/4$  and standard deviation  $\sqrt{23}/4$  ( $\kappa_{3,3}^{(1)} = \kappa_{3,3}^{(2)} \approx 7.32$  and  $\kappa_4^{(1)} \approx 19.56$ ).
5. *Chi-Square Distribution*:  $\varepsilon_j$  is generated from the chi-square distribution with 2 degrees of freedom, standardized by mean 2 and standard deviation 2 ( $\kappa_{3,3}^{(1)} = \kappa_{3,3}^{(2)} = 12$  and  $\kappa_4^{(1)} = 36$ ).
6. *Log-Normal Distribution*:  $\varepsilon_j$  is generated from the lognormal distribution such that  $\log \varepsilon_j \sim N(0, 1/2)$ , standardized by mean  $e^{1/4}$  and standard deviation  $\sqrt{e^{1/2}(e^{1/2} - 1)}$  ( $\kappa_{3,3}^{(1)} = \kappa_{3,3}^{(2)} \approx 17.64$  and  $\kappa_4^{(1)} \approx 111.06$ ).

The skew-Laplace distribution was proposed by Balakrishnan and Ambagaspitiya (1994) (for the probability density function, see e.g., Yanagihara & Yuan, 2005). The distributions in 1, 2, and 3 are symmetric, and distributions in 4, 5, and 6 are skewed.

A sample of size 20 is generated from  $\mathbf{y} = \Sigma^{*1/2}\boldsymbol{\varepsilon}$ . The three candidate models are:

$$\text{Model 1, } M_1 : \mathbf{y}_1, \dots, \mathbf{y}_{20} \sim i.i.d. N_6(\boldsymbol{\mu}, \sigma^2 \mathbf{I}_6),$$

$$\text{Model 2, } M_2 : \mathbf{y}_1, \dots, \mathbf{y}_{20} \sim i.i.d. N_6(\boldsymbol{\mu}, (\sigma^2 - \rho)\mathbf{I}_6 + \rho \mathbf{1}_6 \mathbf{1}'_6),$$

$$\text{Model 3, } M_3 : \mathbf{y}_1, \dots, \mathbf{y}_{20} \sim i.i.d. N_6(\boldsymbol{\mu}, \text{diag}(\sigma_1^2, \sigma_2^2, \sigma_3^2, \sigma_4^2, \sigma_5^2, \sigma_6^2)).$$

Because the sample size  $n$  ( $= 20$ ) is rather small compared with the dimension  $p$  ( $= 6$ ), the saturated model, i.e.,  $\mathbf{y}_1, \dots, \mathbf{y}_{20} \sim i.i.d. N_6(\boldsymbol{\mu}, \boldsymbol{\Sigma})$ , is not considered here. Since  $\boldsymbol{\Sigma}^* \neq \boldsymbol{\Sigma}(\boldsymbol{\xi})$  for any  $\boldsymbol{\xi}$  in any of the candidate models, all the candidate models are misspecified. We use the K-L discrepancy to select the best model among the three candidates. For each of the candidate models and distributions, results of Appendix A.3 imply that  $R_1 > 0$  and  $R_2 > 0$ . Thus,  $R_{\text{PD}} \in \mathcal{G}_\lambda$  in all three models.

The number of replications is chosen as  $N_r = 10,000$ . The following quantities are evaluated at each replication:  $\text{CV}(\lambda)$  with  $\lambda = 0.00, 0.01, 0.02, \dots, 0.98, 0.99, 1.00$ ;  $\text{CCV} = \text{CV}(a_n)$  with  $a_n = \sqrt{n/(n+1)} = \sqrt{20/21}$ ; AIC; TIC; and EIC using  $B = 1,000$  nested resamples. For each of the  $N_r$   $\hat{\boldsymbol{\theta}}$ 's,  $R = \sum_{i=1}^{20} \psi(\mathbf{u}_i | \hat{\boldsymbol{\theta}})$  with  $\mathbf{u}_1, \dots, \mathbf{u}_{20}$  being simulated from  $\mathbf{u} = \Sigma^{*1/2}\boldsymbol{\varepsilon}$  is also obtained, where  $\mathbf{u}_i$  are independent of  $\mathbf{y}_1, \dots, \mathbf{y}_{20}$ . The average of  $R$  across the  $N_r$  replications,  $\bar{R}$ , is regarded as the risk  $R_{\text{PD}}$ . Let  $\overline{\text{IC}}$  be the average of any of the above criteria; the relative bias and relative root mean square error (RMSE) of the criterion are evaluated by

$$\text{Relative Bias} = \frac{\bar{R} - \overline{\text{IC}}}{|\bar{R}|} \times 100, \quad \text{Relative RMSE} = \frac{\sqrt{\sum_{l=1}^{N_r} (\bar{R} - \text{IC}_l)^2 / N_r}}{|\bar{R}|} \times 100.$$

The smallest IC at each replication for a given model is recorded, as are its frequencies among the 10,000 replications.

Table 1 contains the risks ( $\bar{R}$ ) of all the candidate models at each true distribution. Model 3 has the smallest risk when the true distribution is normal and uniform; model 2 becomes the best when the true distribution is Laplace, skew-Laplace, chi-square, or log-normal.

Insert Table 1 around here

Figures 1 and 2 contain the plots of relative biases and RMSEs of  $\text{CV}(\lambda)$  against  $\lambda$ , respectively. Figure 3 contains the frequencies of the model being selected by  $\text{CV}(\lambda)$ . The plots in Figure 1 clearly show that there is an  $\lambda_0$  which makes  $\text{CV}(\lambda)$  an unbiased estimator of  $R_{\text{PD}}$ . In all the figures, the optimal  $\lambda_0$  is close to 1.0 or approximately  $1 - 1/(2n) = 39/40$ .

The bias approaches 0 as  $\lambda$  moves towards  $\lambda_0$ , and departs from 0 as  $\lambda$  moves away from  $\lambda_0$ . Larger biases of  $CV(\lambda)$  are associated with more unknown parameters or larger multivariate kurtosis of the true distribution ( $\kappa_4^{(1)}$ ). Comparing the plots for Laplace and skew-Laplace distributions, we may notice that the sizes of multivariate skewnesses  $\kappa_{3,3}^{(1)}$  and  $\kappa_{3,3}^{(2)}$  have little effect on the bias of  $CV(\lambda)$ . Similar to Figure 1, the plots in Figure 2 clearly show that, regardless of the model and distribution, there exists an  $\lambda_M \in (0, 1)$  such that  $CV(\lambda_M)$  has the smallest RMSE. Furthermore, Figure 3 shows that  $CV(\lambda)$  tends to choose model 3 for smaller  $\lambda$  and model 2 for larger  $\lambda$ .

Insert Figures 1, 2, and 3 around here

Table 2 contains the relative biases, RMSEs, and the frequencies of each of the models being selected by AIC, TIC, EIC, CV, and CCV criteria. The table clearly shows that the CCV criterion has the smallest bias among all the criteria. Moreover, the CCV criterion not only improves the bias of the CV criterion, but also its RMSE. The biases of AIC and TIC are greater than those of EIC, CV, and CCV criteria. In particular, AIC has a very large bias when  $\kappa_4^{(1)}$  is large. RMSE of EIC tends to be smaller than that of the CV criterion, although the bias of EIC tends to be greater than that of the CV criterion. Comparing Tables 1 and 2, CV and CCV select the model with the smallest risk most often. But AIC and TIC select model 3 most often while the best model changes with the true distribution. Notice that the frequency of choosing the best model by each criterion varies when the true distribution changes. Table 3 contains the average frequencies of choosing the best model by each criterion across all the true models. Among the 5 criteria, CCV chooses the best model most frequently; EIC and CV also work well.

In addition to the results reported above, several other models were also studied and similar results were obtained. While the frequency of choosing the best model by each criterion changes with the true model/distribution, the best criterion is mostly among EIC, CV and CCV.

Insert Tables 2 and 3 around here

## 6. Conclusion

In this paper, we defined the class of cross-validatory model selection criterion  $CV(\lambda)$  ( $0 \leq \lambda \leq 1$ ), which includes the sample discrepancy function and the ordinary CV criterion

as special cases.  $\text{CV}(\lambda)$  is an increasing function of  $\lambda$ . In particular, under proper conditions, there exists an  $\lambda_0 \in [0, 1]$  such that  $\text{CV}(\lambda_0)$  is unbiased for  $R_{\text{PD}}$ . Because  $R_1 = O(n^{-2})$  and  $R_2 = \gamma_1 + O(n^{-2})$  with  $\gamma_1 > 0$ ,  $|R_1|$  tends to be smaller than  $|R_2|$ . Thus,  $R_{\text{PD}} \in \mathcal{G}_\lambda$  in most cases. From the viewpoint of second-order asymptotics for the bias,  $\lambda = 1 - 1/(2n) + O(n^{-2})$  is optimal. We found that  $\lambda = \sqrt{n/(n+1)}$  worked well empirically. In particular, without estimating any higher-order cumulants, such a  $\lambda$  reduces the bias in CCV to  $O(n^{-2})$ . Such a result is especially valuable with small samples, where any criterion involving higher-order cumulants will inevitably perform poorly. The Monte Carlo results in the previous section verify the merit of CCV. In addition to the CCV criterion, other second-order bias-corrected criteria also exist. Those other criteria were generally obtained under specified models and distributions. The CCV criterion here is obtained under the general assumption, and can be applied broadly.

The aim of the CCV criterion is to minimize the bias in estimating  $R_{\text{PD}}$ . More important theme is to have a criterion that selects the model with the smallest risk. An unbiased estimator of  $R_{\text{RD}}$  does not necessarily lead to the model with the smallest  $R_{\text{PD}}$  being selected most frequently. Fortunately, the merits of least bias and selecting the best model both occur most frequently with CCV. Thus, we recommend the use of the CCV criterion for general model selection.

## Acknowledgements

Hirokazu Yanagihara's research was supported by the Ministry of Education, Science, Sports and Culture, Grant-in-Aid for Young Scientists (B), #17700274, 2005–2007. Ke-Hai Yuan's research was supported by NSF grant DMS-0437167 and the James McKeen Cattell Fund.

## Appendix

### A.1. Proof of Theorem 1

PROOF OF PROPERTY 1: We omit the proof because it is easy to verify.

PROOF OF PROPERTY 2: Let  $\Psi_i(\boldsymbol{\theta}|\mathbf{Y}, \lambda) = \Psi(\boldsymbol{\theta}|\mathbf{Y}, \mathbf{1}_n - \lambda e_i)$ , and  $\lambda_1 < \lambda_2$ . Because  $\hat{\boldsymbol{\theta}}_i(\lambda)$  minimizes  $\Psi_i(\boldsymbol{\theta}|\mathbf{Y}, \lambda)$ , there exist  $\Psi_i(\hat{\boldsymbol{\theta}}(\lambda_1)|\mathbf{Y}, \lambda_1) \leq \Psi_i(\hat{\boldsymbol{\theta}}(\lambda_2)|\mathbf{Y}, \lambda_1)$  and  $\Psi_i(\hat{\boldsymbol{\theta}}(\lambda_2)|\mathbf{Y}, \lambda_2) \leq$

$\Psi_i(\hat{\boldsymbol{\theta}}_i(\lambda_1)|\mathbf{Y}, \lambda_2)$ . By using these relations we obtain

$$\begin{aligned}
\Psi_i(\hat{\boldsymbol{\theta}}_i(\lambda_1)|\mathbf{Y}) - \lambda_1\psi(\mathbf{y}_i|\hat{\boldsymbol{\theta}}_i(\lambda_1)) &= \Psi_i(\hat{\boldsymbol{\theta}}_i(\lambda_1)|\mathbf{Y}, \lambda_1) \\
&\leq \Psi_i(\hat{\boldsymbol{\theta}}_i(\lambda_2)|\mathbf{Y}, \lambda_1) \\
&= \Psi_i(\hat{\boldsymbol{\theta}}_i(\lambda_2)|\mathbf{Y}) - \lambda_1\psi(\mathbf{y}_i|\hat{\boldsymbol{\theta}}_i(\lambda_2)) \\
&= \Psi_i(\hat{\boldsymbol{\theta}}_i(\lambda_2)|\mathbf{Y}, \lambda_2) + (\lambda_2 - \lambda_1)\psi(\mathbf{y}_i|\hat{\boldsymbol{\theta}}_i(\lambda_2)) \\
&\leq \Psi_i(\hat{\boldsymbol{\theta}}_i(\lambda_1)|\mathbf{Y}, \lambda_2) + (\lambda_2 - \lambda_1)\psi(\mathbf{y}_i|\hat{\boldsymbol{\theta}}_i(\lambda_2)) \\
&= \Psi_i(\hat{\boldsymbol{\theta}}_i(\lambda_1)|\mathbf{Y}) - \lambda_2\psi(\mathbf{y}_i|\hat{\boldsymbol{\theta}}_i(\lambda_1)) + (\lambda_2 - \lambda_1)\psi(\mathbf{y}_i|\hat{\boldsymbol{\theta}}_i(\lambda_2)).
\end{aligned}$$

Thus,

$$\psi(\mathbf{y}_i|\hat{\boldsymbol{\theta}}_i(\lambda_1)) \leq \psi(\mathbf{y}_i|\hat{\boldsymbol{\theta}}_i(\lambda_2)). \quad (\text{A1})$$

It follows from (A1) that

$$\text{CV}(\lambda_1) = \sum_{i=1}^n \psi(\mathbf{y}_i|\hat{\boldsymbol{\theta}}_i(\lambda_1)) \leq \sum_{i=1}^n \psi(\mathbf{y}_i|\hat{\boldsymbol{\theta}}_i(\lambda_2)) = \text{CV}(\lambda_2).$$

Consequently,  $\text{CV}(\lambda)$  is an increasing function of  $\lambda$ .

**PROOF OF PROPERTY 3:** Because  $\hat{\boldsymbol{\theta}}_{[-i]}$  minimizes  $\sum_{j \neq i}^n \psi(\mathbf{y}_j|\boldsymbol{\theta})$ , there exists

$$\sum_{j \neq i}^n \psi(\mathbf{y}_j|\hat{\boldsymbol{\theta}}_{[-i]}) \leq \sum_{j \neq i}^n \psi(\mathbf{y}_j|\hat{\boldsymbol{\theta}}).$$

Thus,

$$\mathbb{E}_{\mathbf{y}}^*[\psi(\mathbf{y}_j|\hat{\boldsymbol{\theta}}_{[-i]})] \leq \mathbb{E}_{\mathbf{y}}^*[\psi(\mathbf{y}_j|\hat{\boldsymbol{\theta}})], \quad (j \neq i). \quad (\text{A2})$$

Let  $\hat{\boldsymbol{\theta}}_n$  be the minimizer of the discrepancy function based on  $\mathbf{y}_1, \dots, \mathbf{y}_n$ , and  $\alpha_n = \mathbb{E}_{\mathbf{y}}^*[\psi(\mathbf{y}_1|\hat{\boldsymbol{\theta}}_n)]$ .

Then  $\alpha_{n-1} = \mathbb{E}_{\mathbf{y}}^*[\psi(\mathbf{y}_j|\hat{\boldsymbol{\theta}}_{[-i]})]$  and  $\alpha_n = \mathbb{E}_{\mathbf{y}}^*[\psi(\mathbf{y}_j|\hat{\boldsymbol{\theta}})]$ . It follows from (A2) that  $\alpha_{n-1} \leq \alpha_n$  for any  $n$ . Thus,  $\alpha_n$  monotonically increases. Let  $\gamma(\boldsymbol{\theta}) = \mathbb{E}_{\mathbf{y}}^*[\psi(\mathbf{y}|\boldsymbol{\theta})]$ . Then  $\lim_{n \rightarrow \infty} \alpha_n = \gamma(\boldsymbol{\theta}_0)$  follows from  $\hat{\boldsymbol{\theta}}_n \xrightarrow{a.s.} \boldsymbol{\theta}_0$ . Therefore,  $\alpha_n$  is bounded and monotonically increases. This directly implies that  $\alpha_n \leq \gamma(\boldsymbol{\theta}_0)$  and

$$\mathbb{E}_{\mathbf{y}}^*[\text{CV}(0)] \leq n\mathbb{E}_{\mathbf{y}}^*[\psi(\mathbf{y}|\boldsymbol{\theta}_0)]. \quad (\text{A3})$$

On the other hand, if  $\boldsymbol{\theta}_0$  is the global minimum of  $\gamma(\boldsymbol{\theta})$ , there must exist  $\gamma(\boldsymbol{\theta}_0) \leq \gamma(\boldsymbol{\theta})$  for any  $\boldsymbol{\theta}$ . Thus,  $\gamma(\boldsymbol{\theta}_0) \leq \mathbb{E}_{\mathbf{y}}^*[\psi(\mathbf{y}_i|\hat{\boldsymbol{\theta}}_{[-i]})]$ , or equivalently

$$n\mathbb{E}_{\mathbf{y}}^*[\psi(\mathbf{y}|\boldsymbol{\theta}_0)] \leq \mathbb{E}_{\mathbf{y}}^*[\text{CV}(1)]. \quad (\text{A4})$$



Equations (A3) and (A4) imply  $n\mathbf{E}_{\mathbf{y}}^*[\psi(\mathbf{y}|\boldsymbol{\theta}_0)] \in \mathcal{G}_\lambda$ .

**PROOF OF PROPERTY 4:** We will first show  $\mathbf{E}_{\mathbf{y}}^*[\text{CV}] \leq R_{\text{PD}}$ , where  $R_{\text{PD}}$  is given by (5). Let  $\hat{\boldsymbol{\theta}}_{\mathbf{U}}$  be the minimizer of  $\Psi(\boldsymbol{\theta}|\mathbf{U})$ . Because  $\mathbf{U}$  and  $\mathbf{Y}$  are identically distributed,  $\mathbf{E}_{\mathbf{y}}^*\mathbf{E}_{\mathbf{u}}^*[\Psi(\hat{\boldsymbol{\theta}}_{\mathbf{U}}|\mathbf{Y})] = \mathbf{E}_{\mathbf{y}}^*\mathbf{E}_{\mathbf{u}}^*[\Psi(\hat{\boldsymbol{\theta}}|\mathbf{U})] = R_{\text{PD}}$ . The property  $\mathbf{E}_{\mathbf{y}}^*[\text{CV}(0)] \leq R_{\text{PD}}$  follows by noticing that  $\text{CV}(0) = \Psi(\hat{\boldsymbol{\theta}}|\mathbf{Y}) \leq \Psi(\hat{\boldsymbol{\theta}}_{\mathbf{U}}|\mathbf{Y})$ . We next show that  $R_{\text{PD}} \leq \mathbf{E}_{\mathbf{y}}^*[\text{CV}(1)]$  when  $R_1 + R_2/2 \geq 0$ , where  $R_1$  and  $R_2$  are given by (11). Notice that  $\text{CV}(1) = \text{CV}$  and  $\hat{\boldsymbol{\theta}}_{[-i]}$  and  $\mathbf{y}_i$  are independent. Because the distribution of  $\mathbf{u}_i$  is identical to that of  $\mathbf{y}_i$ ,  $\mathbf{E}_{\mathbf{y}}^*[\psi(\mathbf{y}_i|\hat{\boldsymbol{\theta}}_{[-i]})] = \mathbf{E}_{\mathbf{y}}^*\mathbf{E}_{\mathbf{u}}^*[\psi(\mathbf{u}_i|\hat{\boldsymbol{\theta}}_{[-i]})]$ . Applying the Taylor expansion at  $\hat{\boldsymbol{\theta}}$ , we obtain

$$\begin{aligned} \sum_{i=1}^n \psi(\mathbf{u}_i|\hat{\boldsymbol{\theta}}_{[-i]}) &= \sum_{i=1}^n \psi(\mathbf{u}_i|\hat{\boldsymbol{\theta}}) + \sum_{i=1}^n \mathbf{g}(\mathbf{u}_i|\hat{\boldsymbol{\theta}})'(\hat{\boldsymbol{\theta}}_{[-i]} - \hat{\boldsymbol{\theta}}) \\ &\quad + \frac{1}{2} \sum_{i=1}^n (\hat{\boldsymbol{\theta}}_{[-i]} - \hat{\boldsymbol{\theta}})' \mathbf{H}(\mathbf{u}_i|\bar{\boldsymbol{\theta}}_i(\delta_i))(\hat{\boldsymbol{\theta}}_{[-i]} - \hat{\boldsymbol{\theta}}), \end{aligned}$$

where  $\bar{\boldsymbol{\theta}}_i(\delta_i)$  is given by (12). Thus,

$$\mathbf{E}_{\mathbf{y}}^*[\text{CV}] = R_{\text{PD}} + R_1 + \frac{1}{2}R_2. \quad (\text{A5})$$

Consequently,  $R_{\text{PD}} \leq \mathbf{E}_{\mathbf{y}}^*[\text{CV}(1)]$  whenever  $R_1 + R_2/2 \geq 0$ .

## A.2. Expansions of Biases of CV and TIC

Let

$$\mathbf{L}(\mathbf{y}|\boldsymbol{\vartheta}) = \left( \frac{\partial}{\partial \boldsymbol{\theta}} \otimes \frac{\partial^2}{\partial \boldsymbol{\theta} \partial \boldsymbol{\theta}'} \right) \psi(\mathbf{y}|\boldsymbol{\theta}) \Big|_{\boldsymbol{\theta}=\boldsymbol{\vartheta}},$$

and

$$\mathbf{K}(\boldsymbol{\theta}) = \mathbf{E}_{\mathbf{y}}^*[\mathbf{L}(\mathbf{y}|\boldsymbol{\theta})], \quad \hat{\mathbf{K}}(\hat{\boldsymbol{\theta}}) = \frac{1}{n} \sum_{i=1}^n \mathbf{L}(\mathbf{y}_i|\hat{\boldsymbol{\theta}}).$$

Because  $\hat{\boldsymbol{\theta}}_i(\lambda)$  is the minimizer of  $\Psi(\boldsymbol{\theta}|\mathbf{Y}, \mathbf{1}_n - \lambda \mathbf{e}_i)$ , there exists

$$\sum_{j=1}^n \mathbf{g}(\mathbf{y}_j|\hat{\boldsymbol{\theta}}_i(\lambda)) = \lambda \mathbf{g}(\mathbf{y}_i|\hat{\boldsymbol{\theta}}_i(\lambda)), \quad (\text{A6})$$

where  $\mathbf{g}(\cdot|\cdot)$  is given by (7). The following stochastic expansion is needed

$$\hat{\boldsymbol{\theta}}_i(\lambda) = \hat{\boldsymbol{\theta}} + \frac{\lambda}{n} \mathbf{z}_{1,i} + \frac{\lambda^2}{n^2} \mathbf{z}_{2,i} + O_p(n^{-3}), \quad (\text{A7})$$

where  $\lambda = O(1)$  and  $\mathbf{z}_{1,i}$  and  $\mathbf{z}_{2,i}$  are to be determined. Applying the Taylor expansion to both sides of (A6) at  $\hat{\boldsymbol{\theta}}$ , replacing  $\hat{\boldsymbol{\theta}}_i(\lambda)$  by (A7), and comparing the  $O(n^{-1})$  and  $O(n^{-2})$

terms in both sides of the resulting equation in sequence, we obtain

$$\mathbf{z}_{1,i} = \hat{\mathbf{J}}(\hat{\boldsymbol{\theta}})^{-1} \mathbf{g}(\mathbf{y}_i | \hat{\boldsymbol{\theta}}), \quad \mathbf{z}_{2,i} = \hat{\mathbf{J}}(\hat{\boldsymbol{\theta}})^{-1} \left\{ \mathbf{H}(\mathbf{y}_i | \hat{\boldsymbol{\theta}}) \mathbf{z}_{1,i} - \frac{1}{2} \hat{\mathbf{K}}(\hat{\boldsymbol{\theta}}) \text{vec}(\mathbf{z}_{1,i} \mathbf{z}'_{1,i}) \right\}, \quad (\text{A8})$$

where  $\mathbf{H}(\cdot | \cdot)$  and  $\hat{\mathbf{J}}(\hat{\boldsymbol{\theta}})$  are given by (7) and (18), respectively. Notice that  $\hat{\boldsymbol{\theta}}_{[-i]} = \hat{\boldsymbol{\theta}}_i(1)$ . Substituting  $\lambda = 1$  into (A7), we obtain the stochastic expansion of  $\hat{\boldsymbol{\theta}}_{[-i]}$  as

$$\hat{\boldsymbol{\theta}}_{[-i]} = \hat{\boldsymbol{\theta}} + \frac{1}{n} \mathbf{z}_{1,i} + \frac{1}{n^2} \mathbf{z}_{2,i} + O_p(n^{-3}). \quad (\text{A9})$$

In order to calculate the asymptotic expansion of the bias of the CV criterion in (6), we first substitute the stochastic expansion of  $\hat{\boldsymbol{\theta}}_{[-i]}$  in (A9) into  $R_1$  and  $R_2$ , where  $R_1$  and  $R_2$  are given by (11); we then use the relation  $\bar{\boldsymbol{\theta}}_i(\delta_i) \xrightarrow{a.s.} \boldsymbol{\theta}_0$ , where  $\bar{\boldsymbol{\theta}}_i(\delta_i)$  and  $\boldsymbol{\theta}_0$  are given by (12) and (9), respectively. These two steps yield

$$\begin{aligned} R_1 &= \frac{1}{n} \sum_{i=1}^n \mathbf{E}_{\mathbf{y}}^* \left[ \mathbf{r}(\hat{\boldsymbol{\theta}})' \left( \mathbf{z}_{1,i} + \frac{1}{n} \mathbf{z}_{2,i} \right) \right] + O(n^{-2}), \\ R_2 &= \frac{1}{n^2} \sum_{i=1}^n \mathbf{E}_{\mathbf{y}}^* \left[ \mathbf{z}'_{1,i} \mathbf{J}(\boldsymbol{\theta}_0) \mathbf{z}_{1,i} \right] + O(n^{-2}), \end{aligned}$$

where  $\mathbf{r}(\cdot)$  and  $\mathbf{J}(\cdot)$  are given by (8). Notice that  $\sum_{i=1}^n \mathbf{g}(\mathbf{y}_i | \hat{\boldsymbol{\theta}}) = \mathbf{0}_q$  due to  $\hat{\boldsymbol{\theta}}$  being the minimizer of  $\Psi(\boldsymbol{\theta} | \mathbf{Y})$ . Thus,

$$\sum_{i=1}^n \mathbf{E}_{\mathbf{y}}^* \left[ \mathbf{r}(\hat{\boldsymbol{\theta}})' \mathbf{z}_{1,i} \right] = \mathbf{E}_{\mathbf{y}}^* \left[ \mathbf{r}(\hat{\boldsymbol{\theta}})' \hat{\mathbf{J}}(\hat{\boldsymbol{\theta}})^{-1} \sum_{i=1}^n \mathbf{g}(\mathbf{y}_i | \hat{\boldsymbol{\theta}}) \right] = 0.$$

Moreover, from  $\hat{\boldsymbol{\theta}} \xrightarrow{a.s.} \boldsymbol{\theta}_0$  and  $\mathbf{r}(\boldsymbol{\theta}_0) = \mathbf{0}_q$ , the second term in the expansion of  $R_1$  is expanded as

$$\frac{1}{n^2} \sum_{i=1}^n \mathbf{E}_{\mathbf{y}}^* \left[ \mathbf{r}(\hat{\boldsymbol{\theta}})' \mathbf{z}_{2,i} \right] = \frac{1}{n^2} \sum_{i=1}^n \mathbf{E}_{\mathbf{y}}^* \left[ \mathbf{r}(\boldsymbol{\theta}_0)' \mathbf{z}_{2,i} \right] + O(n^{-2}) = O(n^{-2}).$$

Consequently,  $R_1 = O(n^{-2})$ . Using  $\hat{\boldsymbol{\theta}} \xrightarrow{a.s.} \boldsymbol{\theta}_0$  and  $\hat{\mathbf{J}}(\hat{\boldsymbol{\theta}}) \xrightarrow{a.s.} \mathbf{J}(\boldsymbol{\theta}_0)$ ,  $R_2$  is expanded as

$$R_2 = \frac{1}{n^2} \sum_{i=1}^n \mathbf{E}_{\mathbf{y}}^* \left[ \mathbf{g}(\mathbf{y}_i | \boldsymbol{\theta}_0)' \mathbf{J}(\boldsymbol{\theta}_0)^{-1} \mathbf{J}(\boldsymbol{\theta}_0) \mathbf{J}(\boldsymbol{\theta}_0)^{-1} \mathbf{g}(\mathbf{y}_i | \boldsymbol{\theta}_0) \right] + O(n^{-2}) = \frac{1}{n} \gamma_1 + O(n^{-2}),$$

where  $\gamma_1$  is given by (13). Substituting the above two results into (A5) yields the following theorem.

**THEOREM A.1.** *When  $|\beta_{abcd}| < \infty$  holds ( $1 \leq a, b, c, d \leq q$ ), the bias of the CV criterion is expanded as*

$$R_{\text{PD}} - \mathbf{E}_{\mathbf{y}}^*[\text{CV}] = -\frac{1}{2n} \gamma_1 + O(n^{-2}). \quad (\text{A10})$$

Using Theorem A.1, we can easily obtain an expansion of the bias of TIC in (17). Applying the Taylor expansion of CV at  $\hat{\boldsymbol{\theta}}$  yields

$$\text{CV} = \Psi(\hat{\boldsymbol{\theta}}|\mathbf{Y}) + C_1 + \frac{1}{n} \left( C_2 + \frac{1}{2}C_3 \right) + O_p(n^{-2}),$$

where  $C_1$ ,  $C_2$ , and  $C_3$  are given by

$$C_1 = \frac{1}{n} \sum_{i=1}^n \mathbf{g}(\mathbf{y}_i|\hat{\boldsymbol{\theta}})' \mathbf{z}_{1,i}, \quad C_2 = \frac{1}{n} \sum_{i=1}^n \mathbf{g}(\mathbf{y}_i|\hat{\boldsymbol{\theta}})' \mathbf{z}_{2,i}, \quad C_3 = \frac{1}{n} \sum_{i=1}^n \mathbf{z}'_{1,i} \mathbf{H}(\mathbf{y}_i|\hat{\boldsymbol{\theta}})' \mathbf{z}_{1,i}. \quad (\text{A11})$$

Notice that  $C_1 = \text{tr}\{\hat{\mathbf{I}}(\hat{\boldsymbol{\theta}})\hat{\mathbf{J}}(\hat{\boldsymbol{\theta}})^{-1}\}$  and  $\text{TIC} = \Psi(\hat{\boldsymbol{\theta}}|\mathbf{Y}) + \text{tr}\{\hat{\mathbf{I}}(\hat{\boldsymbol{\theta}})\hat{\mathbf{J}}(\hat{\boldsymbol{\theta}})^{-1}\}$ . Thus,

$$\mathbf{E}_{\mathbf{y}}^*[\text{TIC}] = \mathbf{E}_{\mathbf{y}}^*[\text{CV}] - \frac{1}{n} \left\{ \mathbf{E}_{\mathbf{y}}^*[C_2] + \frac{1}{2}\mathbf{E}_{\mathbf{y}}^*[C_3] \right\} + O(n^{-2}).$$

By using  $\hat{\boldsymbol{\theta}} \xrightarrow{a.s.} \boldsymbol{\theta}_0$ ,  $\hat{\mathbf{J}}(\hat{\boldsymbol{\theta}}) \xrightarrow{a.s.} \mathbf{J}(\boldsymbol{\theta}_0)$ , and  $\hat{\mathbf{K}}(\hat{\boldsymbol{\theta}}) \xrightarrow{a.s.} \mathbf{K}(\boldsymbol{\theta}_0)$ , we obtain

$$\begin{aligned} \mathbf{E}_{\mathbf{y}}^*[C_2] &= \frac{1}{n} \sum_{i=1}^n \left\{ \mathbf{E}_{\mathbf{y}}^* \left[ \mathbf{g}(\mathbf{y}_i|\hat{\boldsymbol{\theta}})' \hat{\mathbf{J}}(\hat{\boldsymbol{\theta}})^{-1} \mathbf{H}(\mathbf{y}_i|\hat{\boldsymbol{\theta}}) \hat{\mathbf{J}}(\hat{\boldsymbol{\theta}})^{-1} \mathbf{g}(\mathbf{y}_i|\hat{\boldsymbol{\theta}}) \right] \right. \\ &\quad \left. - \frac{1}{2} \mathbf{E}_{\mathbf{y}}^* \left[ \mathbf{g}(\mathbf{y}_i|\hat{\boldsymbol{\theta}})' \hat{\mathbf{J}}(\hat{\boldsymbol{\theta}})^{-1} \hat{\mathbf{K}}(\hat{\boldsymbol{\theta}}) \left\{ [\hat{\mathbf{J}}(\hat{\boldsymbol{\theta}})^{-1} \mathbf{g}(\mathbf{y}_i|\hat{\boldsymbol{\theta}})] \otimes [\hat{\mathbf{J}}(\hat{\boldsymbol{\theta}})^{-1} \mathbf{g}(\mathbf{y}_i|\hat{\boldsymbol{\theta}})] \right\} \right] \right\} \\ &= \gamma_2 - \frac{1}{2}\gamma_3 + O(n^{-1}), \\ \mathbf{E}_{\mathbf{y}}^*[C_3] &= \frac{1}{n} \sum_{i=1}^n \mathbf{E}_{\mathbf{y}}^* \left[ \mathbf{g}(\mathbf{y}_i|\hat{\boldsymbol{\theta}})' \hat{\mathbf{J}}(\hat{\boldsymbol{\theta}})^{-1} \mathbf{H}(\mathbf{y}_i|\hat{\boldsymbol{\theta}}) \hat{\mathbf{J}}(\hat{\boldsymbol{\theta}})^{-1} \mathbf{g}(\mathbf{y}_i|\hat{\boldsymbol{\theta}}) \right] \\ &= \gamma_2 + O(n^{-1}), \end{aligned}$$

where

$$\begin{aligned} \gamma_2 &= \mathbf{E}_{\mathbf{y}}^* \left[ \mathbf{g}(\mathbf{y}|\boldsymbol{\theta}_0)' \mathbf{J}(\boldsymbol{\theta}_0)^{-1} \mathbf{H}(\mathbf{y}|\boldsymbol{\theta}_0) \mathbf{J}(\boldsymbol{\theta}_0)^{-1} \mathbf{g}(\mathbf{y}|\boldsymbol{\theta}_0) \right], \\ \gamma_3 &= \mathbf{E}_{\mathbf{y}}^* \left[ \mathbf{g}(\mathbf{y}|\boldsymbol{\theta}_0)' \mathbf{J}(\boldsymbol{\theta}_0)^{-1} \mathbf{K}(\boldsymbol{\theta}_0) \left\{ [\mathbf{J}(\boldsymbol{\theta}_0)^{-1} \mathbf{g}(\mathbf{y}|\boldsymbol{\theta}_0)] \otimes [\mathbf{J}(\boldsymbol{\theta}_0)^{-1} \mathbf{g}(\mathbf{y}|\boldsymbol{\theta}_0)] \right\} \right]. \end{aligned}$$

Thus,

$$\mathbf{E}_{\mathbf{y}}^*[\text{TIC}] = \mathbf{E}_{\mathbf{y}}^*[\text{CV}] - \frac{1}{2n} (3\gamma_2 - \gamma_3) + O(n^{-2}). \quad (\text{A12})$$

Equations (A10) and (A12) lead to the following theorem.

**THEOREM A.2.** *When  $|\beta_{abcd}| < \infty$  holds ( $1 \leq a, b, c, d \leq q$ ), the bias of TIC is expanded as*

$$R_{\text{PD}} - \mathbf{E}_{\mathbf{y}}^*[\text{TIC}] = -\frac{1}{2n} (\gamma_1 - 3\gamma_2 + \gamma_3) + O(n^{-2}).$$

### A.3. Proof of Example 1

The discrepancy function corresponding to the multivariate normal distribution is

$$\psi(\mathbf{y}|\boldsymbol{\theta}) = \frac{1}{2} \left\{ p \log(2\pi) + \log |\boldsymbol{\Sigma}| + (\mathbf{y} - \boldsymbol{\mu})' \boldsymbol{\Sigma}^{-1} (\mathbf{y} - \boldsymbol{\mu}) \right\},$$

where  $\boldsymbol{\theta} = (\boldsymbol{\mu}', \text{vech}(\boldsymbol{\Sigma})')$  with  $\text{vech}(\mathbf{A})$  being the vector of stacking the distinct elements of a symmetric matrix  $\mathbf{A}$  columnwise. Let  $\mathbf{D}_p$  be the duplication matrix such that  $\text{vec}(\mathbf{A}) = \mathbf{D}_p \text{vech}(\mathbf{A})$  (see Magnus & Neudecker, 1999, p. 48). Then, the corresponding  $\mathbf{r}(\boldsymbol{\theta})$  in (8) is given by

$$\mathbf{r}(\boldsymbol{\theta}) = \frac{1}{2} \begin{pmatrix} \mathbf{I}_p & \mathbf{D}_p' \mathbf{D}_p \end{pmatrix} \begin{pmatrix} 2\boldsymbol{\Sigma}^{-1}(\boldsymbol{\mu} - \boldsymbol{\mu}^*) \\ \text{vech}(\boldsymbol{\Sigma}^{-1} - \boldsymbol{\Sigma}^{-1}\{\boldsymbol{\Sigma}^* + (\boldsymbol{\mu} - \boldsymbol{\mu}^*)(\boldsymbol{\mu} - \boldsymbol{\mu}^*)'\}\boldsymbol{\Sigma}^{-1}) \end{pmatrix}.$$

It is well known that the MLE of  $\boldsymbol{\theta}$  is  $\hat{\boldsymbol{\theta}} = (\bar{\mathbf{y}}', \text{vech}(\mathbf{S})')$ , where  $\bar{\mathbf{y}}$  and  $\mathbf{S}$  are the sample mean and covariance matrix given by (23). On the other hand, the  $i$ th jackknife estimator of  $\boldsymbol{\theta}$  is  $\hat{\boldsymbol{\theta}}_{[-i]} = (\bar{\mathbf{y}}'_{[-i]}, \text{vech}(\mathbf{S}_{[-i]}'))'$ , where  $\bar{\mathbf{y}}_{[-i]} = (n-1)^{-1} \sum_{j \neq i}^n \mathbf{y}_j$  and  $\mathbf{S}_{[-i]} = (n-1)^{-1} \sum_{j \neq i}^n (\mathbf{y}_j - \bar{\mathbf{y}}_{[-i]})(\mathbf{y}_j - \bar{\mathbf{y}}_{[-i]})'$ . Fujikoshi *et al.* (2003) gives

$$\bar{\mathbf{y}}_{[-i]} = \bar{\mathbf{y}} - \frac{1}{n-1}(\mathbf{y}_i - \bar{\mathbf{y}}), \quad \mathbf{S}_{[-i]} = \frac{n}{n-1} \left\{ \mathbf{S} - \frac{1}{n-1}(\mathbf{y}_i - \bar{\mathbf{y}})(\mathbf{y}_i - \bar{\mathbf{y}})' \right\}.$$

Therefore,  $\hat{\boldsymbol{\theta}}_{[-i]} - \hat{\boldsymbol{\theta}}$  becomes

$$\hat{\boldsymbol{\theta}}_{[-i]} - \hat{\boldsymbol{\theta}} = \frac{1}{n-1} \begin{pmatrix} -(\mathbf{y}_i - \bar{\mathbf{y}}) \\ \text{vech}(\mathbf{S}) - \frac{n}{n-1} \text{vech}((\mathbf{y}_i - \bar{\mathbf{y}})(\mathbf{y}_i - \bar{\mathbf{y}})') \end{pmatrix}.$$

Notice that  $\text{vech}(\mathbf{A})' \mathbf{D}_p' \mathbf{D}_p \text{vech}(\mathbf{B}) = \text{tr}(\mathbf{B}'\mathbf{A}) = \text{tr}(\mathbf{A}'\mathbf{B})$  for symmetric matrices  $\mathbf{A}$  and  $\mathbf{B}$ . It follows from the definition of  $R_1$  in (11) that

$$R_1 = \frac{n}{2(n-1)^2} \left\{ \mathbb{E}_{\mathbf{y}}^* [\text{tr}(\boldsymbol{\Sigma}^* \mathbf{S}^{-1})] + \mathbb{E}_{\mathbf{y}}^* [(\bar{\mathbf{y}} - \boldsymbol{\mu}^*)' \mathbf{S}^{-1} (\bar{\mathbf{y}} - \boldsymbol{\mu}^*)] - p \right\}.$$

Jensen's inequality implies  $\mathbb{E}_{\mathbf{y}}^* [\text{tr}(\boldsymbol{\Sigma}^* \mathbf{S}^{-1})] \geq np/(n-1) > p$ , and thus  $R_1 > 0$ .

### A.4. Proof of Theorem 2

It follows from (A9) and (A7) that

$$\hat{\boldsymbol{\theta}}_i(\lambda) = \hat{\boldsymbol{\theta}}_{[-i]} + \frac{1}{n}(\lambda - 1)\mathbf{z}_{1,i} + \frac{1}{n^2}(\lambda^2 - 1)\mathbf{z}_{2,i} + O_p(n^{-3}), \quad (\lambda = O(1)), \quad (\text{A13})$$

where  $\mathbf{z}_{1,i}$  and  $\mathbf{z}_{2,i}$  are given by (A8). Using the expansion (A13) after applying the Taylor expansion of  $\text{CV}(\lambda)$  at  $\hat{\boldsymbol{\theta}}_{[-i]}$  yields

$$\text{CV}(\lambda) = \text{CV} + (\lambda - 1)C_1 + \frac{1}{n} \left\{ (\lambda^2 - 1)C_2 + \frac{1}{2}(\lambda - 1)^2 C_3 \right\} + O_p(n^{-3}),$$

where  $C_1$ ,  $C_2$  and  $C_3$  are given by (A11). Notice that  $\mathbb{E}_{\mathbf{y}}^*[C_1] = O(1)$ ,  $\mathbb{E}_{\mathbf{y}}^*[C_2] = O(1)$  and  $\mathbb{E}_{\mathbf{y}}^*[C_3] = O(1)$ . Using the expansion of  $\mathbb{E}_{\mathbf{y}}^*[\text{CV}]$  in (A10), we obtain

$$\begin{aligned} \mathbb{E}_{\mathbf{y}}^*[\text{CV}(\lambda)] &= R_{\text{PD}} + (\lambda - 1)\mathbb{E}_{\mathbf{y}}^*[C_1] \\ &\quad + \frac{1}{2n} \left\{ \gamma_1 + 2(\lambda^2 - 1)\mathbb{E}_{\mathbf{y}}^*[C_2] + (\lambda - 1)^2 \mathbb{E}_{\mathbf{y}}^*[C_3] \right\} + O(n^{-2}), \end{aligned}$$

where  $\gamma_1$  is given by (13). The first equation in (14) follows by noticing that  $\mathbb{E}_{\mathbf{y}}^*[C_1] = \gamma_1 + O(n^{-1})$ . If  $\lambda = 1 + O(n^{-1})$ , then  $\lambda - 1 = O(n^{-1})$  and  $\lambda^2 - 1 = O(n^{-1})$ . Consequently,

$$\mathbb{E}_{\mathbf{y}}^*[\text{CV}(\lambda)] = R_{\text{PD}} + \frac{1}{2n} \left\{ \gamma_1 + 2n(\lambda - 1)\mathbb{E}_{\mathbf{y}}^*[C_1] \right\} + O(n^{-2}),$$

and from which the second equation in (14) follows.

### A.5. Derivation of Redefining EIC

Notice that the  $b$ th bootstrap resample  $\mathbf{y}_{b,i}^*$  ( $i = 1, \dots, n$ ) is one of  $\mathbf{y}_1, \dots, \mathbf{y}_n$ . Let  $\mathbf{d}_b = (d_{b1}, d_{b2}, \dots, d_{bn})$  with  $d_{bi}$  equal to the number of times  $\mathbf{y}_i$  appears in the  $b$ th bootstrap resample. Then

$$\sum_{i=1}^n \psi(\mathbf{y}_{b,i}^* | \boldsymbol{\theta}) = \sum_{i=1}^n d_{bi} \psi(\mathbf{y}_i | \boldsymbol{\theta}) = \Psi(\boldsymbol{\theta} | \mathbf{Y}, \mathbf{d}_b).$$

Thus,  $\hat{\boldsymbol{\theta}}_b^* = \hat{\boldsymbol{\theta}}(\mathbf{d}_b)$ . Consequently,

$$\begin{aligned} \sum_{i=1}^n \psi(\mathbf{y}_i | \hat{\boldsymbol{\theta}}_b^*) - \sum_{i=1}^n \psi(\mathbf{y}_{b,i}^* | \hat{\boldsymbol{\theta}}_b^*) &= \sum_{i=1}^n \psi(\mathbf{y}_i | \hat{\boldsymbol{\theta}}(\mathbf{d}_b)) - \sum_{i=1}^n d_{bi} \psi(\mathbf{y}_i | \hat{\boldsymbol{\theta}}(\mathbf{d}_b)) \\ &= \sum_{i=1}^n (1 - d_{bi}) \psi(\mathbf{y}_i | \hat{\boldsymbol{\theta}}(\mathbf{d}_b)) \\ &= \Psi(\hat{\boldsymbol{\theta}}(\mathbf{d}_b) | \mathbf{Y}, \mathbf{1}_n - \mathbf{d}_b). \end{aligned} \tag{A14}$$

Substituting (A14) into (19) yields equation (20). The distribution property  $\mathbf{d}_b \sim \text{Multi}_n(n; 1/n, \dots, 1/n)$  is the definition of the bootstrap sampling.

### A.6. Proof of Corollary 1

For the discrepancy function given by (22), we obtain by direct calculation

$$\begin{aligned}\frac{\partial}{\partial \boldsymbol{\mu}} \Psi(\boldsymbol{\theta} | \mathbf{Y}, \mathbf{1}_n - \lambda \mathbf{e}_i) &= -2 \boldsymbol{\Sigma}(\boldsymbol{\xi})^{-1} \{(n - \lambda) \boldsymbol{\mu} - (n \bar{\mathbf{y}} - \lambda \mathbf{y}_i)\}, \\ \frac{\partial}{\partial \boldsymbol{\xi}} \Psi(\boldsymbol{\theta} | \mathbf{Y}, \mathbf{1}_n - \lambda \mathbf{e}_i) &= \frac{\partial}{\partial \boldsymbol{\xi}} (n - \lambda) F(\boldsymbol{\xi} | \mathbf{M}_i(\boldsymbol{\mu}, \lambda)),\end{aligned}$$

where  $F(\boldsymbol{\xi} | \cdot)$  is given by (24) and

$$\mathbf{M}_i(\boldsymbol{\mu}, \lambda) = \frac{1}{n - \lambda} \left\{ \sum_{j=1}^n (\mathbf{y}_j - \boldsymbol{\mu})(\mathbf{y}_j - \boldsymbol{\mu})' - \lambda (\mathbf{y}_i - \boldsymbol{\mu})(\mathbf{y}_i - \boldsymbol{\mu})' \right\}.$$

Denote  $\hat{\boldsymbol{\theta}}_i(\lambda) = (\hat{\boldsymbol{\mu}}_i(\lambda)', \hat{\boldsymbol{\xi}}_i(\lambda)')'$ . Solving the equation  $\partial \Psi(\boldsymbol{\theta} | \mathbf{Y}, \mathbf{1}_n - \lambda \mathbf{e}_i) / \partial \boldsymbol{\mu} = \mathbf{0}_p$  leads to

$$\hat{\boldsymbol{\mu}}_i(\lambda) = \bar{\mathbf{y}} - \frac{\lambda}{n - \lambda} (\mathbf{y}_i - \bar{\mathbf{y}}).$$

Notice that

$$\mathbf{y}_j - \hat{\boldsymbol{\mu}}_i(\lambda) = \begin{cases} \frac{n}{n - \lambda} (\mathbf{y}_i - \bar{\mathbf{y}}) & (j = i) \\ \mathbf{y}_j - \bar{\mathbf{y}} + \frac{\lambda}{n - \lambda} (\mathbf{y}_i - \bar{\mathbf{y}}) & (j \neq i) \end{cases}. \quad (\text{A15})$$

It follows from (A15) that  $\mathbf{M}_i(\hat{\boldsymbol{\mu}}_i(\lambda), \lambda) = \mathbf{S}_i(\lambda)$ , where  $\mathbf{S}_i(\lambda)$  is given by (27). Equation (22) implies

$$\text{CV}(\lambda) = np \log(2\pi) + \sum_{i=1}^n \left\{ \log |\boldsymbol{\Sigma}(\hat{\boldsymbol{\xi}}_i(\lambda))| + (\mathbf{y}_i - \hat{\boldsymbol{\mu}}_i(\lambda))' \boldsymbol{\Sigma}(\hat{\boldsymbol{\xi}}_i(\lambda))^{-1} (\mathbf{y}_i - \hat{\boldsymbol{\mu}}_i(\lambda)) \right\},$$

where  $\hat{\boldsymbol{\xi}}_i(\lambda)$  is given by (26). Substituting (A15) into the above equation yields (25).

### A.7. Derivation of EIC in Corollary 2

Notice that  $\hat{\boldsymbol{\theta}}(\mathbf{d}_b)$  is the minimizer of  $\Psi(\boldsymbol{\theta} | \mathbf{Y}, \mathbf{d}_b)$  and  $\mathbf{1}'_n \mathbf{d}_b = n$ . With the discrepancy function given by (22), by direct calculations we obtain

$$\begin{aligned}\frac{\partial}{\partial \boldsymbol{\mu}} \Psi(\boldsymbol{\theta} | \mathbf{Y}, \mathbf{d}_b) &= -2 \boldsymbol{\Sigma}(\boldsymbol{\xi})^{-1} (\mathbf{Y}' \mathbf{d}_b - n \boldsymbol{\mu}), \\ \frac{\partial}{\partial \boldsymbol{\xi}} \Psi(\boldsymbol{\theta} | \mathbf{Y}, \mathbf{d}_b) &= \frac{\partial}{\partial \boldsymbol{\xi}} (n - \lambda) F(\boldsymbol{\xi} | \mathbf{M}(\boldsymbol{\mu}, \mathbf{d}_b)),\end{aligned}$$

where  $F(\boldsymbol{\xi} | \mathbf{S})$  is given by (24) and

$$\mathbf{M}(\boldsymbol{\mu}, \mathbf{d}_b) = \frac{1}{n} (\mathbf{Y} - \mathbf{1}_n \boldsymbol{\mu}')' \text{diag}(\mathbf{d}_b) (\mathbf{Y} - \mathbf{1}_n \boldsymbol{\mu}').$$

Denote  $\hat{\boldsymbol{\theta}}(\mathbf{d}_b) = (\hat{\boldsymbol{\mu}}(\mathbf{d}_b)', \hat{\boldsymbol{\xi}}(\mathbf{d}_b)')'$ . Solving the equation  $\partial \Psi(\boldsymbol{\theta} | \mathbf{Y}, \mathbf{d}_b) / \partial \boldsymbol{\mu} = \mathbf{0}_p$  leads to

$$\hat{\boldsymbol{\mu}}(\mathbf{d}_b) = \frac{1}{n} \mathbf{Y}' \mathbf{d}_b. \quad (\text{A16})$$

Substituting (A16) into  $\mathbf{M}(\boldsymbol{\mu}, \mathbf{d}_b)$  yields  $\mathbf{M}_i(\hat{\boldsymbol{\mu}}(\mathbf{d}_b), \mathbf{d}_b) = \mathbf{S}(\mathbf{d}_b)$ , where  $\mathbf{S}(\mathbf{d}_b)$  is given by (32). Notice that the EIC under the candidate model  $M$  in (21) is

$$\text{EIC} = nF(\hat{\boldsymbol{\xi}}|\mathbf{S}) + np \log(2\pi) + \frac{n}{B} \sum_{i=1}^n \text{tr} \left\{ \mathbf{M}(\hat{\boldsymbol{\mu}}(\mathbf{d}_b), \mathbf{1}_n - \mathbf{d}_b) \boldsymbol{\Sigma}(\hat{\boldsymbol{\xi}}(\mathbf{d}_b))^{-1} \right\},$$

where  $\hat{\boldsymbol{\xi}}$  is given by (28). Substituting (A16) into  $\mathbf{M}(\boldsymbol{\mu}, \mathbf{1}_n - \mathbf{d}_b)$  yields  $\mathbf{M}(\hat{\boldsymbol{\mu}}(\mathbf{d}_b), \mathbf{1}_n - \mathbf{d}_b) = \mathbf{V}(\mathbf{d}_b)$  in (31), which further leads to (29).

## References

- [1] Akaike, H. (1973). Information theory and an extension of the maximum likelihood principle. In *2nd. International Symposium on Information Theory* (Eds. B. N. Petrov & F. Csáki), pp. 267–281. Akadémiai Kiadó, Budapest.
- [2] Akaike, H. (1974). A new look at the statistical model identification. *IEEE Trans. Automatic Control*, **AC-19**, 716–723.
- [3] Akaike, H. (1987). Factor analysis and AIC. *Psychometrika*, **52**, 317–332.
- [4] Balakrishnan, N., & Ambagaspitiya, R. S. (1994). On skew Laplace distribution. *Technical Report, Department of Mathematics & Statistics, McMaster University*, Hamilton, Ontario, Canada.
- [5] Basu, A., Harris, I. R., Hjort, N. L., & Jones, M. C. (1998). Robust and efficient estimation by minimising a density power divergence. *Biometrika*, **85**, 549–559.
- [6] Bentler, P. M. (1995). *EQS Structural Equation Program Manual*. Multivariate Software, Inc, Encino, CA.
- [7] Bozdogan, H. (1987). Model selection and Akaike’s information criterion (AIC): the general theory and its analytical extensions. *Psychometrika*, **52**, 345–370.
- [8] Fujikoshi, Y., Noguchi, T., Ohtaki, M., & Yanagihara, H. (2003). Corrected versions of cross-validation criteria for selecting multivariate regression and growth curve models. *Ann. Inst. Statist. Math.*, **55**, 537–553.
- [9] Fujisawa, H., & Eguchi, S. (2006). Robust estimation in the normal mixture model. *J. Statist. Plann. Inference*, **136**, 3989–4011.

- [10] Hall, P. (1987). Edgeworth expansion for student's  $t$  statistic under minimal moment condition. *Ann. Probab.*, **15**, 920–931.
- [11] Ishiguro, M., Sakamoto, Y., & Kitagawa, G. (1997). Bootstrapping log likelihood and EIC, an extension of AIC. *Ann. Inst. Stat. Math.*, **49**, 411–434.
- [12] Konishi, S. (1999). Statistical model evaluation and information criteria. In *Multivariate Analysis, Design of Experiments, and Survey Sampling* (Ed. S. Ghosh), pp. 369–399. Marcel Dekker, New York.
- [13] Kullback, S., & Leibler, R. A. (1951). On information and sufficiency. *Ann. Math. Statistics*, **22**, 79–86.
- [14] Lee, S.-Y., & Kontoghiorghes, E. J. (Eds.). (2007). *Handbook of Latent Variable and Related Models*. Elsevier, Amsterdam.
- [15] Lehmann, E. L., & Casella, G. (1998). *Theory of Point Estimation* (2nd. edition). Springer-Verlag, New York.
- [16] Magnus, J. R., & Neudecker, H. (1999). *Matrix Differential Calculus with Applications in Statistics and Econometrics* (revised edition). John Wiley & Sons, New York.
- [17] Mardia, K. V. (1970). Measures of multivariate skewness and kurtosis with applications. *Biometrika*, **57**, 519–530.
- [18] Schwarz, G. (1978). Estimating the dimension of a model. *Ann. Statist.*, **6**, 461–464.
- [19] Stone, M. (1974). Cross-validatory choice and assessment of statistical predictions. *J. Roy. Statist. Soc. Ser. B*, **36**, 111–147.
- [20] Stone, M. (1977). An asymptotic equivalence of choice of model by cross-validation and Akaike's criterion. *J. Roy. Statist. Soc. Ser. B*, **39**, 44–47.
- [21] Takeuchi, K. (1976). Distribution of information statistics and criteria for adequacy of models. *Math. Sci.*, **153**, 12–18 (in Japanese).
- [22] White, H. (1982). Maximum likelihood estimation of misspecified models, *Econometrica*, **50**, 1–25.



- [23] Yanagihara, H. (2005). 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**, 133–157 (in Japanese).
- [24] Yanagihara, H. (2007). A family of estimators for multivariate kurtosis in a nonnormal linear regression model. *J. Multivariate Anal.*, **98**, 1–29.
- [25] Yanagihara, H., & Yuan, K.-H. (2005). Four improved statistics for contrasting means by correcting skewness and kurtosis. *British J. Math. Statist. Psych.*, **58**, 209–237.
- [26] Yanagihara, H., Tonda, T., & Matsumoto, C. (2006). Bias correction of cross-validation criterion based on Kullback-Leibler information under a general condition. *J. Multivariate Anal.*, **97**, 1965–1975.
- [27] Yanagihara, H., Kamo, K., & Tonda, T. (2006). Second-order bias-corrected AIC in multivariate normal linear models under nonnormality. *TR No. 06-01, Statistical Research Group, Hiroshima University, Hiroshima, Japan.*
- [28] Yuan, K.-H., & Bentler, P. M. (2007). Structural equation modeling. In *Handbook of Statistics 26: Psychometrics* (Eds. C. R. Rao & S. Sinharay), pp. 297–358. North-Holland, Amsterdam.

TABLE 1. Risk of each candidate model

True Distribution	Model 1	Model 2	Model 3
1	466.7	464.7	461.5
2	468.9	467.1	470.6
3	466.0	463.9	457.6
4	469.1	467.2	470.7
5	471.0	469.2	480.3
6	475.7	474.0	493.8

TABLE 2. Relative biases, RMSEs, and selection frequencies of five information criteria

True Dist.	Criterion	Model 1			Model 2			Model 3		
		Bias	(RMSE	Freq.)	Bias	(RMSE	Freq.)	Bias	(RMSE	Freq.)
1	AIC	0.52	(4.3	0.0)	1.13	(4.4	43.1)	0.99	(4.2	56.9)
	TIC	-0.51	(4.7	1.6)	0.12	(4.8	35.7)	0.62	(4.4	62.7)
	EIC	0.05	(4.3	11.4)	0.08	(4.4	31.6)	0.16	(4.2	57.0)
	CV	-0.08	(4.3	13.5)	-0.10	(4.4	31.8)	-0.16	(4.2	54.6)
	CCV	0.02	(4.3	12.4)	0.02	(4.4	31.5)	0.05	(4.2	56.1)
2	AIC	1.40	(6.3	0.0)	2.04	(6.6	37.2)	4.21	(7.0	62.8)
	TIC	-2.71	(8.8	1.7)	-2.04	(8.8	26.3)	2.25	(6.6	72.0)
	EIC	0.30	(6.5	19.6)	0.36	(6.7	38.6)	1.08	(6.4	41.7)
	CV	-0.04	(6.7	24.4)	-0.07	(6.9	40.2)	-0.19	(7.0	35.4)
	CCV	0.09	(6.7	22.8)	0.09	(6.9	40.0)	0.22	(6.8	37.2)
3	AIC	0.16	(3.1	0.0)	0.74	(3.2	47.0)	-0.50	(3.2	53.0)
	TIC	0.68	(3.1	1.2)	1.29	(3.3	43.8)	0.11	(3.2	55.1)
	EIC	-0.02	(3.0	4.4)	-0.02	(3.1	22.7)	-0.20	(3.1	72.8)
	CV	-0.08	(3.0	3.7)	-0.11	(3.1	20.9)	-0.15	(3.1	75.3)
	CCV	0.01	(3.0	3.3)	0.00	(3.1	20.7)	-0.01	(3.1	76.0)
4	AIC	1.39	(6.4	0.0)	2.02	(6.6	37.8)	4.16	(6.9	62.2)
	TIC	-2.71	(9.6	1.8)	-2.08	(9.5	26.8)	2.37	(6.8	71.5)
	EIC	0.27	(6.8	19.7)	0.30	(6.9	38.5)	1.02	(6.7	41.9)
	CV	-0.11	(7.1	24.7)	-0.15	(7.3	39.3)	-0.33	(7.6	36.0)
	CCV	0.03	(7.0	22.9)	0.02	(7.2	39.2)	0.12	(7.3	37.8)
5	AIC	2.06	( 7.7	0.0)	2.74	( 8.0	32.2)	7.12	( 9.6	67.8)
	TIC	-4.34	(12.8	1.6)	-3.69	(12.7	19.6)	4.46	( 8.8	78.8)
	EIC	0.37	( 8.4	24.2)	0.46	( 8.5	41.3)	1.73	( 8.7	34.5)
	CV	-0.25	( 9.1	29.4)	-0.25	( 9.4	41.1)	-0.55	(10.9	29.5)
	CCV	-0.07	( 9.0	27.2)	-0.04	( 9.2	41.2)	0.26	(10.0	31.6)
6	AIC	4.07	(10.5	0.0)	4.75	(10.9	27.1)	11.65	(13.8	72.9)
	TIC	-6.10	(20.0	1.1)	-5.45	(19.9	13.4)	8.05	(12.4	85.5)
	EIC	1.22	(12.9	25.6)	1.29	(13.2	42.1)	3.72	(14.4	32.3)
	CV	-0.21	(16.5	30.6)	-0.26	(16.9	42.5)	-0.77	(21.9	26.9)
	CCV	0.13	(15.5	28.7)	0.13	(15.8	42.3)	1.17	(15.9	29.0)

TABLE 3. Averages of frequencies of choosing the model having the smallest risk

Criterion	AIC	TIC	EIC	CV	CCV
Average of Frequencies (%)	40.7	34.0	48.4	48.8	49.1

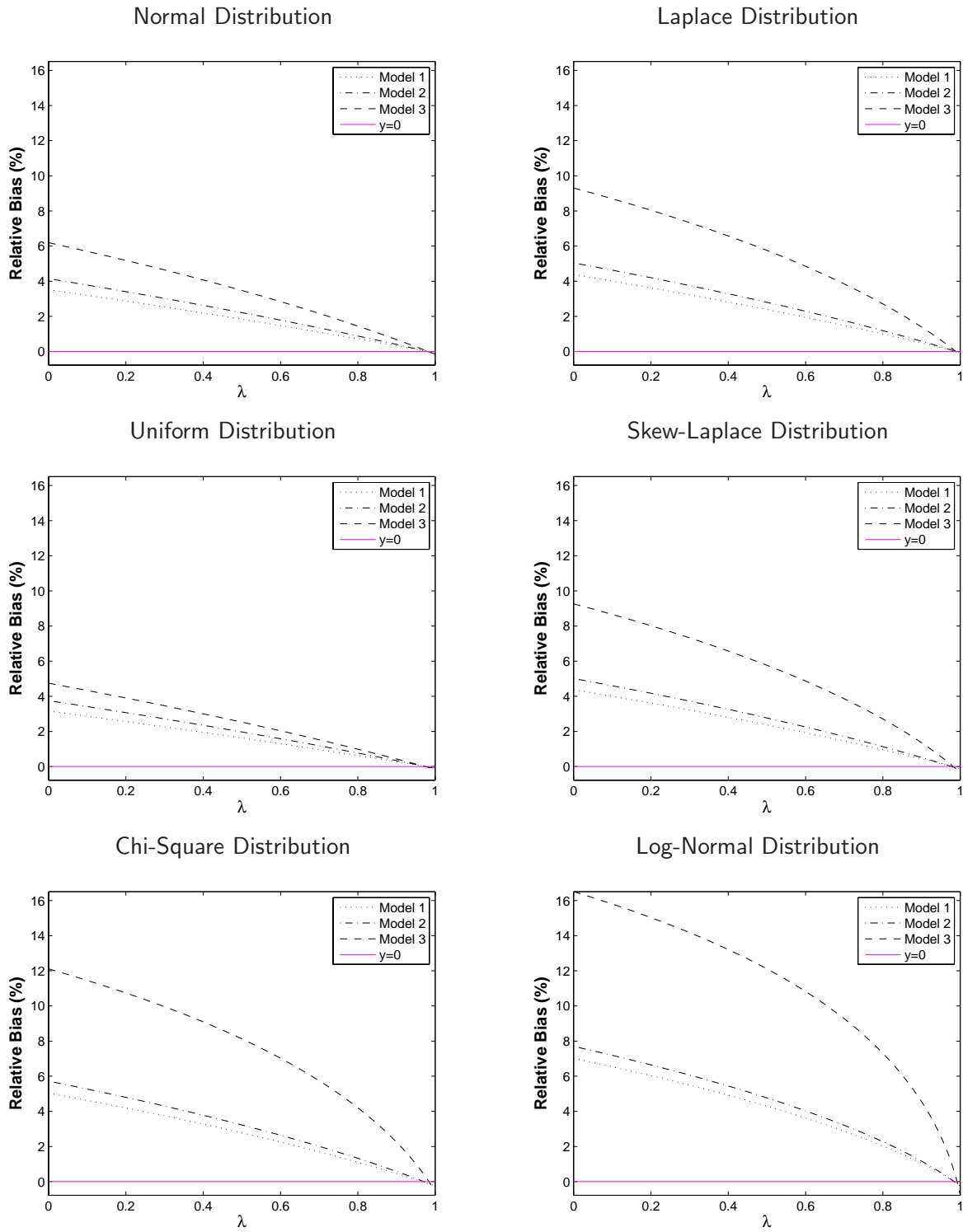


FIGURE 1. Relative biases of cross-validators model selection criteria

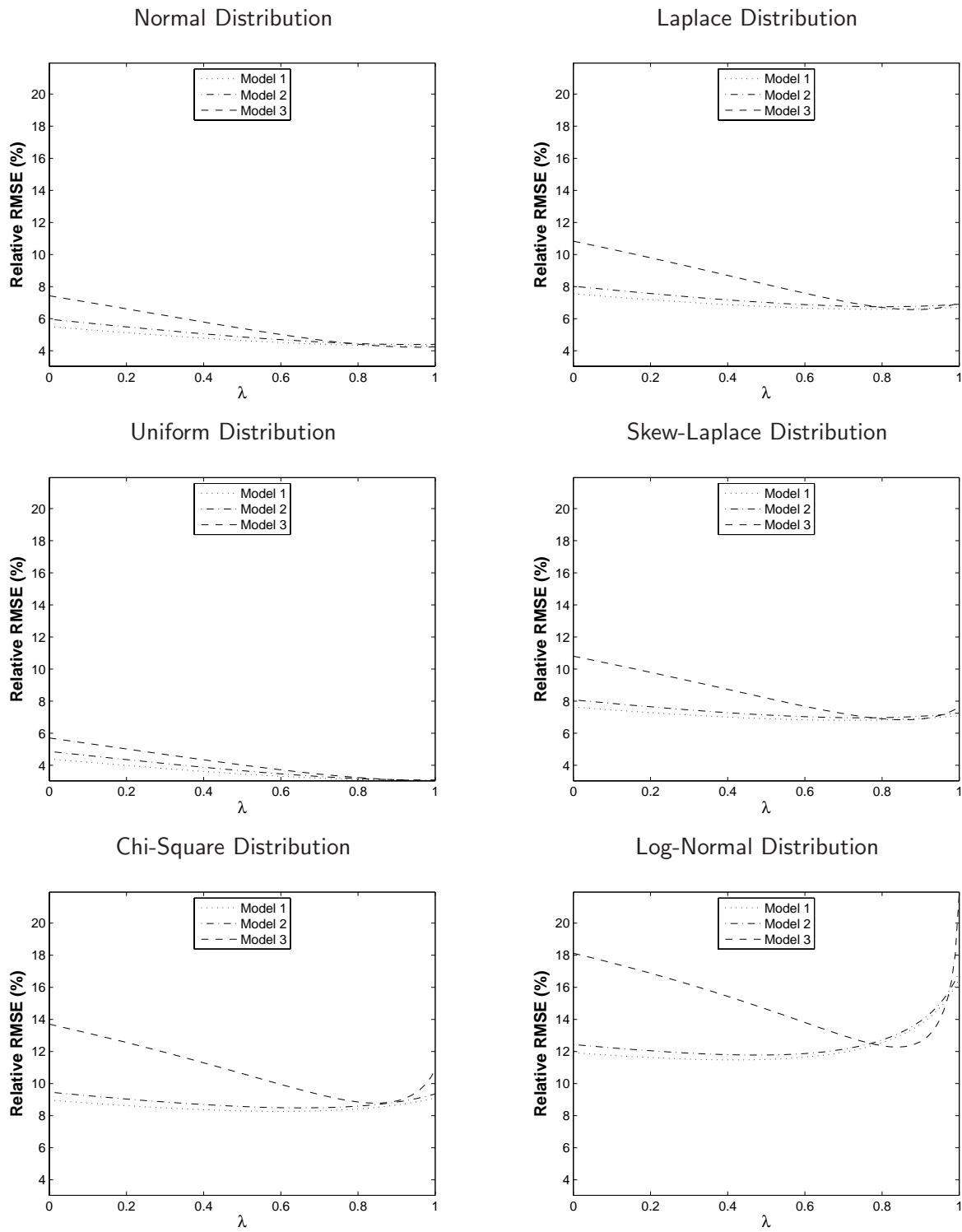


FIGURE 2. Relative RMSEs of cross-validated model selection criteria

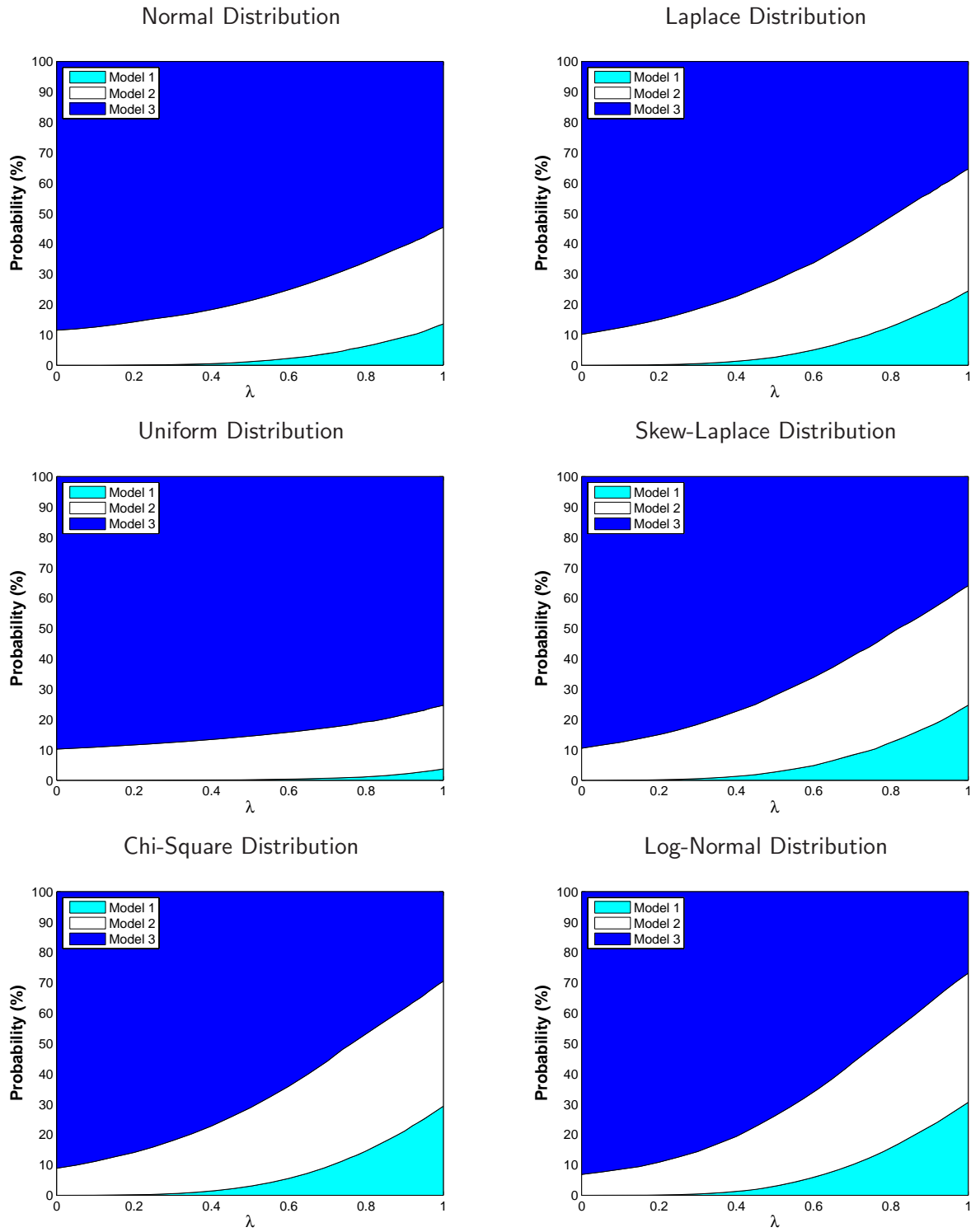


FIGURE 3. Selection frequencies of cross-validators model selection criteria