

博士論文

Variable selection problem in mixed effects models  
with application to small area estimation

(混合効果モデルにおける変数選択問題と  
小地域推定への応用)

Yuki KAWAKUBO

川久保友超



# Abstract

This thesis studies the variable selection problem in mixed effects models, especially the method of information criteria is focused on. The conditional AIC proposed by Vaida and Blanchard (2005) is more appropriate for the focus on clusters than the conventional AIC is. This is useful in problems involving prediction of random effects such as small area estimation. Concerning the conditional AIC and the related fields, several problems are addressed and some new results are obtained. Firstly, the conditional AIC is modified for the underspecified model, which does not include the true model. Secondly, the variable selection problem in linear mixed model under covariate shift situation is considered. It is also shown that considering the covariate shift situation is meaningful in small area estimation problem. Thirdly, the conditional AIC in nonlinear mixed models based on natural exponential family is derived. Lastly, some variants of the AIC and the conditional AIC are proposed and their properties are discussed, which are also related to the Bayesian procedure as well as frequentists' methods.



# Contents

<b>1</b>	<b>Introduction</b>	<b>9</b>
<b>2</b>	<b>Mixed effects models and conditional AIC</b>	<b>11</b>
2.1	Mixed effects models . . . . .	11
2.1.1	General . . . . .	11
2.1.2	Linear mixed model . . . . .	11
2.2	Setup of variable selection problem . . . . .	12
2.2.1	Candidate models, full model and true model . . . . .	12
2.2.2	Overspecified and underspecified . . . . .	12
2.3	Marginal and conditional AIC . . . . .	13
2.3.1	Marginal AIC . . . . .	13
2.3.2	Conditional AIC . . . . .	13
<b>3</b>	<b>Modified conditional AIC in linear mixed models</b>	<b>15</b>
3.1	Motivation . . . . .	15
3.2	Setup of problem . . . . .	16
3.3	Modification of cAIC . . . . .	16
3.3.1	Conditional Akaike information in linear mixed models . . . . .	16
3.3.2	Evaluation of the bias of cAIC . . . . .	17
3.3.3	Estimation of the bias . . . . .	19
3.4	Simulations . . . . .	21
3.5	Real data example . . . . .	23
3.6	Proofs . . . . .	25
3.6.1	Derivation of (3.6) . . . . .	25
3.6.2	Proof of Theorem 3.1 . . . . .	26
3.6.3	Proof of Lemma 3.1 . . . . .	28
3.6.4	Proof of Lemma 3.2 . . . . .	28
3.6.5	Proof of Lemma 3.3 . . . . .	28
3.6.6	Proof of Lemma 3.4 . . . . .	29
<b>4</b>	<b>Conditional AIC under covariate shift with application to small area prediction</b>	<b>31</b>
4.1	Motivation . . . . .	31
4.2	Covariate shift conditional AIC . . . . .	32
4.2.1	Observed model . . . . .	32
4.2.2	Predictive model . . . . .	32
4.2.3	Conditional Akaike information . . . . .	33
4.2.4	Criterion for overspecified model . . . . .	33

4.3	Modification of the criterion . . . . .	34
4.3.1	Drawback of overspecified model assumption . . . . .	34
4.3.2	Evaluation of cAI . . . . .	35
4.3.3	Estimation of cAI . . . . .	36
4.4	Application to small area prediction . . . . .	40
4.5	Simulations . . . . .	41
4.5.1	Simulations of measuring the biases of estimating the true cAI by the criteria . . . . .	41
4.5.2	Simulations of predicting finite population mean . . . . .	44
4.6	Proofs . . . . .	45
4.6.1	Proof of Theorem 4.1 . . . . .	45
4.6.2	Proof of Theorem of 4.2 . . . . .	46
4.6.3	Proof of Lemma 4.1 . . . . .	47
4.6.4	Proof of Lemma 4.2 . . . . .	47
4.6.5	Proof of Lemma 4.3 . . . . .	47
4.6.6	Proof of Lemma 4.4 . . . . .	47
4.6.7	Proof of Lemma 4.5 . . . . .	48
4.6.8	Proof of Theorem 4.3and Theorem 4.4 . . . . .	48
<b>5</b>	<b>Conditional AIC in mixed effects models based on natural exponential family</b>	<b>49</b>
5.1	Motivation . . . . .	49
5.2	Model and conditional AIC . . . . .	50
5.2.1	Mixed effects models based on natural exponential family . . . . .	50
5.2.2	Variable selection problem in nonlinear mixed model . . . . .	52
5.2.3	Conditional AIC in nonlinear mixed model . . . . .	52
5.3	Approximation and estimation of penalty term . . . . .	53
5.3.1	Decomposition of penalty term . . . . .	53
5.3.2	Analytical method for the case of large $n_i$ . . . . .	53
5.3.3	Method for constant $n_i$ by using numerical integration and differentiation . . . . .	58
5.3.4	Numerical method for constant $n_i$ based on parametric bootstrap . . . . .	61
5.4	Simulations . . . . .	62
5.5	Some results of analytical calculations . . . . .	63
5.5.1	Stochastic Expansion of $\hat{\eta}$ . . . . .	63
5.5.2	Expressions of $\mathbf{J}_r$ and $\mathbf{K}_r$ . . . . .	65
5.5.3	Expression of $\partial \mathbf{d}_{ir}^T / \partial \boldsymbol{\eta}$ . . . . .	66
5.5.4	Expression of $\partial \boldsymbol{\Sigma}_i^{-1} / \partial \boldsymbol{\eta}$ . . . . .	67
5.5.5	Approximations of $\xi_{ri}$ . . . . .	68
5.6	Proofs . . . . .	69
5.6.1	Proof of Lemma 5.3 . . . . .	69
5.6.2	Proof of Lemmas 5.4and 5.5 . . . . .	69
<b>6</b>	<b>A variant of AIC using Bayesian marginal likelihood</b>	<b>71</b>
6.1	Motivation . . . . .	71
6.2	Proposed criteria . . . . .	73
6.2.1	Variable selection criteria for linear regression model . . . . .	73
6.2.2	Extension to the case of unknown covariance . . . . .	76
6.3	Consistency of the criteria . . . . .	77
6.4	Simulations . . . . .	77
6.5	Discussion . . . . .	85

6.6	Derivations of the criteria . . . . .	86
6.6.1	Derivation of $IC_{\pi,1}$ in (6.8) . . . . .	86
6.6.2	Derivation of $IC_{\pi,2}$ in (6.9) . . . . .	87
6.6.3	Derivation of $IC_r$ in (6.11) . . . . .	87
6.7	Proof of Theorem 6.1 . . . . .	87
<b>7</b>	<b>Variants of conditional AIC in linear mixed models</b>	<b>91</b>
7.1	Motivation . . . . .	91
7.2	Variants of conditional Akaike information . . . . .	92
7.2.1	Setup of variable selection . . . . .	92
7.2.2	Conditional Kullback–Leibler risk . . . . .	92
7.3	Predictive information criterion . . . . .	93
7.3.1	Bayesian predictive density in linear mixed model . . . . .	93
7.3.2	Derivation of PIC in linear mixed model . . . . .	94
7.3.3	Another PIC putting prior on regression coefficients . . . . .	95
7.4	Conditional AIC variant based on Bayesian marginal likelihood . . . . .	98
7.4.1	Conditional KL risk of predictive density based on Bayesian marginal likelihood . . . . .	98
7.4.2	Case of normal prior . . . . .	98
7.4.3	Conditional RIC . . . . .	100
7.5	Simulations . . . . .	102
7.6	Proofs . . . . .	104
7.6.1	Proof of Proposition 7.1 . . . . .	104
7.6.2	Proof of Proposition 7.2 . . . . .	104





# Chapter 1

## Introduction

This thesis studies the variable selection problem in mixed effects models, especially the method of information criteria is focused on. Vaida and Blanchard (2005) introduced the conditional Akaike information, which is related to the expected Kullback–Leibler divergence based on the conditional likelihood given random effects. The conditional Akaike information (cAI) and the corresponding criterion, called the conditional AIC (cAIC), are more appropriate for the focus on clusters than the conventional AIC (or marginal AIC, mAIC), which is based on the marginal likelihood integrating out the random effects. This is useful in problems involving prediction of random effects such as small area estimation. First of all, we here introduce a class of mixed effects models and its application to small area estimation. Next, we review the variable selection problem in mixed effects models.

The class of mixed effects models has been studied for a long time from both theoretical and applied aspects. Especially, the linear mixed model and the best linear unbiased predictor (BLUP) introduced by Henderson (1950) provide flexible framework for modeling several types of data sets, whose applications are longitudinal data analysis in biostatistics, panel data analysis in econometrics, small area estimation in official statistics and others. The small area estimation problem is how to produce reliable estimates of some characteristic of interest for areas with small sample sizes. Model based estimator in small area estimation problem using mixed effects models can ‘borrow information’ from neighboring areas, which results in stable estimation of small area parameter. Datta and Ghosh (2012), Pfeiffermann (2013) and Rao and Molina (2015) give good reviews about small area estimation.

There are several methods of variable selection in mixed effects models, which include the information criteria such as AIC or BIC, shrinkage methods such as LASSO, the Bayesian procedure, the Fence methods (Jiang et al., 2008), and others. Müller et al. (2013) is a good review of variable selection problem in linear mixed models. In this thesis, we focus on the information criteria, especially the cAIC. Since Vaida and Blanchard (2005), variable selection procedures using the cAIC have been developed and the properties of the cAIC have been discussed. Liang et al. (2008) proposed a different bias correction, who take into account estimation of the unknown parameters included in the covariance matrix of the vector of the random effects. It is noted that their bias correction is closely related to the generalized degrees of freedom of Ye (1998), while Vaida and Blanchard (2005) pointed out their bias correction is the same as the effective degrees of freedom of Hodges and Sargent (2001). Greven and Kneib (2010) derived analytical representation of the bias correction of Liang et al. (2008) and also applied to selecting the random effects. Srivastava and Kubokawa (2010) proposed other versions of the cAIC by changing the estimators of the regression coefficients and the variance parameter. Kubokawa (2011) derived the cAIC with a general covariance matrices of the vectors of the random effects

and the error terms. He also proposed a conditional version of Mallows'  $C_p$  (Mallows, 1973). The cAIC has been also applied to the variable selection problem in the generalized linear mixed models (GLMM). Donohue et al. (2011), Yu and Yau (2012) and Yu et al. (2013) derived the cAIC in the GLMM under the assumption that the sample size in each cluster goes to infinity. Saefken et al. (2014) proposed an exact unbiased estimator of the cAI, which is also justified for finite sample case, in a special case of the GLMM, Poisson mixed regression. As a related procedure to the cAIC, Zhang et al. (2014) proposed a model averaging method based on the prediction risk relative to the conditional Mallows' criterion. This procedure is an extension of the Mallows model averaging of Hansen (2007) to the linear mixed model and Zhang et al. (2014) also proved the model averaging estimator is asymptotically loss efficient under some regularity condition.

In this thesis, some problems are considered and new results are obtained. Firstly, the conditional AIC is modified for the underspecified model, which does not include the true model. Most of the Akaike-type information criteria put the assumption that the candidate model includes the true model, which we call overspecified assumption. Due to the assumption, the cAIC has large bias for estimating the cAI. This problem is considered in Chapter 3. Secondly, the variable selection problem in linear mixed model under covariate shift situation is considered. We also show that considering covariate shift situation is meaningful in small area estimation problem. We discuss this problem in Chapter 4. Thirdly, we derived the cAIC as a variable selection problem in a class of mixed effects models based on natural exponential family, which includes nonlinear mixed models such as Poisson-gamma model and binomial-beta model in Chapter 5. Lastly, some variants of the AIC and the cAIC are proposed and their properties are discussed, which are also related to the Bayesian procedure as well as frequentists' methods. We develop a variant of marginal AIC in Chapter 6 taking linear regression model as an example. In Chapter 7, variants of the cAIC are considered.

## Chapter 2

# Mixed effects models and conditional AIC

In this chapter, mixed effects models are introduced and the setup of the variable selection problem is explained. The mAIC and cAIC are also briefly summarized.

### 2.1 Mixed effects models

#### 2.1.1 General

Let  $\mathbf{y}$  be an observable random vector, and let  $\boldsymbol{\theta}$  be an unobservable random vector. We treat continuous or discrete cases for  $\mathbf{y}$  and  $\boldsymbol{\theta}$ . The conditional probability density (or mass) function of  $\mathbf{y}$  given  $\boldsymbol{\theta}$  is  $f(\mathbf{y}|\boldsymbol{\theta}, \boldsymbol{\eta})$  for a vector of unknown parameters  $\boldsymbol{\eta}$ , and the the probability density (or mass) function of  $\boldsymbol{\theta}$  is  $p(\boldsymbol{\theta}|\boldsymbol{\eta})$ , namely,

$$\begin{aligned}\mathbf{y}|\boldsymbol{\theta} &\sim f(\mathbf{y}|\boldsymbol{\theta}, \boldsymbol{\eta}), \\ \boldsymbol{\theta} &\sim p(\boldsymbol{\theta}|\boldsymbol{\eta}).\end{aligned}\tag{2.1}$$

When  $\boldsymbol{\theta}$  denotes random effects, this expresses general mixed effects model. Because this model can be interpreted as a Bayesian model, we also use the terminology used in Bayes statistics. The marginal likelihood function of  $\mathbf{y}$  and the conditional (or posterior) density function of  $\boldsymbol{\theta}$  given  $\mathbf{y}$  are

$$\begin{aligned}m(\mathbf{y}|\boldsymbol{\eta}) &= \int f(\mathbf{y}|\boldsymbol{\theta}, \boldsymbol{\eta})p(\boldsymbol{\theta}|\boldsymbol{\eta})d\boldsymbol{\theta}, \\ p(\boldsymbol{\theta}|\mathbf{y}, \boldsymbol{\eta}) &= f(\mathbf{y}|\boldsymbol{\theta}, \boldsymbol{\eta})p(\boldsymbol{\theta}|\boldsymbol{\eta})/m(\mathbf{y}|\boldsymbol{\eta}).\end{aligned}$$

#### 2.1.2 Linear mixed model

One of the most important examples of the mixed effects model is the linear mixed model, which is given as the following general form:

$$\mathbf{y} = \mathbf{X}\boldsymbol{\beta} + \mathbf{Z}\mathbf{b} + \boldsymbol{\varepsilon},\tag{2.2}$$

where  $\mathbf{y}$  is an  $n \times 1$  observation vector of the response variables,  $\mathbf{X}$  and  $\mathbf{Z}$  are  $n \times p$  and  $n \times q$  matrices of covariates,  $\boldsymbol{\beta}$  is a  $p \times 1$  vector of unknown regression coefficients,  $\mathbf{b}$  is a  $q \times 1$  vector of random effects, and  $\boldsymbol{\varepsilon}$  is an  $n \times 1$  vector of random errors. It is common to assume that  $\mathbf{b}$  and  $\boldsymbol{\varepsilon}$  are mutually independent and  $\mathbf{b} \sim \mathcal{N}_q(\mathbf{0}, \mathbf{Q})$ ,  $\boldsymbol{\varepsilon} \sim \mathcal{N}_n(\mathbf{0}, \mathbf{R})$ .

The empirical best linear unbiased predictor (EBLUP) of the linear combination of  $\boldsymbol{\beta}$  and  $\mathbf{b}$ , which is of the form  $\mathbf{c}^\top \boldsymbol{\beta} + \mathbf{d}^\top \mathbf{b}$  for a  $p \times 1$  vector  $\mathbf{c}$  and a  $q \times 1$  vector  $\mathbf{d}$ , is  $\mathbf{c}^\top \hat{\boldsymbol{\beta}} + \mathbf{d}^\top \hat{\mathbf{b}}$ , where

$$\begin{aligned}\hat{\boldsymbol{\beta}} &= (\mathbf{X}^\top \boldsymbol{\Sigma}^{-1} \mathbf{X})^{-1} \mathbf{X}^\top \boldsymbol{\Sigma}^{-1} \mathbf{y}, \\ \hat{\mathbf{b}} &= \mathbf{Q} \mathbf{Z}^\top \boldsymbol{\Sigma}^{-1} (\mathbf{y} - \mathbf{X} \hat{\boldsymbol{\beta}}),\end{aligned}$$

where  $\boldsymbol{\Sigma} = \mathbf{Z} \mathbf{Q} \mathbf{Z}^\top + \mathbf{R}$ .

## 2.2 Setup of variable selection problem

### 2.2.1 Candidate models, full model and true model

We focus on the variable selection of the fixed effects. To this end, we clarify candidate models, the true model, and the full model and use the following notations throughout the paper. We take a linear mixed model (2.2) for example.

First, we consider the collection of candidate models as follows. Let  $n \times p_\omega$  matrix  $\mathbf{X}(\omega)$  consist of all the explanatory variables and assume that  $\text{rank}(\mathbf{X}(\omega)) = p_\omega$ . In order to define candidate models by the index set  $j$ , suppose that  $j$  denotes a subset of  $\omega = \{1, \dots, p_\omega\}$  containing  $p_j$  elements, *i.e.*  $p_j = \#(j)$ , and  $\mathbf{X}(j)$  consists of  $p_j$  columns of  $\mathbf{X}(\omega)$  indexed by the elements of  $j$ . We define the index set by  $\mathcal{J} = \mathcal{P}(\omega)$ , namely the power set of  $\omega$ , where we call  $\omega$  the full model. Then we consider the candidate model  $j$  expressed as

$$\mathbf{y} = \mathbf{X}(j) \boldsymbol{\beta}_j + \mathbf{Z} \mathbf{b}_j + \boldsymbol{\varepsilon}_j,$$

where  $\boldsymbol{\beta}_j$  is a  $p_j \times 1$  vector of regression coefficients,  $\mathbf{b}_j$  is a  $q \times 1$  vector of random effects,  $\boldsymbol{\varepsilon}_j$  is an  $n \times 1$  vector of random errors.

Second, we assume that the true model exists in the collection of the candidate models  $\mathcal{P}(\omega)$ , which is denoted by  $j_*$ . It is noted that the dimension of the true model is  $p_{j_*}$ , which is abbreviated to  $p_*$ . Under the assumption, the true model can be written by using the full model design matrix, namely the true mean of  $\mathbf{y}$  can be expressed as

$$E(\mathbf{y}) = \mathbf{X}(\omega) \boldsymbol{\beta}_*,$$

where  $\boldsymbol{\beta}_*$  is  $p_\omega \times 1$  vector of regression coefficients, whose  $p_\omega - p_*$  components are exactly 0 and the rest of components are not 0. Moreover, when the true model is included by the candidate model, the true mean of  $\mathbf{y}$  can be also expressed as

$$E(\mathbf{y}) = \mathbf{X}(j) \boldsymbol{\beta}_j^*,$$

where  $\boldsymbol{\beta}_j^*$  is  $p_j \times 1$  vector of regression coefficients, whose  $p_j - p_*$  components are exactly 0 and the rest of components are not 0. It is common to assume that the true model is included by the candidate model for the derivation of the Akaike-type information criteria.

### 2.2.2 Overspecified and underspecified

We introduce the terms overspecified and underspecified models. Candidate model  $j$  is overspecified if  $\mathbf{X}(\omega) \boldsymbol{\beta}_* \in \mathcal{R}[\mathbf{X}(j)]$ , which means that  $\mathbf{X}(\omega) \boldsymbol{\beta}_*$  is in the column space of  $\mathbf{X}(j)$  following Fujikoshi and Satoh (1997) or Kawakubo and Kubokawa (2014). The set of overspecified models are denoted by  $\mathcal{J}_+ = \{j \in \mathcal{J} | j_* \subseteq j\}$ . On the other hand, candidate model  $j$  is called underspecified when  $\mathbf{X}(\omega) \boldsymbol{\beta}_* \notin \mathcal{R}[\mathbf{X}(j)]$ . The set of underspecified models are denoted by  $\mathcal{J}_- = \mathcal{J} \setminus \mathcal{J}_+$ .

## 2.3 Marginal and conditional AIC

### 2.3.1 Marginal AIC

The marginal AIC (mAIC) is related to the expected Kullback–Leibler (KL) divergence based on the marginal likelihood defined as

$$\int \left[ \int \log \left\{ \frac{m(\tilde{\mathbf{y}}|\boldsymbol{\eta})}{m(\tilde{\mathbf{y}}|\hat{\boldsymbol{\eta}})} \right\} m(\tilde{\mathbf{y}}|\boldsymbol{\eta}) d\tilde{\mathbf{y}} \right] m(\mathbf{y}|\boldsymbol{\eta}) d\mathbf{y}, \quad (2.3)$$

where  $\tilde{\mathbf{y}}$  is an independent replication of  $\mathbf{y}$  and  $\hat{\boldsymbol{\eta}}$  is the maximum likelihood estimator of  $\boldsymbol{\eta}$ . The marginal AIC (mAIC) is an (asymptotically) unbiased estimator of the following marginal Akaike information (mAI):

$$\text{mAI} = \iint -2 \log \{m(\tilde{\mathbf{y}}|\hat{\boldsymbol{\eta}})\} m(\tilde{\mathbf{y}}|\boldsymbol{\eta}) m(\mathbf{y}|\boldsymbol{\eta}) d\tilde{\mathbf{y}} d\mathbf{y},$$

which is a part of (2.3) (multiplied by 2). Then, the mAIC is

$$\text{mAIC} = -2 \log \{m(\mathbf{y}|\hat{\boldsymbol{\eta}})\} + \Delta_{\text{mAIC}},$$

where  $\Delta_{\text{mAIC}}$  is bias correction (or penalty) term, which is given by

$$\Delta_{\text{mAIC}} = \text{mAI} - E[-2 \log \{m(\mathbf{y}|\hat{\boldsymbol{\eta}})\}].$$

Akaike (1973, 1974) proposed that the bias correction converges to  $2p$  where  $p$  is the dimension of the candidate model.

### 2.3.2 Conditional AIC

The marginal AIC is not appropriate for the focus on the prediction of specific clusters or random effects. Taking this point into account, Vaida and Blanchard (2005) considered the expected KL divergence based on the conditional density, which is given by

$$\iint \left[ \int \log \left\{ \frac{f(\tilde{\mathbf{y}}|\boldsymbol{\theta}, \boldsymbol{\eta})}{f(\tilde{\mathbf{y}}|\hat{\boldsymbol{\theta}}, \hat{\boldsymbol{\eta}})} \right\} f(\tilde{\mathbf{y}}|\boldsymbol{\theta}, \boldsymbol{\eta}) d\tilde{\mathbf{y}} \right] f(\mathbf{y}|\boldsymbol{\theta}, \boldsymbol{\eta}) p(\boldsymbol{\theta}|\boldsymbol{\eta}) d\mathbf{y} d\boldsymbol{\theta}, \quad (2.4)$$

where  $\tilde{\mathbf{y}}$  is an independent replication of  $\mathbf{y}$  given  $\boldsymbol{\theta}$ , and  $\hat{\boldsymbol{\theta}}$  is some predictor of  $\boldsymbol{\theta}$ . The conditional AIC (cAIC) is an (asymptotically) unbiased estimator of the following conditional Akaike information (cAI):

$$\text{cAI} = \iiint -2 \log \{f(\tilde{\mathbf{y}}|\hat{\boldsymbol{\theta}}, \hat{\boldsymbol{\eta}})\} f(\tilde{\mathbf{y}}|\boldsymbol{\theta}, \boldsymbol{\eta}) f(\mathbf{y}|\boldsymbol{\theta}, \boldsymbol{\eta}) \pi(\boldsymbol{\theta}|\boldsymbol{\eta}) d\tilde{\mathbf{y}} d\mathbf{y} d\boldsymbol{\theta},$$

which is a part of (2.4) (multiplied by 2). Then, the cAIC is

$$\text{cAIC} = -2 \log \{f(\mathbf{y}|\hat{\boldsymbol{\theta}}, \hat{\boldsymbol{\eta}})\} + \Delta_{\text{cAIC}},$$

where  $\Delta_{\text{cAIC}}$  is bias correction (or penalty) term, which is given as

$$\Delta_{\text{cAIC}} = \text{cAI} - E[-2 \log \{f(\mathbf{y}|\hat{\boldsymbol{\theta}}, \hat{\boldsymbol{\eta}})\}].$$



## Chapter 3

# Modified conditional AIC in linear mixed models

A weak point of cAIC is that it is derived as an unbiased estimator of conditional Akaike information (cAI) in the overspecified case, namely in the case that candidate models include the true model. This results in larger biases in the underspecified case that the true model is not included in candidate models. In this chapter, we derive the modified cAIC (McAIC) to cover both the underspecified and overspecified cases, and investigate properties of McAIC. It is numerically shown that McAIC has less biases and less prediction errors than cAIC. This chapter is based on Kawakubo and Kubokawa (2014).

### 3.1 Motivation

The cAIC by Vaida and Blanchard (2005) is derived under the condition that the candidate model includes the true model. This assumption is called the overspecified assumption. On the other hand, the underspecified case means that a candidate model does not include the true model. Thus, we have the following questions:

- (I) Is cAIC appropriate as an estimator of cAI in the underspecified case ?
- (II) Can one extend cAIC to a procedure useful for both the under- and over-specified cases ?

For the query (I), it is noted that the cAIC is not an asymptotically unbiased estimator of cAI in the underspecified case. In fact, cAIC has large biases in the underspecified case as illustrated in Table 3.1. Thus, the drawback of cAIC gives a motivation for addressing the query (II).

In this chapter, we derive an asymptotically unbiased estimator of cAI in both under- and over-specified cases. This procedure is here called the modified conditional AIC (McAIC). The setup of the problem is explained in Section 3.2. In Section 3.3, we derive the McAIC as an asymptotically unbiased estimator of cAI in both under- and over-specified cases. This approach was used by Fujikoshi and Satoh (1997) to modify AIC and Mallows'  $C_p$  in multivariate linear regression models. The performance of McAIC is investigated numerically by simulation in Section 3.4, and it is shown that McAIC and the corresponding model averaging procedure are better than cAIC in terms of the prediction error. In Section 3.5, we apply the McAIC to estimate small area land prices. All the proofs are given in Section 3.6.

## 3.2 Setup of problem

We focus on the problem of selecting explanatory variables in linear mixed model, whose notations are given in Section 2.2.

First, we explain about the collection of candidate models. We consider the candidate model  $j$ , which is given as follows:

$$\mathbf{y} = \mathbf{X}(j)\boldsymbol{\beta}_j + \mathbf{Z}\mathbf{b}_j + \boldsymbol{\varepsilon}_j, \quad (3.1)$$

where  $\mathbf{y}$  is an  $n \times 1$  observation vector of the response variables,  $\mathbf{X}(j)$  and  $\mathbf{Z}$  are  $n \times p_j$  and  $n \times q$  matrices of covariates,  $\boldsymbol{\beta}_j$  is a  $p_j \times 1$  vector of regression coefficients,  $\mathbf{b}_j$  is a  $q \times 1$  vector of random effects, and  $\boldsymbol{\varepsilon}_j$  is an  $n \times 1$  vector of random errors. We here assume that  $\mathbf{b}_j$  and  $\boldsymbol{\varepsilon}_j$  are mutually independent and that  $\mathbf{b}_j \sim \mathcal{N}_r(\mathbf{0}, \sigma_j^2 \mathbf{G})$ ,  $\boldsymbol{\varepsilon}_j \sim \mathcal{N}_n(\mathbf{0}, \sigma_j^2 \mathbf{I}_n)$  for a common unknown parameter  $\sigma_j^2$  and a known matrix  $\mathbf{G}$ . It is important to point out that the random effects part remains the same over the models  $\{j\}$ . This means that we here consider the problem of selecting only the explanatory variables of the fixed effects. The conditional density function of  $\mathbf{y}$  given  $\mathbf{b}_j$  for model  $j$  is denoted by  $f(\mathbf{y}|\mathbf{b}_j, \boldsymbol{\eta}_j)$ , where  $\boldsymbol{\eta}_j$  is the vector of unknown parameters, namely  $\boldsymbol{\eta}_j = (\boldsymbol{\beta}_j^\top, \sigma_j^2)^\top$ . The density function of  $\mathbf{b}_j$  is denoted by  $p(\mathbf{b}_j|\boldsymbol{\eta}_j)$ .

Second, we assume that the data are generated from the true model which is given by

$$\mathbf{y} = \mathbf{X}(\omega)\boldsymbol{\beta}_* + \mathbf{Z}\mathbf{b}_* + \boldsymbol{\varepsilon}_*$$

for  $\mathbf{b}_* \sim \mathcal{N}_q(\mathbf{0}, \sigma_*^2 \mathbf{G})$ ,  $\boldsymbol{\varepsilon}_* \sim \mathcal{N}_n(\mathbf{0}, \sigma_*^2 \mathbf{I}_n)$ . Thus the marginal distribution of  $\mathbf{y}$  is

$$\mathbf{y} \sim \mathcal{N}_n(\mathbf{X}(\omega)\boldsymbol{\beta}_*, \sigma_*^2 \boldsymbol{\Sigma}), \quad (3.2)$$

where  $\boldsymbol{\Sigma} = \mathbf{Z}\mathbf{G}\mathbf{Z}^\top + \mathbf{I}_n$ . For the true model, the conditional density function of  $\mathbf{y}$  given  $\mathbf{b}_*$  and the density function of  $\mathbf{b}_*$  are denoted by  $f(\mathbf{y}|\mathbf{b}_*, \boldsymbol{\eta}_*)$  and  $p(\mathbf{b}_*|\boldsymbol{\eta}_*)$ , respectively, where  $\boldsymbol{\eta}_* = (\boldsymbol{\beta}_*^\top, \sigma_*^2)^\top$ .

Third, we assume that the collection of candidate models includes both underspecified and overspecified models, and that the full model  $\omega$  includes the true model. This means that the set of overspecified models  $\mathcal{J}_+$  is not empty set.

## 3.3 Modification of cAIC

### 3.3.1 Conditional Akaike information in linear mixed models

We begin with introducing the conditional Akaike information (cAI) in linear mixed models, which was proposed by Vaida and Blanchard (2005). The cAI is the estimand of the cAIC, and is related to the expected Kullback–Leibler (KL) divergence based on the conditional likelihood. The cAI for the setup explained in Section 3.2 is

$$\begin{aligned} \text{cAI} &= \iiint -2 \log\{f(\tilde{\mathbf{y}}|\hat{\mathbf{b}}_j, \hat{\boldsymbol{\eta}}_j)\} f(\tilde{\mathbf{y}}|\mathbf{b}_*, \boldsymbol{\eta}_*) f(\mathbf{y}|\mathbf{b}_*, \boldsymbol{\eta}_*) p(\mathbf{b}_*|\boldsymbol{\eta}_*) d\tilde{\mathbf{y}} d\mathbf{y} d\mathbf{b} \\ &= E^{(\mathbf{y}, \mathbf{b}_*)} E^{\tilde{\mathbf{y}}|\mathbf{b}_*} \left[ n \log(2\pi\hat{\sigma}_j^2) + \|\tilde{\mathbf{y}} - \mathbf{X}(j)\hat{\boldsymbol{\beta}}_j - \mathbf{Z}\hat{\mathbf{b}}_j\|^2 / \hat{\sigma}_j^2 \right], \end{aligned}$$

where  $\tilde{\mathbf{y}}$  is an independent replication of  $\mathbf{y}$  given  $\mathbf{b}$  and  $E^{(\mathbf{y}, \mathbf{b}_*)}$  and  $E^{\tilde{\mathbf{y}}|\mathbf{b}_*}$  denote the expectation with respect to the joint distribution of  $(\mathbf{y}, \mathbf{b}_*)$  and the conditional distribution of  $\tilde{\mathbf{y}}$  given  $\mathbf{b}_*$ . The cAI measures the risk of plug-in predictive density  $f(\tilde{\mathbf{y}}|\hat{\mathbf{b}}_j, \hat{\boldsymbol{\eta}}_j)$ , where  $\hat{\boldsymbol{\eta}}_j$  is the maximum likelihood estimator of  $\boldsymbol{\eta}_j = (\boldsymbol{\beta}_j^\top, \sigma_j^2)^\top$  based on the candidate model  $j$  given as

$$\begin{aligned} \hat{\boldsymbol{\beta}}_j &= (\mathbf{X}(j)^\top \boldsymbol{\Sigma}^{-1} \mathbf{X}(j))^{-1} \mathbf{X}(j)^\top \boldsymbol{\Sigma}^{-1} \mathbf{y}, \\ \hat{\sigma}_j^2 &= (\mathbf{y} - \mathbf{X}(j)\hat{\boldsymbol{\beta}}_j)^\top \boldsymbol{\Sigma}^{-1} (\mathbf{y} - \mathbf{X}(j)\hat{\boldsymbol{\beta}}_j) / n, \end{aligned} \quad (3.3)$$



and  $\hat{\mathbf{b}}_j$  is the empirical Bayes estimator for quadratic loss of  $\mathbf{b}_j$  given as

$$\hat{\mathbf{b}}_j = \mathbf{G}\mathbf{Z}^\top \boldsymbol{\Sigma}^{-1}(\mathbf{y} - \mathbf{X}(j)\hat{\boldsymbol{\beta}}_j).$$

Then the cAIC is a bias corrected unbiased estimator of cAI, which is given as

$$\text{cAIC} = -2 \log\{f(\mathbf{y}|\hat{\mathbf{b}}_j, \hat{\boldsymbol{\eta}}_j)\} + \Delta_{\text{cAI}},$$

where

$$\Delta_{\text{cAI}} = \text{cAI} - E^{(\mathbf{y}, \mathbf{b}_*)}[-2 \log\{f(\mathbf{y}|\hat{\mathbf{b}}_j, \hat{\boldsymbol{\eta}}_j)\}], \quad (3.4)$$

which is called the bias correction (or penatly) term.

### 3.3.2 Evaluation of the bias of cAIC

Vaida and Blanchard (2005) evaluated the bias correction in (3.4) under the assumption that the candidate model  $j$  is overspecified. Then the cAIC is the exact unbiased estimator of cAI when the candidate model  $j$  is overspecified. However, when the candidate model  $j$  is underspecified, the cAIC has large bias for estimating the cAI. Thus we evaluate the bias correction in (3.4) both for overspecified and underspecified case.

First, it follows that

$$-2 \log f(\mathbf{y}|\hat{\mathbf{b}}_j, \hat{\boldsymbol{\eta}}_j) = n \log(2\pi\hat{\sigma}_j^2) + \|\mathbf{y} - \mathbf{X}(j)\hat{\boldsymbol{\beta}}_j - \mathbf{Z}\hat{\mathbf{b}}_j\|^2 / \hat{\sigma}_j^2. \quad (3.5)$$

Then, as shown in Section 3.6, the bias correction can be expressed as

$$\Delta_{\text{cAI}} = E \left[ \left\{ (2n - \text{tr}[\boldsymbol{\Sigma}^{-1}])\sigma_*^2 - \mathbf{u}^\top \boldsymbol{\Sigma}^{-2} \mathbf{u} + 2\mathbf{u}^\top \boldsymbol{\Sigma}^{-2} (\mathbf{X}(j)\hat{\boldsymbol{\beta}}_j - \mathbf{X}(\omega)\boldsymbol{\beta}_*) \right\} / \hat{\sigma}_j^2 \right], \quad (3.6)$$

where expectation is taken with respect to the joint distribution of  $(\mathbf{y}, \mathbf{b}_*)$  and  $\mathbf{u} = \mathbf{y} - \mathbf{X}(\omega)\boldsymbol{\beta}_*$ .

It is important to note that the distribution of  $\hat{\sigma}_j^2$  for the underspecified case is different from that for the overspecified case. Thus, we need to clarify the distribution of  $\hat{\sigma}_j^2$ . To this end,  $n\hat{\sigma}_j^2$  is decomposed as

$$\begin{aligned} n\hat{\sigma}_j^2 &= \sigma_*^2 \{ \mathbf{z}^\top (\mathbf{I}_n - \mathbf{M}_\omega) \mathbf{z} + \mathbf{z}^\top (\mathbf{M}_\omega - \mathbf{M}_j) \mathbf{z} \} \\ &= \sigma_*^2 (K_0 + K_1), \end{aligned}$$

where  $K_0 = \mathbf{z}^\top (\mathbf{I}_n - \mathbf{M}_\omega) \mathbf{z}$ ,  $K_1 = \mathbf{z}^\top (\mathbf{M}_\omega - \mathbf{M}_j) \mathbf{z}$ ,

$$\begin{aligned} \mathbf{z} &= \boldsymbol{\Sigma}^{-1/2} \mathbf{y} / \sigma_*, \\ \mathbf{M}_\omega &= \boldsymbol{\Sigma}^{-1/2} \mathbf{X}(\omega) (\mathbf{X}(\omega)^\top \boldsymbol{\Sigma}^{-1} \mathbf{X}(\omega))^{-1} \mathbf{X}(\omega)^\top \boldsymbol{\Sigma}^{-1/2}, \\ \mathbf{M}_j &= \boldsymbol{\Sigma}^{-1/2} \mathbf{X}(j) (\mathbf{X}(j)^\top \boldsymbol{\Sigma}^{-1} \mathbf{X}(j))^{-1} \mathbf{X}(j)^\top \boldsymbol{\Sigma}^{-1/2}. \end{aligned}$$

Note that  $\mathbf{M}_j$  and  $\mathbf{M}_\omega$  are symmetric and idempotent. Let  $\mathbf{v} = \boldsymbol{\Sigma}^{-1/2} \mathbf{u} / \sigma_*$  and  $\boldsymbol{\xi} = \boldsymbol{\Sigma}^{-1/2} \mathbf{X}(\omega) \boldsymbol{\beta}_* / \sigma_*$ . Then, it is seen that  $\mathbf{M}_\omega \boldsymbol{\xi} = \boldsymbol{\xi}$  and

$$\mathbf{M}_j \boldsymbol{\xi} \begin{cases} = \boldsymbol{\xi} & \text{if } j \in \mathcal{J}_+, \\ \neq \boldsymbol{\xi} & \text{if } j \in \mathcal{J}_-, \end{cases}$$

since  $\mathbf{X}_* \boldsymbol{\beta}_* \in \mathcal{R}[\mathbf{X}_j]$  if  $j \in \mathcal{J}_+$ . Thus  $K_0$  can be rewritten as

$$K_0 = (\boldsymbol{\xi} + \mathbf{v})^\top (\mathbf{I}_N - \mathbf{M}_\omega) (\boldsymbol{\xi} + \mathbf{v}) = \mathbf{v}^\top (\mathbf{I}_N - \mathbf{M}_\omega) \mathbf{v}, \quad (3.7)$$

so that  $K_0$  follows a chi-squared distribution with  $n - p_\omega$  degrees of freedom, denoted by

$$K_0 \sim \chi_{n-p_\omega}^2.$$

Also,  $K_1$  can be rewritten as

$$\begin{aligned} K_1 &= \mathbf{v}^\top (\mathbf{M}_\omega - \mathbf{M}_j) \mathbf{v} + 2\xi^\top (\mathbf{M}_\omega - \mathbf{M}_j) \mathbf{v} + \xi^\top (\mathbf{M}_\omega - \mathbf{M}_j) \xi \\ &= \mathbf{v}^\top (\mathbf{M}_\omega - \mathbf{M}_j) \mathbf{v} + 2L + n\delta, \end{aligned} \quad (3.8)$$

where

$$\begin{aligned} L &= \xi^\top (\mathbf{M}_\omega - \mathbf{M}_j) \mathbf{v}, \\ \delta &= \xi^\top (\mathbf{M}_\omega - \mathbf{M}_j) \xi / n. \end{aligned} \quad (3.9)$$

In the overspecified case, we have  $K_1 \sim \chi_{p_\omega - p_j}^2$  since  $\mathbf{M}_\omega \xi = \mathbf{M}_j \xi = \xi$ . In the underspecified case,  $K_1$  follows a noncentral chi-squared distribution with  $p_\omega - p_j$  degrees of freedom and with the noncentrality parameter  $N\delta$ , denoted by  $K_1 \sim \chi_{p_\omega - p_j}^2(n\delta)$ . Thus,

$$K_1 \sim \begin{cases} \chi_{p_\omega - p_j}^2 & \text{if } j \in \mathcal{J}_+, \\ \chi_{p_\omega - p_j}^2(n\delta) & \text{if } j \in \mathcal{J}_-. \end{cases}$$

Since  $\mathbf{u}^\top \Sigma^{-2} \mathbf{u} = \sigma_*^2 \mathbf{v}^\top \Sigma^{-1} \mathbf{v}$  and

$$\mathbf{u}^\top \Sigma^{-2} (\mathbf{X}(j) \hat{\boldsymbol{\beta}}_j - \mathbf{X}(\omega) \boldsymbol{\beta}_*) = \sigma_*^2 \{ \mathbf{v}^\top \Sigma^{-1} \mathbf{M}_j \mathbf{v} - \xi^\top (\mathbf{M}_\omega - \mathbf{M}_j) \Sigma^{-1} \mathbf{v} \},$$

we can rewrite (3.6) as

$$\Delta_{\text{cAI}} = n \cdot E [(K_0 + K_1)^{-1} \{ (2n - \text{tr}[\Sigma^{-1}]) - \mathbf{v}^\top \Sigma^{-1} \mathbf{v} + 2\mathbf{v}^\top \Sigma^{-1} \mathbf{M}_j \mathbf{v} - 2\xi^\top (\mathbf{M}_\omega - \mathbf{M}_j) \Sigma^{-1} \mathbf{v} \}]. \quad (3.10)$$

Although  $K_0 + K_1$  has a central chi-squared distribution in the overspecified case, it has a noncentral chi-squared distribution in the underspecified case. Thus, we need to approximate the bias  $\Delta_{\text{cAI}}$ . To this end, we assume the following conditions:

- (A1)  $\xi^\top (\mathbf{M}_\omega - \mathbf{M}_j) \xi = O(n)$ , which is the non-centrality parameter of  $K_1$ .
- (A2)  $\xi^\top (\mathbf{M}_\omega - \mathbf{M}_j) \Sigma^{-1} (\mathbf{M}_\omega - \mathbf{M}_j) \xi = O(n)$ .

The condition (A1) is equivalent to  $\delta = O(1)$  given in (3.9). It is also noted that the condition (A2) is satisfied by (A1) if the maximum eigenvalue of  $\Sigma^{-1}$  is uniformly bounded. Under these assumptions, we can get the following theorem, which will be proved in Section 3.6.

**Theorem 3.1** *In the overspecified case, the bias correction of cAIC is provided by the exact expression  $\Delta_{\text{cAI}} = B^*$ , where*

$$B^* = 2n \times \left\{ \frac{n - \text{tr}[\Sigma^{-1}] + \text{tr}[\mathbf{P}_j]}{n - p_j - 2} + \frac{\text{tr}[\Sigma^{-1}] - \text{tr}[\mathbf{P}_j]}{(n - p_j - 2)(n - p_j)} \right\}, \quad (3.11)$$

for  $\mathbf{P}_j = \Sigma^{-1} \mathbf{X}(j) (\mathbf{X}(j)^\top \Sigma^{-1} \mathbf{X}(j))^{-1} \mathbf{X}(j)^\top \Sigma^{-1}$ . In the underspecified case, the bias correction of cAIC is approximated as

$$\Delta_{\text{cAI}} = B^* + B_1 + B_2 + B_3 + O(n^{-1}), \quad (3.12)$$

where  $B_1$ ,  $B_2$  and  $B_3$  are defined as

$$B_1 = \frac{2n(\lambda - 1)}{n - p_j - 2}(n - \text{tr}[\boldsymbol{\Sigma}^{-1}]), \quad (3.13)$$

$$B_2 = 2p_j\lambda(\lambda - 1) - 4\lambda(\lambda - 1)^2 + 2\text{tr}[\mathbf{P}_j](\lambda - 1) + 2(\lambda - 1)\text{tr}[\boldsymbol{\Sigma}^{-1}] \times \frac{2\lambda^2 - (p_j + 1)\lambda + 1}{n}, \quad (3.14)$$

$$B_3 = \frac{4\lambda^2}{n}\boldsymbol{\xi}^\top(\mathbf{M}_\omega - \mathbf{M}_j)\boldsymbol{\Sigma}^{-1}(\mathbf{M}_\omega - \mathbf{M}_j)\boldsymbol{\xi}, \quad (3.15)$$

where  $\lambda = 1/(1 + \delta)$ .

It is noted that in the overspecified case the bias  $B^*$  given in (3.11) is identical to that in Vaida and Blanchard (2005). It is also noted that  $B_1 = B_2 = B_3 = 0$  in the overspecified case, since  $\lambda = 1$  and  $\mathbf{M}_j\boldsymbol{\xi} = \boldsymbol{\xi}$ .

### 3.3.3 Estimation of the bias

We now derive an asymptotically unbiased estimator of the bias  $\Delta_{\text{cAI}}$ . It follows from Theorem 3.1 that it is sufficient to estimate  $B_1$ ,  $B_2$  and  $B_3$  because  $B_*$  does not include any unknown parameters. Since  $B_1$  and  $B_2$  are linear functions of  $\lambda$ ,  $\lambda^2$  and  $\lambda^3$ , we begin by estimating these polynomials of  $\lambda$ .

Let us define  $\hat{\lambda}$ ,  $\widehat{\lambda^2}$  and  $\widehat{\lambda^3}$  as

$$\hat{\lambda} = \frac{n - p_j}{n - p_\omega} \frac{\hat{\sigma}_\omega^2}{\hat{\sigma}_j^2}, \quad (3.16)$$

$$\widehat{\lambda^2} = \frac{(n - p_j)(n - p_j + 2)}{(n - p_\omega)(n - p_\omega + 2)} \left( \frac{\hat{\sigma}_\omega^2}{\hat{\sigma}_j^2} \right)^2, \quad (3.17)$$

$$\widehat{\lambda^3} = \frac{(n - p_j)(n - p_j + 2)(n - p_j + 4)}{(n - p_\omega)(n - p_\omega + 2)(n - p_\omega + 4)} \left( \frac{\hat{\sigma}_\omega^2}{\hat{\sigma}_j^2} \right)^3. \quad (3.18)$$

In the overspecified case, it is noted that  $n\hat{\sigma}_\omega^2 = \sigma_*^2 K_0 \sim \sigma_*^2 \chi_{n-p_\omega}^2$ ,  $K_1 \sim \chi_{p_\omega-p_j}^2$  and  $n\hat{\sigma}_j^2 = \sigma_*^2(K_0 + K_1)$ , so that

$$\frac{\hat{\sigma}_\omega^2}{\hat{\sigma}_j^2} \sim \text{Be} \left( \frac{n - p_\omega}{2}, \frac{p_\omega - p_j}{2} \right),$$

where  $\text{Be}(\cdot, \cdot)$  denotes the beta distribution. This implies that  $E[\hat{\lambda}] = E[\widehat{\lambda^2}] = E[\widehat{\lambda^3}] = 1$  in the overspecified case. In the underspecified case, on the other hand, it follows that  $E[(\hat{\sigma}_\omega^2/\hat{\sigma}_j^2)^k] \rightarrow \lambda^k$  as  $n \rightarrow \infty$  for  $k = 1, 2, 3$ , where the brief proof is given in Section 3.6.

**Lemma 3.1** *In the overspecified case,  $E[\hat{\lambda}] = E[\widehat{\lambda^2}] = E[\widehat{\lambda^3}] = 1$ . In the underspecified case,  $\hat{\lambda}$ ,  $\widehat{\lambda^2}$  and  $\widehat{\lambda^3}$  are asymptotically unbiased estimators of  $\lambda$ ,  $\lambda^2$  and  $\lambda^3$ , respectively.*

Using estimators (3.16), (3.17) and (3.18), we can estimate  $B_1$  and  $B_2$  in (3.13) and (3.14). However, because of  $B_1 = O(n)$ , a naive estimator that just substitutes  $\hat{\lambda}$  for  $\lambda$  in  $B_1$  has a bias

with order  $O(1)$ . Then  $E[\hat{\lambda}]$  can be expanded up to  $O(n^{-1})$  as

$$\begin{aligned} E[\hat{\lambda}] &= \frac{n - p_j}{n - p_\omega} E \left[ \frac{K_0}{K_0 + K_1} \right] \\ &= \lambda + \frac{-2\lambda^2(\lambda - 1) + p_j\lambda(\lambda - 1)}{n} + O(n^{-2}), \end{aligned} \quad (3.19)$$

where the proof is given in (3.33) in Section 3.6.

**Lemma 3.2** *Consider the following estimator for  $B_1$ :*

$$\widehat{B}_1 = \frac{2n(n - \text{tr}[\Sigma^{-1}])}{n - p_j - 2} \left\{ \hat{\lambda} - 1 + \frac{2(\widehat{\lambda}^3 - \widehat{\lambda}^2) - p_j(\widehat{\lambda}^2 - \hat{\lambda})}{n} \right\}. \quad (3.20)$$

Then, in the overspecified case,  $E[\widehat{B}_1] = 0$ , and in the underspecified case,  $E[\widehat{B}_1] = B_1 + O(n^{-1})$ .

We next obtain an estimator of  $\boldsymbol{\xi}^\top (\mathbf{M}_\omega - \mathbf{M}_j) \Sigma^{-1} (\mathbf{M}_\omega - \mathbf{M}_j) \boldsymbol{\xi}$  which is a part of  $B_3$ . Define  $\tilde{\sigma}_j^2$  by

$$\tilde{\sigma}_j^2 = (\mathbf{y} - \mathbf{X}(j)\widehat{\boldsymbol{\beta}}_j)^\top \Sigma^{-2} (\mathbf{y} - \mathbf{X}(j)\widehat{\boldsymbol{\beta}}_j).$$

From the fact that  $\tilde{\sigma}_j^2 = \sigma_*^2 (\mathbf{v} + \boldsymbol{\xi})^\top (\mathbf{I}_n - \mathbf{M}_j) \Sigma^{-1} (\mathbf{I}_n - \mathbf{M}_j) (\mathbf{v} + \boldsymbol{\xi})$ , it follows that

$$E[\tilde{\sigma}_j^2] = \sigma_*^2 \{ \text{tr}[\Sigma^{-1} - \mathbf{P}_j] + \boldsymbol{\xi}^\top (\mathbf{M}_\omega - \mathbf{M}_j) \Sigma^{-1} (\mathbf{M}_\omega - \mathbf{M}_j) \boldsymbol{\xi} \}.$$

Hence an estimator of  $\boldsymbol{\xi}^\top (\mathbf{M}_\omega - \mathbf{M}_j) \Sigma^{-1} (\mathbf{M}_\omega - \mathbf{M}_j) \boldsymbol{\xi}$  is given by

$$\tilde{\sigma}_j^2 / \hat{\sigma}_\omega^2 - \text{tr}[\Sigma^{-1} - \mathbf{P}_j].$$

**Lemma 3.3** *Consider the following estimator for  $B_3$ :*

$$\widetilde{B}_3 = \frac{4}{n} \left( \frac{\hat{\sigma}_\omega^2}{\hat{\sigma}_j^2} \right)^2 \times \left\{ \frac{\tilde{\sigma}_j^2}{\hat{\sigma}_\omega^2} - \text{tr}[\Sigma^{-1} - \mathbf{P}_j] \right\}.$$

Then in the overspecified case,  $E[\widetilde{B}_3] = O(n^{-1})$ . In the underspecified case,  $E[\widetilde{B}_3] = B_3 + O(n^{-1})$ .

Lemma 3.3 implies that in both overspecified and underspecified cases,  $\widetilde{B}_3$  is an asymptotically unbiased estimator of  $B_3$  up to  $O(1)$ , but  $\widetilde{B}_3$  has an  $O(n^{-1})$  bias that cannot be negligible for overspecified models. Since the cAIC by Vaida and Blanchard (2005) is an exact unbiased estimator of cAI, we want to adjust  $\widetilde{B}_3$  so that the adjusted estimator can have a bias with order  $O(n^{-2})$  in the overspecified case.

**Lemma 3.4** *For  $B_3$ , consider the following estimator as a higher order unbiased estimator than  $\widetilde{B}_3$ :*

$$\widehat{B}_3 = \widetilde{B}_3 - \frac{4p_\omega \cdot \text{tr}[\Sigma^{-1} - \mathbf{P}_j]}{n^2} + \frac{8\text{tr}[\mathbf{P}_\omega - \mathbf{P}_j]}{n^2}, \quad (3.21)$$

where  $\mathbf{P}_\omega = \Sigma^{-1} \mathbf{X}(\omega) (\mathbf{X}(\omega)^\top \Sigma^{-1} \mathbf{X}(\omega))^{-1} \mathbf{X}(\omega)^\top \Sigma^{-1}$ . Then, in the overspecified case,  $E[\widehat{B}_3] = O(n^{-2})$ . In the underspecified case,  $E[\widehat{B}_3] = B_3 + O(n^{-1})$ .

Using Lemmas 3.1, 3.2 and 3.4, we can estimate the bias correction  $\Delta_{\text{cAI}}$  by the estimator

$$\widehat{\Delta}_{\text{cAI}} = B^* + \widehat{B}_1 + \widehat{B}_2 + \widehat{B}_3. \quad (3.22)$$

The bias correction estimator can be used not only for overspecified models, but also for under-specified models. Thus, we get the modified conditional Akaike information criterion (McAIC) given as

$$\text{McAIC} = -2 \log f(\mathbf{y} | \widehat{\mathbf{b}}_j, \widehat{\boldsymbol{\beta}}_j, \widehat{\sigma}_j^2) + \widehat{\Delta}_{\text{cAI}}. \quad (3.23)$$

**Theorem 3.2** *In the overspecified case, it follows that*

$$E[\widehat{\Delta}_{\text{cAI}}] = \Delta_{\text{cAI}} + O(n^{-2}) \quad \text{and} \quad E[\text{McAIC}] = \text{cAI} + O(n^{-2}).$$

*In the underspecified case, it follows that*

$$E[\widehat{\Delta}_{\text{cAI}}] = \Delta_{\text{cAI}} + O(n^{-1}) \quad \text{and} \quad E[\text{McAIC}] = \text{cAI} + O(n^{-1}).$$

**Remark 3.1** In the derivation of McAIC, we assume that the covariance matrix of  $\mathbf{b}_*$  is  $\sigma_*^2 \mathbf{G}$  for a known matrix  $\mathbf{G}$ . This setup seems restrictive, since  $\sigma_*^2 \mathbf{G}$  involves some unknown variance components in most linear mixed models. For example, we consider the nested error regression model which will be treated in the next section for simulation. In this model,  $\mathbf{G}$  is a function of  $\psi_* = \tau_*^2 / \sigma_*^2$  where  $\tau_*^2$  is a variance component of random effects. We then propose to use a consistent estimator  $\widehat{\psi} = \widehat{\psi}(\mathbf{y})$  which satisfies  $\widehat{\psi}(\mathbf{y} + \mathbf{X}\boldsymbol{\alpha}) = \widehat{\psi}(\mathbf{y})$  for any  $p$ -dimensional vector  $\boldsymbol{\alpha}$ , and replace  $\mathbf{G}(\psi_*)$  in (3.23) by its plug-in estimator  $\mathbf{G}(\widehat{\psi})$ . For the location invariant property of  $\widehat{\psi}$ , the influence by replacement may be limited as long as we consider the problem of selecting only explanatory variables. Our recommendation is to estimate  $\psi_*$  by the full model. Estimating nuisance parameters by the full model is a methodology similar to Mallows'  $C_p$ .

## 3.4 Simulations

In this section, we investigate the behaviors of the suggested criterion McAIC by simulation through two kinds of experiments.

**[the first experiment]** In the first experiment, we compare the performance of the criteria cAIC and McAIC by measuring the biases of estimating cAI and the relative frequency of selecting the true model. It is important to note that Akaike-type criteria may be useful in getting the "best predictive model", but are not very suited for selecting the 'correct model' unlike the Bayesian information criterion (BIC) by Schwarz (1978). However, it has been reported that the performance of selecting the true model improves by modification or bias correction of Akaike-type criteria in Hurvich and Tsai (1989), Fujikoshi and Satoh (1997) and others. Thus it is meaningful to investigate the performance of selecting the true model, and we handle the experiment to measure the prediction error next.

We consider a class of the nested candidate models  $j_\alpha = \{1, \dots, \alpha\}$  for  $\alpha = 1, \dots, p_\omega$ . The observed vector  $\mathbf{y}$  is generated by the true model (3.2) with  $\mathbf{Z} = \text{diag}(\mathbf{Z}_1, \dots, \mathbf{Z}_q)$ ,  $\mathbf{Z}_i = \mathbf{1}_r$  for  $i = 1, \dots, q$  and  $\mathbf{G} = \mathbf{I}_q$ , where  $\mathbf{1}_r$  denotes an  $r \times 1$  vector of ones and  $r = n/q$ . Note that this model is known as the nested error regression model (NERM), and that  $q$  and  $r$  denote the number of clusters (or areas) and the sample size in each cluster, respectively. Let  $\mathbf{X}(\omega)$  be generated as  $\text{vec}(\mathbf{X}(\omega)^\text{T}) \sim \mathcal{N}(\mathbf{0}, \mathbf{I}_n \otimes \boldsymbol{\Sigma}_x)$  with  $\boldsymbol{\Sigma}_x = 0.9\mathbf{I}_{p_\omega} + 0.1\mathbf{J}_{p_\omega}$  where  $\mathbf{J}_{p_\omega} = \mathbf{1}_{p_\omega} \mathbf{1}_{p_\omega}^\text{T}$ . The true coefficient vector  $\boldsymbol{\beta}_*$  is given by  $\boldsymbol{\beta}_* = (\beta_1, \dots, \beta_{p_*}, 0, \dots, 0)^\text{T}$  and  $\beta_l$  is generated as

$\beta_l = 2 \times ((-1)^l / (l + 0.7)) \times U(1, 2)$ ,  $1 \leq l \leq p_*$  for a uniform random variable  $U(1, 2)$  on the interval  $(1, 2)$ .

We here handle the case that  $p_\omega = 7$ ,  $p_* = 5$ ,  $n = 50$  and  $q = 10$ . The true values of cAI in each model are calculated from (3.27) in Section 3.6 based on 10,000 replications. The biases and the rates of selecting each model are computed as their averages based on 1,000 replications. The results are shown in Table 3.1. In each row, the several values, which are explained below, for each corresponding model  $j_\alpha$  are reported. In the second column of the table, the true values of cAI in each model are reported. In the third and fourth columns, the averages of biases of estimating cAI by cAIC and McAIC are reported, respectively. In the fifth and sixth columns, the rates of selecting each model by cAIC and McAIC are reported, respectively.

model	cAI	bias		selection rates	
		cAIC	McAIC	cAIC	McAIC
<b>pattern (a): <math>\sigma_*^2 = 1</math></b>					
$j_1$	217.81	16.234	-0.63453	0	0
$j_2$	184.14	12.170	-0.55772	0	0
$j_3$	167.68	7.2864	-0.25035	0.005	0.025
$j_4$	163.66	4.7962	-0.15384	0.050	0.087
$j_5$	158.84	-0.16459	-0.13939	0.778	0.812
$j_6$	160.59	-0.22447	-0.14562	0.117	0.062
$j_7$	162.47	-0.23391	-0.15563	0.050	0.014
<b>pattern (b): <math>\sigma_*^2 = 0.5</math></b>					
$j_1$	210.62	18.465	-0.53668	0	0
$j_2$	168.35	16.200	-0.47607	0	0
$j_3$	142.90	11.564	-0.20047	0	0.001
$j_4$	134.82	8.4066	-0.083391	0.004	0.014
$j_5$	124.18	-0.16459	-0.13939	0.824	0.907
$j_6$	125.94	-0.22447	-0.14562	0.122	0.063
$j_7$	127.81	-0.23391	-0.15563	0.050	0.015

Table 3.1: Biases of estimating cAI by cAIC and McAIC, and the rates of selecting each model by cAIC and McAIC

Table 3.1 shows that although the conventional cAIC has large biases for underspecified models, namely model  $j_1$  to  $j_4$ , our proposed McAIC has smaller biases for both underspecified and overspecified models. Especially, because cAIC overestimates the cAI for underspecified cases, cAIC does not tend to select smaller models. The McAIC procedures also have better performances of selecting the true model than the conventional cAIC procedures do. The fact that McAIC can estimate with small biases for each model may imply that this criterion provides an appropriate weight vector for model averaging methods. We will check this hypothesis in the next experiment.

**[the second experiment]** In the second experiment, we investigate the prediction error of the best model chosen by cAIC and McAIC and that of model averaged predictor based on cAIC and McAIC. We here handle the case of unknown  $\mathbf{G}$  and consider the model class which consists of all subsets of  $\omega = \{1, \dots, p_\omega\}$ . The other setup is the same as that in the first experiment except for  $\mathbf{G} = \mathbf{G}(\psi_*) = \psi_* \mathbf{I}_q$  where  $\psi_* = \tau_*^2 / \sigma_*^2$ , namely  $\mathbf{b}_* \sim \mathcal{N}(\mathbf{0}, \tau_*^2 \mathbf{I}_q)$ .  $\sigma_*^2$  and  $\tau_*^2$  are estimated by unbiased estimators proposed by Prasad and Rao (1990) under the full model,

which are explained as follows. Let  $S = \mathbf{y}^T \{ \mathbf{I}_n - \mathbf{X}(\omega)(\mathbf{X}(\omega)^T \mathbf{X}(\omega))^{-1} \mathbf{X}(\omega)^T \} \mathbf{y}$  and  $S_1 = \mathbf{y}^T \{ \mathbf{E} - \mathbf{E} \mathbf{X}(\omega)(\mathbf{X}(\omega)^T \mathbf{E} \mathbf{X}(\omega))^{-1} \mathbf{X}(\omega)^T \mathbf{E} \} \mathbf{y}$  where  $\mathbf{E} = \text{diag}(\mathbf{E}_1, \dots, \mathbf{E}_q)$ ,  $\mathbf{E}_i = \mathbf{I}_r - r^{-1} \mathbf{J}_r$  for  $i = 1, \dots, q$ . Then, the Prasad–Rao estimators of  $\sigma_*^2$  and  $\tau_*^2$  are given by

$$\hat{\sigma}^{2\text{PR}} = S_1 / (n - q - p_\omega) \quad \text{and} \quad \hat{\tau}^{2\text{PR}} = \{ S - (n - p_\omega) \hat{\sigma}^{2\text{PR}} \} / n^*, \quad (3.24)$$

where  $n^* = n - \text{tr}[\mathbf{Z}^T \mathbf{X}(\omega)(\mathbf{X}(\omega)^T \mathbf{X}(\omega))^{-1} \mathbf{X}(\omega)^T \mathbf{Z}]$ . We estimate  $\psi_*$  by  $\hat{\psi} = \hat{\tau}^{2\text{PR}} / \hat{\sigma}^{2\text{PR}}$  and use the plug-in value of  $\mathbf{G}(\hat{\psi})$  for cAIC and McAIC. This  $\hat{\psi}$  satisfies the properties of consistency and location invariance mentioned in the Remark 3.1. It is important to point out that we use the ML estimator  $\hat{\sigma}_j^2$  or  $\hat{\sigma}_\omega^2$  in cAIC or McAIC given in (3.3), though we use  $\hat{\sigma}^{2\text{PR}}$  to estimate  $\psi_*$  by substituting  $\Sigma(\hat{\psi})$  for  $\Sigma$ . We measure the performance of cAIC and McAIC via  $\|\hat{\mathbf{y}}_j - \mathbf{X}(\omega)\boldsymbol{\beta}_* - \mathbf{Z}\mathbf{b}_*\|^2/n$  for  $\hat{\mathbf{y}}_j = \mathbf{X}(j)\hat{\boldsymbol{\beta}}_j + \mathbf{Z}\hat{\mathbf{b}}_j$ , which is here called the prediction error because  $\hat{\mathbf{b}}_j$  is a predictor of  $\mathbf{b}_*$ . The prediction errors are given as averages based on 1,000 replications.

In addition to the procedures that select the best model by the cAIC and McAIC, we consider a model averaging method. The aim of model averaging is to predict a future value by a weighted mean of the predictors for the candidate models. The weighting functions are important in the model averaging method, and we use optimal weights suggested in Burnham and Anderson (2002). In the context of McAIC, the weight is defined as follows. Let  $\text{McAIC}_j$  denote the value of McAIC in the model  $j$  and let  $\text{McAIC}_{\min}$  be the minimum McAIC value. Also, let  $\overline{\text{McAIC}}_j = \text{McAIC}_j - \text{McAIC}_{\min}$ . Then the weight is defined by

$$w_j = \frac{\exp(-\overline{\text{McAIC}}_j/2)}{\sum_l \exp(-\overline{\text{McAIC}}_l/2)}. \quad (3.25)$$

Based on the weights given in (3.25), we can obtain a model averaged predictor

$$\hat{\mathbf{y}}_{\text{Ave}} = \sum_j w_j \hat{\mathbf{y}}_j,$$

where  $\hat{\mathbf{y}}_j$  is the predictor based on model  $j$ , and the summation is taken over all the candidate models. We call this method ‘Smoothed McAIC (S-McAIC)’. A similar method based on cAIC is called ‘Smoothed cAIC (S-cAIC)’.

Table 3.2 reports the prediction errors for the best model selected by cAIC and McAIC and for the model averaged predictors based on S-cAIC and S-McAIC for the six cases in which  $n$ ,  $\sigma_*^2$  and  $\tau_*^2$  take different values. We set  $p_\omega = 5$  and  $p_* = 3$ . The values in parentheses are the improvement over the prediction error by the cAIC procedure expressed in percentage. From the table, it can be seen that McAIC and the corresponding averaging procedure S-McAIC are better than cAIC and S-cAIC. Also, it is revealed that the prediction errors get smaller as the sample size is larger. This implies that the information criteria can estimate the cAI more accurately for the large sample size.

### 3.5 Real data example

We apply the variable selection procedures cAIC and McAIC to the posted land price data along the Keikyū train line, which connects the suburbs in Kanagawa prefecture to the Tokyo metropolitan area. This data set was used by Kubokawa and Nagashima (2012), who studied parametric bootstrap methods in the linear mixed models.

	cAIC	McAIC	S-cAIC	S-McAIC
<b>pattern (a)</b>	0.23198	0.22628	0.22581	0.22240
$N = 50, \sigma_*^2 = 1.0, \tau_*^2 = 0.5$		(2.46)	(2.66)	(4.13)
<b>pattern (b)</b>	0.13107	0.12818	0.12837	0.12629
$N = 50, \sigma_*^2 = 0.5, \tau_*^2 = 1.0$		(2.20)	(2.06)	(3.65)
<b>pattern (c)</b>	0.12516	0.12231	0.12234	0.12040
$N = 50, \sigma_*^2 = 0.5, \tau_*^2 = 0.5$		(2.28)	(2.25)	(3.80)
<b>pattern (d)</b>	0.15561	0.15344	0.15091	0.14930
$N = 80, \sigma_*^2 = 1.0, \tau_*^2 = 0.5$		(1.39)	(3.02)	(4.06)
<b>pattern (e)</b>	0.083838	0.082674	0.081682	0.080744
$N = 80, \sigma_*^2 = 0.5, \tau_*^2 = 1.0$		(1.39)	(2.57)	(3.69)
<b>pattern (f)</b>	0.081538	0.080418	0.079368	0.078477
$N = 80, \sigma_*^2 = 0.5, \tau_*^2 = 0.5$		(1.37)	(2.66)	(3.75)

Table 3.2: Prediction errors of the predictors based on cAIC, McAIC, S-cAIC and S-McAIC, and improvement over cAIC procedure

We analyze the land price data in 2001 with covariates for 47 stations which we consider as small areas and let  $q = 47$ . For the  $i$ th small area, there are data of  $n_i$  land spots, and the total sample size is  $n = \sum_{i=1}^q n_i = 189$ . The land price (Yen in hundreds of thousands) per  $m^2$  of the  $k$  spot in the  $i$ th small area is denoted by  $y_{ik}$ ,  $\text{TRN}_i$  is the time to take by train from the station  $i$  to the Tokyo station around 9:00 in the morning,  $\text{DST}_{ik}$  is the geographical distance from the spot  $k$  to the nearby station  $i$ ,  $\text{FOOT}_{ik}$  is the time to take on foot from the spot  $k$  to the nearby station  $i$  and  $\text{FAR}_{ik}$  denotes the floor-area ratio of the spot  $k$ . As explanatory variables, we consider nine variables  $\text{FAR}_{ik}$ ,  $\text{TRN}_i$ ,  $\text{TRN}_i^2$ ,  $\text{DST}_{ik}$ ,  $\text{DST}_{ik}^2$ ,  $\text{FOOT}_{ik}^2$ ,  $\text{TRN}_i \times \text{DST}_{ik}$  and  $\text{TRN}_i \times \text{FOOT}_{ik}$ , which are denoted by  $x_1, \dots, x_9$  and  $x_0$  denotes constant term.

When one wants to estimate the mean land prices of each small area, naive estimators such as sample mean may have large variabilities because the sample sizes in each small area are small. Then we use some model based small area estimation technique, which introduce random effects and borrow the strength from the information of other areas. For the details, see Rao and Molina (2015). One example is the nested error regression model (NERM), which we handled in the second experiment of the previous section. We here employ NERM to estimate the mean land prices in the places around each station. In order to specify the model, or to select the explanatory variables from 10 covariates  $x_0, \dots, x_9$ , we use the variable selection criteria cAIC and McAIC. The procedure is that regressors which minimizes the information criteria are added to the model known as the forward stepwise selection. The unknown parameter  $\psi_* = \tau_*^2 / \sigma_*^2$  is estimated by  $\hat{\tau}^{2\text{PR}} / \hat{\sigma}^{2\text{PR}}$  given in (3.24) by the full model.

Table 3.3 reports the values of cAIC and McAIC in each model. Both criteria select the same model  $j = \{1, 0, 2, 3\}$ , namely,

$$y_{ik} = \beta_0 + \text{FAR}_i \beta_1 + \text{TRN}_i \beta_2 + (\text{TRN}_i)^2 \beta_3 + b_i + \varepsilon_{ik},$$

where the parameters are estimated by  $\hat{\sigma}^{2\text{PR}} = 0.46382$ ,  $\hat{\tau}^{2\text{PR}} = 0.08199$ , i.e.,  $\hat{\psi} = 0.17722$  and  $(\hat{\beta}_0, \hat{\beta}_1, \hat{\beta}_2, \hat{\beta}_3) = (5.0823, 6.3790 \times 10^{-3}, -1.0561 \times 10^{-1}, 6.4938 \times 10^{-4})$ ,  $\hat{\sigma}_\omega^2 = 0.43889$ ,  $\hat{\sigma}_j^2 = 0.44465$ . This result demonstrates that the most important factor to decide land prices is the floor area ratio, which is very natural, and that the land prices are decreasing as a quadratic function of the time to take by train to the metropolitan area. It is interesting to see that the



model	cAIC	McAIC
{1}	531.89	512.08
{1, 0}	482.63	467.98
{1, 0, 2}	419.18	417.31
{1, 0, 2, 3}	410.07*	410.92*
{1, 0, 2, 3, 4}	412.21	413.01
{1, 0, 2, 3, 4, 5}	413.09	413.79
{1, 0, 2, 3, 4, 5, 8}	415.02	415.54
{1, 0, 2, 3, 4, 5, 8, 6}	417.14	417.60
{1, 0, 2, 3, 4, 5, 8, 6, 9}	419.02	419.23
$\omega$	420.97	420.98

Table 3.3: The values of cAIC and McAIC in each model in the example of posted land price data

values of cAIC is bigger than that of McAIC for the models {1}, {1, 0} and {1, 0, 2} and both criteria take similar values for the models bigger than {1, 0, 2, 3}, which we take as the ‘true’ model. This result seems to consistent with the demonstration of the first numerical experiment in the previous section, in which cAIC tends to overestimate cAI for the underspecified models.

## 3.6 Proofs

### 3.6.1 Derivation of (3.6)

First compute the expectation with respect to  $\tilde{\mathbf{y}} \sim f(\tilde{\mathbf{y}}|\mathbf{b}_*, \boldsymbol{\eta}_*)$  in cAI. Then, cAI can be written as

$$\begin{aligned} \text{cAI} &= E[n \log(2\pi\hat{\sigma}_j^2) + n\sigma_*^2/\hat{\sigma}_j^2] \\ &\quad + E\left[\|(\mathbf{X}(j)\hat{\boldsymbol{\beta}}_j - \mathbf{X}(\omega)\boldsymbol{\beta}_*) + \mathbf{Z}(\hat{\mathbf{b}}_j - \mathbf{b}_*)\|^2/\hat{\sigma}_j^2\right]. \end{aligned} \quad (3.26)$$

Note that  $\mathbf{b}_*|\mathbf{y} \sim \mathcal{N}(\mathbf{G}\mathbf{Z}^T\boldsymbol{\Sigma}^{-1}\mathbf{u}, \sigma_*^2(\mathbf{G} - \mathbf{G}\mathbf{Z}^T\boldsymbol{\Sigma}^{-1}\mathbf{Z}\mathbf{G}))$  and that

$$\mathbf{X}(j)\hat{\boldsymbol{\beta}}_j - \mathbf{X}(\omega)\boldsymbol{\beta}_* + \mathbf{Z}(\hat{\mathbf{b}}_j - \mathbf{b}_*) = \boldsymbol{\Sigma}^{-1}(\mathbf{X}(j)\hat{\boldsymbol{\beta}}_j - \mathbf{X}(\omega)\boldsymbol{\beta}_*) - \mathbf{Z}(\mathbf{b}_* - \mathbf{G}\mathbf{Z}^T\boldsymbol{\Sigma}^{-1}\mathbf{u}).$$

Taking the expectation with respect to  $\mathbf{b}_*|\mathbf{y} \sim \pi(\mathbf{b}_*|\mathbf{y}, \boldsymbol{\eta}_*)$  in (3.26), we rewrite the cAI as

$$\begin{aligned} \text{cAI} &= E[n \log(2\pi\hat{\sigma}_j^2) + (2n - \text{tr}[\boldsymbol{\Sigma}^{-1}])\sigma_*^2/\hat{\sigma}_j^2] \\ &\quad + E[(\mathbf{X}(j)\hat{\boldsymbol{\beta}}_j - \mathbf{X}(\omega)\boldsymbol{\beta}_*)^T \boldsymbol{\Sigma}^{-2} (\mathbf{X}(j)\hat{\boldsymbol{\beta}}_j - \mathbf{X}(\omega)\boldsymbol{\beta}_*)/\hat{\sigma}_j^2]. \end{aligned} \quad (3.27)$$

Next, for a part of  $-2 \log f(\mathbf{y}|\hat{\mathbf{b}}_j, \hat{\boldsymbol{\eta}}_j)$  in (3.5), it is noted that

$$\begin{aligned} &\mathbf{y} - \mathbf{X}(j)\hat{\boldsymbol{\beta}}_j - \mathbf{Z}\hat{\mathbf{b}}_j \\ &= \mathbf{y} - \mathbf{X}(\omega)\boldsymbol{\beta}_* - (\mathbf{X}(j)\hat{\boldsymbol{\beta}}_j - \mathbf{X}(\omega)\boldsymbol{\beta}_*) - \mathbf{Z}\mathbf{G}\mathbf{Z}^T\boldsymbol{\Sigma}^{-1} \left\{ \mathbf{y} - \mathbf{X}(\omega)\boldsymbol{\beta}_* - (\mathbf{X}(j)\hat{\boldsymbol{\beta}}_j - \mathbf{X}(\omega)\boldsymbol{\beta}_*) \right\} \\ &= \boldsymbol{\Sigma}^{-1}\mathbf{u} - \boldsymbol{\Sigma}^{-1}(\mathbf{X}(j)\hat{\boldsymbol{\beta}}_j - \mathbf{X}(\omega)\boldsymbol{\beta}_*). \end{aligned} \quad (3.28)$$

Thus, from (3.27) and (3.28), we can see that  $\Delta_{\text{cAI}} = \text{cAI} - E[-2 \log f(\mathbf{y}|\hat{\mathbf{b}}_j, \hat{\boldsymbol{\eta}}_j)]$  is expressed as (3.6).  $\square$

### 3.6.2 Proof of Theorem 3.1

For (3.10), we decompose  $\Delta_{\text{cAI}}$  as

$$\Delta_{\text{cAI}} = b_1 + b_2 + b_3 + b_4,$$

where  $b_1 = nE[(2n - \text{tr } \Sigma^{-1})/(K_0 + K_1)]$ ,  $b_2 = -nE[\mathbf{v}^T \Sigma^{-1} \mathbf{v}/(K_0 + K_1)]$ ,  $b_3 = 2nE[\mathbf{v}^T \Sigma^{-1} \mathbf{M}_j \mathbf{v}/(K_0 + K_1)]$  and  $b_4 = -2nE[\boldsymbol{\xi}^T (\mathbf{M}_\omega - \mathbf{M}_j) \Sigma^{-1} \mathbf{v}/(K_0 + K_1)]$ .

We begin with expanding  $(K_0 + K_1)^{-1}$  up to  $O_p(n^{-2})$ . Let  $\mathbf{v} = \mathbf{v}^T (\mathbf{I}_n - \mathbf{M}_\omega) \mathbf{v} + \mathbf{v}^T (\mathbf{M}_\omega - \mathbf{M}_j) \mathbf{v}$ . Then,  $K_0 + K_1 = K + 2L + n\delta$ ,  $K \sim \chi_{n-p_j}^2$ ,  $K = O_p(n)$  and  $L = O_p(n^{1/2})$ , so that we can write  $(2L + n\delta)/K = \delta + D$  and  $D = O_p(n^{-1/2})$ . Thus, it follows that

$$\begin{aligned} (K_0 + K_1)^{-1} &= (K + 2L + n\delta)^{-1} = K^{-1}(1 + \delta + D)^{-1} = \frac{\lambda}{K}(1 + D\lambda)^{-1} \\ &= \frac{\lambda}{K} \left\{ 1 - D\lambda + (D\lambda)^2 + O_p(n^{-3/2}) \right\}. \end{aligned} \quad (3.29)$$

Since  $\delta\lambda = 1 - \lambda$ , it is seen that

$$D\lambda = \left(1 - \frac{n}{K}\right)(\lambda - 1) + \frac{2L\lambda}{K}.$$

Let  $A = W/\{(m - p_j) - 1\}$ , which is of  $O_p(n^{-1/2})$ . then,

$$\begin{aligned} \frac{1}{K} &= \frac{1}{n - p_j} \left\{ 1 - A + A^2 + O_p(n^{-3/2}) \right\} \\ &= \frac{1}{n} \left( 1 - A + A^2 + \frac{p_j}{n} \right) + O_p(n^{-5/2}), \\ 1 - \frac{n}{K} &= A - A^2 - \frac{p_j}{n} + O_p(n^{-3/2}). \end{aligned} \quad (3.30)$$

Hence,  $(K_0 + K_1)^{-1}$  can be evaluated as

$$(K_0 + K_1)^{-1} = \frac{\lambda}{n} \left\{ 1 - \left(A + \frac{2L}{n}\right)\lambda + A^2\lambda^2 + \frac{p_j\lambda + 4AL\lambda^2}{n} + \frac{4L^2\lambda^2}{n^2} \right\} + O_p(n^{-5/2}). \quad (3.31)$$

For any function  $q(\cdot)$ , we have  $E[q(\mathbf{v}^T \mathbf{G} \mathbf{v})L] = 0$  since  $q(\mathbf{v}^T \mathbf{G} \mathbf{v})L$  is an odd function of  $\mathbf{v}$ . Also,  $E[A] = 0$ ,  $E[A^2] = 2/(n - p_j)$ ,  $E[L^2] = n\delta$ . Hence, it is observed that

$$E[(K_0 + K_1)^{-1}] = \frac{\lambda}{n} \left\{ 1 + \frac{-2\lambda^2 + (p_j + 4)\lambda}{n} \right\} + O(n^{-3}). \quad (3.32)$$

Using the expansions (3.31) and (3.32), we can evaluate  $b_1$ ,  $b_2$ ,  $b_3$  and  $b_4$ , respectively.

First,  $b_1$  can be evaluated as

$$\begin{aligned} b_1 &= n(2n - \text{tr } [\Sigma^{-1}]) \times \frac{\lambda}{n} \times \left\{ 1 + \frac{-2\lambda^2 + (p_j + 4)\lambda}{n} \right\} + O(n^{-1}) \\ &= b_1^* + (2n - \text{tr } [\Sigma^{-1}])(\lambda - 1) \left\{ \frac{n}{n - p_j - 2} + \frac{-2\lambda^2 + (p_j + 2)\lambda}{n} \right\} + O(n^{-1}), \end{aligned}$$

where

$$b_1^* = \frac{n(2n - \text{tr } [\Sigma^{-1}])}{n - p_j - 2},$$

which is the exact  $b_1$  for overspecified models.

Next note that  $\text{tr}[\boldsymbol{\Sigma}^{-1}] = O(n)$ ,  $\mathbf{v}^\top \boldsymbol{\Sigma}^{-1} \mathbf{v} - \text{tr}[\boldsymbol{\Sigma}^{-1}] = O_p(n^{1/2})$ . Then,  $b_2$  is evaluated as

$$\begin{aligned} b_2 &= -nE \left[ \left\{ \text{tr}[\boldsymbol{\Sigma}^{-1}] + (\mathbf{v}^\top \boldsymbol{\Sigma}^{-1} \mathbf{v} - \text{tr}[\boldsymbol{\Sigma}^{-1}]) \right\} \right. \\ &\quad \times \left. \frac{\lambda}{n} \left\{ 1 - \left(A + \frac{2L}{n}\right)\lambda + A^2\lambda^2 + \frac{p_j\lambda + 4AL\lambda^2}{n} + \frac{4L^2\lambda^2}{n^2} \right\} \right] + O(n^{-1}) \\ &= -\lambda \text{tr}[\boldsymbol{\Sigma}^{-1}] \left\{ 1 + \frac{-2\lambda^2 + (p_j + 4)\lambda}{n} \right\} + \lambda^2 E[\mathbf{v}^\top \boldsymbol{\Sigma}^{-1} \mathbf{v} A] + O(n^{-1}). \end{aligned}$$

From the second order moment of quadratic forms of standard normal random vectors, it follows that

$$E[\mathbf{v}^\top \boldsymbol{\Sigma}^{-1} \mathbf{v} A] = \frac{2}{n - p_j} \left\{ \text{tr}[\boldsymbol{\Sigma}^{-1}] - \text{tr}[\boldsymbol{\Sigma}^{-1} \mathbf{M}_j] \right\}.$$

Using this equality, we can evaluate  $b_2$  as

$$\begin{aligned} b_2 &= -\lambda \text{tr}[\boldsymbol{\Sigma}^{-1}] \left\{ 1 + \frac{-2\lambda^2 + (p_j + 2)\lambda}{n} \right\} + O(n^{-1}) \\ &= b_2^* - \text{tr}[\boldsymbol{\Sigma}^{-1}] (\lambda - 1) \left\{ \frac{n}{n - p_j - 2} + \frac{-2\lambda^2 + p_j\lambda - 2}{n} \right\} + O(n^{-1}), \end{aligned}$$

where

$$b_2^* = -n \times \left\{ \frac{\text{tr}[\boldsymbol{\Sigma}^{-1}]}{n - p_j - 2} - \frac{2\text{tr}[\boldsymbol{\Sigma}^{-1}] - 2\text{tr}[\boldsymbol{\Sigma}^{-1} \mathbf{M}_j]}{(n - p_j)(n - p_j - 2)} \right\},$$

which is the exact  $b_2$  for overspecified models.

As for  $b_3$ , it can be decomposed as

$$\frac{\mathbf{v}^\top \boldsymbol{\Sigma}^{-1} \mathbf{M}_j \mathbf{v}}{K_0 + K_1} = \frac{\mathbf{v}^\top \mathbf{M}_j \boldsymbol{\Sigma}^{-1} \mathbf{M}_j \mathbf{v}}{K_0 + K_1} + \frac{\mathbf{v}^\top (\mathbf{I}_n - \mathbf{M}_j) \boldsymbol{\Sigma}^{-1} \mathbf{M}_j \mathbf{v}}{K_0 + K_1}.$$

Since  $\mathbf{M}_j \mathbf{v}$  is independent of  $(\mathbf{I}_n - \mathbf{M}_j) \mathbf{v}$  and  $K_0 + K_1$ , and  $E[\mathbf{M}_j \mathbf{v}] = 0$ , it follows that

$$E \left[ \frac{\mathbf{v}^\top (\mathbf{I}_n - \mathbf{M}_j) \boldsymbol{\Sigma}^{-1} \mathbf{M}_j \mathbf{v}}{K_0 + K_1} \right] = 0.$$

Further, because  $\mathbf{v}^\top \mathbf{M}_j \boldsymbol{\Sigma}^{-1} \mathbf{M}_j \mathbf{v}$  and  $K_0 + K_1$  are mutually independent,  $b_3$  is evaluated as

$$\begin{aligned} b_3 &= 2n \times E[\mathbf{v}^\top \mathbf{M}_j \boldsymbol{\Sigma}^{-1} \mathbf{M}_j \mathbf{v}] \times E[(K_0 + K_1)^{-1}] \\ &= b_3^* + 2\text{tr}[\boldsymbol{\Sigma}^{-1} \mathbf{M}_j] (\lambda - 1) + O(n^{-1}), \end{aligned}$$

where

$$b_3^* = \frac{2n \text{tr}[\boldsymbol{\Sigma}^{-1} \mathbf{M}_j]}{n - p_j - 2},$$

which is the exact  $b_3$  for overspecified models.

Finally, we evaluate  $b_4$ . Note that  $\boldsymbol{\xi}^\top (\mathbf{M}_\omega - \mathbf{M}_j) \boldsymbol{\Sigma}^{-1} \mathbf{v} = O_p(n^{1/2})$  from the assumption (A2). Then,  $b_4$  can be expanded as

$$\begin{aligned} b_4 &= -2n \times E \left[ \boldsymbol{\xi}^\top (\mathbf{M}_\omega - \mathbf{M}_j) \boldsymbol{\Sigma}^{-1} \mathbf{v} \times \frac{\lambda}{n} \left\{ 1 - \left(A + \frac{2L}{n}\right)\lambda \right\} \right] + O(n^{-1}) \\ &= \frac{4\lambda^2}{n} \boldsymbol{\xi}^\top (\mathbf{M}_\omega - \mathbf{M}_j) \boldsymbol{\Sigma}^{-1} (\mathbf{M}_\omega - \mathbf{M}_j) \boldsymbol{\xi} + O(n^{-1}). \end{aligned}$$

Combining the evaluations of  $b_1$ ,  $b_2$ ,  $b_3$  and  $b_4$  yields the result in (3.12), where  $B^*$  is defined by  $B^* = b_1^* + b_2^* + b_3^*$ .  $\square$

### 3.6.3 Proof of Lemma 3.1

It follows from  $K_0 = n + O_p(n^{1/2})$  and (3.31) that

$$E \left[ \left( \frac{\hat{\sigma}_\omega^2}{\hat{\sigma}_j^2} \right)^k \right] = E \left[ \left( \frac{K_0}{K_0 + K_1} \right)^k \right] \rightarrow \lambda^k \quad (n \rightarrow \infty),$$

which proves lemma 3.1.  $\square$

### 3.6.4 Proof of Lemma 3.2

In the overspecified case, it follows from Lemma 3.1 that  $E[\hat{\lambda}] = E[\widehat{\lambda^2}] = E[\widehat{\lambda^3}] = 1$ , so that  $E[\widehat{B}_1] = 0$ . In the underspecified case, we shall check (3.19). Using the expansion (3.31) of  $(K_0 + K_1)^{-1}$ , we can approximate  $E[\hat{\lambda}]$  as

$$\begin{aligned} E[\hat{\lambda}] &= \frac{n - p_j}{n - p_\omega} \times E \left[ \frac{K_0}{K_0 + K_1} \right] \\ &= \frac{n - p_j}{n - p_\omega} \frac{\lambda}{n} \times E \left[ \{(n - p_\omega) + \mathbf{v}^T(\mathbf{I}_n - \mathbf{M}_\omega)\mathbf{v} - (n - p_\omega)\} \right. \\ &\quad \left. \times \left\{ 1 - \left( A + \frac{2L}{n} \right) \lambda + A^2 \lambda^2 + \frac{p_j \lambda + 4AL\lambda^2}{n} + \frac{4L^2 \lambda^2}{n^2} \right\} \right] + O(n^{-2}). \end{aligned} \quad (3.33)$$

Evaluating (3.33) up to  $O(n^{-1})$ , we can get (3.19) and Lemma 3.2.  $\square$

### 3.6.5 Proof of Lemma 3.3

Let  $c_1 = \text{tr}[\Sigma^{-1}(\mathbf{I}_n - \mathbf{M}_j)]$ ,  $c_2 = \text{tr}[\Sigma^{-1}(\mathbf{I}_n - \mathbf{M}_\omega)]$  and

$$\begin{aligned} D_1 &= \mathbf{v}^T(\mathbf{I}_n - \mathbf{M}_j)\Sigma^{-1}(\mathbf{I}_n - \mathbf{M}_j)\mathbf{v}, \\ D_2 &= 2\xi^T(\mathbf{I}_n - \mathbf{M}_j)\Sigma^{-1}(\mathbf{I}_n - \mathbf{M}_j)\mathbf{v}, \\ D_3 &= \xi^T(\mathbf{M}_\omega - \mathbf{M}_j)\Sigma^{-1}(\mathbf{M}_\omega - \mathbf{M}_j)\xi. \end{aligned} \quad (3.34)$$

Since  $\tilde{\sigma}_j^2 = \sigma_*^2(D_1 + D_2 + D_3)$ , we can rewrite  $\widetilde{B}_3$  as

$$\begin{aligned} \widetilde{B}_3 &= \frac{4K_0(D_1 + D_2 + D_3)}{(K_0 + K_1)^2} - \frac{4c_1}{n} \frac{K_0^2}{(K_0 + K_1)^2} \\ &= \widetilde{B}_{31} - \widetilde{B}_{32}. \quad (\text{say}) \end{aligned}$$

From the expansion (3.31) of  $(K_0 + K_1)^{-1}$ , it follows that

$$\begin{aligned} E[\widetilde{B}_{31}] &= \frac{4\lambda^2}{n^2} \times \{E[K_0 D_1] + E[K_0 D_3]\} + O(n^{-1}) \\ &= \frac{4c_1 \lambda^2}{n} + \frac{4\lambda^2 D_3}{n} + O(n^{-1}), \\ E[\widetilde{B}_{32}] &= \frac{4c_1}{n} \times \frac{\lambda^2}{n^2} \times E[K_0^2] + O(n^{-1}) \\ &= \frac{4c_1 \lambda^2}{n} + O(n^{-1}), \end{aligned}$$

which proves Lemma 3.3.  $\square$

### 3.6.6 Proof of Lemma 3.4

It is noted that the adjustment term of  $\widehat{B}_3$  is of order  $O(n^{-1})$  from (3.21). Then it follows from Lemma 3.3 that  $E[\widehat{B}_3] = B_3 + O(n^{-1})$ . Thus it is sufficient to evaluate  $E[\widetilde{B}_3]$  up to  $O(n^{-1})$  for overspecified models.

In the overspecified case, it is noted that  $n\hat{\sigma}_j^2 = \sigma_*^2 K$  and  $\tilde{\sigma}_j^2 = \sigma_*^2 D_1$  for  $D_1$  given in (3.34). Then,

$$\begin{aligned}\widetilde{B}_3 &= 4 \times \frac{K_0 \times D_1}{K^2} - \frac{4c_1}{n} \left( \frac{K_0}{K} \right)^2 \\ &= \widetilde{B}_{33} - \widetilde{B}_{34}. \quad (\text{say})\end{aligned}$$

From (3.30),  $K^{-2}$  is expanded as

$$\frac{1}{K^2} = \frac{1}{n^2} \left( 1 - 2A + 3A^2 + \frac{2p_j}{n} \right) + O_p(n^{-7/2}).$$

Thus,  $E[\widetilde{B}_{33}]$  is written as

$$\begin{aligned}E[\widetilde{B}_{33}] &= \frac{4}{n^2} E \left[ \{ (n - p_\omega) + K_0 - (n - p_\omega) \} \{ c_1 + (D_1 - c_1) \} \right. \\ &\quad \left. \times \left\{ 1 - 2A + 3A^2 + \frac{2p_j}{n} \right\} \right] + O(n^{-2}) \\ &= \frac{4}{n^2} \times \left\{ (n - p_\omega) c_1 \left( 1 + 3E[A^2] + \frac{2p_j}{n} \right) \right. \\ &\quad \left. - 2c_1 E[K_0 A] - 2(n - p_\omega) E[D_1 A] + E[K_0 D_1] - c_1 (n - p_\omega) \right\} + O(n^{-2}) \\ &= \frac{4(n - p_\omega + 2p_j)c_1}{n^2} + \frac{8(c_2 - c_1)}{n^2} + O(n^{-2}),\end{aligned}\tag{3.35}$$

since  $K_0 - (n - p_\omega) = O_p(n^{1/2})$ ,  $c_1 = O(n)$  and  $D_1 - c_1 = O_p(n^{1/2})$ . Noting that  $K_0/K \sim Be((n - p_\omega)/2, (p_\omega - p_j)/2)$ , we can evaluate  $E[\widetilde{B}_{34}]$  as

$$\begin{aligned}E[\widetilde{B}_{34}] &= \frac{4c_1}{n} \frac{(n - p_\omega)(n - p_\omega + 2)}{(n - p_j)(n - p_j + 2)} \\ &= \frac{4(n - 2p_\omega + 2p_j)c_1}{n^2} + O(n^{-2}).\end{aligned}\tag{3.36}$$

Combining (3.35) and (3.36) gives

$$E[\widetilde{B}_3] = \frac{4c_1 p_\omega}{n^2} + \frac{8(c_2 - c_1)}{n^2} + O(n^{-2}),$$

which shows Lemma 3.4. □



## Chapter 4

# Conditional AIC under covariate shift with application to small area prediction

In this chapter, we consider the problem of selecting explanatory variables of fixed effects in linear mixed models under covariate shift, which is the situation that the values of covariates in the predictive model are different from those in the observed model. We construct a variable selection criterion based on the conditional Akaike information introduced by Vaida and Blanchard (2005) and the proposed criterion is generalization of the conditional AIC in terms of covariate shift. We especially focus on covariate shift in small area prediction and show usefulness of the proposed criterion.

### 4.1 Motivation

For prediction problems, it is often the case that the values of covariates in the predictive model are different from those in the observed model, which we call covariate shift. However, even when the information about the covariates in the predictive model can be used, most Akaike-type criteria do not use it. This is because most criteria put the assumption that the predictive model is the same as the observed model. We explain what this means in the context of the conditional AIC.

Vaida and Blanchard (2005) proposed the conditional AIC (cAIC) as an unbiased estimator of the cAI, which is given by

$$\text{cAI} = \iiint -2 \log\{f(\tilde{\mathbf{y}}|\hat{\boldsymbol{\theta}}, \hat{\boldsymbol{\eta}})\} f(\tilde{\mathbf{y}}|\boldsymbol{\theta}, \boldsymbol{\eta}) f(\mathbf{y}|\boldsymbol{\theta}, \boldsymbol{\eta}) p(\boldsymbol{\theta}|\boldsymbol{\eta}) d\tilde{\mathbf{y}} d\mathbf{y} d\boldsymbol{\theta}, \quad (4.1)$$

in general, where the notations are the same as those in Chapter 2. We call the model in which  $\mathbf{y}$  is the vector of the response variables the ‘observed model’, and call the model in which  $\tilde{\mathbf{y}}$  is the vector of the response variables the ‘predictive model’. The cAI in (4.1) assumes that the conditional density of  $\mathbf{y}$  given  $\boldsymbol{\theta}$  and that of  $\tilde{\mathbf{y}}$  given  $\boldsymbol{\theta}$  are the same and both of them are denoted by  $f(\cdot|\boldsymbol{\theta}, \boldsymbol{\eta})$ , which implies that the observed model and the predictive model are the same. However, under covariate shift, the conditional density of  $\tilde{\mathbf{y}}$  given  $\boldsymbol{\theta}$  should be different from that of  $\mathbf{y}$  given  $\boldsymbol{\theta}$  and is denoted by  $g(\tilde{\mathbf{y}}|\boldsymbol{\theta}, \boldsymbol{\eta})$ . Then, we redefine the cAI under covariate shift as follows:

$$\text{cAI} = \iiint -2 \log\{g(\tilde{\mathbf{y}}|\hat{\boldsymbol{\theta}}, \hat{\boldsymbol{\eta}})\} g(\tilde{\mathbf{y}}|\boldsymbol{\theta}, \boldsymbol{\eta}) f(\mathbf{y}|\boldsymbol{\theta}, \boldsymbol{\eta}) p(\boldsymbol{\theta}|\boldsymbol{\eta}) d\tilde{\mathbf{y}} d\mathbf{y} d\boldsymbol{\theta}.$$

In the next section, we define the cAI under covariate shift in linear mixed model.

## 4.2 Covariate shift conditional AIC

### 4.2.1 Observed model

The candidate observed model  $j$  is the linear mixed model

$$\mathbf{y} = \mathbf{X}(j)\boldsymbol{\beta}_j + \mathbf{Z}\mathbf{b}_j + \boldsymbol{\varepsilon}_j, \quad (4.2)$$

where  $\mathbf{y}$  is an  $n \times 1$  observation vector of response variables,  $\mathbf{X}(j)$  and  $\mathbf{Z}$  are  $n \times p_j$  and  $n \times q$  matrices of covariates, respectively,  $\boldsymbol{\beta}_j$  is a  $p_j \times 1$  vector of regression coefficients,  $\mathbf{b}_j$  is a  $q \times 1$  vector of random effects, and  $\boldsymbol{\varepsilon}_j$  is an  $n \times 1$  vector of random errors. Let  $\mathbf{b}_j$  and  $\boldsymbol{\varepsilon}_j$  be mutually independent and  $\mathbf{b}_j \sim \mathcal{N}_q(\mathbf{0}, \sigma_j^2 \mathbf{G})$ ,  $\boldsymbol{\varepsilon}_j \sim \mathcal{N}_n(\mathbf{0}, \sigma_j^2 \mathbf{R})$ , where  $\mathbf{G}$  and  $\mathbf{R}$  are  $q \times q$  and  $n \times n$  positive definite matrices and  $\sigma_j^2$  is a scalar. We assume that  $\mathbf{G}$  and  $\mathbf{R}$  are known and  $\sigma_j^2$  is unknown. The conditional density function of  $\mathbf{y}$  given  $\mathbf{b}_j$  and the density function of  $\mathbf{b}_j$  for the model  $j$  are denoted by  $f(\mathbf{y}|\mathbf{b}_j, \boldsymbol{\beta}_j, \sigma_j^2)$  and  $p(\mathbf{b}_j|\sigma_j^2)$ , respectively.

The true observed model  $j_*$  is

$$\mathbf{y} = \mathbf{X}(\omega)\boldsymbol{\beta}_* + \mathbf{Z}\mathbf{b}_* + \boldsymbol{\varepsilon}_*,$$

where  $\mathbf{b}_* \sim \mathcal{N}_q(\mathbf{0}, \sigma_*^2 \mathbf{G})$ ,  $\boldsymbol{\varepsilon}_* \sim \mathcal{N}_n(\mathbf{0}, \sigma_*^2 \mathbf{R})$ . Note that  $\mathbf{X}(\omega)$  is  $n \times p_\omega$  matrix of covariates for the full model  $\omega$  and that  $\boldsymbol{\beta}_*$  is  $p_\omega \times 1$  vector of regression coefficients, whose  $p_\omega - p_j$  components are exactly 0 and the rest of components are not 0, as explained in Section 2.2. Then the marginal distribution of  $\mathbf{y}$  is

$$\mathbf{y} \sim \mathcal{N}_n(\mathbf{X}(\omega)\boldsymbol{\beta}_*, \sigma_*^2 \boldsymbol{\Sigma}),$$

where  $\boldsymbol{\Sigma} = \mathbf{Z}\mathbf{G}\mathbf{Z}^T + \mathbf{R}$ . For the true model, the conditional density function of  $\mathbf{y}$  given  $\mathbf{b}_*$  and the density function of  $\mathbf{b}_*$  are denoted by  $f(\mathbf{y}|\mathbf{b}_*, \boldsymbol{\beta}_*, \sigma_*^2)$  and  $p(\mathbf{b}_*|\sigma_*^2)$ , respectively.

### 4.2.2 Predictive model

The candidate predictive model  $j$  is the linear mixed model which has the same regression coefficients  $\boldsymbol{\beta}_j$  and random effects  $\mathbf{b}_j$  as in the candidate observed model  $j$ , but different covariates, namely

$$\tilde{\mathbf{y}} = \tilde{\mathbf{X}}(j)\boldsymbol{\beta}_j + \tilde{\mathbf{Z}}\mathbf{b}_j + \tilde{\boldsymbol{\varepsilon}}_j, \quad (4.3)$$

where  $\tilde{\mathbf{y}}$  is  $m \times 1$  random vector of the target of prediction,  $\tilde{\mathbf{X}}(j)$  and  $\tilde{\mathbf{Z}}$  are  $m \times p_j$  and  $m \times q$  matrices of covariates whose columns correspond to those of  $\mathbf{X}(j)$  and  $\mathbf{Z}$ , respectively, and  $\tilde{\boldsymbol{\varepsilon}}_j$  is  $m \times 1$  vector of random errors, which is independent of  $\mathbf{b}_j$  and  $\boldsymbol{\varepsilon}_j$  and is distributed as  $\tilde{\boldsymbol{\varepsilon}}_j \sim \mathcal{N}_m(\mathbf{0}, \sigma_j^2 \tilde{\mathbf{R}})$ , where  $\tilde{\mathbf{R}}$  is a known  $m \times m$  positive definite matrix. We assume that we know the values of  $\tilde{\mathbf{X}}(j)$  and  $\tilde{\mathbf{Z}}$  in the predictive model and that they are not necessarily the same as those of  $\mathbf{X}(j)$  and  $\mathbf{Z}$  in the observed model. We call this situation covariate shift. The conditional density function of  $\tilde{\mathbf{y}}$  given  $\mathbf{b}_j$  for the model  $j$  is denoted by  $g(\tilde{\mathbf{y}}|\mathbf{b}_j, \boldsymbol{\beta}_j, \sigma_j^2)$ .

The true predictive model  $j_*$  is

$$\tilde{\mathbf{y}} = \tilde{\mathbf{X}}(\omega)\boldsymbol{\beta}_* + \tilde{\mathbf{Z}}\mathbf{b}_* + \tilde{\boldsymbol{\varepsilon}}_*,$$

where  $\tilde{\mathbf{X}}(\omega)$  is  $m \times p_\omega$  matrix of covariates and  $\tilde{\boldsymbol{\varepsilon}}_* \sim \mathcal{N}_m(\mathbf{0}, \sigma_*^2 \tilde{\mathbf{R}})$ . Then the marginal distribution of  $\tilde{\mathbf{y}}$  is

$$\tilde{\mathbf{y}} \sim \mathcal{N}_m(\tilde{\mathbf{X}}(\omega)\boldsymbol{\beta}_*, \sigma_*^2 \tilde{\boldsymbol{\Sigma}}),$$

where  $\tilde{\boldsymbol{\Sigma}} = \tilde{\mathbf{Z}}\mathbf{G}\tilde{\mathbf{Z}}^T + \tilde{\mathbf{R}}$ . For the true model, the conditional density function of  $\tilde{\mathbf{y}}$  given  $\mathbf{b}_*$  is denoted by  $g(\tilde{\mathbf{y}}|\mathbf{b}_*, \boldsymbol{\beta}_*, \sigma_*^2)$ .



### 4.2.3 Conditional Akaike information

The conditional Akaike information (cAI) measures the prediction risk of the plug-in predictive density  $g(\tilde{\mathbf{y}}|\hat{\mathbf{b}}_j, \hat{\boldsymbol{\beta}}_j, \hat{\sigma}_j^2)$ , where  $\hat{\boldsymbol{\beta}}_j$  and  $\hat{\sigma}_j^2$  are maximum likelihood estimators of  $\boldsymbol{\beta}_j$  and  $\sigma_j^2$ , which are given as

$$\begin{aligned}\hat{\boldsymbol{\beta}}_j &= (\mathbf{X}(j)^\top \boldsymbol{\Sigma}^{-1} \mathbf{X}(j))^{-1} \mathbf{X}(j)^\top \boldsymbol{\Sigma}^{-1} \mathbf{y}, \\ \hat{\sigma}_j^2 &= (\mathbf{y} - \mathbf{X}(j) \hat{\boldsymbol{\beta}}_j)^\top \boldsymbol{\Sigma}^{-1} (\mathbf{y} - \mathbf{X}(j) \hat{\boldsymbol{\beta}}_j) / n,\end{aligned}$$

and  $\hat{\mathbf{b}}_j$  is empirical Bayes estimator of  $\mathbf{b}_j$  for quadratic loss, which is given by

$$\hat{\mathbf{b}}_j = \mathbf{G} \mathbf{Z}^\top \boldsymbol{\Sigma}^{-1} (\mathbf{y} - \mathbf{X}(j) \hat{\boldsymbol{\beta}}_j).$$

Then, the cAI under covariate shift is

$$\begin{aligned}\text{cAI} &= E^{(\mathbf{y}, \mathbf{b}_*)} E^{\tilde{\mathbf{y}}|\mathbf{b}_*} \left[ -2 \log \{g(\tilde{\mathbf{y}}|\hat{\mathbf{b}}_j, \hat{\boldsymbol{\beta}}_j, \hat{\sigma}_j^2)\} \right] \\ &= E^{(\mathbf{y}, \mathbf{b}_*)} E^{\tilde{\mathbf{y}}|\mathbf{b}_*} \left[ m \log(2\pi \hat{\sigma}_j^2) + \log |\tilde{\mathbf{R}}| + (\tilde{\mathbf{y}} - \tilde{\mathbf{X}}(j) \hat{\boldsymbol{\beta}}_j - \tilde{\mathbf{Z}} \hat{\mathbf{b}}_j)^\top \tilde{\mathbf{R}}^{-1} (\tilde{\mathbf{y}} - \tilde{\mathbf{X}}(j) \hat{\boldsymbol{\beta}}_j - \tilde{\mathbf{Z}} \hat{\mathbf{b}}_j) / \hat{\sigma}_j^2 \right],\end{aligned}$$

where  $E^{(\mathbf{y}, \mathbf{b}_*)}$  and  $E^{\tilde{\mathbf{y}}|\mathbf{b}_*}$  denote expectation with respect to the joint distribution of  $(\mathbf{y}, \mathbf{b}_*) \sim f(\mathbf{y}|\mathbf{b}_*, \boldsymbol{\beta}_*, \sigma_*^2) p(\mathbf{b}_*|\sigma_*^2)$  and the conditional distribution of  $\tilde{\mathbf{y}}$  given  $\mathbf{b}_*$ , namely  $\tilde{\mathbf{y}}|\mathbf{b}_* \sim g(\tilde{\mathbf{y}}|\mathbf{b}_*, \boldsymbol{\beta}_*, \sigma_*^2)$ , respectively. Taking expectation with respect to  $\tilde{\mathbf{y}}|\mathbf{b}_* \sim g(\tilde{\mathbf{y}}|\mathbf{b}_*, \boldsymbol{\beta}_*, \sigma_*^2)$  and  $\mathbf{b}_*|\mathbf{y} \sim \mathcal{N}_q(\tilde{\mathbf{b}}_*, \sigma_*^2 (\mathbf{G} - \mathbf{G} \mathbf{Z}^\top \boldsymbol{\Sigma}^{-1} \mathbf{Z} \mathbf{G}))$  for  $\tilde{\mathbf{b}}_* = \mathbf{G} \mathbf{Z}^\top \boldsymbol{\Sigma}^{-1} (\mathbf{y} - \mathbf{X}(\omega) \boldsymbol{\beta}_*)$ , we can obtain

$$\text{cAI} = E \left[ m \log(2\pi \hat{\sigma}_j^2) + \log |\tilde{\mathbf{R}}| + \text{tr}(\tilde{\mathbf{R}}^{-1} \boldsymbol{\Lambda}) \cdot \sigma_*^2 / \hat{\sigma}_j^2 + \mathbf{a}^\top \tilde{\mathbf{R}}^{-1} \mathbf{a} / \hat{\sigma}_j^2 \right], \quad (4.4)$$

where

$$\begin{aligned}\boldsymbol{\Lambda} &= \tilde{\boldsymbol{\Sigma}} - \tilde{\mathbf{Z}} \mathbf{G} \mathbf{Z}^\top \boldsymbol{\Sigma}^{-1} \mathbf{Z} \mathbf{G} \tilde{\mathbf{Z}}^\top, \\ \mathbf{a} &= (\tilde{\mathbf{X}}(j) \hat{\boldsymbol{\beta}}_j - \tilde{\mathbf{X}}(\omega) \boldsymbol{\beta}_*) - \tilde{\mathbf{Z}} \mathbf{G} \mathbf{Z}^\top \boldsymbol{\Sigma}^{-1} (\mathbf{X}(j) \hat{\boldsymbol{\beta}}_j - \mathbf{X}(\omega) \boldsymbol{\beta}_*).\end{aligned} \quad (4.5)$$

### 4.2.4 Criterion for overspecified model

In this subsection, we propose the covariate shift cAIC (CScAIC) as an unbiased estimator of cAI in (4.4) under the assumption that the candidate model is overspecified. When the candidate model  $j$  is overspecified, the true mean vector of  $\mathbf{y}$  and  $\tilde{\mathbf{y}}$  can be expressed as

$$E(\mathbf{y}) = \mathbf{X}(j) \boldsymbol{\beta}_j^*, \quad \text{and} \quad E(\tilde{\mathbf{y}}) = \tilde{\mathbf{X}}(j) \boldsymbol{\beta}_j^*$$

where  $\boldsymbol{\beta}_j^*$  is  $p_j \times 1$  vector, whose  $p_j - p_*$  components are exactly 0 and the rest of components are not 0. In this subsection, we hereafter abbreviate  $\mathbf{X}(j)$  to  $\mathbf{X}$ ,  $\tilde{\mathbf{X}}(j)$  to  $\tilde{\mathbf{X}}$  and  $\boldsymbol{\beta}_j^*$  to  $\boldsymbol{\beta}$  for notational convenience. Then,  $\mathbf{a}$  in cAI given by (4.5) can be reduced to

$$\mathbf{a} = \mathbf{A}(\hat{\boldsymbol{\beta}}_j - \boldsymbol{\beta}),$$

where  $\mathbf{A} = \tilde{\mathbf{X}} - \tilde{\mathbf{Z}} \mathbf{G} \mathbf{Z}^\top \boldsymbol{\Sigma}^{-1} \mathbf{X}$ . Furthermore,  $n \hat{\sigma}_j^2 / \sigma_*^2$  follows chi-squared distribution with  $n - p_j$  degrees of freedom, denoted by  $n \hat{\sigma}_j^2 / \sigma_*^2 \sim \chi_{n-p_j}^2$ , which implies that

$$E[(\hat{\boldsymbol{\beta}}_j - \boldsymbol{\beta})(\hat{\boldsymbol{\beta}}_j - \boldsymbol{\beta})^\top] = \sigma_*^2 (\mathbf{X}^\top \boldsymbol{\Sigma}^{-1} \mathbf{X})^{-1},$$

when the candidate model  $j$  is overspecified. Noting that  $\hat{\beta}_j$  and  $\hat{\sigma}_j^2$  are mutually independent, we can evaluate the cAI in (4.4) as

$$\text{cAI} = E[m \log(2\pi\hat{\sigma}_j^2)] + \log |\tilde{\mathbf{R}}| + \frac{n}{n - p_j - 2} \left\{ \text{tr} [\tilde{\mathbf{R}}^{-1} \mathbf{\Lambda}] + \text{tr} [\tilde{\mathbf{R}}^{-1} \mathbf{A}(\mathbf{X}^T \mathbf{\Sigma}^{-1} \mathbf{X})^{-1} \mathbf{A}^T] \right\}. \quad (4.6)$$

We propose the CScAIC as an unbiased estimator of cAI as follows:

$$\text{CScAIC} = l(j) + \Delta_{\text{CS}}, \quad (4.7)$$

where  $l(j)$  is likelihood (or goodness of fit) part and  $\Delta_{\text{CS}}$  is bias correction, which are given by

$$l(j) = m \log(2\pi\hat{\sigma}_j^2) + \log |\tilde{\mathbf{R}}| + (\mathbf{y} - \mathbf{X}\hat{\beta}_j - \mathbf{Z}\hat{\mathbf{b}}_j)^T \tilde{\mathbf{R}}^{-1} (\mathbf{y} - \mathbf{X}\hat{\beta}_j - \mathbf{Z}\hat{\mathbf{b}}_j) / \hat{\sigma}_j^2, \quad (4.8)$$

$$\begin{aligned} \Delta_{\text{CS}} &= \frac{n}{n - p_j - 2} \left\{ \text{tr} [\tilde{\mathbf{R}}^{-1} \mathbf{\Lambda}] + \text{tr} [\tilde{\mathbf{R}}^{-1} \mathbf{A}(\mathbf{X}^T \mathbf{\Sigma}^{-1} \mathbf{X})^{-1} \mathbf{A}^T] \right\} \\ &\quad + \frac{n}{n - p_j} \left\{ -\text{tr} [\mathbf{R}\mathbf{\Sigma}^{-1}] + \text{tr} [\mathbf{R}\mathbf{P}] \right\}, \end{aligned} \quad (4.9)$$

where  $\mathbf{P} = \mathbf{\Sigma}^{-1} \mathbf{X}(\mathbf{X}^T \mathbf{\Sigma}^{-1} \mathbf{X})^{-1} \mathbf{X}^T \mathbf{\Sigma}^{-1}$ . Then we obtain the following theorem.

**Theorem 4.1** *When the candidate model is overspecified, the covariate shift cAIC (CScAIC) in (4.7) is an unbiased estimator of the cAI in (4.4).*

Next corollary shows that our CScAIC includes the cAIC by Vaida and Blanchard (2005) as special case.

**Corollary 4.1** *Suppose that covariate shift does not occur, namely  $\tilde{\mathbf{X}} = \mathbf{X}$ ,  $\tilde{\mathbf{Z}} = \mathbf{Z}$  and  $n = m$ . In addition, let the covariance matrix of  $\varepsilon$  and  $\tilde{\varepsilon}$  be both  $\sigma^2 \mathbf{I}_n$ , namely  $\mathbf{R} = \tilde{\mathbf{R}} = \mathbf{I}_n$ . Then the bias correction of the CScAIC in (4.9) is reduced to*

$$\Delta_{\text{CS}} = \frac{2n^2}{n - p_j - 2} + \frac{2n(n - p_j - 1)}{(n - p_j)(n - p_j - 2)} \left\{ -\text{tr} [\mathbf{\Sigma}^{-1}] + \text{tr} [\mathbf{\Sigma}^{-1} \mathbf{X}(\mathbf{X}^T \mathbf{\Sigma}^{-1} \mathbf{X})^{-1} \mathbf{X}^T \mathbf{\Sigma}^{-1}] \right\}, \quad (4.10)$$

which are identical to bias correction of the cAIC by Theorem 2 in Vaida and Blanchard (2005).

## 4.3 Modification of the criterion

### 4.3.1 Drawback of overspecified model assumption

Most of the Akaike-type criteria are derived under the assumption that ‘the candidate model includes the true model’, namely overspecified assumption. Although the assumption is too strong, the influence is restrictive in practical use. This is because the likelihood part of the criterion is a naive estimator of the risk function, namely the cAI in the context of the cAIC.

However, under the covariate shift situation, the likelihood part of the CScAIC is not a good estimator of the cAI. As a result, the CScAIC in (4.7) has large biases for estimating the cAI of the underspecified models as illustrated in simulations in Section 4.5.

Thus we evaluate and estimate the cAI directly both for the overspecified models and underspecified models in the following subsections.

### 4.3.2 Evaluation of cAI

We evaluate the cAI in (4.4) both for overspecified model and for underspecified model. We assume that the full model  $\omega$  is overspecified, namely the collection of the overspecified models  $\mathcal{J}_+$  is not empty set. We also assume that the size of response variable in predictive model  $m$  is of order  $O(n)$ .

When the candidate model  $j$  is overspecified,  $n\hat{\sigma}_j^2/\sigma_*^2$  follows the chi-squared distribution. However, for the underspecified model, this is not true. We decompose  $n\hat{\sigma}_j^2/\sigma_*^2$  as

$$\begin{aligned} n\hat{\sigma}_j^2/\sigma_*^2 &= \mathbf{z}^\top(\mathbf{I}_n - \mathbf{M}_\omega)\mathbf{z} + \mathbf{z}^\top(\mathbf{M}_\omega - \mathbf{M}_j)\mathbf{z} \\ &= K_0 + K_1 \quad (\text{say}), \end{aligned}$$

where

$$\begin{aligned} \mathbf{z} &= \boldsymbol{\Sigma}^{-1/2}\mathbf{y}/\sigma_*, \\ \mathbf{M}_\omega &= \mathbf{W}_\omega(\mathbf{W}_\omega^\top\mathbf{W}_\omega)^{-1}\mathbf{W}_\omega^\top = \boldsymbol{\Sigma}^{-1/2}\mathbf{X}(\omega)(\mathbf{X}(\omega)^\top\boldsymbol{\Sigma}^{-1}\mathbf{X}(\omega))^{-1}\mathbf{X}(\omega)^\top\boldsymbol{\Sigma}^{-1/2}, \\ \mathbf{M}_j &= \mathbf{W}_j(\mathbf{W}_j^\top\mathbf{W}_j)^{-1}\mathbf{W}_j^\top = \boldsymbol{\Sigma}^{-1/2}\mathbf{X}(j)(\mathbf{X}(j)^\top\boldsymbol{\Sigma}^{-1}\mathbf{X}(j))^{-1}\mathbf{X}(j)^\top\boldsymbol{\Sigma}^{-1/2}, \end{aligned}$$

for  $\mathbf{W}_\omega = \boldsymbol{\Sigma}^{-1/2}\mathbf{X}(\omega)/\sigma_*$  and  $\mathbf{W}_j = \boldsymbol{\Sigma}^{-1/2}\mathbf{X}(j)/\sigma_*$ . Note that  $\mathbf{M}_j$  and  $\mathbf{M}_\omega$  are projection matrices, namely symmetric and idempotent. Let  $\mathbf{v} = \boldsymbol{\Sigma}^{-1/2}\mathbf{u}/\sigma_*$  and  $\boldsymbol{\xi} = \mathbf{W}_\omega\boldsymbol{\beta}_*$ . Then, it can be seen that  $\mathbf{M}_\omega\boldsymbol{\xi} = \boldsymbol{\xi}$  and

$$\mathbf{M}_j\boldsymbol{\xi} \begin{cases} = \boldsymbol{\xi} & \text{if } j \in \mathcal{J}_+, \\ \neq \boldsymbol{\xi} & \text{if } j \in \mathcal{J}_-, \end{cases}$$

since  $\mathbf{X}(\omega)\boldsymbol{\beta}_* \in \mathcal{R}[\mathbf{X}(j)]$  if  $j \in \mathcal{J}_+$ . Thus  $K_0$  can be rewritten as

$$K_0 = (\boldsymbol{\xi} + \mathbf{v})^\top(\mathbf{I}_n - \mathbf{M}_\omega)(\boldsymbol{\xi} + \mathbf{v}) = \mathbf{v}^\top(\mathbf{I}_n - \mathbf{M}_\omega)\mathbf{v},$$

so that  $K_0 \sim \chi_{n-p_\omega}^2$ . Also,  $K_1$  can be rewritten as

$$\begin{aligned} K_1 &= \mathbf{v}^\top(\mathbf{M}_\omega - \mathbf{M}_j)\mathbf{v} + 2\boldsymbol{\xi}^\top(\mathbf{M}_\omega - \mathbf{M}_j)\mathbf{v} + \boldsymbol{\xi}^\top(\mathbf{M}_\omega - \mathbf{M}_j)\boldsymbol{\xi} \\ &= \mathbf{v}^\top(\mathbf{M}_\omega - \mathbf{M}_j)\mathbf{v} + 2L + n\delta, \end{aligned}$$

where

$$\begin{aligned} L &= \boldsymbol{\xi}^\top(\mathbf{M}_\omega - \mathbf{M}_j)\mathbf{v}, \\ \delta &= \boldsymbol{\xi}^\top(\mathbf{M}_\omega - \mathbf{M}_j)\boldsymbol{\xi}/n. \end{aligned}$$

In the overspecified case, we have  $K_1 \sim \chi_{p_\omega-p_j}^2$  since  $\mathbf{M}_\omega\boldsymbol{\xi} = \mathbf{M}_j\boldsymbol{\xi} = \boldsymbol{\xi}$ . In the underspecified case,  $K_1$  follows a non-central chi-squared distribution with  $p_\omega - p_j$  degrees of freedom and with the noncentrality parameter  $n\delta$ , denoted by  $K_1 \sim \chi_{p_\omega-p_j}^2(n\delta)$ . Thus,

$$K_1 \sim \begin{cases} \chi_{p_\omega-p_j}^2 & \text{if } j \in \mathcal{J}_+, \\ \chi_{p_\omega-p_j}^2(n\delta) & \text{if } j \in \mathcal{J}_-. \end{cases}$$

Because  $\hat{\boldsymbol{\beta}}_j$  and  $\hat{\sigma}_j^2$  are mutually independent, the cAI in (4.4) can be rewritten as

$$\text{cAI} = E[m \log(2\pi\hat{\sigma}_j^2)] + \log|\tilde{\mathbf{R}}| + n \cdot E[(K_0 + K_1)^{-1}] \{ \text{tr}(\tilde{\mathbf{R}}^{-1}\boldsymbol{\Lambda}) + E[\mathbf{a}^\top\tilde{\mathbf{R}}^{-1}\mathbf{a}/\sigma_*^2] \}. \quad (4.11)$$

Evaluating  $E[(K_0 + K_1)^{-1}]$  and  $E[\mathbf{a}^\top\tilde{\mathbf{R}}^{-1}\mathbf{a}/\sigma_*^2]$ , we can obtain the following theorem.

**Theorem 4.2** For the overspecified case, it follows that  $\text{cAI} = E[m \log(2\pi\hat{\sigma}_j^2)] + \log|\tilde{\mathbf{R}}| + R^*$ , where

$$R^* = \frac{n\gamma}{n - p_j - 2},$$

for  $\gamma = \text{tr}(\tilde{\mathbf{R}}^{-1}\mathbf{\Lambda}) + \text{tr}[\tilde{\mathbf{R}}^{-1}\mathbf{A}(\mathbf{X}(j)^\top \mathbf{\Sigma}^{-1} \mathbf{X}(j))^{-1} \mathbf{A}^\top]$  and  $\mathbf{A} = \tilde{\mathbf{X}}(j) - \tilde{\mathbf{Z}}\mathbf{G}\mathbf{Z}^\top \mathbf{\Sigma}^{-1} \mathbf{X}(j)$ . For the underspecified case, cAI is approximated as

$$\text{cAI} = E[m \log(2\pi\hat{\sigma}_j^2)] + \log|\tilde{\mathbf{R}}| + R^* + R_1 + R_2 + R_3 + R_4 + O(n^{-1}), \quad (4.12)$$

where

$$\begin{aligned} R_1 &= \gamma(\lambda - 1), \\ R_2 &= \gamma \cdot n^{-1} \{-2\lambda^3 + (p_j + 4)\lambda^2 - (p_j + 2)\}, \\ R_3 &= \lambda \cdot \beta_*^\top \mathbf{B}^\top \tilde{\mathbf{R}}^{-1} \mathbf{B} \beta_* / \sigma_*^2, \end{aligned}$$

and

$$R_4 = n^{-1} \{-2\lambda^3 + (p_j + 4)\lambda^2\} \times \beta_*^\top \mathbf{B}^\top \tilde{\mathbf{R}}^{-1} \mathbf{B} \beta_* / \sigma_*^2,$$

for  $\lambda = 1/(1 + \delta)$ ,

$$\begin{aligned} \mathbf{B} &= \{\tilde{\mathbf{P}}_j \mathbf{X}(\omega) - \tilde{\mathbf{X}}(\omega) + \tilde{\mathbf{Z}}\mathbf{G}\mathbf{Z}^\top (\mathbf{P}_\omega - \mathbf{P}_j) \mathbf{X}(\omega)\}, \\ \mathbf{P}_j &= \mathbf{\Sigma}^{-1} \mathbf{X}(j) (\mathbf{X}(j)^\top \mathbf{\Sigma}^{-1} \mathbf{X}(j))^{-1} \mathbf{X}(j)^\top \mathbf{\Sigma}^{-1}, \\ \mathbf{P}_\omega &= \mathbf{\Sigma}^{-1} \mathbf{X}(\omega) (\mathbf{X}(\omega)^\top \mathbf{\Sigma}^{-1} \mathbf{X}(\omega))^{-1} \mathbf{X}(\omega)^\top \mathbf{\Sigma}^{-1}, \end{aligned}$$

and

$$\tilde{\mathbf{P}}_j = \tilde{\mathbf{X}}(j) (\mathbf{X}(j)^\top \mathbf{\Sigma}^{-1} \mathbf{X}(j))^{-1} \mathbf{X}(j)^\top \mathbf{\Sigma}^{-1}.$$

When the candidate model  $j$  is overspecified, it follows that  $R_1, R_2, R_3$  and  $R_4$  are exactly 0.

### 4.3.3 Estimation of cAI

Because the approximation of cAI in (4.12) includes unknown parameters, we have to provide an estimator of cAI for practical use. Firstly, we obtain estimators of  $R_1$  and  $R_2$ , which are polynomials of  $\lambda$ . We define  $\hat{\lambda}$ ,  $\widehat{\lambda^2}$  and  $\widehat{\lambda^3}$  as

$$\begin{aligned} \hat{\lambda} &= \frac{n - p_j}{n - p_\omega} \frac{\hat{\sigma}_\omega^2}{\hat{\sigma}_j^2}, \\ \widehat{\lambda^2} &= \frac{(n - p_j)(n - p_j + 2)}{(n - p_\omega)(n - p_\omega + 2)} \left( \frac{\hat{\sigma}_\omega^2}{\hat{\sigma}_j^2} \right)^2, \end{aligned}$$

and

$$\widehat{\lambda^3} = \frac{(n - p_j)(n - p_j + 2)(n - p_j + 4)}{(n - p_\omega)(n - p_\omega + 2)(n - p_\omega + 4)} \left( \frac{\hat{\sigma}_\omega^2}{\hat{\sigma}_j^2} \right)^3.$$

When  $j \in \mathcal{J}_+$ , it follows that  $n\hat{\sigma}_\omega^2/\sigma_*^2 = K_0$ ,  $n\hat{\sigma}_j^2/\sigma_*^2 = K_0 + K_1$ , and that  $K_0$  and  $K_1$  are mutually independent and distribute as  $K_0 \sim \chi_{n-p_\omega}^2$ ,  $K_1 \sim \chi_{p_\omega-p_j}^2$ . Thus,

$$\frac{\hat{\sigma}_\omega^2}{\hat{\sigma}_j^2} \sim \text{Be} \left( \frac{n - p_\omega}{2}, \frac{p_\omega - p_j}{2} \right),$$

where  $\text{Be}(\cdot, \cdot)$  denotes the beta distribution. This implies that  $E(\hat{\lambda}) = E(\widehat{\lambda^2}) = E(\widehat{\lambda^3}) = 1$  for the overspecified case. For the underspecified case, on the other hand, it follows that  $E[(\hat{\sigma}_\omega^2/\hat{\sigma}_j^2)^k] = \lambda^k + O(n^{-1})$  as  $n \rightarrow \infty$  for  $k = 1, 2, 3$ . Then we can obtain an estimator of  $R_2$  in the approximation of cAI given by (4.12), which is given as follows:

$$\widehat{R}_2 = \gamma \cdot \frac{-2\widehat{\lambda^3} + (p_j + 4)\widehat{\lambda^2} - (p_j + 2)}{n}. \quad (4.13)$$

Noting that  $\gamma = O(n)$ , we can get the following lemma.

**Lemma 4.1** *When the candidate model  $j$  is underspecified,  $\widehat{R}_2$  in (4.13) is an asymptotically unbiased estimator of  $R_2$  whose bias is of order  $O(n^{-1})$ , namely*

$$E(\widehat{R}_2) = R_2 + O(n^{-1}).$$

*When the candidate model  $j$  is overspecified, it follows that  $E(\widehat{R}_2) = 0$ .*

Because  $R_1$  is of order  $O(n)$ , we have to estimate  $\lambda$  with higher order accuracy in order to obtain an estimator of  $R_1$  whose bias is of order  $O(n^{-1})$  for the underspecified case. To this end, we expand  $E(\hat{\lambda})$  up to  $O(n^{-1})$  as

$$\begin{aligned} E(\hat{\lambda}) &= \frac{n - p_j}{n - p_\omega} \cdot E\left(\frac{K_0}{K_0 + K_1}\right) \\ &= \lambda + \frac{-2\lambda^3 + (p_j + 2)\lambda^2 - p_j\lambda}{n} + O(n^{-2}). \end{aligned}$$

Then we can obtain an estimator of  $R_1$  given as

$$\widehat{R}_1 = \gamma \cdot \left\{ \hat{\lambda} - \frac{-2\widehat{\lambda^3} + (p_j + 2)\widehat{\lambda^2} - p_j\hat{\lambda}}{n} - 1 \right\}. \quad (4.14)$$

**Lemma 4.2** *When the candidate model  $j$  is underspecified,  $\widehat{R}_1$  in (4.14) is an asymptotically unbiased estimator of  $R_1$  whose bias is of order  $O(n^{-1})$ , namely*

$$E(\widehat{R}_1) = R_1 + O(n^{-1}).$$

*When the candidate model  $j$  is overspecified, it follows that  $E(\widehat{R}_1) = 0$ .*

We next consider estimation procedures of  $R_3$  and  $R_4$ , which are complex functions of unknown parameters. We see  $R_3$  and  $R_4$  as functions of  $\boldsymbol{\eta}_* = (\boldsymbol{\beta}_*^\top, \sigma_*^2)^\top$ , namely  $R_3 = R_3(\boldsymbol{\eta}_*)$ ,  $R_4 = R_4(\boldsymbol{\eta}_*)$  and substitute their unbiased estimators  $\tilde{\boldsymbol{\eta}} = (\tilde{\boldsymbol{\beta}}^\top, \tilde{\sigma}^2)^\top$ , which are given by

$$\begin{aligned} \tilde{\boldsymbol{\beta}} &= \widehat{\boldsymbol{\beta}}_\omega = (\mathbf{X}(\omega)^\top \boldsymbol{\Sigma}^{-1} \mathbf{X}(\omega))^{-1} \mathbf{X}(\omega)^\top \boldsymbol{\Sigma}^{-1} \mathbf{y}, \\ \tilde{\sigma}^2 &= (\mathbf{y} - \mathbf{X}(\omega)\tilde{\boldsymbol{\beta}})^\top \boldsymbol{\Sigma}^{-1} (\mathbf{y} - \mathbf{X}(\omega)\tilde{\boldsymbol{\beta}}) / (n - p_\omega). \end{aligned}$$

Then, plug-in estimators of  $R_3$  and  $R_4$  are

$$\begin{aligned} \widetilde{R}_3 &= R_3(\tilde{\boldsymbol{\eta}}) = \tilde{\lambda} \cdot \tilde{\boldsymbol{\beta}}^\top \mathbf{B}^\top \tilde{\mathbf{R}}^{-1} \mathbf{B} \tilde{\boldsymbol{\beta}} / \tilde{\sigma}^2, \\ \widetilde{R}_4 &= R_4(\tilde{\boldsymbol{\eta}}) = n^{-1} \{-2\tilde{\lambda}^3 + (p_j + 4)\tilde{\lambda}^2\} \times \tilde{\boldsymbol{\beta}}^\top \mathbf{B}^\top \tilde{\mathbf{R}}^{-1} \mathbf{B} \tilde{\boldsymbol{\beta}} / \tilde{\sigma}^2, \end{aligned} \quad (4.15)$$

where  $\tilde{\lambda} = 1/(1 + \tilde{\delta})$  for

$$\tilde{\delta} = \tilde{\boldsymbol{\beta}}^T \mathbf{X}(\omega)^T (\mathbf{P}_\omega - \mathbf{P}_j) \mathbf{X}(\omega) \tilde{\boldsymbol{\beta}} / (n\tilde{\sigma}^2).$$

Because  $R_3$  is of order  $O(n)$ , plug-in estimator  $\widetilde{R}_3$  has second order bias. Then we correct the bias by analytical method based on Taylor series expansions. We can see that expectation of the plug-in estimator  $R_3(\tilde{\boldsymbol{\eta}})$  is expanded as

$$E[R_3(\tilde{\boldsymbol{\eta}})] = R_3(\boldsymbol{\eta}_*) + B_1(\boldsymbol{\eta}_*) + B_2(\boldsymbol{\eta}_*) + O(n^{-2}), \quad (4.16)$$

where  $B_1(\boldsymbol{\eta}_*)$  and  $B_2(\boldsymbol{\eta}_*)$  are second- and third-order biases of  $R_3(\tilde{\boldsymbol{\eta}})$ , respectively, namely  $B_1(\boldsymbol{\eta}_*) = O(1)$  and  $B_2(\boldsymbol{\eta}_*) = O(n^{-1})$ . Because  $\tilde{\boldsymbol{\beta}}$  and  $\tilde{\sigma}^2$  are independent, it follows that

$$B_1(\boldsymbol{\eta}_*) = \frac{1}{2} \cdot \text{tr} \left[ \frac{\partial^2 R_3(\boldsymbol{\eta}_*)}{\partial \boldsymbol{\beta}_*^T \partial \boldsymbol{\beta}_*^T} E[(\tilde{\boldsymbol{\beta}} - \boldsymbol{\beta}_*)(\tilde{\boldsymbol{\beta}} - \boldsymbol{\beta}_*)^T] \right] + \frac{1}{2} \frac{\partial^2 R_3(\boldsymbol{\eta}_*)}{(\partial \sigma_*^2)^2} E[(\tilde{\sigma}^2 - \sigma_*^2)^2], \quad (4.17)$$

where  $E[(\tilde{\boldsymbol{\beta}} - \boldsymbol{\beta}_*)(\tilde{\boldsymbol{\beta}} - \boldsymbol{\beta}_*)^T] = \sigma_*^2 (\mathbf{X}(\omega)^T \boldsymbol{\Sigma}^{-1} \mathbf{X}(\omega))^{-1}$  and  $E[(\tilde{\sigma}^2 - \sigma_*^2)^2] = 2(\sigma_*^2)^2 / (n - p_\omega)$ . Second order partial derivatives of  $R_3$  are given by the following lemma.

**Lemma 4.3** *Second order partial derivative of  $R_3(\boldsymbol{\eta}_*)$  with respect to  $\boldsymbol{\beta}_*$  is*

$$\begin{aligned} \frac{\partial^2 R_3(\boldsymbol{\eta}_*)}{\partial \boldsymbol{\beta}_*^T \partial \boldsymbol{\beta}_*^T} &= \frac{\boldsymbol{\beta}_*^T \mathbf{B}^T \tilde{\mathbf{R}}^{-1} \mathbf{B} \boldsymbol{\beta}_*}{\sigma_*^2} \times \left\{ -\frac{2\mathbf{C}}{n\sigma_*^2(1+\delta)^2} + \frac{8\mathbf{C}\boldsymbol{\beta}_* \boldsymbol{\beta}_*^T \mathbf{C}}{n^2(\sigma_*^2)^2(1+\delta)^3} \right\} \\ &\quad - \frac{4\mathbf{B}^T \tilde{\mathbf{R}}^{-1} \mathbf{B} \boldsymbol{\beta}_* \boldsymbol{\beta}_*^T \mathbf{C} + 4\mathbf{C}\boldsymbol{\beta}_* \boldsymbol{\beta}_*^T \mathbf{B}^T \tilde{\mathbf{R}}^{-1} \mathbf{B}}{n(\sigma_*^2)^2(1+\delta)^2} + 2\lambda \cdot \mathbf{B}^T \tilde{\mathbf{R}}^{-1} \mathbf{B} / \sigma_*^2, \end{aligned}$$

where  $\mathbf{C} = \mathbf{X}(\omega)^T (\mathbf{P}_\omega - \mathbf{P}_j) \mathbf{X}(\omega)$ . Second order partial derivative of  $R_3(\boldsymbol{\eta}_*)$  with respect to  $\sigma_*^2$  is

$$\begin{aligned} \frac{\partial^2 R_3(\boldsymbol{\eta}_*)}{(\partial \sigma_*^2)^2} &= \frac{\boldsymbol{\beta}_*^T \mathbf{B}^T \tilde{\mathbf{R}}^{-1} \mathbf{B} \boldsymbol{\beta}_*}{\sigma_*^2} \times \left\{ -\frac{2\boldsymbol{\beta}_*^T \mathbf{C} \boldsymbol{\beta}_*}{n(\sigma_*^2)^3(1+\delta)^2} + \frac{2(\boldsymbol{\beta}_*^T \mathbf{C} \boldsymbol{\beta}_*)^2}{n^2(\sigma_*^2)^4(1+\delta)^3} \right\} \\ &\quad - \frac{2\boldsymbol{\beta}_*^T \mathbf{B}^T \tilde{\mathbf{R}}^{-1} \mathbf{B} \boldsymbol{\beta}_* \boldsymbol{\beta}_*^T \mathbf{C} \boldsymbol{\beta}_*}{n(\sigma_*^2)^4(1+\delta)^2} + 2\lambda \cdot \boldsymbol{\beta}_*^T \mathbf{B}^T \tilde{\mathbf{R}}^{-1} \mathbf{B} \boldsymbol{\beta}_* / (\sigma_*^2)^3. \end{aligned}$$

When the candidate model  $j$  is overspecified, second order bias  $B_1(\boldsymbol{\eta}_*)$  can be simplified to

$$B_1(\boldsymbol{\eta}_*) = \text{tr} [\mathbf{B}^T \tilde{\mathbf{R}}^{-1} \mathbf{B} (\mathbf{X}(\omega)^T \boldsymbol{\Sigma}^{-1} \mathbf{X}(\omega))^{-1}],$$

because  $\mathbf{C}\boldsymbol{\beta}_* = \mathbf{B}\boldsymbol{\beta}_* = \mathbf{0}$  and  $\lambda = 1$ , which implies that  $(\partial^2 R_3(\boldsymbol{\eta}_*)) / (\partial \boldsymbol{\beta}_*^T \partial \boldsymbol{\beta}_*^T) = 2\mathbf{B}^T \tilde{\mathbf{R}}^{-1} \mathbf{B} / \sigma_*^2$  and that  $(\partial^2 R_3(\boldsymbol{\eta}_*)) / (\partial \sigma_*^2)^2 = 0$ . However, one cannot know which candidate models are overspecified. Then, we propose the following bias corrected estimator of  $R_3$ :

$$\widetilde{\widetilde{R}}_3 = R_3(\tilde{\boldsymbol{\eta}}) - B_1(\tilde{\boldsymbol{\eta}}). \quad (4.18)$$

**Lemma 4.4** *Both for the cases where the candidate model  $j$  is overspecified and where  $j$  is underspecified,  $\widetilde{\widetilde{R}}_3$  and  $\widetilde{R}_4$  in (4.18) and (4.15) are asymptotically unbiased estimators of  $R_3$  and  $R_4$ , whose biases are of order  $O(n^{-1})$ , namely*

$$E(\widetilde{\widetilde{R}}_3) = R_3 + O(n^{-1}), \quad \text{and} \quad E(\widetilde{R}_4) = R_4 + O(n^{-1}).$$

Using  $\widehat{R}_1$ ,  $\widehat{R}_2$ ,  $\widetilde{R}_3$  and  $\widetilde{R}_4$  given by (4.14), (4.13), (4.18) and (4.15), respectively, we can construct an estimator of cAI as follows:

$$\widehat{\text{cAI}} = m \log(2\pi\hat{\sigma}_j^2) + \log|\widetilde{\mathbf{R}}| + R^* + \widehat{R}_1 + \widehat{R}_2 + \widetilde{R}_3 + \widetilde{R}_4. \quad (4.19)$$

**Theorem 4.3** *Both for the cases where the candidate model  $j$  is overspecified and where  $j$  is underspecified,  $\widehat{\text{cAI}}$  in (4.19) is a second-order asymptotically unbiased estimator of cAI, namely*

$$E(\widehat{\text{cAI}}) = \text{cAI} + O(n^{-1}).$$

Because the CScAIC in Section 4.2 is an exact unbiased estimator of the cAI when the candidate model is overspecified, the performance of  $\widehat{\text{cAI}}$  in (4.19) is not as good as that of the CScAIC for the overspecified model when  $n$  is small. Thus we should improve the estimator of  $R_3$  and  $R_4$  to remove the biases which are of order  $O(n^{-1})$ . To this end, we adopt parametric bootstrap method. Bootstrap sample  $\mathbf{y}^\dagger$  is generated by

$$\mathbf{y}^\dagger = \mathbf{X}(\omega)\widetilde{\boldsymbol{\beta}} + \mathbf{Z}\mathbf{b}^\dagger + \boldsymbol{\varepsilon}^\dagger,$$

where  $\mathbf{b}^\dagger$  and  $\boldsymbol{\varepsilon}^\dagger$  are generated by the following distribution:

$$\mathbf{b}^\dagger \sim \mathcal{N}(\mathbf{0}, \tilde{\sigma}^2 \mathbf{G}), \quad \text{and} \quad \boldsymbol{\varepsilon}^\dagger \sim \mathcal{N}(\mathbf{0}, \tilde{\sigma}^2 \mathbf{I}_n).$$

Then, we use the following estimator of  $R_4$ :

$$\widehat{R}_4 = 2R_4(\tilde{\boldsymbol{\eta}}) - E_{\tilde{\boldsymbol{\eta}}}[R_4(\tilde{\boldsymbol{\eta}}^\dagger)] \quad (4.20)$$

where  $E_{\tilde{\boldsymbol{\eta}}}$  denotes expectation with respect to bootstrap distribution and  $\tilde{\boldsymbol{\eta}}^\dagger = ((\tilde{\boldsymbol{\beta}}^\dagger)^\top, \tilde{\sigma}^{2\dagger})^\top$  is

$$\begin{aligned} \tilde{\boldsymbol{\beta}}^\dagger &= (\mathbf{X}(\omega)^\top \boldsymbol{\Sigma}^{-1} \mathbf{X}(\omega))^{-1} \mathbf{X}(\omega) \boldsymbol{\Sigma}^{-1} \mathbf{y}^\dagger, \\ \tilde{\sigma}^{2\dagger} &= (\mathbf{y}^\dagger - \mathbf{X}(\omega)\tilde{\boldsymbol{\beta}}^\dagger)^\top \boldsymbol{\Sigma}^{-1} (\mathbf{y}^\dagger - \mathbf{X}(\omega)\tilde{\boldsymbol{\beta}}^\dagger) / (n - p_\omega). \end{aligned}$$

As for  $R_3$ , it follows from (4.16) that

$$E_{\tilde{\boldsymbol{\eta}}}[R_3(\tilde{\boldsymbol{\eta}}^\dagger)] = R_3(\tilde{\boldsymbol{\eta}}) + B_1(\tilde{\boldsymbol{\eta}}) + B_2(\tilde{\boldsymbol{\eta}}) + O_p(n^{-2}).$$

However,  $B_1(\tilde{\boldsymbol{\eta}})$  has bias which is of order  $O(n^{-1})$ , namely

$$E[B_1(\tilde{\boldsymbol{\eta}})] = B_1(\boldsymbol{\eta}_*) + B_{11}(\boldsymbol{\eta}_*) + O(n^{-2}),$$

where  $B_{11}(\boldsymbol{\eta}_*) = O(n^{-1})$ . Because this bias is not negligible when one wants to estimate  $R_3$  with third order accuracy, we estimate this bias by bootstrap method as follows:

$$\widehat{B}_{11} = E_{\tilde{\boldsymbol{\eta}}}[B_1(\tilde{\boldsymbol{\eta}}^\dagger)] - B_1(\tilde{\boldsymbol{\eta}}).$$

Then we obtain an estimator of  $R_3$ , which is given as

$$\begin{aligned} \widehat{R}_3 &= 2R_3(\tilde{\boldsymbol{\eta}}) - E_{\tilde{\boldsymbol{\eta}}}[R_3(\tilde{\boldsymbol{\eta}}^\dagger)] + \widehat{B}_{11} \\ &= 2R_3(\tilde{\boldsymbol{\eta}}) - E_{\tilde{\boldsymbol{\eta}}}[R_3(\tilde{\boldsymbol{\eta}}^\dagger)] + E_{\tilde{\boldsymbol{\eta}}}[B_1(\tilde{\boldsymbol{\eta}}^\dagger)] - B_1(\tilde{\boldsymbol{\eta}}). \end{aligned} \quad (4.21)$$

**Lemma 4.5** *Both for the cases where the candidate model  $j$  is overspecified and where  $j$  is underspecified,  $\widehat{R}_3$  and  $\widehat{R}_4$  in (4.21) and (4.20) are asymptotically unbiased estimators of  $R_3$  and  $R_4$ , whose biases are of order  $O(n^{-2})$ , namely*

$$E(\widehat{R}_3) = R_3 + O(n^{-2}), \quad \text{and} \quad E(\widehat{R}_4) = R_4 + O(n^{-2}).$$

Using  $\widehat{R}_1$ ,  $\widehat{R}_2$ ,  $\widehat{R}_3$  and  $\widehat{R}_4$  given by (4.14), (4.13), (4.21) and (4.20), we can obtain an estimator of cAI as follows:

$$\widehat{\text{cAI}}^\dagger = m \log(2\pi\hat{\sigma}_j^2) + \log|\widetilde{\mathbf{R}}| + R^* + \widehat{R}_1 + \widehat{R}_2 + \widehat{R}_3 + \widehat{R}_4, \quad (4.22)$$

which improves  $\widehat{\text{cAI}}$  in unbiasedness.

**Theorem 4.4** *When the candidate model  $j$  is overspecified,  $\widehat{\text{cAI}}^\dagger$  in (4.22) is a third order asymptotically unbiased estimator of cAI, namely*

$$E(\widehat{\text{cAI}}^\dagger) = \text{cAI} + O(n^{-2}).$$

*When the candidate model  $j$  is underspecified,  $\widehat{\text{cAI}}^\dagger$  is a second order asymptotically unbiased estimator of cAI, namely*

$$E(\widehat{\text{cAI}}^\dagger) = \text{cAI} + O(n^{-1}).$$

## 4.4 Application to small area prediction

A typical example of the covariate shift situation appears in small area prediction problem. The model for small area prediction supposes that the observed small area data have the finite population which has the super-population model with random effects, one of which is the well-known nested error regression model (NERM) proposed by Battese et al. (1988).

Let  $Y_{ik}$  and  $\mathbf{x}_{ik}(j)$  denote the value of a characteristic of interest and its  $p_j$ -dimensional auxiliary variable for the  $k$ th unit of the  $i$ th area for  $i = 1, \dots, q$  and  $k = 1, \dots, N_i$ . Note that  $\mathbf{x}_{ik}(j)$  is subvector of  $\mathbf{x}_{ik}(\omega)$ , which is the vector of the explanatory variables in the full model  $\omega$ , and we hereafter abbreviate the model index  $j$  and write  $\mathbf{x}_{ik}$  instead of  $\mathbf{x}_{ik}(j)$ ,  $p$  instead of  $p_j$  and others. Then, the NERM is

$$Y_{ik} = \mathbf{x}_{ik}^\top \boldsymbol{\beta} + b_i + \varepsilon_{ik} \quad (i = 1, \dots, q; k = 1, \dots, N_i), \quad (4.23)$$

where  $\boldsymbol{\beta}$  is a  $p \times 1$  vector of regression coefficients,  $b_i$  is a random effect for the  $i$ th area and  $b_i$ 's and  $\varepsilon_{ik}$ 's are mutually independently distributed as  $b_i \sim \mathcal{N}(0, \tau^2)$  and  $\varepsilon_{ik} \sim \mathcal{N}(0, \sigma^2)$ . We consider the situation that only  $n_i$  values of the  $Y_{ik}$ 's are observed through some sampling procedure. We define the number of the unobserved variables in the  $i$ th area by  $N_i - n_i = r_i$  and let  $n = n_1 + \dots + n_q$ ,  $r = r_1 + \dots + r_q$ . Suppose, without loss of generality, the first  $n_i$  elements of  $\{Y_{i1}, \dots, Y_{i, N_i}\}$  are observed, which are denoted by  $y_1, \dots, y_{i, n_i}$ , and  $Y_{i, n_i+1}, \dots, Y_{i, N_i}$  are unobserved. Then the observed model is defined as

$$y_{ik} = \mathbf{x}_{ik}^\top \boldsymbol{\beta} + b_i + \varepsilon_{ik} \quad (i = 1, \dots, q; k = 1, \dots, n_i), \quad (4.24)$$

which corresponds to (4.2) with  $\mathbf{y} = (\mathbf{y}_1^\top, \dots, \mathbf{y}_q^\top)^\top$  for  $\mathbf{y}_i = (y_{i1}, \dots, y_{i, n_i})^\top$ ,  $\mathbf{X} = (\mathbf{X}_1^\top, \dots, \mathbf{X}_q^\top)^\top$  for  $\mathbf{X}_i = (\mathbf{x}_{i1}, \dots, \mathbf{x}_{i, n_i})^\top$ ,  $\mathbf{Z} = \text{diag}(\mathbf{Z}_1, \dots, \mathbf{Z}_q)$  for  $\mathbf{Z}_i = \mathbf{1}_{n_i}$ ,  $\mathbf{G} = \psi \mathbf{I}_q$  and  $\mathbf{R} = \mathbf{I}_n$ , where  $\mathbf{1}_{n_i}$  denotes an  $n_i \times 1$  vector of ones and  $\psi = \tau^2/\sigma^2$ . In the derivation of our proposed criteria,



we have assumed that the covariance matrix of  $\mathbf{b}$  is  $\sigma^2 \mathbf{G}$  for a known matrix  $\mathbf{G}$ . However in the NERM,  $\mathbf{G}$  includes the parameter  $\psi$ , which is usually unknown and has to be estimated. In this case, we propose that  $\mathbf{G}$  in the bias correction should be replaced with its plug-in estimator  $\mathbf{G}(\hat{\psi})$ . The influence caused by the replacement may be limited because  $\psi$  is the nuisance parameter when one is interested in selecting only explanatory variables. Kawakubo and Kubokawa (2014) discussed the problem in their Remark 3.1.

We consider two types of predictive models. The first one can be used in the situation where all  $\mathbf{x}_{ik}$ 's are available. Then the predictive model, which we call the 'unit level predictive model', is defined by

$$Y_{ik} = \mathbf{x}_{ik}^T \boldsymbol{\beta} + b_i + \varepsilon_{ik} \quad (i = 1, \dots, q; k = n_i + 1, \dots, N_i), \quad (4.25)$$

which corresponds to (4.3) with  $\tilde{\mathbf{y}} = (\tilde{y}_1^T, \dots, \tilde{y}_q^T)^T$  for  $\tilde{y}_i = (Y_{i,n_i+1}, \dots, Y_{i,N_i})^T$ ,  $\tilde{\mathbf{X}} = (\tilde{\mathbf{X}}_1^T, \dots, \tilde{\mathbf{X}}_q^T)^T$  for  $\tilde{\mathbf{X}}_i = (\mathbf{x}_{i,n_i+1}, \dots, \mathbf{x}_{i,N_i}^T)^T$ ,  $\tilde{\mathbf{Z}} = \text{diag}(\tilde{\mathbf{Z}}_1, \dots, \tilde{\mathbf{Z}}_q)$  for  $\tilde{\mathbf{Z}}_i = \mathbf{1}_{r_i}$ ,  $\tilde{\mathbf{R}} = \mathbf{I}_r$ . Note that  $m = r$ .

In the problem of small area prediction, we often encounter the situation where all  $\mathbf{x}_{ik}$ 's are not observed but the area mean  $\bar{\mathbf{x}}_i = N_i^{-1} \sum_{k=1}^{N_i} \mathbf{x}_{ik}$  is known and we are interested in predicting  $\bar{Y}_i$ , which is the mean of finite population  $\{Y_{i1}, \dots, Y_{iN_i}\}$ , by using the value of  $\bar{\mathbf{x}}_i$ . Then the second type of predictive model, which we call the 'area level predictive model', can be defined as

$$\bar{Y}_{i(u)} = \bar{\mathbf{x}}_{i(u)}^T \boldsymbol{\beta} + b_i + \bar{\varepsilon}_{i(u)} \quad (i = 1, \dots, q), \quad (4.26)$$

where  $\bar{Y}_{i(u)} = r_i^{-1} \sum_{k=n_i+1}^{N_i} Y_{ik}$ , the mean of unobserved variables,  $\bar{\mathbf{x}}_{i(u)} = r_i^{-1} \sum_{k=n_i+1}^{N_i} \mathbf{x}_{ik}$ , calculated from  $\bar{\mathbf{x}}_i$  and  $(\mathbf{x}_{i1}, \dots, \mathbf{x}_{in_i})$ , and  $\bar{\varepsilon}_{i(u)} = r_i^{-1} \sum_{k=n_i+1}^{N_i} \varepsilon_{ik}$  distributed as  $\mathcal{N}(0, \sigma^2/r_i)$ . The model (4.26) corresponds to (4.3) with  $\tilde{\mathbf{y}} = (\bar{Y}_{1(u)}, \dots, \bar{Y}_{q(u)})^T$ ,  $\tilde{\mathbf{X}} = (\bar{\mathbf{x}}_{1(u)}, \dots, \bar{\mathbf{x}}_{q(u)})^T$ ,  $\tilde{\mathbf{Z}} = \mathbf{I}_q$  and  $\tilde{\mathbf{R}} = \text{diag}(\bar{R}_1, \dots, \bar{R}_q)$  for  $\bar{R}_i = 1/r_i$ . Note that  $m = q$ .

After selecting explanatory variables with our proposed criteria, we predict  $\bar{Y}_{i(u)}$  by the empirical best linear unbiased predictor  $\hat{\bar{Y}}_{i(u)} = \bar{\mathbf{x}}_{i(u)}^T \hat{\boldsymbol{\beta}} + \hat{b}_i$  and obtain a predictor of the finite population mean  $\bar{Y}_i$ , which is given as

$$\hat{\bar{Y}}_i = \frac{1}{N_i} \left\{ \sum_{k=1}^{n_i} y_{ik} + r_i \hat{\bar{Y}}_{i(u)} \right\}. \quad (4.27)$$

Thus, covariate shift appears in standard models for small area prediction and the proposed criterion is important and useful in such a situation.

## 4.5 Simulations

### 4.5.1 Simulations of measuring the biases of estimating the true cAI by the criteria

In this subsection, we compare the performance of the criteria by measuring the biases of estimating the cAI. We consider a class of the nested candidate models  $j_\alpha = \{1, \dots, \alpha\}$  for  $\alpha = 1, \dots, p_\omega$  where  $p_\omega = 7$ . The true observed model is the NERM in (4.24) with  $\sigma^2 = \tau^2 = 1$  and  $n_i = 3$  for  $i = 1, \dots, q$ . We consider the unit level predictive model (4.25) for the first experiment and the area level predictive model (4.26) for the second experiment. The explanatory variables in the full model  $\mathbf{x}_{ik}(\omega)$ 's ( $i = 1, \dots, q; k = 1, \dots, N_i$ ) are independently generated by  $\mathcal{N}(\mathbf{0}, \boldsymbol{\Sigma}_x)$  where  $\boldsymbol{\Sigma}_x = 0.9 \mathbf{I}_{p_\omega} + 0.1 \mathbf{J}_{p_\omega}$  for  $\mathbf{J}_{p_\omega} = \mathbf{1}_{p_\omega} \mathbf{1}_{p_\omega}^T$ . The true coefficient vector  $\boldsymbol{\beta}_*$  is  $\boldsymbol{\beta}_* = (\beta_1, \dots, \beta_{p_*}, 0, 0)$

for  $p_* = 5$  and  $\beta_l$ 's ( $1 \leq l \leq p_*$ ) are generated by  $\beta_l = 2 \times ((-1)^l / (l + 0.7)) \times U(1, 2)$  for a uniform random variable  $U(1, 2)$  on the interval  $(1, 2)$ . The values of the explanatory variables  $\mathbf{x}_{ik}$ 's and the vector of regression coefficients  $\beta_*$  are fixed through simulations.

Table 4.1: Relative biases of estimating cAI by CScAIC,  $\widehat{\text{cAI}}$  and  $\widehat{\text{cAI}}^\dagger$  for unit level predictive model

model	true value	relative bias		
	cAI	CScAIC	$\widehat{\text{cAI}}$	$\widehat{\text{cAI}}^\dagger$
pattern (a): $q = 10$				
$j_1$	206.38	-33.439	-0.33371	-0.080972
$j_2$	152.81	-18.840	-0.2414	-0.16414
$j_3$	140.79	-18.556	-0.23026	-0.46789
$j_4$	132.61	-11.451	0.26445	-0.19514
$j_5$	116.51	-0.0019291	1.5979	0.49514
$j_6$	122.46	-0.050686	0.77756	0.15429
$j_7$	128.88	0.0086256	0.09468	0.09468
pattern (b): $q = 15$				
$j_1$	233.35	-0.22534	0.11255	0.1159
$j_2$	189.54	9.431	0.24145	0.21667
$j_3$	177.28	14.197	0.42246	0.37347
$j_4$	163.62	-0.76563	0.32597	0.02934
$j_5$	152.94	0.13627	0.8566	0.25115
$j_6$	156.73	0.068668	0.53817	0.15002
$j_7$	161.65	0.083897	0.015869	0.015869
pattern (c): $q = 20$				
$j_1$	299.12	4.3775	0.084838	0.084654
$j_2$	252.60	6.1677	0.24072	0.23581
$j_3$	250.27	-5.1825	-0.016634	0.010911
$j_4$	208.25	2.6115	0.36013	0.23929
$j_5$	197.52	0.25682	0.53321	0.26101
$j_6$	200.38	0.24977	0.40855	0.25647
$j_7$	203.57	0.22713	0.21272	0.21272

Table 4.1 and 4.2 report the true values of the cAI and the relative biases of estimating the cAI by the criteria CScAIC in (4.7),  $\widehat{\text{cAI}}$  in (4.19) and  $\widehat{\text{cAI}}^\dagger$  in (4.22), for the experiment using the unit level predictive model and for the experiment using the area level predictive model, respectively. We handle the cases where the number of the areas  $q = 10, 15, 20$ . The true values of cAI in each candidate model are calculated based on (4.4) with 10000 Monte Carlo iterations. The relative biases of estimating the cAI by the criteria is defined as

$$100 \times \frac{E[\text{IC}] - \text{cAI}}{\text{cAI}},$$

where  $\text{IC} = \text{CScAIC}, \widehat{\text{cAI}}, \widehat{\text{cAI}}^\dagger$  and expectation is computed based on 1000 replications. From the tables, we can see the following facts. Firstly, the CScAIC has large biases for underspecified

Table 4.2: Relative biases of estimating cAI by CScAIC,  $\widehat{\text{cAI}}$  and  $\widehat{\text{cAI}}^\dagger$  for area level predictive model

model	true value	relative bias		
	cAI	CScAIC	$\widehat{\text{cAI}}$	$\widehat{\text{cAI}}^\dagger$
pattern (a): $q = 10$				
$j_1$	61.095	-10.429	-0.27157	-0.21359
$j_2$	46.635	5.4222	0.40292	-0.055653
$j_3$	50.744	-10.285	0.36735	-0.23857
$j_4$	47.323	-0.77484	1.3412	0.36198
$j_5$	45.735	-0.21108	2.2751	0.60073
$j_6$	49.452	-0.32466	0.82334	0.017969
$j_7$	52.805	-0.29278	-0.082743	-0.082743
pattern (b): $q = 15$				
$j_1$	95.393	-5.4716	0.12331	0.13930
$j_2$	70.056	17.521	0.39905	0.36599
$j_3$	66.412	21.039	0.59611	0.53872
$j_4$	61.310	5.174	0.58453	0.10853
$j_5$	60.532	0.23723	1.0853	0.25515
$j_6$	63.109	0.13276	0.50306	0.096935
$j_7$	65.379	0.16648	-0.0017143	-0.0017143
pattern (c): $q = 20$				
$j_1$	98.841	21.017	0.18161	0.17565
$j_2$	94.83	10.059	0.31607	0.30894
$j_3$	88.346	5.6464	0.13835	0.11302
$j_4$	78.383	7.8734	0.46942	0.27593
$j_5$	77.794	0.18930	0.54202	0.17009
$j_6$	79.277	0.18908	0.41365	0.18534
$j_7$	81.216	0.15395	0.11783	0.11783

models, namely  $j_1$ ,  $j_2$ ,  $j_3$  and  $j_4$ , while the modified estimators of the cAI,  $\widehat{\text{cAI}}$  and  $\widehat{\text{cAI}}^\dagger$  have smaller biases for both overspecified and underspecified models. Secondly,  $\widehat{\text{cAI}}^\dagger$  can estimate the cAI more unbiasedly than  $\widehat{\text{cAI}}$  can for the case of small sample size because  $\widehat{\text{cAI}}^\dagger$  is third order asymptotically unbiased estimator of the cAI. However, the relative biases of  $\widehat{\text{cAI}}$ , which is the second order asymptotically unbiased estimator of the cAI, gets smaller as the sample size is larger and the difference in performance between  $\widehat{\text{cAI}}$  and  $\widehat{\text{cAI}}^\dagger$  is not very important.

#### 4.5.2 Simulations of predicting finite population mean

In this subsection, we investigate the numerical performance of the small area prediction problem explained in Section 4.4. We consider the model class which consists of all subsets of the full model  $\omega = \{1, \dots, p_\omega\}$  for  $p_\omega = 5$ . The true observed model is the NERM in (4.23) with  $\sigma^2 = \tau^2 = 1$  and  $n_i = 3$  for  $i = 1, \dots, q$ , but we consider that  $\tau^2$  is unknown, namely  $\mathbf{G} = \psi \mathbf{I}_q$  for  $\psi = \tau^2/\sigma^2$  and  $\psi$  is the unknown parameter. Estimating procedure of  $\psi$  is the same as the one introduced in Section 3.4. We consider both the unit level predictive model (4.25) and the area level predictive model (4.26). The explanatory variables  $\mathbf{x}_{ik}$ 's for  $k = 1, \dots, n_i$ , namely the observed model, are independently generated by  $\mathcal{N}(4\mathbf{1}_{p_\omega}, \boldsymbol{\Sigma}_x)$  where  $\boldsymbol{\Sigma}_x = 0.9\mathbf{I}_{p_\omega} + 0.1\mathbf{J}_{p_\omega}$ . On the other hand, the explanatory variables  $\mathbf{x}_{ik}$ 's for  $k = n_i + 1, \dots, N_i$ , namely the predictive model, are independently generated by  $\mathcal{N}(a\mathbf{1}_{p_\omega}, \boldsymbol{\Sigma}_x)$  for  $a = 2, 4, 6$ . Under this setting, we measure the simulated prediction error of the best model selected by the conventional cAIC of Vaida and Blanchard (2005) and  $\widehat{\text{cAI}}$ 's based on unit level predictive model and area level predictive model. The prediction error is measured by quadratic loss

$$\sum_{i=1}^q (\widehat{Y}_i - \bar{Y})^2,$$

where  $\widehat{Y}$  is given by (4.27). The prediction errors are given as the averages based on 1000 replications.

Table 4.3: Prediction errors based on Vaida and Blanchard (2005)'s cAIC (VB),  $\widehat{\text{cAI}}$  using unit level predictive model (unit) and  $\widehat{\text{cAI}}$  using area level predictive model (area). The values in parentheses are the improvement over VB expressed in percentage.

	VB	unit	area
a=2	0.16966	0.16925 (0.24)	0.16887 (0.47)
a=4	0.15419	0.15382 (0.24)	0.15322 (0.63)
a=6	0.15346	0.15316 (0.20)	0.15257 (0.58)

Table 4.3 reports the prediction errors of the best model selected by the cAIC of Vaida and Blanchard (2005) denoted by 'VB',  $\widehat{\text{cAI}}$  using unit level predictive model denoted by 'unit', and  $\widehat{\text{cAI}}$  using area level predictive model denoted by 'area'. The values in parentheses are the improvement over the prediction error based on the cAIC expressed in percentage. It can be seen that our proposed criteria are better than the cAIC in the sense that predicting the small

area mean. It is valuable to point out that the prediction error of the mean of finite population can be improved by using our proposed criteria, which motivates us to use them for variable selection in small area prediction of the finite population.

## 4.6 Proofs

Firstly, we introduce the following lemma, which was shown in Section A.2 of Srivastava and Kubokawa (2010).

**Lemma 4.6** *Assume that  $\mathbf{C}$  is an  $n \times n$  symmetric matrix,  $\mathbf{M}$  is an idempotent matrix of rank  $p$  and that  $\mathbf{v} \sim \mathcal{N}(\mathbf{0}, \mathbf{I}_n)$ . Then,*

$$E \left[ \frac{\mathbf{v}^\top \mathbf{C} \mathbf{v}}{\mathbf{v}^\top (\mathbf{I}_n - \mathbf{M}) \mathbf{v}} \right] = \frac{\text{tr}(\mathbf{C})}{n - p - 2} - \frac{2 \text{tr}[\mathbf{C}(\mathbf{I}_n - \mathbf{M})]}{(n - p)(n - p - 2)}.$$

### 4.6.1 Proof of Theorem 4.1

Because the cAI can be evaluated as (4.6) for the overspecified case, it suffices to show that

$$\begin{aligned} & E[(\mathbf{y} - \mathbf{X}\hat{\boldsymbol{\beta}}_j - \mathbf{Z}\hat{\mathbf{b}}_j)^\top \tilde{\mathbf{R}}^{-1} (\mathbf{y} - \mathbf{X}\hat{\boldsymbol{\beta}}_j - \mathbf{Z}\hat{\mathbf{b}}_j) / \hat{\sigma}_j^2] \\ &= \frac{n}{n - p_j} \{ \text{tr}[\mathbf{R}\boldsymbol{\Sigma}^{-1}] - \text{tr}[\mathbf{R}\boldsymbol{\Sigma}^{-1} \mathbf{X}(\mathbf{X}^\top \boldsymbol{\Sigma}^{-1} \mathbf{X})^{-1} \mathbf{X}^\top \boldsymbol{\Sigma}^{-1}] \}. \end{aligned} \quad (4.28)$$

It follows that

$$\begin{aligned} & \mathbf{y} - \mathbf{X}\hat{\boldsymbol{\beta}}_j - \mathbf{Z}\hat{\mathbf{b}}_j \\ &= \mathbf{y} - \mathbf{X}\hat{\boldsymbol{\beta}}_j - \mathbf{Z}\mathbf{G}\mathbf{Z}^\top \boldsymbol{\Sigma}^{-1} (\mathbf{y} - \mathbf{X}\hat{\boldsymbol{\beta}}_j) \\ &= \mathbf{R}\boldsymbol{\Sigma}^{-1} \left\{ \mathbf{y} - \mathbf{X}\boldsymbol{\beta} - \mathbf{X}(\hat{\boldsymbol{\beta}}_j - \boldsymbol{\beta}) \right\} \\ &= \mathbf{R}\boldsymbol{\Sigma}^{-1} (\mathbf{y} - \mathbf{X}\boldsymbol{\beta}) - \mathbf{R}\boldsymbol{\Sigma}^{-1} \mathbf{X}(\hat{\boldsymbol{\beta}}_j - \boldsymbol{\beta}). \end{aligned}$$

Then, we can see that

$$\begin{aligned} & E[(\mathbf{y} - \mathbf{X}\hat{\boldsymbol{\beta}}_j - \mathbf{Z}\hat{\mathbf{b}}_j)^\top \tilde{\mathbf{R}}^{-1} (\mathbf{y} - \mathbf{X}\hat{\boldsymbol{\beta}}_j - \mathbf{Z}\hat{\mathbf{b}}_j) / \hat{\sigma}_j^2] \\ &= E[(\mathbf{y} - \mathbf{X}\boldsymbol{\beta})^\top \boldsymbol{\Sigma}^{-1} \mathbf{R}\boldsymbol{\Sigma}^{-1} (\mathbf{y} - \mathbf{X}\boldsymbol{\beta}) / \hat{\sigma}_j^2] + E[(\hat{\boldsymbol{\beta}}_j - \boldsymbol{\beta})^\top \mathbf{X}^\top \boldsymbol{\Sigma}^{-1} \mathbf{R}\boldsymbol{\Sigma}^{-1} \mathbf{X}(\hat{\boldsymbol{\beta}}_j - \boldsymbol{\beta}) / \hat{\sigma}_j^2] \\ &\quad - 2E[(\mathbf{y} - \mathbf{X}\boldsymbol{\beta})^\top \boldsymbol{\Sigma}^{-1} \mathbf{R}\boldsymbol{\Sigma}^{-1} \mathbf{X}(\hat{\boldsymbol{\beta}}_j - \boldsymbol{\beta}) / \hat{\sigma}_j^2]. \end{aligned}$$

We define  $\mathbf{v} = \boldsymbol{\Sigma}^{-1/2}(\mathbf{y} - \mathbf{X}\boldsymbol{\beta})/\sigma_*$  and  $\mathbf{M} = \boldsymbol{\Sigma}^{-1/2} \mathbf{X}(\mathbf{X}^\top \boldsymbol{\Sigma}^{-1} \mathbf{X})^{-1} \mathbf{X}^\top \boldsymbol{\Sigma}^{-1/2}$ , which are the same notations as those in Section 4.3. Then, we can rewrite the equation above as

$$\begin{aligned} & E[(\mathbf{y} - \mathbf{X}\hat{\boldsymbol{\beta}}_j - \mathbf{Z}\hat{\mathbf{b}}_j)^\top \tilde{\mathbf{R}}^{-1} (\mathbf{y} - \mathbf{X}\hat{\boldsymbol{\beta}}_j - \mathbf{Z}\hat{\mathbf{b}}_j) / \hat{\sigma}_j^2] \\ &= nE \left[ \frac{\mathbf{v}^\top \boldsymbol{\Sigma}^{-1/2} \mathbf{R}\boldsymbol{\Sigma}^{-1/2} \mathbf{v}}{\mathbf{v}^\top (\mathbf{I}_n - \mathbf{M}) \mathbf{v}} \right] + nE \left[ \frac{\mathbf{v}^\top \mathbf{M}\boldsymbol{\Sigma}^{-1/2} \mathbf{R}\boldsymbol{\Sigma}^{-1/2} \mathbf{M}\mathbf{v}}{\mathbf{v}^\top (\mathbf{I}_n - \mathbf{M}) \mathbf{v}} \right] - 2nE \left[ \frac{\mathbf{v}^\top \boldsymbol{\Sigma}^{-1/2} \mathbf{R}\boldsymbol{\Sigma}^{-1/2} \mathbf{M}\mathbf{v}}{\mathbf{v}^\top (\mathbf{I}_n - \mathbf{M}) \mathbf{v}} \right] \\ &= I_1 + I_2 - 2I_3 \quad (\text{say}). \end{aligned}$$

It follows from Lemma 4.6 that

$$\begin{aligned} I_1 &= n \times \left\{ \frac{\text{tr}[\mathbf{R}\boldsymbol{\Sigma}^{-1}]}{n-p_j-2} - \frac{2\text{tr}[\boldsymbol{\Sigma}^{-1/2}\mathbf{R}\boldsymbol{\Sigma}^{-1/2}(\mathbf{I}_n - \mathbf{M})]}{(n-p_j)(n-p_j-2)} \right\} \\ &= n \times \left\{ \frac{\text{tr}[\mathbf{R}\boldsymbol{\Sigma}^{-1}]}{n-p_j-2} - \frac{2\text{tr}[\mathbf{R}\boldsymbol{\Sigma}^{-1}] - 2\text{tr}[\mathbf{R}\mathbf{P}]}{(n-p_j)(n-p_j-2)} \right\}. \end{aligned} \quad (4.29)$$

As for  $I_2$ , because  $\mathbf{v}^\top \mathbf{M}\boldsymbol{\Sigma}^{-1/2}\mathbf{R}\boldsymbol{\Sigma}^{-1/2}\mathbf{v}$  is independent of  $\mathbf{v}^\top(\mathbf{I}_n - \mathbf{M})\mathbf{v}$ , we can obtain

$$I_2 = \frac{n}{n-p_j-2} \cdot \text{tr}[\mathbf{R}\mathbf{P}].$$

To evaluate  $I_3$ , we rewrite

$$I_3 = nE \left[ \frac{\mathbf{v}^\top \mathbf{M}\boldsymbol{\Sigma}^{-1/2}\mathbf{R}\boldsymbol{\Sigma}^{-1/2}\mathbf{M}\mathbf{v}}{\mathbf{v}^\top(\mathbf{I}_n - \mathbf{M})\mathbf{v}} \right] + nE \left[ \frac{\mathbf{v}^\top(\mathbf{I}_n - \mathbf{M})\boldsymbol{\Sigma}^{-1/2}\mathbf{R}\boldsymbol{\Sigma}^{-1/2}\mathbf{M}\mathbf{v}}{\mathbf{v}^\top(\mathbf{I}_n - \mathbf{M})\mathbf{v}} \right]. \quad (4.30)$$

Because  $\mathbf{M}\mathbf{v}$  is independent of  $(\mathbf{I}_n - \mathbf{M})\mathbf{v}$  and  $E[\mathbf{M}\mathbf{v}] = \mathbf{0}$ , the second term of the right hand side of the above equation is 0. Then we get  $I_3 = I_2$  and (4.28) form (4.29) and (4.30).  $\square$

#### 4.6.2 Proof of Theorem of 4.2

For the overspecified case, the cAI is evaluated as (4.6), which is the same expression as Theorem 4.2. Thus we show that the cAI in (4.11) is evaluated as (4.12) for the underspecified case.

Because  $E[(K_0 + K_1)^{-1}]$  is evaluated as (3.32) in Chapter 3, it suffices to evaluate  $E[\mathbf{a}^\top \tilde{\mathbf{R}}^{-1} \mathbf{a} / \sigma_*^2]$ . Let  $\mathbf{u} = \mathbf{y} - \mathbf{X}(\omega)\boldsymbol{\beta}_*$ . Then, we can rewrite  $\tilde{\mathbf{X}}(j)\hat{\boldsymbol{\beta}}_j - \tilde{\mathbf{X}}(\omega)\boldsymbol{\beta}_*$  in  $\mathbf{a}$  as

$$\begin{aligned} \tilde{\mathbf{X}}(j)\hat{\boldsymbol{\beta}}_j - \tilde{\mathbf{X}}(\omega)\boldsymbol{\beta}_* &= \tilde{\mathbf{P}}_j(\mathbf{X}(\omega)\boldsymbol{\beta}_* + \mathbf{u}) - \tilde{\mathbf{X}}(\omega)\boldsymbol{\beta}_* \\ &= (\tilde{\mathbf{P}}_j\mathbf{X}(\omega) - \tilde{\mathbf{X}}(\omega))\boldsymbol{\beta}_* + \sigma_*\tilde{\mathbf{P}}_j\boldsymbol{\Sigma}^{1/2}\mathbf{v}. \end{aligned}$$

Next, we can rewrite  $\mathbf{X}(j)\hat{\boldsymbol{\beta}}_j - \mathbf{X}(\omega)\boldsymbol{\beta}_*$  in  $\mathbf{a}$  as

$$\begin{aligned} \mathbf{X}(j)\hat{\boldsymbol{\beta}}_j - \mathbf{X}(\omega)\boldsymbol{\beta}_* &= \sigma_*\boldsymbol{\Sigma}^{1/2}\{\mathbf{M}_j(\mathbf{W}_\omega\boldsymbol{\beta}_* + \mathbf{v}) - \mathbf{W}_\omega\boldsymbol{\beta}_*\} \\ &= \sigma_*\boldsymbol{\Sigma}^{1/2}\{-(\mathbf{M}_\omega - \mathbf{M}_j)\mathbf{W}_\omega\boldsymbol{\beta}_* + \mathbf{M}_j\mathbf{v}\} \\ &= -\boldsymbol{\Sigma}(\mathbf{P}_\omega - \mathbf{P}_j)\mathbf{X}(\omega)\boldsymbol{\beta}_* + \sigma_*\boldsymbol{\Sigma}^{1/2}\mathbf{M}_j\mathbf{v}. \end{aligned}$$

Then, we obtain

$$\mathbf{a} = \mathbf{B}\boldsymbol{\beta}_* + \sigma_*(\tilde{\mathbf{P}}_j\boldsymbol{\Sigma}^{1/2} - \tilde{\mathbf{Z}}\mathbf{G}\mathbf{Z}^\top\boldsymbol{\Sigma}^{-1/2}\mathbf{M}_j)\mathbf{v}.$$

Moreover, it follows that

$$\begin{aligned} \tilde{\mathbf{P}}_j\boldsymbol{\Sigma}^{1/2} - \tilde{\mathbf{Z}}\mathbf{G}\mathbf{Z}^\top\boldsymbol{\Sigma}^{-1/2}\mathbf{M}_j &= (\tilde{\mathbf{X}}(j) - \tilde{\mathbf{Z}}\mathbf{G}\mathbf{Z}^\top\boldsymbol{\Sigma}^{-1}\mathbf{X}(j))(\mathbf{X}(j)^\top\boldsymbol{\Sigma}^{-1}\mathbf{X}(j))^{-1}\mathbf{X}(j)^\top\boldsymbol{\Sigma}^{-1/2} \\ &= \mathbf{A}(\mathbf{X}(j)^\top\boldsymbol{\Sigma}^{-1}\mathbf{X}(j))^{-1}\mathbf{X}(j)^\top\boldsymbol{\Sigma}^{-1/2} \end{aligned}$$

Thus,  $E[\mathbf{a}^\top \tilde{\mathbf{R}}^{-1} \mathbf{a} / \sigma_*^2]$  can be evaluated as

$$E[\mathbf{a}^\top \tilde{\mathbf{R}}^{-1} \mathbf{a} / \sigma_*^2] = \text{tr}[\tilde{\mathbf{R}}^{-1}\mathbf{A}(\mathbf{X}(j)^\top\boldsymbol{\Sigma}^{-1}\mathbf{X}(j))^{-1}\mathbf{A}^\top] + \boldsymbol{\beta}_*^\top \mathbf{B}^\top \tilde{\mathbf{R}}^{-1} \mathbf{B}\boldsymbol{\beta}_* / \sigma_*^2. \quad (4.31)$$

It follows from (3.32) and (4.31) that

$$\begin{aligned} &n \cdot E[(K_0 + K_1)^{-1}]\{\text{tr}(\tilde{\mathbf{R}}^{-1}\boldsymbol{\Lambda}) + E[\mathbf{a}^\top \tilde{\mathbf{R}}^{-1} \mathbf{a} / \sigma_*^2]\} \\ &= (\gamma + \boldsymbol{\beta}_*^\top \mathbf{B}^\top \tilde{\mathbf{R}}^{-1} \mathbf{B}\boldsymbol{\beta}_* / \sigma_*^2) \times \left\{ \lambda + \frac{-2\lambda^3 + (p_j + 4)\lambda^2}{n} \right\} + O(n^{-1}), \end{aligned}$$

which shows that the cAI in (4.11) is approximated to (4.12).  $\square$

### 4.6.3 Proof of Lemma 4.1

When the candidate model  $j$  is overspecified,  $E(\hat{\lambda}) = E(\widehat{\lambda^2}) = E(\widehat{\lambda^3}) = 1$  because  $\hat{\sigma}_\omega^2/\hat{\sigma}_j^2 \sim \text{Be}((n - p_\omega)/2, (p_\omega - p_j)/2)$ . Thus, it follows that  $E(\widehat{R_2}) = 0$ .

When the candidate model  $j$  is underspecified, it follows that

$$E \left[ \left( \frac{\hat{\sigma}_\omega^2}{\hat{\sigma}_j^2} \right)^k \right] = E \left[ \left( \frac{K_0}{K_0 + K_1} \right)^k \right] = \lambda^k + O(n^{-1}), \quad (4.32)$$

for  $k = 1, 2, 3$  because  $(K_0 + K_1)^{-1}$  is expanded as (3.31) and  $K_0 = \mathbf{v}^\top (\mathbf{I}_n - \mathbf{M}_\omega) \mathbf{v}$ . Thus, it follows that  $\widehat{R_2} = R_2 + O(n^{-1})$ .  $\square$

### 4.6.4 Proof of Lemma 4.2

When the candidate model  $j$  is overspecified,  $E(\widehat{R_1}) = 0$  because  $E(\hat{\lambda}) = E(\widehat{\lambda^2}) = E(\widehat{\lambda^3}) = 1$ .

When the candidate model  $j$  is underspecified, it follows from (3.33) that

$$E(\hat{\lambda}) = \lambda + \frac{-\lambda^3 + (p_j + 2)\lambda^2 - p_j\lambda}{n} + O(n^{-2}). \quad (4.33)$$

Thus it follows from (4.32) and (4.33) that  $E(\widehat{R_1}) = R_1 + O(n^{-1})$ .  $\square$

### 4.6.5 Proof of Lemma 4.3

Firstly, note that

$$R_3(\boldsymbol{\eta}_*) = \lambda \cdot \boldsymbol{\beta}_*^\top \mathbf{B}^\top \widetilde{\mathbf{R}}^{-1} \mathbf{B} \boldsymbol{\beta}_* / \sigma_*^2,$$

where  $\lambda = 1/(1 + \delta)$  for

$$\delta = \boldsymbol{\beta}_*^\top \mathbf{X}(\omega)^\top (\mathbf{P}_\omega - \mathbf{P}_j) \mathbf{X}(\omega) \boldsymbol{\beta}_* / (n\sigma_*^2).$$

Then, we can see that

$$\frac{\partial R_3(\boldsymbol{\eta}_*)}{\partial \boldsymbol{\eta}_*} = \frac{\partial \lambda}{\partial \delta} \cdot \frac{\partial \delta}{\partial \boldsymbol{\eta}_*}$$

and that  $\partial \lambda / \partial \delta = -(1 + \delta)^{-2}$ . After some calculations, we can obtain Lemma 4.3.  $\square$

### 4.6.6 Proof of Lemma 4.4

Firstly, note that  $E[R_3(\widetilde{\boldsymbol{\eta}})]$  is expanded as

$$E[R_3(\widetilde{\boldsymbol{\eta}})] = R_3(\boldsymbol{\eta}_*) + B_1(\boldsymbol{\eta}_*) + O(n^{-1}),$$

where  $B_1(\boldsymbol{\eta}_*)$  is given as (4.17). Because  $B_1(\boldsymbol{\eta}_*) = O(1)$ , it follows that  $B_1(\widetilde{\boldsymbol{\eta}}) = B_1(\boldsymbol{\eta}_*) + O(n^{-1})$ , which shows that

$$E[\widetilde{R_3}] = R_3 + O(n^{-1}).$$

In the same way, we can obtain  $E[\widetilde{R_4}] = E[R_4(\widetilde{\boldsymbol{\eta}})] = R_4(\boldsymbol{\eta}_*) + O(n^{-1})$ .  $\square$

### 4.6.7 Proof of Lemma 4.5

It follows from (4.21) that

$$E[\widehat{R}_3] = E \left[ 2R_3(\tilde{\boldsymbol{\eta}}) - E_{\tilde{\boldsymbol{\eta}}}[R_3(\tilde{\boldsymbol{\eta}}^\dagger)] + E_{\tilde{\boldsymbol{\eta}}}[B_1(\tilde{\boldsymbol{\eta}}^\dagger)] - B_1(\tilde{\boldsymbol{\eta}}) \right].$$

Because  $E[R_3(\tilde{\boldsymbol{\eta}})]$  is expanded as  $E[R_3(\tilde{\boldsymbol{\eta}})] = R_3(\boldsymbol{\eta}_*) + B_1(\boldsymbol{\eta}_*) + B_2(\boldsymbol{\eta}_*) + O(n^{-2})$ , we can see that

$$\begin{aligned} E \left[ 2R_3(\tilde{\boldsymbol{\eta}}) - E_{\tilde{\boldsymbol{\eta}}}[R_3(\tilde{\boldsymbol{\eta}}^\dagger)] \right] &= 2 \{R_3(\boldsymbol{\eta}_*) + B_1(\boldsymbol{\eta}_*) + B_2(\boldsymbol{\eta}_*)\} - E[R_3(\tilde{\boldsymbol{\eta}}) + B_1(\tilde{\boldsymbol{\eta}}) + B_2(\tilde{\boldsymbol{\eta}})] + O(n^{-2}) \\ &= R_3(\boldsymbol{\eta}_*) + B_1(\boldsymbol{\eta}_*) + B_2(\boldsymbol{\eta}_*) - E[B_1(\tilde{\boldsymbol{\eta}}) + B_2(\tilde{\boldsymbol{\eta}})] + O(n^{-2}). \end{aligned}$$

Moreover, because  $E[B_1(\tilde{\boldsymbol{\eta}})] = B_1(\boldsymbol{\eta}_*) + B_{11}(\boldsymbol{\eta}_*) + O(n^{-2})$  and  $E[B_2(\tilde{\boldsymbol{\eta}})] = B_2(\boldsymbol{\eta}_*) + O(n^{-2})$ , the equation above can be rewritten as

$$E \left[ 2R_3(\tilde{\boldsymbol{\eta}}) - E_{\tilde{\boldsymbol{\eta}}}[R_3(\tilde{\boldsymbol{\eta}}^\dagger)] \right] = R_3(\boldsymbol{\eta}_*) - B_{11}(\boldsymbol{\eta}_*) + O(n^{-2}). \quad (4.34)$$

Next, it is seen that

$$\begin{aligned} E \left[ E_{\tilde{\boldsymbol{\eta}}}[B_1(\tilde{\boldsymbol{\eta}}^\dagger)] - B_1(\tilde{\boldsymbol{\eta}}) \right] &= E[B_1(\tilde{\boldsymbol{\eta}}) + B_{11}(\tilde{\boldsymbol{\eta}})] - \{B_1(\boldsymbol{\eta}_*) + B_{11}(\boldsymbol{\eta}_*)\} + O(n^{-2}) \\ &= B_{11}(\boldsymbol{\eta}_*) + O(n^{-2}). \end{aligned} \quad (4.35)$$

Thus, it follows from (4.34) and (4.35) that

$$E[\widehat{R}_3] = R_3(\boldsymbol{\eta}_*) + O(n^{-2}).$$

Similarly, we can show that  $E[\widehat{R}_4] = R_4(\boldsymbol{\eta}_*) + O(n^{-2})$ . □

### 4.6.8 Proof of Theorem 4.3 and Theorem 4.4

From Theorem 4.2 and Lemmas 4.1–4.5, we can easily show the theorems. □



## Chapter 5

# Conditional AIC in mixed effects models based on natural exponential family

In this chapter, we consider the variable selection problem for the class of mixed effects models based on natural exponential family, which includes useful nonlinear mixed models, Poisson-gamma model and binomial-beta model. We construct the conditional AIC in the models and show the usefulness for variable selection.

### 5.1 Motivation

For variable selection problem in mixed effects models, Vaida and Blanchard (2005) introduced the conditional Akaike information (cAI), which is related to the expected Kullback–Leibler divergence based on the conditional likelihood given random effects. The cAI and the resulting information criterion conditional AIC (cAIC) are appropriate when one is interested in predicting the random effects. Vaida and Blanchard (2005) proposed the cAIC for variable selection criterion in normal linear mixed model. Since then, the cAIC has been studied for various models, which include generalized linear mixed model (GLMM) as well as linear mixed model.

However, variable selection problem in nonlinear mixed model, for example Poisson-gamma model or binomial-beta model, has not been considered well. Although the conventional AIC (or marginal AIC, mAIC), which is based on the marginal likelihood integrating out the random effects, can be used as a variable selection criterion, the mAIC is not appropriate for predicting the random effects. Then, we consider the cAI for mixed effects models based on natural exponential family and construct the cAIC as an asymptotically unbiased estimator of the cAI.

The rest of this chapter is organized as follows. In Section 5.2, we explain about the mixed effects models based on natural exponential family and define the cAI for the class of the models. In Section 5.3, we evaluate and estimate the bias correction and construct the cAIC in three ways. The numerical performance of the proposed criteria is investigated by simulations in Section 5.4. Section 5.5 shows some results of analytical calculations and Section 5.6 gives proofs of lemmas and theorems.

## 5.2 Model and conditional AIC

### 5.2.1 Mixed effects models based on natural exponential family

Let  $y_1, \dots, y_m$  be mutually independent random variables where the conditional distribution of  $y_i$  given  $\theta_i$  belongs to the following natural exponential family (NEF):

$$y_i|\theta_i \sim f(y_i|\theta_i) = \exp[n_i(\theta_i y_i - \psi(\theta_i)) + c(y_i, n_i)], \quad (i = 1, \dots, m), \quad (5.1)$$

where  $\theta_i$  is the natural parameter,  $n_i$  is a known scale parameter and  $\psi(\cdot)$  and  $c(\cdot, \cdot)$  are functions specific to each distribution. As an example of the random variable  $y_i$  with probability density (or mass) function (5.1), we consider the following situation.

For  $i = 1, \dots, m$ , which denote clusters (or areas in the context of small area estimation), let the random variables  $Z_{i1}, \dots, Z_{i,n_i}$  be independent and identically distributed given  $\theta_i$  with the common distribution belonging to the one parameter exponential family, namely,

$$P(Z_{ij} \in A) = \int_A \exp\{\theta_i z - \psi(\theta_i)\} dF(z),$$

with  $F$  a Stieltjes measure on  $\mathbb{R}$ . When we define  $y_i = (Z_{i1} + \dots + Z_{i,n_i})/n_i$ , the conditional distribution of  $y_i$  given  $\theta_i$  has the density function of the form (5.1). This situation is also considered in Ghosh and Maiti (2008), who developed empirical Bayes (EB) confidence intervals for population means in each small area with NEF distributions. In their work, the result is based on an asymptotic theory in the sense that  $n_i \rightarrow \infty$  and  $m \rightarrow \infty$ . We also consider the same setup to derive a criterion in Section 5.3.2.

From the property of NEF, the mean of  $y_i$  given  $\theta_i$  is

$$\mu_i = E(y_i|\theta_i) = \psi'(\theta_i),$$

where  $\psi'(\cdot)$  denotes the derivative of  $\psi(\cdot)$ . Define  $Q(\mu_i) = \psi''(\theta_i)$ , where  $\psi''(\cdot)$  is the second derivative of  $\psi(\cdot)$ , then the variance of  $y_i$  given  $\theta_i$  is

$$V(y_i|\theta_i) = \frac{\psi''(\theta_i)}{n_i} = \frac{Q(\mu_i)}{n_i}.$$

$Q(\cdot)$  is called the variance function and we hereafter assume that  $Q(\cdot)$  is a quadratic function, namely  $Q(x) = v_0 + v_1 x + v_2 x^2$  for known constants  $v_0$ ,  $v_1$  and  $v_2$ , which are not simultaneously zero. The family of such distributions is called natural exponential family with quadratic variance function (NEF-QVF), many properties of which are studied by Morris (1982, 1983).

As the random cluster (area) effect, we consider the conjugate prior for  $\theta_i$  with the probability density function

$$p(\theta_i|\lambda, m_i) = \exp[\lambda(m_i \theta_i - \psi(\theta_i))] C(\lambda, m_i), \quad (5.2)$$

where  $\lambda$  is an unknown scalar hyperparameter and  $C(\cdot, \cdot)$  is a function specific to each distribution. The mean and variance of  $\mu_i = E(y_i|\theta_i)$  are

$$E(\mu_i|\lambda, m_i) = m_i \quad \text{and} \quad V(\mu_i|\lambda, m_i) = \frac{Q(m_i)}{\lambda - v_2},$$

respectively. When  $p \times 1$  auxiliary variable  $\mathbf{x}_i$  is available, one uses  $\mathbf{x}_i^T \boldsymbol{\beta}$  as a predictor of  $m_i$ , where  $\boldsymbol{\beta}$  is an unknown  $p \times 1$  vector of regression coefficients. We here consider the following link function  $h(\cdot)$  between  $m_i$  and  $\mathbf{x}_i^T \boldsymbol{\beta}$ :

$$m_i = h^{-1}(\mathbf{x}_i^T \boldsymbol{\beta}).$$

If  $h^{-1}(\cdot) = \psi'(\cdot)$ , the link function  $h(\cdot)$  is called canonical. We henceforth consider the canonical link. We define the vector of the unknown parameters in this model as  $\boldsymbol{\eta} = (\boldsymbol{\beta}^T, \lambda)^T$ .

The posterior density function of  $\theta_i$  given  $y_i$ , and the marginal density (or mass) function of  $y_i$  are

$$p(\theta_i|y_i, \lambda, m_i) = \exp[(n_i + \lambda)(\hat{\mu}_i^B \theta_i - \psi(\theta_i))]C(n_i + \lambda, \hat{\mu}_i^B), \quad (5.3)$$

$$m(y_i|\lambda, m_i) = \frac{C(\lambda, m_i)}{C(n_i + \lambda, \hat{\mu}_i^B)} \exp[c(y_i, n_i)], \quad (5.4)$$

respectively, where  $\hat{\mu}_i^B$  is the Bayes estimator of  $\mu_i$  under quadratic loss, namely

$$\hat{\mu}_i^B = \hat{\mu}_i(y_i, \boldsymbol{\eta}) = \frac{n_i y_i + \lambda m_i}{n_i + \lambda}, \quad m_i = \psi'(\mathbf{x}_i^T \boldsymbol{\beta}).$$

We can get the empirical Bayes (EB) estimator of  $\mu_i$  by substituting  $\hat{\boldsymbol{\eta}} = (\hat{\boldsymbol{\beta}}^T, \hat{\lambda})^T$  for  $\boldsymbol{\eta}$  in  $\hat{\mu}_i(y_i, \boldsymbol{\eta})$  as follows:

$$\hat{\mu}_i^{\text{EB}} = \hat{\mu}_i(y_i, \hat{\boldsymbol{\eta}}) = \frac{n_i y_i + \hat{\lambda} \hat{m}_i}{n_i + \hat{\lambda}}, \quad \hat{m}_i = \psi'(\mathbf{x}_i^T \hat{\boldsymbol{\beta}}),$$

where  $\hat{\boldsymbol{\beta}}$  and  $\hat{\lambda}$  are some estimators of  $\boldsymbol{\beta}$  and  $\lambda$  based on the marginal distribution of  $\mathbf{y} = (y_1, \dots, y_m)^T$ .

The model explained above is an example of general mixed effects model (2.1) for  $\mathbf{y}$  and  $\boldsymbol{\theta} = (\theta_1, \dots, \theta_m)^T$  and is called nonlinear mixed model unless  $h(\cdot)$  is the identity function, which is the case of normal linear mixed model. Nonlinear mixed model is useful in the context of small area estimation and is used by Ghosh and Maiti (2004, 2008), Lohr and Rao (2009) and others. We provide two useful examples belonging to the nonlinear mixed model, one for the count and the other for binary data sets.

**[1] Poisson-gamma mixture model.** Let  $z_1, \dots, z_m$  be mutually independent random variables having

$$z_i|\mu_i \sim \text{Po}(n_i \mu_i) \quad \text{and} \quad \mu_i \sim \text{Ga}(\lambda m_i, \lambda^{-1}),$$

where  $\mu_1, \dots, \mu_m$  are mutually independent,  $\text{Po}(\mu)$  denotes Poisson distribution with mean  $\mu$ , and  $\text{Ga}(a, b)$  denotes gamma distribution with density function

$$f(x) = \frac{1}{\Gamma(a) b^a} x^{a-1} \exp(-x/b), \quad x > 0,$$

where  $\Gamma(\cdot)$  denotes gamma function. Let  $y_i = z_i/n_i$ . Then the probability mass function of  $y_i$  given  $\mu_i$  is (5.1) with  $\theta_i = \log(\mu_i)$ ,  $\psi(\cdot) = \exp(\cdot)$  and the probability density function of  $\theta_i$  is (5.2). In this model, canonical link function is  $h(\cdot) = \psi'^{-1}(\cdot) = \log(\cdot)$  and the quadratic variance function is  $Q(x) = x$ , namely  $v_1 = 1$ ,  $v_0 = v_2 = 0$ .

**[2] binomial-beta mixture model.** Let  $z_1, \dots, z_m$  be mutually independent random variables having

$$z_i|\mu_i \sim \text{Bin}(n_i, \mu_i) \quad \text{and} \quad \mu_i \sim \text{Beta}(\lambda m_i, \lambda(1 - m_i)),$$

where  $\mu_1, \dots, \mu_m$  are mutually independent,  $\text{Bin}(n, p)$  denotes binomial distribution, and  $\text{Beta}(a, b)$  denotes beta distribution. Let  $y_i = z_i/n_i$ . Then the probability mass function of  $y_i$  given  $\mu_i$  is (5.1) with  $\theta_i = \text{logit}(\mu_i)$  and  $\psi(\theta_i) = \log(1 + \exp(\theta_i))$  and the probability density function of  $\theta_i$  is (5.2). In this model, canonical link function is  $h(\cdot) = \psi'^{-1}(\cdot) = \text{logit}(\cdot)$  and the quadratic variance function is  $Q(x) = x - x^2$ , namely  $v_0 = 0$ ,  $v_1 = 1$ ,  $v_2 = -1$ .

### 5.2.2 Variable selection problem in nonlinear mixed model

For the variable selection problem in the nonlinear mixed model, the marginal AIC (mAIC) can be easily used, because the marginal likelihood (5.4) is obtained analytically. However, when one wants to predict the random area effect  $\theta_i$ , the mAIC based on the marginal likelihood is not appropriate (Vaida and Blanchard, 2005). Then we derive the cAIC for nonlinear mixed model.

In the same way as the previous chapters, we define candidate models by the index set  $j$ , which is a subset of  $\omega = \{1, \dots, p_\omega\}$ . Then, we define  $\mathbf{X}(j) = (\mathbf{x}_1(j), \dots, \mathbf{x}_m(j))^T$ , where  $\mathbf{x}_i(j)$  is a  $p_j \times 1$  vector for  $p_j = \#(j)$ . Let  $j_*$  denote the true model and the dimension of the true model be  $p_{j_*}$ , which is abbreviated to  $p_*$ . To derive the criterion, we here assume that the candidate model  $j$  is overspecified. Under the assumption, the mean of the true model can be expressed as

$$m_i = \psi'(\mathbf{x}_i(j)^T \boldsymbol{\beta}_j^*),$$

where  $\boldsymbol{\beta}_j^*$  is  $p_j \times 1$  vector of regression coefficients, whose  $p_j - p_*$  components are exactly 0 and the rest of components are not 0. Thus we henceforth abbreviate  $\mathbf{x}_i(j)$  to  $\mathbf{x}_i$ ,  $\boldsymbol{\beta}_j^*$  to  $\boldsymbol{\beta}$  for notational convenience.

### 5.2.3 Conditional AIC in nonlinear mixed model

Vaida and Blanchard (2005) proposed the cAIC as an (asymptotically) unbiased estimator of a part of the prediction risk of the plug-in predictive density relative to the expected Kullback–Leibler divergence, which is called the conditional Akaike Information (cAI). Firstly, let us define the conditional Akaike Information (cAI) for the nonlinear mixed model. Let  $\tilde{y}_i$  be an independent replication of  $y_i$  given  $\theta_i$ . Then the logarithm of the plug-in predictive density is

$$\log\{f(\tilde{y}_i|\hat{\theta}_i^{\text{EB}})\} = n_i(\hat{\theta}_i^{\text{EB}}\tilde{y}_i - \psi(\hat{\theta}_i^{\text{EB}})) + c(\tilde{y}_i, n_i),$$

where  $\hat{\theta}_i^{\text{EB}} = \psi'^{-1}(\hat{\mu}_i^{\text{EB}}) = h(\hat{\mu}_i^{\text{EB}})$ . Because the second term of the equation above, namely  $c(\tilde{y}_i, n_i)$ , is irrelevant to the candidate model, we define the cAI in the nonlinear mixed model as

$$\begin{aligned} \text{cAI} &= -2E^{(\mathbf{y}, \boldsymbol{\theta})} E^{\tilde{\mathbf{y}}|\boldsymbol{\theta}} \left[ \sum_{i=1}^m n_i(\hat{\theta}_i^{\text{EB}}\tilde{y}_i - \psi(\hat{\theta}_i^{\text{EB}})) \right] \\ &= -2E^{(\mathbf{y}, \boldsymbol{\theta})} \left[ \sum_{i=1}^m n_i(\hat{\theta}_i^{\text{EB}}\mu_i - \psi(\hat{\theta}_i^{\text{EB}})) \right], \end{aligned} \quad (5.5)$$

where  $E^{(\mathbf{y}, \boldsymbol{\theta})}$  and  $E^{\tilde{\mathbf{y}}|\boldsymbol{\theta}}$  denote expectation with respect to the joint distribution of  $(\mathbf{y}, \boldsymbol{\theta})$  and the conditional distribution of  $\tilde{\mathbf{y}}$  given  $\boldsymbol{\theta}$  for  $\tilde{\mathbf{y}} = (\tilde{y}_1, \dots, \tilde{y}_m)^T$ , respectively. When we estimate the cAI by  $-2 \sum_{i=1}^m n_i(\hat{\theta}_i^{\text{EB}}y_i - \psi(\hat{\theta}_i^{\text{EB}}))$ , the bias is

$$\begin{aligned} E \left[ -2 \sum_{i=1}^m n_i(\hat{\theta}_i^{\text{EB}}y_i - \psi(\hat{\theta}_i^{\text{EB}})) \right] - \text{cAI} &= -2E \left[ \sum_{i=1}^m n_i\hat{\theta}_i^{\text{EB}}(y_i - \mu_i) \right] \\ &= -2B \quad (\text{say}). \end{aligned} \quad (5.6)$$

Then we propose the cAIC for the nonlinear mixed model as bias corrected estimator of the cAI, which is given by

$$\text{cAIC} = -2 \sum_{i=1}^m n_i(\hat{\theta}_i^{\text{EB}}y_i - \psi(\hat{\theta}_i^{\text{EB}})) + 2\hat{B},$$

where  $\widehat{B}$  is an asymptotically unbiased estimator of  $B$ , which we call the bias correction term or the penalty term. In the next section, we give an asymptotic approximation and an estimator of the penalty term  $B$ .

## 5.3 Approximation and estimation of penalty term

### 5.3.1 Decomposition of penalty term

Because it is difficult to evaluate the penalty term  $B$  exactly, we give second-order approximation of  $B$  for large  $m$ . Taylor series expansion of  $\hat{\theta}_i^{\text{EB}} = h\{\widehat{\mu}_i(y_i, \widehat{\boldsymbol{\eta}})\}$  around  $\widehat{\boldsymbol{\eta}} = \boldsymbol{\eta}$  gives

$$h(\widehat{\mu}_i^{\text{EB}}) = h(\widehat{\mu}_i^{\text{B}}) + \frac{\partial h(\widehat{\mu}_i^{\text{B}})}{\partial \boldsymbol{\eta}^{\text{T}}}(\widehat{\boldsymbol{\eta}} - \boldsymbol{\eta}) + \frac{1}{2}(\widehat{\boldsymbol{\eta}} - \boldsymbol{\eta})^{\text{T}} \frac{\partial^2 h(\widehat{\mu}_i^{\text{B}})}{\partial \boldsymbol{\eta} \partial \boldsymbol{\eta}^{\text{T}}}(\widehat{\boldsymbol{\eta}} - \boldsymbol{\eta}) + o_p(m^{-1}).$$

Then the penalty term  $B$  in (5.6) can be expanded as

$$\begin{aligned} B &= E \left[ \sum_{i=1}^m n_i h(\widehat{\mu}_i^{\text{B}})(y_i - \mu_i) \right] + E \left[ \sum_{i=1}^m n_i (y_i - \mu_i) \frac{\partial h(\widehat{\mu}_i^{\text{B}})}{\partial \boldsymbol{\eta}^{\text{T}}}(\widehat{\boldsymbol{\eta}} - \boldsymbol{\eta}) \right], \\ &\quad + \frac{1}{2} E \left[ \sum_{i=1}^m n_i (y_i - \mu_i) (\widehat{\boldsymbol{\eta}} - \boldsymbol{\eta})^{\text{T}} \frac{\partial^2 h(\widehat{\mu}_i^{\text{B}})}{\partial \boldsymbol{\eta} \partial \boldsymbol{\eta}^{\text{T}}}(\widehat{\boldsymbol{\eta}} - \boldsymbol{\eta}) \right] + o(1) \\ &= B_1(\boldsymbol{\eta}) + B_2(\boldsymbol{\eta}) + 2^{-1} B_3(\boldsymbol{\eta}) + o(1), \quad (\text{say}) \end{aligned} \quad (5.7)$$

where  $B_1 = O(m)$ ,  $B_2 = O(1)$  and  $B_3 = O(1)$ . We evaluate each term of (5.7) in the following subsections by three different methods. To this end, we consider the following conditions:

(C1) The number of clusters  $m$  goes to infinity, and for each  $i = 1, \dots, m$ ,  $n_i$  is sufficiently large, which is of order  $n_i = O(m^{1/2+\delta_i})$  for some positive  $\delta_i > 0$ .

(C2) The number of clusters  $m$  goes to infinity, and for each  $i = 1, \dots, m$ ,  $n_i = O(1)$ .

In the next subsection, we evaluate and estimate  $B_1$ ,  $B_2$  and  $B_3$  by analytical method under the condition (C1). In Section 5.3.3 and 5.3.4, we consider the case of constant  $n_i$ , namely the condition (C2). In Section 5.3.3, we propose a method using numerical integration and differentiation, and in Section 5.3.4 we give a numerical method based on parametric bootstrap.

### 5.3.2 Analytical method for the case of large $n_i$

Firstly, we evaluate  $B_1$ , which is given by

$$B_1(\boldsymbol{\eta}) = E \left[ \sum_{i=1}^m n_i h(\widehat{\mu}_i^{\text{B}})(y_i - \mu_i) \right]. \quad (5.8)$$

Though the Bayes estimator  $\widehat{\mu}_i^{\text{B}} = (n_i y_i + \lambda \mu_i) / (n_i + \lambda)$  is written as the linear function of  $y_i$ , the function  $h(\cdot)$ , which is the link between the mean parameter and the natural parameter, is nonlinear for most of the members of the natural exponential family except for the normal distribution. Thus it is difficult to evaluate  $B_1$  exactly. We here use closeness of  $\widehat{\mu}_i^{\text{B}}$  and  $\mu_i$  for large  $n_i$ , and expand  $h(\widehat{\mu}_i^{\text{B}}) = h(\mu_i + \widehat{\mu}_i^{\text{B}} - \mu_i)$  in (5.8). Then  $h(\widehat{\mu}_i^{\text{B}})$  is approximated by polynomials of the linear function of random variables  $y_i$  and  $\mu_i$ , and  $B_1$  can be asymptotically approximated.

We give two examples of approximation of  $B_1$ , one is the Poisson-gamma mixture model, where  $h(\cdot)$  is the logarithmic function, the other is the binomial-beta mixture model, where  $h(\cdot)$  is the logit function. For the Poisson-gamma mixture model,  $h(\hat{\mu}_i^B) = \log(\hat{\mu}_i^B)$  is expanded as

$$\begin{aligned} h(\hat{\mu}_i^B) &= \log \left\{ \mu_i \left( 1 + \frac{\hat{\mu}_i^B - \mu_i}{\mu_i} \right) \right\} \\ &= \theta_i + \frac{\hat{\mu}_i^B - \mu_i}{\mu_i} - \frac{1}{2} \left( \frac{\hat{\mu}_i^B - \mu_i}{\mu_i} \right)^2 + \frac{1}{3} \left( \frac{\hat{\mu}_i^B - \mu_i}{\mu_i} \right)^3 + \sum_{k=4}^{\infty} \frac{(-1)^{k+1}}{k} \left( \frac{\hat{\mu}_i^B - \mu_i}{\mu_i} \right)^k. \end{aligned}$$

Because  $E \left[ \left\{ \frac{\hat{\mu}_i^B - \mu_i}{\mu_i} \right\}^r (y_i - \mu_i) \right] = O(n_i^{-3})$  for  $r \geq 4$ , if  $n_i = O(m^{1/2+\delta_i})$  for some positive  $\delta_i > 0$ , it follows that

$$\begin{aligned} B_1(\boldsymbol{\eta}) &= \sum_{i=1}^m n_i E[\theta_i (y_i - \mu_i)] + \sum_{i=1}^m n_i E \left[ \frac{\hat{\mu}_i^B - \mu_i}{\mu_i} (y_i - \mu_i) \right] - \frac{1}{2} \sum_{i=1}^m n_i E \left[ \left( \frac{\hat{\mu}_i^B - \mu_i}{\mu_i} \right)^2 (y_i - \mu_i) \right] \\ &\quad + \frac{1}{3} \sum_{i=1}^m n_i E \left[ \left( \frac{\hat{\mu}_i^B - \mu_i}{\mu_i} \right)^3 (y_i - \mu_i) \right] + o(1). \end{aligned}$$

We can evaluate each term in the equation above exactly, noting that

$$\hat{\mu}_i^B - \mu_i = \frac{n_i}{n_i + \lambda} (y_i - \mu_i) - \frac{\lambda}{n_i + \lambda} (\mu_i - m_i).$$

After some calculations, we can obtain the following lemma.

**Lemma 5.1** *Under the condition (C1), for the Poisson-gamma mixture model,  $B_1$  in (5.8) is approximated up to second-order as*

$$B_1(\boldsymbol{\eta}) = B_{11}(\boldsymbol{\eta}) + o(1),$$

where

$$B_{11}(\boldsymbol{\eta}) = m - \sum_{i=1}^m \frac{2\lambda^2 m_i - \lambda}{2n_i(\lambda m_i - 1)}. \quad (5.9)$$

For the binomial-beta mixture model, where  $h(\cdot)$  is the logit function,  $h(\hat{\mu}_i^B) = \log\{\hat{\mu}_i^B/(1 - \hat{\mu}_i^B)\}$  is expanded as

$$\begin{aligned} h(\hat{\mu}_i^B) &= \log \left\{ \mu_i \left( 1 + \frac{\hat{\mu}_i^B - \mu_i}{\mu_i} \right) \right\} - \log \left\{ (1 - \mu_i) \left( 1 - \frac{\hat{\mu}_i^B - \mu_i}{1 - \mu_i} \right) \right\} \\ &= \theta_i + \log \left( 1 + \frac{\hat{\mu}_i^B - \mu_i}{\mu_i} \right) - \log \left( 1 - \frac{\hat{\mu}_i^B - \mu_i}{1 - \mu_i} \right) \\ &= \theta_i + (C_i + D_i) - \frac{1}{2}(C_i^2 - D_i^2) + \frac{1}{3}(C_i^3 + D_i^3) + \sum_{k=4}^{\infty} \frac{1}{k} \left\{ (-1)^{k+1} C_i^k + D_i^k \right\}, \end{aligned}$$

where  $C_i = (\hat{\mu}_i^B - \mu_i)/\mu_i$ ,  $D_i = (\hat{\mu}_i^B - \mu_i)/(1 - \mu_i)$ . If  $n_i = O(m^{1/2+\delta_i})$ ,  $\delta_i > 0$ , it follows that

$$\begin{aligned} B_1 &= \sum_{i=1}^m n_i E[\theta_i (y_i - \mu_i)] + \sum_{i=1}^m n_i E[(C_i + D_i)(y_i - \mu_i)] - \frac{1}{2} \sum_{i=1}^m n_i E[(C_i^2 - D_i^2)(y_i - \mu_i)] \\ &\quad + \frac{1}{3} \sum_{i=1}^m n_i E[(C_i^3 + D_i^3)(y_i - \mu_i)] + o(1), \end{aligned}$$

where each term in the equation above can be evaluated exactly. After some calculations, we get the following lemma.

**Lemma 5.2** *Under the condition (C1), for the binomial-beta mixture model,  $B_1$  in (5.8) is approximated up to second-order as*

$$B_1(\boldsymbol{\eta}) = B_{11}(\boldsymbol{\eta}) + o(1),$$

where

$$B_{11}(\boldsymbol{\eta}) = m - \sum_{i=1}^m \frac{\lambda + 1}{n_i} + \sum_{i=1}^m \frac{4\lambda^2 m_i (1 - m_i) - (\lambda + 2)(\lambda - 1)}{2n_i(\lambda m_i - 1)\{\lambda(1 - m_i) - 1\}}. \quad (5.10)$$

$B_{11}(\hat{\boldsymbol{\eta}})$ , where  $B_{11}(\boldsymbol{\eta})$  is given by (5.9) or (5.10), is expanded as

$$B_{11}(\hat{\boldsymbol{\eta}}) = B_{11}(\boldsymbol{\eta}) + \frac{\partial B_{11}(\boldsymbol{\eta})}{\partial \boldsymbol{\eta}^T} (\hat{\boldsymbol{\eta}} - \boldsymbol{\eta}) + \frac{1}{2} \text{tr} \left[ \frac{\partial^2 B_{11}(\boldsymbol{\eta})}{\partial \boldsymbol{\eta} \partial \boldsymbol{\eta}^T} (\hat{\boldsymbol{\eta}} - \boldsymbol{\eta})(\hat{\boldsymbol{\eta}} - \boldsymbol{\eta})^T \right] + o_p(1),$$

then we propose a bias corrected estimator of  $B_1$  given by

$$\hat{B}_1 = B_{11}(\hat{\boldsymbol{\eta}}) - B_{12}(\hat{\boldsymbol{\eta}}) - B_{13}(\hat{\boldsymbol{\eta}}), \quad (5.11)$$

where

$$B_{12}(\boldsymbol{\eta}) = \frac{\partial B_{11}(\boldsymbol{\eta})}{\partial \boldsymbol{\eta}^T} E(\hat{\boldsymbol{\eta}} - \boldsymbol{\eta}),$$

$$B_{13}(\boldsymbol{\eta}) = \frac{1}{2} \text{tr} \left[ \frac{\partial^2 B_{11}(\boldsymbol{\eta})}{\partial \boldsymbol{\eta} \partial \boldsymbol{\eta}^T} E[(\hat{\boldsymbol{\eta}} - \boldsymbol{\eta})(\hat{\boldsymbol{\eta}} - \boldsymbol{\eta})^T] \right].$$

When  $\boldsymbol{\eta}$  is estimated by the estimating equation (5.12) explained later,  $E(\hat{\boldsymbol{\eta}} - \boldsymbol{\eta})$  and  $E[(\hat{\boldsymbol{\eta}} - \boldsymbol{\eta})(\hat{\boldsymbol{\eta}} - \boldsymbol{\eta})^T]$  are analytically approximated by (5.34) and (5.31). For the Poisson-gamma mixture model,  $\partial B_{11}(\boldsymbol{\eta})/\partial \boldsymbol{\beta}$  and  $(\partial^2 B_{11}(\boldsymbol{\eta})) / (\partial \boldsymbol{\eta} \partial \boldsymbol{\eta}^T)$  are

$$\frac{\partial B_{11}(\boldsymbol{\eta})}{\partial \boldsymbol{\beta}} = \sum_{i=1}^m \frac{\lambda^2 m_i}{2n_i(\lambda m_i - 1)^2} \mathbf{x}_i, \quad \frac{\partial B_{11}(\boldsymbol{\eta})}{\partial \lambda} = \sum_{i=1}^m \frac{1}{n_i} \left\{ \frac{1}{2(\lambda m_i - 1)^2} - 1 \right\},$$

and

$$\frac{\partial^2 B_{11}(\boldsymbol{\eta})}{\partial \boldsymbol{\beta} \partial \boldsymbol{\beta}^T} = - \sum_{i=1}^m \frac{\lambda^2 m_i (\lambda m_i + 1)}{2n_i (\lambda m_i - 1)^3} \mathbf{x}_i \mathbf{x}_i^T, \quad \frac{\partial^2 B_{11}(\boldsymbol{\eta})}{\partial \boldsymbol{\beta} \partial \lambda} = - \sum_{i=1}^m \frac{\lambda m_i}{n_i (\lambda m_i - 1)^3} \mathbf{x}_i,$$

$$\frac{\partial^2 B_{11}(\boldsymbol{\eta})}{\partial \lambda^2} = - \sum_{i=1}^m \frac{m_i}{n_i (\lambda m_i - 1)^3}.$$

For the binomial-beta mixture model, it follows that

$$\frac{\partial B_{11}(\boldsymbol{\eta})}{\partial \boldsymbol{\beta}} = \sum_{i=1}^m \frac{\lambda^2 (\lambda - 1)(\lambda - 2) m_i (1 - m_i)(1 - 2m_i)}{2n_i (\lambda m_i - 1)^2 \{\lambda(1 - m_i) - 1\}^2} \mathbf{x}_i,$$

$$\frac{\partial B_{11}(\boldsymbol{\eta})}{\partial \lambda} = - \sum_{i=1}^m \frac{1}{n_i} + \sum_{i=1}^m \frac{-\lambda(3\lambda - 4) m_i (1 - m_i) + (\lambda - 1)^2}{2n_i (\lambda m_i - 1)^2 \{\lambda(1 - m_i) - 1\}^2},$$

and

$$\begin{aligned}\frac{\partial^2 B_{11}(\boldsymbol{\eta})}{\partial \boldsymbol{\beta} \partial \boldsymbol{\beta}^T} &= \lambda^2 (\lambda - 1) (\lambda - 2) \sum_{i=1}^m \left\{ \frac{\{1 - 6m_i(1 - m_i)\}}{2n_i(\lambda m_i - 1)^2 \{\lambda(1 - m_i) - 1\}^2} \right. \\ &\quad \left. - \frac{\lambda^2 m_i(1 - m_i)(1 - 2m_i)^2}{n_i(\lambda m_i - 1)^3 \{\lambda(1 - m_i) - 1\}^3} \right\} m_i(1 - m_i) \mathbf{x}_i \mathbf{x}_i^T, \\ \frac{\partial^2 B_{11}(\boldsymbol{\eta})}{\partial \boldsymbol{\beta} \partial \lambda} &= - \sum_{i=1}^m \frac{\lambda(3\lambda - 4)(1 - 2m_i)}{2n_i(\lambda m_i - 1)^2 \{\lambda(1 - m_i) - 1\}^2} \\ &\quad + \sum_{i=1}^m \frac{\lambda^2(1 - 2m_i)\{\lambda(3\lambda - 4)m_i(1 - m_i) - (\lambda - 1)^2\}}{n_i(\lambda m_i - 1)^3 \{\lambda(1 - m_i) - 1\}^3} \\ \frac{\partial^2 B_{11}(\boldsymbol{\eta})}{\partial \lambda^2} &= \frac{\{2(\lambda - 1) - m_i(1 - m_i)(6\lambda - 4)\}}{2n_i(\lambda m_i - 1)^2 \{\lambda(1 - m_i) - 1\}^2} \\ &\quad - \frac{\{2m_i(1 - m_i)\lambda - 1\}\{(\lambda - 1)^2 - \lambda(3\lambda - 4)m_i(1 - m_i)\}}{n_i(\lambda m_i - 1)^3 \{\lambda(1 - m_i) - 1\}^3}.\end{aligned}$$

Next, we give an asymptotic approximation of  $B_2$  and  $B_3$  analytically when the hyperparameter  $\boldsymbol{\eta}$  is estimated by the estimating equation suggested by Godambe and Thompson (1989). We define  $\mathbf{g}_i = (g_{1i}, g_{2i})^T$  for  $g_{1i} = y_i - m_i$  and  $g_{2i} = (y_i - m_i)^2 - \phi_i Q(m_i)$  and

$$\begin{aligned}\mathbf{D}_i^T &= E \left( - \frac{\partial \mathbf{g}_i^T}{\partial \boldsymbol{\eta}} \right) \\ &= Q(m_i) \begin{bmatrix} \mathbf{x}_i & Q'(m_i) \phi_i \mathbf{x}_i \\ 0 & -(1 + v_2/n_i)(\lambda - v_2)^{-2} \end{bmatrix}, \\ \boldsymbol{\Sigma}_i &= \mathbf{Cov}(\mathbf{g}_i) = \begin{bmatrix} \mu_{2i} & \mu_{3i} \\ \mu_{3i} & \mu_{4i} - \mu_{2i}^2 \end{bmatrix},\end{aligned}$$

where  $\mu_{ri} = E[(y_i - m_i)^r]$  and  $\phi_i = (\lambda/n_i + 1)/(\lambda - v_2)$ . Following Ghosh and Maiti (2004), the exact expressions of  $\mu_{ri}$ 's for  $r = 2, 3, 4$ , are

$$\begin{aligned}\mu_{2i} &= \phi_i Q(m_i), \quad \mu_{3i} = \frac{Q(m_i)Q'(m_i)(\lambda/n_i + 1)(\lambda/n_i + 2)}{(\lambda - v_2)(\lambda - 2v_2)}, \\ \mu_{4i} &= (d_i + 1)(2d_i + 1)(3d_i + 1)E[(\mu_i - m_i)^4] + 6n_i^{-1}Q'(m_i)(d_i + 1)(2d_i + 1)E[(\mu_i - m_i)^3] \\ &\quad + n_i^{-2}(d_i + 1)[7\{Q'(m_i)\}^2 + 2n_i(4d_i + 3)Q(m_i)]E[(\mu_i - m_i)^2] \\ &\quad + n_i^{-3}Q(m_i)[n_i(2d_i + 3)Q(m_i) + \{Q'(m_i)\}^2],\end{aligned}$$

where  $d_i = v_2/n_i$  and  $E[(\mu_i - m_i)^2] = Q(m_i)/(\lambda - v_2)$ ,  $E[(\mu_i - m_i)^3] = 2Q(m_i)Q'(m_i)/\{(\lambda - v_2)(\lambda - 2v_2)\}$  and

$$E[(\mu_i - m_i)^4] = \frac{3Q(m_i)[(\lambda - 2v_2)Q(m_i) + 2\{Q'(m_i)\}^2]}{(\lambda - v_2)(\lambda - 2v_2)(\lambda - 3v_2)}.$$

Then, Ghosh and Maiti (2004) derived the estimating equation:

$$\mathbf{s}(\boldsymbol{\eta}) = \mathbf{0} \quad \text{for} \quad \mathbf{s}(\boldsymbol{\eta}) = \sum_{i=1}^m \mathbf{D}_i^T \boldsymbol{\Sigma}_i^{-1} \mathbf{g}_i. \quad (5.12)$$

This method is also used in Ghosh and Maiti (2008) and Kubokawa et al. (2014). The estimating equation by Godambe and Thompson (1989) is an extension of quasi-likelihood methods



proposed by Wedderburn (1974) and  $\mathbf{s}(\boldsymbol{\eta})$  is the ‘extended quasi-score function’. In this context, ‘quasi-Fisher information’ is given by

$$\begin{aligned} E(\mathbf{s}\mathbf{s}^\top) &= \sum_{i=1}^m \mathbf{D}_i^\top \boldsymbol{\Sigma}_i^{-1} E(\mathbf{g}_i \mathbf{g}_i^\top) \boldsymbol{\Sigma}_i^{-1} \mathbf{D}_i = \sum_{i=1}^m \mathbf{D}_i^\top \boldsymbol{\Sigma}_i^{-1} \mathbf{D}_i \\ &= \mathbf{U} \quad (\text{say}), \end{aligned}$$

and the asymptotic variance of  $\widehat{\boldsymbol{\eta}}$ , which is the solution of  $\mathbf{s}(\boldsymbol{\eta}) = \mathbf{0}$ , is  $E(\widehat{\boldsymbol{\eta}}\widehat{\boldsymbol{\eta}}^\top) = \mathbf{U}^{-1} + o(m^{-1})$ . Derivation of the asymptotic bias and variance of  $\widehat{\boldsymbol{\eta}}$  are shown in Section 5.5.1 which gives a stochastic expansion of  $\widehat{\boldsymbol{\eta}}$ .

Following Ghosh and Maiti (2004), we define

$$\begin{aligned} \mathbf{J}_r &= \text{Cov}\left(\mathbf{s}, \frac{\partial s_r}{\partial \boldsymbol{\eta}}\right), \quad \mathbf{K}_r = E\left(\frac{\partial^2 s_r}{\partial \boldsymbol{\eta} \partial \boldsymbol{\eta}^\top}\right), \quad r = 1, \dots, p+1 \\ \mathbf{a}^\top &= [\text{tr}(\mathbf{U}^{-1} \mathbf{J}_1), \dots, \text{tr}(\mathbf{U}^{-1} \mathbf{J}_{p+1})], \quad \mathbf{b}^\top = [\text{tr}(\mathbf{U}^{-1} \mathbf{K}_1), \dots, \text{tr}(\mathbf{U}^{-1} \mathbf{K}_{p+1})], \end{aligned}$$

where  $\mathbf{s} = (s_1, \dots, s_{p+1})^\top$  and the expressions of  $\mathbf{J}_r$  and  $\mathbf{K}_r$  are given by (5.37) and (5.38), respectively. Then the asymptotic bias of  $\widehat{\boldsymbol{\eta}}$  is  $E(\widehat{\boldsymbol{\eta}} - \boldsymbol{\eta}) = \mathbf{U}^{-1}(\mathbf{a} + \mathbf{b}/2) + o(m^{-1})$ . Furthermore, we decompose  $\mathbf{U}^{-1}$  as

$$\mathbf{U}^{-1} = \begin{bmatrix} \mathbf{U}_1 \\ \mathbf{U}_2 \end{bmatrix}, \quad \mathbf{U}^{-1} = \begin{bmatrix} \mathbf{U}_{11} & \mathbf{U}_{12} \\ \mathbf{U}_{12}^\top & U_{22} \end{bmatrix},$$

where  $\mathbf{U}_1$  and  $\mathbf{U}_2$  are  $(p, p+1)$  and  $(1, p+1)$  matrices,  $\mathbf{U}_{11}$ ,  $\mathbf{U}_{12}$  are  $(p, p)$  and  $(p, 1)$  matrices, respectively, and  $U_{22}$  is a scalar. Now we can evaluate  $B_2$  as the following lemma.

**Lemma 5.3** *Under the condition (C1),  $B_2$  in (5.7) is approximated as*

$$B_2(\boldsymbol{\eta}) = B_{21}(\boldsymbol{\eta}) + B_{22}(\boldsymbol{\eta}) + o(1),$$

when  $\boldsymbol{\eta}$  is estimated by the estimating equation (5.12).  $B_{21}$  and  $B_{22}$  are

$$\begin{aligned} B_{21}(\boldsymbol{\eta}) &= \sum_{i=1}^m \frac{n_i \lambda^2}{(n_i + \lambda)^2} Q(m_i) \mathbf{x}_i^\top \mathbf{U}_1 \mathbf{D}_i^\top \boldsymbol{\Sigma}_i^{-1} \begin{bmatrix} \xi_{2i} \\ \xi_{3i} - \phi_i Q(m_i) \xi_{1i} \end{bmatrix} \\ &\quad - \sum_{i=1}^m \frac{n_i^2 \lambda}{(n_i + \lambda)^3} \mathbf{U}_2 \mathbf{D}_i^\top \boldsymbol{\Sigma}_i^{-1} \begin{bmatrix} \xi_{3i} \\ \xi_{4i} - \phi_i Q(m_i) \xi_{2i} \end{bmatrix}, \\ B_{22}(\boldsymbol{\eta}) &= \sum_{i=1}^m \frac{n_i \lambda^2}{(n_i + \lambda)^2} Q(m_i) \xi_{1i} \mathbf{x}_i^\top \mathbf{U}_1 (\mathbf{a} + 2^{-1} \mathbf{b}) - \sum_{i=1}^m \frac{n_i^2 \lambda}{(n_i + \lambda)^3} \xi_{2i} \mathbf{U}_2 (\mathbf{a} + 2^{-1} \mathbf{b}), \end{aligned}$$

where  $\xi_{ri} = E[h'(\widehat{\mu}_i^B)(y_i - m_i)^r]$  for  $r = 1, \dots, 4$  and their approximations in the case of  $n_i = O(m^{1/2+\delta_i})$ ,  $\delta_i > 0$  are shown in Section 5.5.5.

By the lemma above, we can estimate  $B_2$  by

$$\widehat{B}_2 = B_{21}(\widehat{\boldsymbol{\eta}}) + B_{22}(\widehat{\boldsymbol{\eta}}), \quad (5.13)$$

which is second-order unbiased.

Next,  $B_3$  can be evaluated as follows.

**Lemma 5.4** Under the condition (C1),  $B_3$  in (5.7) is approximated as

$$B_3(\boldsymbol{\eta}) = B_{31}(\boldsymbol{\eta}) + 2B_{32}(\boldsymbol{\eta}) + o(1),$$

when  $\boldsymbol{\eta}$  is estimated by the estimating equation (5.12).  $B_{31}$  and  $B_{32}$  are

$$B_{31}(\boldsymbol{\eta}) = \sum_{i=1}^m \frac{n_i \lambda^2}{(n_i + \lambda)^2} Q(m_i) Q'(m_i) \xi_{1i} \mathbf{x}_i^T \mathbf{U}_{11} \mathbf{x}_i,$$

$$B_{32}(\boldsymbol{\eta}) = \sum_{i=1}^m \frac{n_i^2 \lambda}{(n_i + \lambda)^3} Q(m_i) \xi_{1i} \mathbf{x}_i^T \mathbf{U}_{12},$$

where  $\xi_{1i} = E[h'(\widehat{\mu}_i^B)(y_i - m_i)]$  and its approximation in the case of  $n_i = O(m^{1/2+\delta_i})$ ,  $\delta_i > 0$  are shown in Section 5.5.5.

By the lemma above, we can estimate  $B_3$  by

$$\widehat{B}_3 = B_{31}(\widehat{\boldsymbol{\eta}}) + 2B_{32}(\widehat{\boldsymbol{\eta}}), \quad (5.14)$$

which is second-order unbiased.

Using  $\widehat{B}_1$ ,  $\widehat{B}_2$  and  $\widehat{B}_3$  given by (5.11), (5.13) and (5.14), we propose the cAIC in nonlinear mixed model as follows:

$$\text{cAIC} = -2 \sum_{i=1}^m n_i (\widehat{\theta}_i^{\text{EB}} y_i - \psi(\widehat{\theta}_i^{\text{EB}})) + 2(\widehat{B}_1 + \widehat{B}_2 + 2^{-1} \widehat{B}_3). \quad (5.15)$$

**Theorem 5.1** Under the condition (C1), the cAIC in (5.15) is a second-order unbiased estimator of cAI in (5.5), namely

$$E(\text{cAIC}) = \text{cAI} + o(1),$$

when  $\boldsymbol{\eta}$  is estimated by the estimating equation (5.12).

### 5.3.3 Method for constant $n_i$ by using numerical integration and differentiation

In the last subsection, we consider the condition of large  $n_i$  to derive the criterion. However, this condition is not appropriate for many cases in real data. The most typical example is small area estimation, where the number of observations in each area is not very large. Then, in this and the next subsection, we propose alternative methods of estimating the penalty term, which does not need the assumption of large  $n_i$ .

The first method is based on stochastic expansion of  $\widehat{\boldsymbol{\eta}}$ , which is the same as the method of the previous subsection. However we cannot approximate  $h(\widehat{\mu}_i^B)$  by polynomials of  $\widehat{\mu}_i^B$  for the case of  $n_i = O(1)$ . Then, we alternatively use numerical integration to evaluate the expectation of complex functions of random variables. We firstly estimate  $B_1$ . By the law of iterated expectations and the fact that  $y_i - \widehat{\mu}_i^B = \lambda(y_i - m_i)/(n_i + \lambda)$ ,  $B_1$  can be rewritten as

$$B_1(\boldsymbol{\eta}) = \sum_{i=1}^m \frac{n_i \lambda}{n_i + \lambda} E[h(\widehat{\mu}_i^B)(y_i - m_i)].$$

Because it is hard to obtain closed form expression of  $B_1(\boldsymbol{\eta})$ , we propose to calculate  $B_1(\widehat{\boldsymbol{\eta}})$  by Monte Carlo simulation, or which can be seen as parametric bootstrap method. Bootstrap sample  $\mathbf{y}^* = (y_1^*, \dots, y_m^*)^T$  is generated by

$$y_i^* \sim f(y_i^* | \theta_i^*), \quad \text{and} \quad \theta_i^* \sim p(\theta_i^* | \widehat{\lambda}, \widehat{m}_i). \quad (5.16)$$

Then  $B_1(\hat{\boldsymbol{\eta}})$  can be written as

$$B_1(\hat{\boldsymbol{\eta}}) = \sum_{i=1}^m \frac{n_i \hat{\lambda}}{n_i + \hat{\lambda}} E_* [h(\hat{\mu}_i^{B*})(y_i^* - \hat{m}_i) \mid \mathbf{y}], \quad (5.17)$$

where  $\hat{\mu}_i^{B*} = (n_i y_i^* + \hat{\lambda} \hat{m}_i) / (n_i + \hat{\lambda})$  and  $E_*$  denotes expectation with respect to the distribution of  $\mathbf{y}^*$  given  $\mathbf{y}$ . Monte Carlo approximation of (5.17) is

$$B_1(\hat{\boldsymbol{\eta}}) \approx \widetilde{B}_1^* = B^{-1} \sum_{b=1}^B \sum_{i=1}^m \frac{n_i \hat{\lambda}}{n_i + \hat{\lambda}} h(\hat{\mu}_i^{B*}(b))(y_i^*(b) - \hat{m}_i), \quad (5.18)$$

for large  $B$ , where  $y_i^*(b)$  is the  $b$ th bootstrap sample based on (5.16) and  $\hat{\mu}_i^{B*}(b)$  is the version of  $\hat{\mu}_i^{B*}$  based on the bootstrap sample.

Because  $B_1(\hat{\boldsymbol{\eta}})$  has second-order bias, we consider bias correction as follows:

$$B_1(\hat{\boldsymbol{\eta}}) - B_{14}(\hat{\boldsymbol{\eta}}) - B_{15}(\hat{\boldsymbol{\eta}}),$$

where

$$B_{14}(\boldsymbol{\eta}) = \frac{\partial B_1(\boldsymbol{\eta})}{\partial \boldsymbol{\eta}^T} \mathbf{U}(\boldsymbol{\eta})^{-1} (\mathbf{a}(\boldsymbol{\eta}) + \mathbf{b}(\boldsymbol{\eta})/2), \quad \text{and} \quad B_{15}(\boldsymbol{\eta}) = \frac{1}{2} \text{tr} \left[ \frac{\partial^2 B_1(\boldsymbol{\eta})}{\partial \boldsymbol{\eta} \partial \boldsymbol{\eta}^T} \mathbf{U}(\boldsymbol{\eta})^{-1} \right],$$

noting that  $E(\hat{\boldsymbol{\eta}} - \boldsymbol{\eta}) = \mathbf{U}^{-1}(\mathbf{a} + \mathbf{b}/2) + o(m^{-1})$  and  $E[(\hat{\boldsymbol{\eta}} - \boldsymbol{\eta})(\hat{\boldsymbol{\eta}} - \boldsymbol{\eta})^T] = \mathbf{U}^{-1} + o(m^{-1})$  when  $\hat{\boldsymbol{\eta}}$  is estimated by the estimating equation (5.12). The values of the first- and second-derivative of  $B_1(\boldsymbol{\eta})$  evaluated at  $\boldsymbol{\eta} = \hat{\boldsymbol{\eta}}$  is calculated by numerical differentiation method. We propose the following procedure:

$$\left. \frac{\partial B_1(\boldsymbol{\eta})}{\partial \eta_k} \right|_{\boldsymbol{\eta}=\hat{\boldsymbol{\eta}}} \approx \frac{\widetilde{B}_1^*(\hat{\boldsymbol{\eta}} + \varepsilon \mathbf{e}_k) - \widetilde{B}_1^*(\hat{\boldsymbol{\eta}} - \varepsilon \mathbf{e}_k)}{2\varepsilon}, \quad (k = 1, \dots, p+1),$$

and

$$\left. \frac{\partial^2 B_1(\boldsymbol{\eta})}{\partial \eta_k \partial \eta_l} \right|_{\boldsymbol{\eta}=\hat{\boldsymbol{\eta}}} \approx \begin{cases} \widetilde{B}_1^{*(kk)} = \{\widetilde{B}_1^*(\hat{\boldsymbol{\eta}} + \varepsilon \mathbf{e}_k) + \widetilde{B}_1^*(\hat{\boldsymbol{\eta}} - \varepsilon \mathbf{e}_k) - 2\widetilde{B}_1^*(\hat{\boldsymbol{\eta}})\} / \varepsilon^2 & (\text{if } k = l), \\ \widetilde{B}_1^{*(kl)} = [\{\widetilde{B}_1^*(\hat{\boldsymbol{\eta}} + \varepsilon(\mathbf{e}_k + \mathbf{e}_l)) + \widetilde{B}_1^*(\hat{\boldsymbol{\eta}} - \varepsilon(\mathbf{e}_k + \mathbf{e}_l)) - 2\widetilde{B}_1^*(\hat{\boldsymbol{\eta}})\} \\ \quad - \varepsilon^2 \{\widetilde{B}_1^{*(kk)} + \widetilde{B}_1^{*(ll)}\}] / (2\varepsilon^2) & (\text{if } k \neq l), \end{cases}$$

for small positive  $\varepsilon$ , where  $\mathbf{e}_k$  is the  $(p+1) \times 1$  vector, with the  $k$ th component equal to 1 and the other components equal to 0. Using these numerical integration and differentiation, we can approximate  $B_{14}(\hat{\boldsymbol{\eta}})$  and  $B_{15}(\hat{\boldsymbol{\eta}})$ , which we call  $\widehat{B}_{14}^*$  and  $\widehat{B}_{15}^*$ . Then we can obtain an estimator of  $B_1$  as follows:

$$\widehat{B}_1^* = \widetilde{B}_1^* - \widehat{B}_{14}^* - \widehat{B}_{15}^*. \quad (5.19)$$

Next we estimate  $B_2$  and  $B_3$ . By Lemma 5.3, we can approximate  $B_2$  as

$$B_2(\boldsymbol{\eta}) = B_{21}(\boldsymbol{\eta}) + B_{22}(\boldsymbol{\eta}) + o(1),$$

which is also valid under the condition (C2). In the previous subsection,  $\xi_{ri} = E[h'(\hat{\mu}_i^B)(y_i - m_i)^r]$  ( $r = 1, \dots, 4$ ) in  $B_{21}$  and  $B_{22}$  are approximated by Taylor series expansion of  $h(\cdot)$  for large  $n_i$ .

However, we now consider the case of  $n_i = O(1)$ . Then we calculate  $\xi_{ri}(\hat{\boldsymbol{\eta}})$  by Monte Carlo approximation based on bootstrap samples (5.16) as follows:

$$\xi_{ri}(\hat{\boldsymbol{\eta}}) \approx B^{-1} \sum_{b=1}^B \sum_{i=1}^m h'(\hat{\mu}_i^{B*}(b))(y_i^*(b) - \hat{m}_i)^r.$$

Calculating  $\xi_{ri}(\hat{\boldsymbol{\eta}})$ 's, we can obtain  $B_{21}(\hat{\boldsymbol{\eta}})$  and  $B_{22}(\hat{\boldsymbol{\eta}})$  based on Monte Carlo approximation, which we call  $\widehat{B}_{21}^*$  and  $\widehat{B}_{22}^*$ , and the following estimator of  $B_2$ :

$$\widehat{B}_2^* = \widehat{B}_{21}^* + \widehat{B}_{22}^*. \quad (5.20)$$

As for  $B_3$ , Lemma 5.4 gave an asymptotic approximation. However, this approximation is based on the condition (C1), we have to modify the evaluation. Then we give the following lemma.

**Lemma 5.5** *Under the condition (C2),  $B_3$  in (5.7) is approximated as*

$$B_3(\boldsymbol{\eta}) = B_{33}(\boldsymbol{\eta}) + 2B_{34}(\boldsymbol{\eta}) + B_{35}(\boldsymbol{\eta}) + o(1),$$

when  $\boldsymbol{\eta}$  is estimated by the estimating equation (5.12).  $B_{33}(\boldsymbol{\eta})$ ,  $B_{34}(\boldsymbol{\eta})$  and  $B_{35}(\boldsymbol{\eta})$  are

$$\begin{aligned} B_{33}(\boldsymbol{\eta}) &= B_{31}(\boldsymbol{\eta}) + \sum_{i=1}^m \frac{n_i \lambda^3}{(n_i + \lambda)^3} \{Q(m_i)\}^2 \nu_{1i} \mathbf{x}_i^T \mathbf{U}_{11} \mathbf{x}_i, \\ B_{34}(\boldsymbol{\eta}) &= B_{32}(\boldsymbol{\eta}) - \sum_{i=1}^m \frac{n_i^2 \lambda^2}{(n_i + \lambda)^4} Q(m_i) \nu_{2i} \mathbf{x}_i^T \mathbf{U}_{12}, \\ B_{35}(\boldsymbol{\eta}) &= \sum_{i=1}^m \left\{ \frac{2n_i^2 \lambda}{(n_i + \lambda)^4} \xi_{2i} + \frac{n_i^3 \lambda}{(n_i + \lambda)^5} \nu_{3i} \right\} \times U_{22}, \end{aligned}$$

where  $\nu_{ri} = E[h''(\hat{\mu}_i^B)(y_i - m_i)^r]$  for  $r = 1, \dots, 3$ .

We calculate  $\nu_{ri}(\hat{\boldsymbol{\eta}})$  in the lemma above by Monte Carlo approximation based on bootstrap samples (5.16) as follows:

$$\nu_{ri}(\hat{\boldsymbol{\eta}}) \approx B^{-1} \sum_{b=1}^B \sum_{i=1}^m h''(\hat{\mu}_i^{B*}(b))(y_i^*(b) - \hat{m}_i)^r$$

Calculating  $\nu_{ri}(\hat{\boldsymbol{\eta}})$ 's, we can obtain  $B_{33}(\hat{\boldsymbol{\eta}})$ ,  $B_{34}(\hat{\boldsymbol{\eta}})$  and  $B_{35}(\hat{\boldsymbol{\eta}})$ , which we call  $\widehat{B}_{33}^*$ ,  $\widehat{B}_{34}^*$  and  $\widehat{B}_{35}^*$ , and the following estimator of  $B_3$ :

$$\widehat{B}_3^* = \widehat{B}_{33}^* + 2\widehat{B}_{34}^* + \widehat{B}_{35}^*. \quad (5.21)$$

Using  $\widehat{B}_1^*$ ,  $\widehat{B}_2^*$  and  $\widehat{B}_3^*$  given by (5.19), (5.20) and (5.21), we propose the following cAIC\* using numerical integration and differentiation:

$$\text{cAIC}^* = -2 \sum_{i=1}^m n_i (\hat{\theta}_i^{\text{EB}} y_i - \psi(\hat{\theta}_i^{\text{EB}})) + 2(\widehat{B}_1^* + \widehat{B}_2^* + 2^{-1} \widehat{B}_3^*). \quad (5.22)$$

**Theorem 5.2** *Under the condition (C2), the cAIC\* in (5.22) is second-order asymptotically unbiased estimator of cAI in (5.5), namely*

$$E(\text{cAIC}^*) = \text{cAI} + o(1),$$

when  $\boldsymbol{\eta}$  is estimated by the estimating equation (5.12).

### 5.3.4 Numerical method for constant $n_i$ based on parametric bootstrap

The second method to estimate the penalty term without the assumption of large  $n_i$  is based on parametric bootstrap. This method does not need stochastic expansion of  $\hat{\boldsymbol{\eta}}$ .

We firstly estimate  $B_1$ . The plug-in estimator  $B_1(\hat{\boldsymbol{\eta}})$  can be calculated by Monte Carlo approximation  $\widetilde{B}_1^*$  given by (5.18). However, this naive estimator has second-order bias. When it is hard to obtain analytical form of asymptotic bias and variance of  $\hat{\boldsymbol{\eta}}$ , we cannot correct the bias of  $B_1(\hat{\boldsymbol{\eta}})$  analytically. Then we propose numerical bias correction based on parametric bootstrap method as follows:

$$2B_1(\hat{\boldsymbol{\eta}}) - E_*[B_1(\hat{\boldsymbol{\eta}}^*) | \mathbf{y}].$$

The second term of the equation above can be written as

$$E_*[B_1(\hat{\boldsymbol{\eta}}^*) | \mathbf{y}] = \sum_{i=1}^m E_* \left[ \frac{n_i \hat{\lambda}^*}{n_i + \hat{\lambda}^*} \cdot E_{**} [h(\hat{\mu}_i^{B**})(y_i^{**} - \hat{m}_i^*) | (\mathbf{y}^*, \mathbf{y})] | \mathbf{y} \right], \quad (5.23)$$

where the distribution of  $y_i^{**}$  is

$$y_i^{**} \sim f(y_i^{**} | \theta_i^{**}), \quad \text{and} \quad \theta_i^{**} \sim p(\theta_i^{**} | \hat{\lambda}^*, \hat{m}_i^*), \quad (5.24)$$

$\hat{\mu}_i^{B**} = (n_i y_i^{**} + \hat{\lambda}^* \hat{m}_i^*) / (n_i + \hat{\lambda}^*)$ , and  $E_{**}$  denotes expectation with respect to the distribution of  $\mathbf{y}^{**} = (y_1^{**}, \dots, y_m^{**})^T$  given  $\mathbf{y}^*$  and  $\mathbf{y}$ . Monte Carlo approximation of (5.23) is

$$E_*[B_1(\hat{\boldsymbol{\eta}}^*) | \mathbf{y}] \approx \widetilde{B}_1^{**} = (BC)^{-1} \sum_{b=1}^B \sum_{c=1}^C \sum_{i=1}^m \frac{n_i \hat{\lambda}^*(b)}{n_i + \hat{\lambda}^*(b)} h\{\hat{\mu}_i^{B**}(c(b))\} \{y_i^{**}(c(b)) - \hat{m}_i^*(b)\},$$

for large  $B$  and  $C$ , where  $y_i^{**}(c(b))$  is the  $c$ th double bootstrap sample from the  $b$ th bootstrap sample based on (5.24), and  $\hat{\mu}_i^{B**}(c(b))$  is the version of  $\hat{\mu}_i^{B**}$  based on the bootstrap sample. Then, we propose the following bias corrected estimator of  $B_1$ :

$$\widehat{B}_1^{**} = 2\widetilde{B}_1^{**} - \widetilde{B}_1^*. \quad (5.25)$$

We next estimate  $B_2$  and  $B_3$ . It is noted that  $B_2(\boldsymbol{\eta})$  and  $B_3(\boldsymbol{\eta})$  can be rewritten as

$$B_2(\boldsymbol{\eta}) = \sum_{i=1}^m \frac{n_i \lambda}{n_i + \lambda} E \left[ (y_i - m_i) \frac{\partial h(\hat{\mu}_i^B)}{\partial \boldsymbol{\eta}^T} (\hat{\boldsymbol{\eta}} - \boldsymbol{\eta}) \right],$$

$$B_3(\boldsymbol{\eta}) = \sum_{i=1}^m \frac{n_i \lambda}{n_i + \lambda} E \left[ (y_i - m_i) (\hat{\boldsymbol{\eta}} - \boldsymbol{\eta})^T \frac{\partial h(\hat{\mu}_i^B)}{\partial \boldsymbol{\eta}^T} (\hat{\boldsymbol{\eta}} - \boldsymbol{\eta}) \right],$$

by the law of iterated expectations and the fact that  $y_i - \hat{\mu}_i^B = \lambda(y_i - m_i) / (n_i + \lambda)$ . Then  $B_2(\hat{\boldsymbol{\eta}})$  and  $B_3(\hat{\boldsymbol{\eta}})$  can be written as

$$B_2(\hat{\boldsymbol{\eta}}) = \sum_{i=1}^m \frac{n_i \hat{\lambda}}{n_i + \hat{\lambda}} E_* \left[ (y_i^* - \hat{m}_i) \frac{\partial h(\hat{\mu}_i^{B*})}{\partial \hat{\boldsymbol{\eta}}^T} (\hat{\boldsymbol{\eta}}^* - \hat{\boldsymbol{\eta}}) | \mathbf{y} \right],$$

$$B_3(\hat{\boldsymbol{\eta}}) = \sum_{i=1}^m \frac{n_i \hat{\lambda}}{n_i + \hat{\lambda}} E_* \left[ (y_i^* - \hat{m}_i) (\hat{\boldsymbol{\eta}}^* - \hat{\boldsymbol{\eta}})^T \frac{\partial^2 h(\hat{\mu}_i^{B*})}{\partial \hat{\boldsymbol{\eta}} \partial \hat{\boldsymbol{\eta}}^T} (\hat{\boldsymbol{\eta}}^* - \hat{\boldsymbol{\eta}}) | \mathbf{y} \right], \quad (5.26)$$

where  $\hat{\boldsymbol{\eta}}^*$  is an estimator of  $\boldsymbol{\eta}$  based on bootstrap sample  $\mathbf{y}^*$ . Note that the exact expressions of  $\partial h(\hat{\mu}_i^{B*}) / \partial \hat{\boldsymbol{\eta}}^T$  and  $\partial^2 h(\hat{\mu}_i^{B*}) / \partial \hat{\boldsymbol{\eta}} \partial \hat{\boldsymbol{\eta}}^T$  are

$$\frac{\partial h(\hat{\mu}_i^{B*})}{\partial \hat{\boldsymbol{\eta}}} = \begin{bmatrix} h'(\hat{\mu}_i^{B*}) \hat{\lambda} (n_i + \hat{\lambda})^{-1} Q(\hat{m}_i) \mathbf{x}_i \\ -h'(\hat{\mu}_i^{B*}) n_i (n_i + \hat{\lambda})^{-2} (y_i^* - \hat{m}_i) \end{bmatrix}$$

and

$$\frac{\partial^2 h(\hat{\mu}_i^{B*})}{\partial \hat{\boldsymbol{\eta}} \partial \hat{\boldsymbol{\eta}}^\top} = \begin{bmatrix} \partial^2 h(\hat{\mu}_i^{B*}) / (\partial \hat{\boldsymbol{\beta}} \partial \hat{\boldsymbol{\beta}}^\top) & \partial^2 h(\hat{\mu}_i^{B*}) / (\partial \hat{\boldsymbol{\beta}} \partial \hat{\lambda}) \\ \partial^2 h(\hat{\mu}_i^{B*}) / (\partial \hat{\lambda} \partial \hat{\boldsymbol{\beta}}^\top) & \partial^2 h(\hat{\mu}_i^{B*}) / (\partial \hat{\lambda})^2 \end{bmatrix}$$

where

$$\begin{aligned} \frac{\partial^2 h(\hat{\mu}_i^{B*})}{\partial \hat{\boldsymbol{\beta}} \partial \hat{\boldsymbol{\beta}}^\top} &= \frac{\hat{\lambda}}{n_i + \hat{\lambda}} Q(\hat{m}_i) \left\{ h''(\hat{\mu}_i^{B*}) \frac{\hat{\lambda}}{n_i + \hat{\lambda}} Q(\hat{m}_i) + h'(\hat{\mu}_i^{B*}) Q'(\hat{m}_i) \right\} \mathbf{x}_i \mathbf{x}_i^\top, \\ \frac{\partial^2 h(\hat{\mu}_i^{B*})}{\partial \hat{\boldsymbol{\beta}} \partial \hat{\lambda}} &= \frac{n_i}{(n_i + \hat{\lambda})^2} Q(\hat{m}_i) \left\{ -h''(\hat{\mu}_i^{B*}) \frac{\hat{\lambda}}{n_i + \hat{\lambda}} (y_i^* - \hat{m}_i) + h'(\hat{\mu}_i^{B*}) \right\} \mathbf{x}_i, \\ \frac{\partial^2 h(\hat{\mu}_i^{B*})}{\partial \hat{\lambda}^2} &= \frac{n_i}{(n_i + \hat{\lambda})^3} (y_i^* - \hat{m}_i) \left\{ h''(\hat{\mu}_i^{B*}) \frac{n_i}{n_i + \hat{\lambda}} (y_i^* - \hat{m}_i) + 2h'(\hat{\mu}_i^{B*}) \right\}. \end{aligned}$$

Monte Carlo approximation of (5.26) is

$$\begin{aligned} B_2(\hat{\boldsymbol{\eta}}) &\approx \widehat{B}_2^{**} = B^{-1} \sum_{b=1}^B \sum_{i=1}^m \frac{n_i \hat{\lambda}}{n_i + \hat{\lambda}} (y_i^*(b) - \hat{m}_i) \frac{\partial h(\hat{\mu}_i^{B*}(b))}{\partial \hat{\boldsymbol{\eta}}^\top} (\hat{\boldsymbol{\eta}}^*(b) - \hat{\boldsymbol{\eta}}), \\ B_3(\hat{\boldsymbol{\eta}}) &\approx \widehat{B}_3^{**} = B^{-1} \sum_{b=1}^B \sum_{i=1}^m \frac{n_i \hat{\lambda}}{n_i + \hat{\lambda}} (y_i^*(b) - \hat{m}_i) (\hat{\boldsymbol{\eta}}^*(b) - \hat{\boldsymbol{\eta}})^\top \frac{\partial^2 h(\hat{\mu}_i^{B*}(b))}{\partial \hat{\boldsymbol{\eta}} \partial \hat{\boldsymbol{\eta}}^\top} (\hat{\boldsymbol{\eta}}^*(b) - \hat{\boldsymbol{\eta}}), \end{aligned} \quad (5.27)$$

for large  $B$ , where  $y_i^*(b)$  is the  $b$ th bootstrap sample based on (5.16) and  $\hat{\mu}_i^{B*}(b)$  and  $\hat{\boldsymbol{\eta}}^*(b)$  are the versions of  $\hat{\mu}_i^{B*}$  and  $\hat{\boldsymbol{\eta}}^*$  based on the bootstrap sample. Because  $B_2(\boldsymbol{\eta})$  and  $B_3(\boldsymbol{\eta})$  are of order  $O(1)$ ,  $\widehat{B}_2^{**}$  and  $\widehat{B}_3^{**}$  are asymptotically unbiased estimators of  $B_2$  and  $B_3$  whose biases are of order  $o(1)$ .

Using  $\widehat{B}_1^{**}$ ,  $\widehat{B}_2^{**}$  and  $\widehat{B}_3^{**}$  given by (5.25) and (5.27), we propose the following cAIC\*\* based on parametric bootstrap:

$$\text{cAIC}^{**} = -2 \sum_{i=1}^m n_i (\hat{\theta}_i^{\text{EB}} y_i - \psi(\hat{\theta}_i^{\text{EB}})) + 2(\widehat{B}_1^{**} + \widehat{B}_2^{**} + 2^{-1} \widehat{B}_3^{**}). \quad (5.28)$$

**Theorem 5.3** *Under the condition (C2), the cAIC\*\* in (5.28) is a second-order asymptotically unbiased estimator of cAI in (5.5), namely*

$$E(\text{cAIC}^{**}) = \text{cAI} + o(1).$$

Though the cAIC\*\* does not depend on the method of estimating  $\boldsymbol{\eta}$ , this requires parametric bootstrap, which is computationally harder than cAIC and cAIC\*. Especially, estimation of  $B_1$  needs double bootstrap method. However, because each step in estimation of  $B_1$  does not include maximization algorithm, namely  $B_1$  does not require estimating the hyperparameter, the computational load of estimating  $B_1$  is the same level as estimating  $B_2$  and  $B_3$ , which only requires single bootstrap.

## 5.4 Simulations

In this section, we compare the numerical performance between the conventional mAIC and the proposed criteria. We handle Poisson-gamma model introduced in Section 5.2. We consider a class of the nested candidate models  $j_\alpha = \{1, \dots, \alpha\}$  for  $\alpha = 1, \dots, p_\omega$ . Let  $\mathbf{x}_i(\omega)$ 's be

independently generated as  $\mathcal{N}(\mathbf{0}, \boldsymbol{\Sigma}_x)$  where  $\boldsymbol{\Sigma}_x = 0.9\mathbf{I}_{p_\omega} + 0.1\mathbf{J}_{p_\omega}$  for  $\mathbf{J}_{p_\omega} = \mathbf{1}_{p_\omega}\mathbf{1}_{p_\omega}^\top$ . The true vector of coefficients  $\boldsymbol{\beta}_*$  is given by  $\boldsymbol{\beta}_* = (1, 1, 1, 0, 0)^\top$ , namely  $p_\omega = 5$  and  $p_* = 3$ . Scale parameters are given by  $\lambda = 4$  and  $n_i = 11$  for  $i = 1, \dots, m$ , where the sample size  $m$  varies with simulations.

We compare the performance between the mAIC

$$\text{mAIC} = -2 \sum_{i=1}^m \log\{m(y_i|\hat{\lambda}, \hat{m}_i)\} + 2(p_j + 1),$$

where  $m(y_i|\lambda, m_i)$  is given in (5.4), and the proposed criteria, based on 500 replications. We handle the cases of  $m = 100$  and  $m = 70$ . Note that the case of  $m = 100$  satisfies the condition (C1) but the case of  $m = 70$  does not.

Table 5.1: The number of selecting each model based on 500 replications and prediction error of the best model selected by the criteria

	the number of selection					prediction error
	$j_1$	$j_2$	$j_3$	$j_4$	$j_5$	
$m = 100$						
mAIC	91	102	211	58	38	7.42962
cAIC	59	98	223	72	48	7.42248
$m = 70$						
mAIC	112	114	190	59	25	5.30810
cAIC	97	163	145	62	33	5.30644

Table 5.1 reports the number of selecting each model and prediction error of the best model selected by the mAIC and the cAIC in (5.15). The prediction error is measured by quadratic loss

$$\sum_{i=1}^m (\hat{\mu}_i^{\text{EB}} - \mu_i)^2,$$

and the values are given as the averages based on 500 replications. For the case of  $m = 100$ , where the condition (C1) is satisfied, the performance of the cAIC is better than that of the mAIC in terms of both selecting the true model and prediction error. For the case of  $m = 70$ , where the condition (C1) is not satisfied, the performance of the cAIC is comparable to that of the mAIC in terms of the prediction error, but is not in terms of selecting the true model. Although we conducted simulations which investigate cAIC\* in (5.22) and cAIC\*\* in (5.28), the performance of the criteria was numerically instable. This is the future work to be resolved.

## 5.5 Some results of analytical calculations

### 5.5.1 Stochastic Expansion of $\hat{\boldsymbol{\eta}}$

We give a stochastic expansion of  $\hat{\boldsymbol{\eta}}$  and derive the asymptotic bias and variance of  $\hat{\boldsymbol{\eta}}$ . Let  $\mathbf{t} = (t_1, \dots, t_{p+1})^\top$  where

$$t_r = (\hat{\boldsymbol{\eta}} - \boldsymbol{\eta})^\top \frac{\partial^2 s_r(\boldsymbol{\eta})}{\partial \boldsymbol{\eta} \partial \boldsymbol{\eta}^\top} (\hat{\boldsymbol{\eta}} - \boldsymbol{\eta}), \quad r = 1, \dots, p+1.$$

It follows from the Taylor series expansion of  $\mathbf{s}(\hat{\boldsymbol{\eta}})$  around  $\hat{\boldsymbol{\eta}} = \boldsymbol{\eta}$  that

$$\mathbf{0} = \mathbf{s}(\boldsymbol{\eta}) + \frac{\partial \mathbf{s}(\boldsymbol{\eta})}{\partial \boldsymbol{\eta}^T} (\hat{\boldsymbol{\eta}} - \boldsymbol{\eta}) + \frac{1}{2} \mathbf{t} + o_p(1),$$

thus we can get the expression

$$\hat{\boldsymbol{\eta}} - \boldsymbol{\eta} = \left( -\frac{\partial \mathbf{s}(\boldsymbol{\eta})}{\partial \boldsymbol{\eta}^T} \right)^{-1} \left\{ \mathbf{s}(\boldsymbol{\eta}) + \frac{1}{2} \mathbf{t} + o_p(1) \right\}. \quad (5.29)$$

It follows that

$$\begin{aligned} -\frac{\partial \mathbf{s}(\boldsymbol{\eta})}{\partial \boldsymbol{\eta}^T} &= \sum_{i=1}^m \left( -\frac{\partial \mathbf{D}_i^T \boldsymbol{\Sigma}_i^{-1}}{\partial \boldsymbol{\eta}^T} \right) (\mathbf{I}_{p+1} \otimes \mathbf{g}_i) + \sum_{i=1}^m \mathbf{D}_i^T \boldsymbol{\Sigma}_i^{-1} \left( -\frac{\partial \mathbf{g}_i}{\partial \boldsymbol{\eta}^T} \right) \\ &(\equiv \mathbf{V}_1 + \mathbf{V}_2), \end{aligned}$$

and  $E(\mathbf{V}_1) = \mathbf{0}$ ,  $\mathbf{V}_1 = O_p(m^{1/2})$ ,  $E(\mathbf{V}_2) = \mathbf{U} = O(m)$ . Let  $\mathbf{W} = \mathbf{V}_1 + \mathbf{V}_2 - \mathbf{U}$ , then noting that  $\mathbf{V}_2 - \mathbf{U} = O_p(m^{1/2})$ ,  $\mathbf{W} = O_p(m^{1/2})$ , we can expand  $\hat{\boldsymbol{\eta}} - \boldsymbol{\eta}$  in (5.29) as

$$\begin{aligned} \hat{\boldsymbol{\eta}} - \boldsymbol{\eta} &= (\mathbf{U} + \mathbf{W})^{-1} \left\{ \mathbf{s} + \frac{1}{2} \mathbf{t} + o_p(1) \right\} \\ &= \left\{ \mathbf{U}^{-1} - \mathbf{U}^{-1} \mathbf{W} \mathbf{U}^{-1} + o_p(m^{-3/2}) \right\} \left\{ \mathbf{s} + \frac{1}{2} \mathbf{t} + o_p(1) \right\} \\ &= \mathbf{U}^{-1} \mathbf{s} - \mathbf{U}^{-1} \mathbf{W} \mathbf{U}^{-1} \mathbf{s} + \frac{1}{2} \mathbf{U}^{-1} \mathbf{t} + o_p(m^{-1}) \\ &\equiv \hat{\boldsymbol{\eta}}^{(1)} + \hat{\boldsymbol{\eta}}^{(2)} + \hat{\boldsymbol{\eta}}^{(3)} + o_p(m^{-1}), \end{aligned} \quad (5.30)$$

where  $\hat{\boldsymbol{\eta}}^{(1)} = O_p(m^{-1/2})$ ,  $\hat{\boldsymbol{\eta}}^{(2)} = O_p(m^{-1})$ ,  $\hat{\boldsymbol{\eta}}^{(3)} = O_p(m^{-1})$ . Now we evaluate each term in (5.30).

First we evaluate the first and the second moment of  $\hat{\boldsymbol{\eta}}^{(1)} = \mathbf{U}^{-1} \sum_{i=1}^m \mathbf{D}_i^T \boldsymbol{\Sigma}_i^{-1} \mathbf{g}_i$ . It is easy to see that  $E(\hat{\boldsymbol{\eta}}^{(1)}) = \mathbf{0}$ . The second moment is

$$\begin{aligned} E[\hat{\boldsymbol{\eta}}^{(1)} (\hat{\boldsymbol{\eta}}^{(1)})^T] &= \sum_{i=1}^m \mathbf{U}^{-1} \mathbf{D}_i^T \boldsymbol{\Sigma}_i^{-1} E[\mathbf{g}_i \mathbf{g}_i^T] \boldsymbol{\Sigma}_i^{-1} \mathbf{D}_i \mathbf{U}^{-1} \\ &= \sum_{i=1}^m \mathbf{U}^{-1} \mathbf{D}_i^T \boldsymbol{\Sigma}_i^{-1} \mathbf{D}_i \mathbf{U}^{-1} \\ &= \mathbf{U}^{-1}, \end{aligned}$$

thus it follows that

$$E[(\hat{\boldsymbol{\eta}} - \boldsymbol{\eta})(\hat{\boldsymbol{\eta}} - \boldsymbol{\eta})^T] = \mathbf{U}^{-1} + o(m^{-1}), \quad (5.31)$$

where  $\mathbf{U}^{-1}$  is the asymptotic variance of  $\hat{\boldsymbol{\eta}}$ .

Next we see  $\hat{\boldsymbol{\eta}}^{(2)}$ . Let  $\mathbf{W} = (\mathbf{w}_1, \dots, \mathbf{w}_{p+1})^T$ , then  $\hat{\boldsymbol{\eta}}^{(2)}$  can be rewritten as

$$\begin{aligned} \hat{\boldsymbol{\eta}}^{(2)} &= \mathbf{U}^{-1} (-\mathbf{W}) \mathbf{U}^{-1} \mathbf{s} \\ &= \mathbf{U}^{-1} \left[ \text{tr} [\mathbf{U}^{-1} \mathbf{s} (-\mathbf{w}_1^T)], \dots, \text{tr} [\mathbf{U}^{-1} \mathbf{s} (-\mathbf{w}_{p+1}^T)] \right]^T. \end{aligned}$$



It is noted that  $E[\mathbf{s}(-\mathbf{w}_r^\top)] = \mathbf{Cov}(\mathbf{s}, -\mathbf{w}_r) = \mathbf{Cov}(\mathbf{s}, \partial s_r / \partial \boldsymbol{\eta}) = \mathbf{J}_r$ , then

$$E(\widehat{\boldsymbol{\eta}}^{(2)}) = \mathbf{U}^{-1} \mathbf{a}. \quad (5.32)$$

Finally we evaluate  $\widehat{\boldsymbol{\eta}}^{(3)} = 2^{-1} \mathbf{U}^{-1} \mathbf{t}$ . It can be seen that

$$\begin{aligned} t_r &= (\widehat{\boldsymbol{\eta}}^{(1)})^\top \frac{\partial^2 s_r}{\partial \boldsymbol{\eta} \partial \boldsymbol{\eta}^\top} \widehat{\boldsymbol{\eta}}^{(1)} + o_p(1) \\ &= \text{tr} \left[ \frac{\partial^2 s_r}{\partial \boldsymbol{\eta} \partial \boldsymbol{\eta}^\top} \mathbf{U}^{-1} \mathbf{s} \mathbf{s}^\top \mathbf{U}^{-1} \right] + o_p(1), \end{aligned}$$

Because  $E(\mathbf{s}) = \mathbf{0}$  and  $y_i$ 's are mutually independent, it follows that

$$\begin{aligned} E(t_r) &= \text{tr} \left[ E \left( \frac{\partial^2 s_r}{\partial \boldsymbol{\eta} \partial \boldsymbol{\eta}^\top} \right) \mathbf{U}^{-1} E(\mathbf{s} \mathbf{s}^\top) \mathbf{U}^{-1} \right] + o(1) \\ &= \text{tr} [\mathbf{K}_r \mathbf{U}^{-1}] + o(1). \end{aligned}$$

Thus we can get

$$E(\widehat{\boldsymbol{\eta}}^{(3)}) = 2^{-1} \mathbf{U}^{-1} \mathbf{b} + o(m^{-1}). \quad (5.33)$$

From (5.32) and (5.33), the bias of  $\widehat{\boldsymbol{\eta}}$  is

$$E(\widehat{\boldsymbol{\eta}} - \boldsymbol{\eta}) = \mathbf{U}^{-1} (\mathbf{a} + 2^{-1} \mathbf{b}) + o(m^{-1}). \quad (5.34)$$

### 5.5.2 Expressions of $\mathbf{J}_r$ and $\mathbf{K}_r$

First we get the expression of  $\mathbf{J}_r$ . Let  $\mathbf{V}_1 = [\mathbf{v}_{11}, \dots, \mathbf{v}_{1,p+1}]^\top$ ,  $\mathbf{V}_2 = [\mathbf{v}_{21}, \dots, \mathbf{v}_{2,p+1}]^\top$ . It follows from  $\mathbf{W} = \mathbf{V}_1 + \mathbf{V}_2 - \mathbf{U}$  that

$$\mathbf{J}_r = \mathbf{Cov}(\mathbf{s}, -\mathbf{w}_r) = \mathbf{Cov}(\mathbf{s}, -\mathbf{v}_{1r}) + \mathbf{Cov}(\mathbf{s}, -\mathbf{v}_{2r}).$$

Let  $\mathbf{D}_i^\top = [\mathbf{d}_{i1}, \dots, \mathbf{d}_{i,p+1}]^\top$ . Then  $\mathbf{Cov}(\mathbf{s}, -\mathbf{v}_{1r})$  can be written as

$$\begin{aligned} \mathbf{Cov}(\mathbf{s}, -\mathbf{v}_{1r}) &= E \left[ \sum_{i,j=1}^m \mathbf{D}_i^\top \boldsymbol{\Sigma}_i^{-1} \mathbf{g}_i \left( \frac{\partial \mathbf{d}_{jr}^\top \boldsymbol{\Sigma}_j^{-1}}{\partial \boldsymbol{\eta}^\top} \right) (\mathbf{I}_{p+1} \otimes \mathbf{g}_j) \right] \\ &= \sum_{i=1}^m \mathbf{D}_i^\top \boldsymbol{\Sigma}_i^{-1} E(\mathbf{g}_i \mathbf{g}_i^\top) \left( \frac{\partial \mathbf{d}_{ir}^\top \boldsymbol{\Sigma}_i^{-1}}{\partial \boldsymbol{\eta}} \right)^\top \\ &= \sum_{i=1}^m \mathbf{D}_i^\top \left\{ \frac{\partial \mathbf{d}_{ir}^\top}{\partial \boldsymbol{\eta}} \boldsymbol{\Sigma}_i^{-1} + (\mathbf{I}_{p+1} \otimes \mathbf{d}_{ir}^\top) \frac{\partial \boldsymbol{\Sigma}_i^{-1}}{\partial \boldsymbol{\eta}} \right\}^\top. \end{aligned} \quad (5.35)$$

$\mathbf{Cov}(\mathbf{s}, -\mathbf{v}_{2r})$  is

$$\begin{aligned} \mathbf{Cov}(\mathbf{s}, -\mathbf{v}_{2r}) &= E \left[ \sum_{i,j=1}^m \mathbf{D}_i^\top \boldsymbol{\Sigma}_i^{-1} \mathbf{g}_i \mathbf{d}_{jr}^\top \boldsymbol{\Sigma}_j^{-1} \frac{\partial \mathbf{g}_j}{\partial \boldsymbol{\eta}^\top} \right] \\ &= E \left[ \sum_{i=1}^m \mathbf{D}_i^\top \boldsymbol{\Sigma}_i^{-1} \mathbf{g}_i \mathbf{d}_{ir}^\top \boldsymbol{\Sigma}_i^{-1} \frac{\partial \mathbf{g}_i}{\partial \boldsymbol{\eta}^\top} \right]. \end{aligned}$$

Here  $\partial \mathbf{g}_i / \partial \boldsymbol{\eta}^T$  is expressed as

$$\frac{\partial \mathbf{g}_i}{\partial \boldsymbol{\eta}^T} = -2Q(m_i)(y_i - m_i) \begin{bmatrix} \mathbf{0} & 0 \\ \mathbf{x}_i^T & 0 \end{bmatrix} - Q(m_i) \begin{bmatrix} \mathbf{x}_i^T & 0 \\ \phi_i Q'(m_i) \mathbf{x}_i^T & -(1 + v_2/n_i)(\lambda - v_2)^{-2} \end{bmatrix},$$

thus it follows that

$$\text{Cov}(\mathbf{s}, -\mathbf{v}_{2r}) = -2 \sum_{i=1}^m Q(m_i) \mathbf{D}_i^T \boldsymbol{\Sigma}_i^{-1} \begin{bmatrix} \mu_{2i} \\ \mu_{3i} \end{bmatrix} \mathbf{d}_{ir}^T \boldsymbol{\Sigma}_i^{-1} \begin{bmatrix} \mathbf{0} & 0 \\ \mathbf{x}_i^T & 0 \end{bmatrix}. \quad (5.36)$$

From (5.35) and (5.36),

$$\mathbf{J}_r = \sum_{i=1}^m \mathbf{D}_i^T \left\{ \frac{\partial \mathbf{d}_{ir}^T}{\partial \boldsymbol{\eta}} \boldsymbol{\Sigma}_i^{-1} + (\mathbf{I}_{p+1} \otimes \mathbf{d}_{ir}^T) \frac{\partial \boldsymbol{\Sigma}_i^{-1}}{\partial \boldsymbol{\eta}} \right\}^T - 2 \sum_{i=1}^m Q(m_i) \mathbf{D}_i^T \boldsymbol{\Sigma}_i^{-1} \begin{bmatrix} \mu_{2i} \\ \mu_{3i} \end{bmatrix} \mathbf{d}_{ir}^T \boldsymbol{\Sigma}_i^{-1} \begin{bmatrix} \mathbf{0} & 0 \\ \mathbf{x}_i^T & 0 \end{bmatrix}, \quad (5.37)$$

where the expressions of  $\partial \mathbf{d}_{ir}^T / \partial \boldsymbol{\eta}$  and  $\partial \boldsymbol{\Sigma}_i^{-1} / \partial \boldsymbol{\eta}$  are given in Section 5.5.3 and 5.5.4.

Next we get the expression of  $\mathbf{K}_r = E[(\partial^2 s_r) / (\partial \boldsymbol{\eta} \partial \boldsymbol{\eta}^T)]$ . It can be seen that

$$\begin{aligned} \mathbf{K}_r &= \sum_{i=1}^m \left\{ 2 \frac{\partial \mathbf{d}_{ir}^T}{\partial \boldsymbol{\eta}} \boldsymbol{\Sigma}_i^{-1} E \left( \frac{\partial \mathbf{g}_i}{\partial \boldsymbol{\eta}^T} \right) + 2 (\mathbf{I}_{p+1} \otimes \mathbf{d}_{ir}^T) \frac{\partial \boldsymbol{\Sigma}_i^{-1}}{\partial \boldsymbol{\eta}} E \left( \frac{\partial \mathbf{g}_i}{\partial \boldsymbol{\eta}^T} \right) + \{ \mathbf{I}_{p+1} \otimes (\mathbf{d}_{ir}^T \boldsymbol{\Sigma}_i^{-1}) \} E \left( \frac{\partial^2 \mathbf{g}_i}{\partial \boldsymbol{\eta} \partial \boldsymbol{\eta}^T} \right) \right\} \\ &= \sum_{i=1}^m \left\{ -2 \frac{\partial \mathbf{d}_{ir}^T}{\partial \boldsymbol{\eta}} \boldsymbol{\Sigma}_i^{-1} \mathbf{D}_i - 2 (\mathbf{I}_{p+1} \otimes \mathbf{d}_{ir}^T) \frac{\partial \boldsymbol{\Sigma}_i^{-1}}{\partial \boldsymbol{\eta}} \mathbf{D}_i + \{ \mathbf{I}_{p+1} \otimes (\mathbf{d}_{ir}^T \boldsymbol{\Sigma}_i^{-1}) \} E \left( \frac{\partial^2 \mathbf{g}_i}{\partial \boldsymbol{\eta} \partial \boldsymbol{\eta}^T} \right) \right\}, \end{aligned} \quad (5.38)$$

where the expressions of  $\partial \mathbf{d}_{ir}^T / \partial \boldsymbol{\eta}$  and  $\partial \boldsymbol{\Sigma}_i^{-1} / \partial \boldsymbol{\eta}$  are given in Section 5.5.3 and 5.5.4, and  $E[(\partial^2 \mathbf{g}_i) / (\partial \boldsymbol{\eta} \partial \boldsymbol{\eta}^T)] = [\mathbf{G}_1^T, \dots, \mathbf{G}_{p+1}^T]^T$  is written as

$$\begin{aligned} &[\mathbf{G}_1^T, \dots, \mathbf{G}_p^T]^T \\ &= -Q(m_i) \mathbf{x}_i \otimes \begin{bmatrix} Q'(m_i) \mathbf{x}_i^T & 0 \\ \{-2Q(m_i) + (Q'(m_i))^2 + 2v_2 \phi_i Q(m_i)\} \mathbf{x}_i^T & -Q'(m_i)(1 + v_2/n_i)(\lambda - v_2)^{-2} \end{bmatrix}, \\ &\mathbf{G}_{p+1} \\ &= -Q(m_i) \begin{bmatrix} \mathbf{0} & 0 \\ -(1 + v_2/n_i)(\lambda - v_2)^{-2} Q'(m_i) \mathbf{x}_i^T & 2(1 + v_2/n_i)(\lambda - v_2)^{-3} \end{bmatrix}. \end{aligned}$$

### 5.5.3 Expression of $\partial \mathbf{d}_{ir}^T / \partial \boldsymbol{\eta}$

First, we give the expression of  $\partial \mathbf{d}_{ir}^T / \partial \boldsymbol{\eta}$  for  $r = 1, \dots, p$ . It follows that

$$\mathbf{d}_{ir}^T = Q(m_i) [x_{ir}, Q'(m_i) \phi_i x_{ir}],$$

then we get

$$\begin{aligned} \frac{\partial \mathbf{d}_{ir}^T}{\partial \boldsymbol{\beta}} &= Q(m_i) x_{ir} [Q'(m_i), \{(Q'(m_i))^2 + 2v_2 Q(m_i)\} \phi_i] \otimes \mathbf{x}_i, \\ \frac{\partial \mathbf{d}_{ir}^T}{\partial \lambda} &= Q(m_i) Q'(m_i) x_{ir} [0, -(1 + v_2/n_i)(\lambda - v_2)^{-2}]. \end{aligned}$$

Next, we give the expression of  $\partial \mathbf{d}_{i,p+1}^T / \partial \boldsymbol{\eta}$ . It follows that

$$\mathbf{d}_{i,p+1}^T = Q(m_i) [0, -(1 + v_2/n_i)(\lambda - v_2)^{-2}],$$

then we can get

$$\begin{aligned}\frac{\partial \mathbf{d}_{i,p+1}^T}{\partial \boldsymbol{\beta}} &= -Q'(m_i)Q(m_i)(1+v_2/n_i)(\lambda-v_2)^{-2}[\mathbf{0}; \mathbf{x}_i], \\ \frac{\partial \mathbf{d}_{i,p+1}^T}{\partial \lambda} &= [0, 2Q(m_i)(1+v_2/n_i)(\lambda-v_2)^{-3}].\end{aligned}$$

#### 5.5.4 Expression of $\partial \boldsymbol{\Sigma}_i^{-1} / \partial \boldsymbol{\eta}$

It follows that

$$\frac{\partial \boldsymbol{\Sigma}_i^{-1}}{\partial \boldsymbol{\eta}} = -(\mathbf{I}_{p+1} \otimes \boldsymbol{\Sigma}_i^{-1}) \frac{\partial \boldsymbol{\Sigma}_i}{\partial \boldsymbol{\eta}} \boldsymbol{\Sigma}_i^{-1},$$

thus it suffices to evaluate  $\partial \boldsymbol{\Sigma}_i / \partial \boldsymbol{\eta}$ .

Let  $\partial \boldsymbol{\Sigma}_i / \partial \boldsymbol{\eta} = [\mathbf{H}_1^T, \dots, \mathbf{H}_{p+1}^T]^T$ . For  $r = 1, \dots, p$ ,

$$\mathbf{H}_r = \begin{bmatrix} \partial \mu_{2i} / \partial \beta_r & \partial \mu_{3i} / \partial \beta_r \\ \partial \mu_{3i} / \partial \beta_r & \partial \mu_{4i} / \partial \beta_r - \partial \mu_{2i}^2 / \partial \beta_r \end{bmatrix},$$

where

$$\begin{aligned}\frac{\partial \mu_{2i}}{\partial \beta_r} &= \phi_i Q(m_i) Q'(m_i) x_{ir}, \\ \frac{\partial \mu_{2i}^2}{\partial \beta_r} &= 2\phi_i^2 \{Q(m_i)\}^2 Q'(m_i) x_{ir}, \\ \frac{\partial \mu_{3i}}{\partial \beta_r} &= \frac{(\lambda/n_i + 1)(\lambda/n_i + 2)}{(\lambda - v_2)(\lambda - 2v_2)} \left[ \{Q'(m_i)\}^2 + 2v_2 Q(m_i) \right] Q(m_i) x_{ir}.\end{aligned}$$

Let  $d_i = v_2/n_i$ , then

$$\begin{aligned}\frac{\partial \mu_{4i}}{\partial \beta_r} &= \left\{ \frac{6(d_i + 1)(2d_i + 1)(3d_i + 1)}{(\lambda - v_2)(\lambda - 2v_2)(\lambda - 3v_2)} + \frac{12(d_i + 1)(2d_i + 1)}{n_i(\lambda - v_2)(\lambda - 2v_2)} + \frac{7(d_i + 1)}{n_i^2(\lambda - v_2)} + \frac{1}{n_i^3} \right\} J_1, \\ &+ \left\{ \frac{3(d_i + 1)(2d_i + 1)(3d_i + 1)}{(\lambda - v_2)(\lambda - 3v_2)} + \frac{2(d_i + 1)(4d_i + 3)}{n_i(\lambda - v_2)} + \frac{2d_i + 3}{n_i^2} \right\} J_2,\end{aligned}$$

where

$$\begin{aligned}J_1 &= \left[ \{Q'(m_i)\}^2 + 4v_2 Q(m_i) \right] Q(m_i) Q'(m_i) x_{ir}, \\ J_2 &= 2 \{Q(m_i)\}^2 Q'(m_i) x_{ir}.\end{aligned}$$

It can be seen that

$$\mathbf{H}_{p+1} = \begin{bmatrix} \partial \mu_{2i} / \partial \lambda & \partial \mu_{3i} / \partial \lambda \\ \partial \mu_{3i} / \partial \lambda & \partial \mu_{4i} / \partial \lambda - \partial \mu_{2i}^2 / \partial \lambda \end{bmatrix},$$

where

$$\begin{aligned}\frac{\partial \mu_{2i}}{\partial \lambda} &= -Q(m_i) \frac{v_2/n_i + 1}{(\lambda - v_2)^2}, \\ \frac{\partial \mu_{2i}^2}{\partial \lambda} &= -2 \{Q(m_i)\}^2 \frac{(\lambda/n_i + 1)(v_2/n_i + 1)}{(\lambda - v_2)^3}, \\ \frac{\partial \mu_{3i}}{\partial \lambda} &= Q(m_i)Q'(m_i) \frac{(2\lambda + 3n_i)(\lambda - v_2)(\lambda - 2v_2) - (\lambda + n_i)(\lambda + 2n_i)(2\lambda - 3v_2)}{n_i^2(\lambda - v_2)^2(\lambda - 2v_2)^2}, \\ \frac{\partial \mu_{4i}}{\partial \lambda} &= 6(d_i + 1)(2d_i + 1)(3d_i + 1)J_3 + \frac{12}{n_i}(d_i + 1)(2d_i + 1)Q(m_i) \{Q'(m_i)\}^2 J_4 \\ &\quad - \left\{ \frac{2(d_i + 1)(4d_i + 3)}{n_i} \{Q(m_i)\}^2 + \frac{7(d_i + 1)}{n_i^2} Q(m_i) \{Q'(m_i)\}^2 \right\} \frac{1}{(\lambda - v_2)^2},\end{aligned}$$

for

$$\begin{aligned}J_3 &= -\{Q(m_i)\}^2 \frac{\lambda - 2v_2}{(\lambda - v_2)^2(\lambda - 3v_2)^2} - Q(m_i) \{Q'(m_i)\}^2 \frac{3\lambda^2 - 12v_2\lambda + 11v_2^2}{(\lambda - v_2)^2(\lambda - 2v_2)^2(\lambda - 3v_2)^2}, \\ J_4 &= -\frac{2\lambda - 3v_2}{(\lambda - v_2)^2(\lambda - 2v_2)^2}.\end{aligned}$$

### 5.5.5 Approximations of $\xi_{ri}$

It suffices to evaluate  $\xi_{ri}$  up to  $O(1)$ , namely the remainder is of order  $O(n_i^{-1})$ . It is noted that  $\xi_{ri} = E[h'(\hat{\mu}_i^B)(y_i - m_i)^r] = E[h'(\mu_i)(y_i - m_i)^r] + O(n_i^{-1})$ . After some calculations, we can get  $\xi_{ri}$ . For the Poisson-gamma mixture model,

$$\begin{aligned}\xi_{1i} &= -\frac{1}{\lambda m_i - 1} + O(n_i^{-1}), \quad \xi_{2i} = \frac{m_i}{\lambda m_i - 1} + O(n_i^{-1}), \\ \xi_{3i} &= -\frac{m_i}{\lambda(\lambda m_i - 1)} + O(n_i^{-1}), \quad \xi_{4i} = \frac{3\lambda m_i^2 - 2m_i}{\lambda^2(\lambda m_i - 1)} + O(n_i^{-1}),\end{aligned}$$

for the binomial-beta mixture model,

$$\begin{aligned}\xi_{1i} &= -\frac{(\lambda - 1)(1 - 2m_i)}{(\lambda m_i - 1)\{\lambda(1 - m_i) - 1\}} + O(n_i^{-1}), \\ \xi_{2i} &= \frac{(\lambda - 2)m_i(1 - m_i)}{(\lambda m_i - 1)\{\lambda(1 - m_i) - 1\}} + O(n_i^{-1}), \\ \xi_{3i} &= -\frac{m_i(1 - m_i)(1 - 2m_i)}{(\lambda m_i - 1)\{\lambda(1 - m_i) - 1\}} + O(n_i^{-1}), \\ \xi_{4i} &= \frac{3\lambda m_i^2(1 - m_i)^2 - 2m_i(1 - m_i)(m_i^2 - m_i + 1)}{(\lambda m_i - 1)\{\lambda(1 - m_i) - 1\}} + O(n_i^{-1}).\end{aligned}$$

## 5.6 Proofs

### 5.6.1 Proof of Lemma 5.3

First,  $B_2$  can be rewritten as

$$\begin{aligned}
B_2 &= \sum_{i=1}^m n_i E \left[ E \left[ (y_i - \mu_i) \frac{\partial h(\hat{\mu}_i^B)}{\partial \boldsymbol{\eta}^T} (\hat{\boldsymbol{\eta}} - \boldsymbol{\eta}) \mid \mathbf{y} \right] \right] \\
&= \sum_{i=1}^m n_i E \left[ (y_i - \hat{\mu}_i^B) \frac{\partial h(\hat{\mu}_i^B)}{\partial \boldsymbol{\eta}^T} (\hat{\boldsymbol{\eta}} - \boldsymbol{\eta}) \right] \\
&= \sum_{i=1}^m \frac{n_i \lambda}{n_i + \lambda} E [\mathbf{f}_i^T (\hat{\boldsymbol{\eta}} - \boldsymbol{\eta})], \tag{5.39}
\end{aligned}$$

where

$$\mathbf{f}_i = (y_i - m_i) \frac{\partial h(\hat{\mu}_i^B)}{\partial \boldsymbol{\eta}} = \left[ \frac{\lambda Q(m_i) h'(\hat{\mu}_i^B) (y_i - m_i)}{n_i + \lambda} \mathbf{x}_i^T, -\frac{n_i h'(\hat{\mu}_i^B) (y_i - m_i)^2}{(n_i + \lambda)^2} \right]^T.$$

Using the stochastic expansion of  $\hat{\boldsymbol{\eta}}$  in Section 5.5.1, we decompose  $B_2$  as

$$\begin{aligned}
B_2 &= \sum_{i=1}^m \frac{n_i \lambda}{n_i + \lambda} E(\mathbf{f}_i^T \hat{\boldsymbol{\eta}}^{(1)}) + \sum_{i=1}^m \frac{n_i \lambda}{n_i + \lambda} E[\mathbf{f}_i^T (\hat{\boldsymbol{\eta}}^{(2)} + \hat{\boldsymbol{\eta}}^{(3)})] + o(1) \\
&\equiv I_1 + I_2 + o(1).
\end{aligned}$$

$E(\mathbf{f}_i^T \hat{\boldsymbol{\eta}}^{(1)})$  is evaluated as

$$\begin{aligned}
E(\mathbf{f}_i^T \hat{\boldsymbol{\eta}}^{(1)}) &= E \left[ \left[ \frac{\lambda Q(m_i) h'(\hat{\mu}_i^B) (y_i - m_i)}{n_i + \lambda} \mathbf{x}_i^T, -\frac{n_i h'(\hat{\mu}_i^B) (y_i - m_i)^2}{(n_i + \lambda)^2} \right] \mathbf{U}^{-1} \sum_{j=1}^m \mathbf{D}_j^T \boldsymbol{\Sigma}_j^{-1} \mathbf{g}_j \right] \\
&= \frac{\lambda}{n_i + \lambda} Q(m_i) \mathbf{x}_i^T \mathbf{U}_1 \mathbf{D}_i^T \boldsymbol{\Sigma}_i^{-1} \begin{bmatrix} \xi_{2i} \\ \xi_{3i} - \phi_i Q(m_i) \xi_{1i} \end{bmatrix} \\
&\quad - \frac{n_i}{(n_i + \lambda)^2} \mathbf{U}_2 \mathbf{D}_i^T \boldsymbol{\Sigma}_i^{-1} \begin{bmatrix} \xi_{3i} \\ \xi_{4i} - \phi_i Q(m_i) \xi_{2i} \end{bmatrix},
\end{aligned}$$

which yields  $I_1 = B_{21}$ . From (5.32) and (5.33),  $E[\mathbf{f}_i^T (\hat{\boldsymbol{\eta}}^{(2)} + \hat{\boldsymbol{\eta}}^{(3)})]$  is evaluated as

$$\begin{aligned}
E[\mathbf{f}_i^T (\hat{\boldsymbol{\eta}}^{(2)} + \hat{\boldsymbol{\eta}}^{(3)})] &= E(\mathbf{f}_i^T) E(\hat{\boldsymbol{\eta}}^{(2)} + \hat{\boldsymbol{\eta}}^{(3)}) + o(m^{-1}) \\
&= \left[ \frac{\lambda}{n_i + \lambda} Q(m_i) \xi_{1i} \mathbf{x}_i^T, -\frac{n_i}{(n_i + \lambda)^2} \xi_{2i} \right] \mathbf{U}^{-1} (\mathbf{a} + 2^{-1} \mathbf{b}) + o(m^{-1}) \\
&= \frac{\lambda}{n_i + \lambda} Q(m_i) \xi_{1i} \mathbf{x}_i^T \mathbf{U}_1 (\mathbf{a} + 2^{-1} \mathbf{b}) - \frac{n_i}{(n_i + \lambda)^2} \xi_{2i} \mathbf{U}_2 (\mathbf{a} + 2^{-1} \mathbf{b}) + o(m^{-1}),
\end{aligned}$$

which yields  $I_2 = B_{22} + o(1)$ . □

### 5.6.2 Proof of Lemmas 5.4 and 5.5

In the same manner as (5.39),  $B_3$  can be rewritten as

$$B_3 = \sum_{i=1}^m \frac{n_i \lambda}{n_i + \lambda} \text{tr} [E[\mathbf{F}_i (\hat{\boldsymbol{\eta}} - \boldsymbol{\eta}) (\hat{\boldsymbol{\eta}} - \boldsymbol{\eta})^T]],$$

where

$$\mathbf{F}_i = (y_i - m_i) \frac{\partial^2 h(\hat{\mu}_i^B)}{\partial \boldsymbol{\eta} \partial \boldsymbol{\eta}^T} = (y_i - m_i) \begin{bmatrix} \partial^2 h(\hat{\mu}_i^B) / (\partial \boldsymbol{\beta} \partial \boldsymbol{\beta}^T) & \partial^2 h(\hat{\mu}_i^B) / (\partial \boldsymbol{\beta} \partial \lambda) \\ \partial^2 h(\hat{\mu}_i^B) / (\partial \lambda \partial \boldsymbol{\beta}^T) & \partial^2 h(\hat{\mu}_i^B) / (\partial \lambda)^2 \end{bmatrix},$$

for

$$\begin{aligned} \frac{\partial^2 h(\hat{\mu}_i^B)}{\partial \boldsymbol{\beta} \partial \boldsymbol{\beta}^T} &= \frac{\lambda}{n_i + \lambda} Q(m_i) Q'(m_i) h'(\hat{\mu}_i^B) \mathbf{x}_i \mathbf{x}_i^T + \left( \frac{\lambda}{n_i + \lambda} \right)^2 \{Q(m_i)\}^2 h''(\hat{\mu}_i^B) \mathbf{x}_i \mathbf{x}_i^T, \\ \frac{\partial^2 h(\hat{\mu}_i^B)}{\partial \boldsymbol{\beta} \partial \lambda} &= \frac{n_i}{(n_i + \lambda)^2} Q(m_i) h'(\hat{\mu}_i^B) \mathbf{x}_i - \frac{n_i \lambda}{(n_i + \lambda)^3} Q(m_i) h''(\hat{\mu}_i^B) \mathbf{x}_i, \\ \frac{\partial^2 h(\hat{\mu}_i^B)}{\partial \lambda \partial \lambda} &= \frac{2n_i}{(n_i + \lambda)^3} h'(\hat{\mu}_i^B) (y_i - m_i) + \frac{n_i^2}{(n_i + \lambda)^4} h''(\hat{\mu}_i^B) (y_i - m_i)^2. \end{aligned}$$

It follows that

$$\begin{aligned} E[\mathbf{F}_i(\hat{\boldsymbol{\eta}} - \boldsymbol{\eta})(\hat{\boldsymbol{\eta}} - \boldsymbol{\eta})^T] &= E(\mathbf{F}_i) E[\hat{\boldsymbol{\eta}}^{(1)}(\hat{\boldsymbol{\eta}}^{(1)})^T] + o(m^{-1}) \\ &= E(\mathbf{F}_i) \mathbf{U}^{-1} + o(m^{-1}), \end{aligned}$$

which yields  $B_3 = B_{31} + 2B_{32} + o(1)$  under the condition (C1), and  $B_3 = B_{33} + 2B_{34} + B_{35} + o(1)$  under the condition (C2), respectively.  $\square$

## Chapter 6

# A variant of AIC using Bayesian marginal likelihood

In this chapter, we propose an information criterion which measures the prediction risk of the predictive density based on the Bayesian marginal likelihood from a frequentist point of view. We derive the criteria for selecting variables in linear regression models by putting the prior on the regression coefficients, and discuss the relationship between the proposed criteria and other related ones. There are three advantages of our method. Firstly, this is a compromise between the frequentist and Bayesian standpoint because it evaluates the frequentist's risk of the Bayesian model. Thus it is less influenced by prior misspecification. Secondly, non-informative improper prior can be also used for constructing the criterion. When the uniform prior is assumed on the regression coefficients, the resulting criterion is identical to the residual information criterion (RIC) of Shi and Tsai (2002). Lastly, the criteria have the consistency property for selecting the true model.

### 6.1 Motivation

The problem of selecting appropriate models has been extensively studied in the literature since Akaike (1973, 1974), who derived so called the Akaike information criterion (AIC). Since the AIC and their variants are based on the risk of the predictive densities with respect to the Kullback–Leibler (KL) divergence, they can select a good model in the light of prediction. It is known, however, that the AIC-type criteria do not have the consistency property, namely, the probability that the criteria select the true model does not converges to 1. Another approach to model selection is Bayesian procedures such as Bayes factors and the Bayesian information criterion (BIC) suggested by Schwarz (1978), both of which are constructed based on the Bayesian marginal likelihood. Bayesian procedures for model selection have the consistency property in some specific models, while they do not select models in terms of prediction. In addition, it is known that Bayes factors do not work for improper prior distributions and that the BIC does not use any specific prior information. In this chapter, we provide a unified framework to derive an information criterion for model selection so that it can produce various information criteria including AIC, BIC and the residual information criterion (RIC) suggested by of Shi and Tsai (2002). Especially, we propose an intermediate criterion between AIC and BIC using the empirical Bayes method.

To explain the unified framework in the general setup, let  $\mathbf{y}$  be an  $n$ -variate observable random vector whose density is  $m(\mathbf{y}|\boldsymbol{\eta})$  for a vector of unknown parameters  $\boldsymbol{\eta}$ . Let  $\hat{m}(\hat{\mathbf{y}}; \mathbf{y})$  be

a predictive density for  $m(\tilde{\mathbf{y}}|\boldsymbol{\eta})$ , where  $\tilde{\mathbf{y}}$  is an independent replication of  $\mathbf{y}$ . We here evaluate the predictive performance of  $\hat{m}(\tilde{\mathbf{y}}; \mathbf{y})$  in terms of the following risk:

$$R(\boldsymbol{\eta}; \hat{m}) = \int \left[ \int \log \left\{ \frac{m(\tilde{\mathbf{y}}|\boldsymbol{\eta})}{\hat{m}(\tilde{\mathbf{y}}; \mathbf{y})} \right\} m(\tilde{\mathbf{y}}|\boldsymbol{\eta}) d\tilde{\mathbf{y}} \right] m(\mathbf{y}|\boldsymbol{\eta}) d\mathbf{y}. \quad (6.1)$$

Since this is interpreted as a risk with respect to the KL divergence, we call it the KL risk. The spirit of AIC suggests that we can provide an information criterion for model selection as an (asymptotically) unbiased estimator of the information

$$\begin{aligned} I(\boldsymbol{\eta}; \hat{m}) &= \iint -2 \log\{\hat{m}(\tilde{\mathbf{y}}; \mathbf{y})\} m(\tilde{\mathbf{y}}|\boldsymbol{\eta}) m(\mathbf{y}|\boldsymbol{\eta}) d\tilde{\mathbf{y}} d\mathbf{y} \\ &= E_{\boldsymbol{\eta}} [-2 \log\{\hat{m}(\tilde{\mathbf{y}}; \mathbf{y})\}], \end{aligned} \quad (6.2)$$

which is a part of (6.1) (multiplied by 2), where  $E_{\boldsymbol{\eta}}$  denotes the expectation with respect to the distribution of  $m(\tilde{\mathbf{y}}, \mathbf{y}|\boldsymbol{\eta}) = m(\tilde{\mathbf{y}}|\boldsymbol{\eta})m(\mathbf{y}|\boldsymbol{\eta})$ . Let  $\Delta = I(\boldsymbol{\eta}; \hat{m}) - E_{\boldsymbol{\eta}}[-2 \log\{\hat{m}(\tilde{\mathbf{y}}; \mathbf{y})\}]$ . Then, the AIC variant based on the predictor  $\hat{m}(\tilde{\mathbf{y}}; \mathbf{y})$  is defined by

$$\text{IC}(\hat{m}) = -2 \log\{\hat{m}(\mathbf{y}; \mathbf{y})\} + \hat{\Delta},$$

where  $\hat{\Delta}$  is an (asymptotically) unbiased estimator of  $\Delta$ .

It is interesting to point out that  $\text{IC}(\hat{m})$  produces AIC and BIC for specific predictors.

(AIC) Put  $\hat{m}(\tilde{\mathbf{y}}; \mathbf{y}) = m(\tilde{\mathbf{y}}|\hat{\boldsymbol{\eta}})$  for the maximum likelihood estimator  $\hat{\boldsymbol{\eta}}$  of  $\boldsymbol{\eta}$ . Then,  $\text{IC}(m(\tilde{\mathbf{y}}|\hat{\boldsymbol{\eta}}))$  is the exact AIC or the corrected AIC suggested by Sugiura (1978) and Hurvich and Tsai (1989), which is approximated by AIC of Akaike (1973, 1974) as  $-2 \log\{m(\mathbf{y}|\hat{\boldsymbol{\eta}})\} + 2 \dim(\boldsymbol{\eta})$ .

(BIC) Put  $\hat{m}(\tilde{\mathbf{y}}; \mathbf{y}) = m_{\pi_0}(\tilde{\mathbf{y}}) = \int m(\tilde{\mathbf{y}}|\boldsymbol{\eta}) \pi_0(\boldsymbol{\eta}) d\boldsymbol{\eta}$  for a proper prior distribution  $\pi_0(\boldsymbol{\eta})$ . Since it can be easily seen that  $I(\boldsymbol{\eta}; m_{\pi_0}) = E_{\boldsymbol{\eta}}[-2 \log\{m_{\pi_0}(\mathbf{y})\}]$ , we have  $\Delta = 0$  in this case, so that  $\text{IC}(m_{\pi_0}) = -2 \log\{m_{\pi_0}(\mathbf{y})\}$ , which is the Bayesian marginal likelihood. It is noted that  $-2 \log\{m_{\pi_0}(\mathbf{y})\}$  is approximated by  $\text{BIC} = -2 \log\{m(\mathbf{y}|\hat{\boldsymbol{\eta}})\} + \log(n) \cdot \dim(\boldsymbol{\eta})$ .

The criterion  $\text{IC}(\hat{m})$  can produce not only the conventional information criteria AIC and BIC, but also various criteria between AIC and BIC. For example, it is supposed that  $\boldsymbol{\eta}$  is divided as  $\boldsymbol{\eta} = (\boldsymbol{\beta}^T, \boldsymbol{\omega}^T)^T$  for a  $p$ -dimensional parameter vector of interest  $\boldsymbol{\beta}$  and an  $r$ -dimensional nuisance parameter vector  $\boldsymbol{\omega}$ . We assume that  $\boldsymbol{\beta}$  has a prior density  $\pi(\boldsymbol{\beta}|\boldsymbol{\lambda}, \boldsymbol{\omega})$  with hyperparameter  $\boldsymbol{\lambda}$ . The model is described as

$$\begin{aligned} \mathbf{y}|\boldsymbol{\beta} &\sim m(\mathbf{y}|\boldsymbol{\beta}, \boldsymbol{\omega}), \\ \boldsymbol{\beta} &\sim \pi(\boldsymbol{\beta}|\boldsymbol{\lambda}, \boldsymbol{\omega}), \end{aligned}$$

and  $\boldsymbol{\omega}$  and  $\boldsymbol{\lambda}$  are estimated by data. Inference based on such a model is called an empirical Bayes procedure. Put  $\hat{m}(\tilde{\mathbf{y}}; \mathbf{y}) = m_{\pi}(\tilde{\mathbf{y}}|\hat{\boldsymbol{\lambda}}, \hat{\boldsymbol{\omega}}) = \int f(\tilde{\mathbf{y}}|\boldsymbol{\beta}, \hat{\boldsymbol{\omega}}) \pi(\boldsymbol{\beta}|\hat{\boldsymbol{\lambda}}, \hat{\boldsymbol{\omega}}) d\boldsymbol{\beta}$  for some estimators  $\hat{\boldsymbol{\lambda}}$  and  $\hat{\boldsymbol{\omega}}$ . Then, the information in (6.2) is

$$I(\boldsymbol{\eta}; m_{\pi}) = \iint -2 \log\{m_{\pi}(\tilde{\mathbf{y}}|\hat{\boldsymbol{\lambda}}, \hat{\boldsymbol{\omega}})\} m(\tilde{\mathbf{y}}|\boldsymbol{\beta}, \boldsymbol{\omega}) m(\mathbf{y}|\boldsymbol{\beta}, \boldsymbol{\omega}) d\tilde{\mathbf{y}} d\mathbf{y}, \quad (6.3)$$

and the resulting information criterion is

$$\text{IC}(m_{\pi}) = -2 \log\{m_{\pi}(\mathbf{y}|\hat{\boldsymbol{\lambda}}, \hat{\boldsymbol{\omega}})\} + \hat{\Delta}, \quad (6.4)$$

where  $\hat{\Delta}$  is an (asymptotically) unbiased estimator of  $\Delta = I(\boldsymbol{\eta}; m_{\pi}) - E_{\boldsymbol{\eta}}[-2 \log\{m_{\pi}(\mathbf{y}|\hat{\boldsymbol{\lambda}}, \hat{\boldsymbol{\omega}})\}]$ .



There are three motivations to consider the information  $I(\boldsymbol{\eta}; m_\pi)$  in (6.3) and the information criterion  $IC(m_\pi)$  in (6.4).

Firstly, it is noted that the Bayesian predictor  $m_\pi(\tilde{\mathbf{y}}|\hat{\boldsymbol{\lambda}}, \hat{\boldsymbol{\omega}})$  is evaluated by the risk  $R(\boldsymbol{\eta}; m_\pi)$  in (6.1), which is based on a frequentist point of view. On the other hand, the Bayesian risk is

$$r(\boldsymbol{\psi}; \hat{m}) = \int R(\boldsymbol{\eta}; \hat{m})\pi(\boldsymbol{\beta}|\boldsymbol{\lambda}, \boldsymbol{\omega})d\boldsymbol{\beta}, \quad (6.5)$$

which measures the prediction error of  $\hat{m}(\tilde{\mathbf{y}}; \mathbf{y})$  under the assumption that the prior information is correct, where  $\boldsymbol{\psi} = (\boldsymbol{\lambda}^\top, \boldsymbol{\omega}^\top)^\top$ . The resulting Bayesian criteria such as PIC (Kitagawa, 1997) and DIC (Spiegelhalter et al., 2002) are sensitive to the prior misspecification, since they depend on the prior information. Because  $R(\boldsymbol{\eta}; m_\pi)$  can measure the prediction error of the Bayesian model from a standpoint of frequentists, however, the resulting criterion  $IC(m_\pi)$  is less influenced by the prior misspecification.

Secondly, we can construct the information criterion  $IC(m_\pi)$  when the prior distribution of  $\boldsymbol{\beta}$  is improper, since the information  $I(\boldsymbol{\eta}; m_\pi)$  in (6.3) can be defined formally for the corresponding improper marginal likelihood. Because the Bayesian risk  $r(\boldsymbol{\psi}; m_\pi)$  does not exist for the improper prior, however, we cannot obtain the corresponding Bayesian criteria and cannot use the Bayesian risk. Objective Bayesians want to avoid informative prior and many non-informative priors are improper. The suggested criterion  $IC(m_\pi)$  can respond to such a request. For example, objective Bayesians assume the uniform improper prior on regression coefficients  $\boldsymbol{\beta}$  in linear regression models. It is interesting to note that the resulting variable selection criterion (6.4) is identical to the residual information criterion (RIC) of Shi and Tsai (2002), which is shown in the next section.

Lastly, this criterion has the consistency property. We derive the criterion for the variable selection problem in general linear regression model and prove that the criterion selects the true model with probability tending to one. The BIC or marginal likelihood are known to have the consistency (Nishii, 1984), while most AIC-type criteria are not consistent. But AIC-type criteria have the property to choose a good model in the sense of minimizing the prediction error (Shibata, 1981; Shao, 1997). Our proposed criterion is consistent for selection of the parameters of interest  $\boldsymbol{\beta}$  and selects a good model in the light of prediction based on the empirical Bayes model.

The rest of this chapter is organized as follows. In Section 6.2, we obtain the information criterion (6.4) in linear regression model with general covariance structure and compare it with other related criteria. In Section 6.3, we prove the consistency of the criteria. In Section 6.4, we investigate the performance of the criteria through simulations. Section 6.5 concludes the chapter with some discussions.

## 6.2 Proposed criteria

### 6.2.1 Variable selection criteria for linear regression model

We consider the linear regression model as the candidate model, which is given as

$$\mathbf{y} = \mathbf{X}(j)\boldsymbol{\beta}_j + \mathbf{u}_j, \quad (6.6)$$

where  $\mathbf{y}$  is an  $n \times 1$  observation vector of the response variables,  $\mathbf{X}(j)$  is an  $n \times p_j$  matrix of the explanatory variables,  $\boldsymbol{\beta}_j$  is a  $p_j \times 1$  vector of the regression coefficients,  $\mathbf{u}_j$  is an  $n \times 1$  vector

of the random errors, and  $j$  is the index set which denotes the candidate model. Throughout the chapter, we assume that  $\mathbf{X}(j)$  has full column rank  $p_j$ . Here, the random error  $\mathbf{u}_j$  has the distribution  $\mathcal{N}_n(\mathbf{0}, \sigma_j^2 \boldsymbol{\Sigma})$ , where  $\sigma_j^2$  is an unknown scalar and  $\boldsymbol{\Sigma}$  is a known positive definite matrix. The linear regression model (6.6) includes the linear mixed model (3.1) in Chapter 3 as a special case.

We consider the problem of selecting the explanatory variables and assume that the true model is included by the candidate model, namely all the candidate model is overspecified. This is the common assumption to obtain the criterion. Under this assumption, the true mean of  $\mathbf{y}$  can be written as

$$E(\mathbf{y}) = \mathbf{X}(j)\boldsymbol{\beta}_j^*,$$

where  $\boldsymbol{\beta}_j^*$  is  $p_j \times 1$  vector whose  $p_j - p_*$  components are exactly 0 and the rest of components are not 0. We hereafter abbreviate the model index  $j$  for notational convenience. We also abbreviate  $\boldsymbol{\beta}_j^*$  as  $\boldsymbol{\beta}$  and write the true variance parameter as  $\sigma^2$ .

We shall construct the variable selection criteria for the regression model (6.6) which is of the form (6.4). We consider the following two situations.

[i] **Normal prior for  $\boldsymbol{\beta}$ .** We first assume the prior distribution of  $\boldsymbol{\beta}$ ,

$$\pi(\boldsymbol{\beta}|\sigma^2) \sim \mathcal{N}(\mathbf{0}, \sigma^2 \mathbf{W}),$$

where  $\mathbf{W}$  is a  $p \times p$  matrix suitably chosen with full rank. Examples of  $\mathbf{W}$  are  $\mathbf{W} = (\lambda \mathbf{X}^T \mathbf{X})^{-1}$  for  $\lambda > 0$  when  $\boldsymbol{\Sigma}$  is identity matrix, which is introduced by Zellner (1986), or more simply  $\mathbf{W} = \lambda^{-1} \mathbf{I}_p$ . Because the likelihood is  $m(\mathbf{y}|\boldsymbol{\beta}, \sigma^2) \sim \mathcal{N}(\mathbf{X}\boldsymbol{\beta}, \sigma^2 \boldsymbol{\Sigma})$ , the marginal likelihood function is

$$\begin{aligned} m_\pi(\mathbf{y}|\sigma^2) &= \int m(\mathbf{y}|\boldsymbol{\beta}, \sigma^2) \pi(\boldsymbol{\beta}|\sigma^2) d\boldsymbol{\beta} \\ &= (2\pi\sigma^2)^{-n/2} \cdot |\boldsymbol{\Sigma}|^{-1/2} \cdot |\mathbf{W}|^{-1/2} \cdot |\mathbf{X}^T \boldsymbol{\Sigma}^{-1} \mathbf{X} + \mathbf{W}^{-1}|^{-1/2} \cdot \exp\{-\mathbf{y}^T \mathbf{A} \mathbf{y} / (2\sigma^2)\}, \end{aligned}$$

where  $\mathbf{A} = \boldsymbol{\Sigma}^{-1} - \boldsymbol{\Sigma}^{-1} \mathbf{X} (\mathbf{X}^T \boldsymbol{\Sigma}^{-1} \mathbf{X} + \mathbf{W}^{-1})^{-1} \mathbf{X}^T \boldsymbol{\Sigma}^{-1}$ . Note that  $\mathbf{A} = (\boldsymbol{\Sigma} + \mathbf{B})^{-1}$  for  $\mathbf{B} = \mathbf{X} \mathbf{W} \mathbf{X}^T$ , namely  $m_\pi(\mathbf{y}|\sigma^2) \sim \mathcal{N}(\mathbf{0}, \sigma^2 (\boldsymbol{\Sigma} + \mathbf{B}))$ . Then we take the predictive density as  $\hat{m}(\tilde{\mathbf{y}}; \mathbf{y}) = m_\pi(\tilde{\mathbf{y}}|\hat{\sigma}^2)$  and the information (6.3) can be written as

$$I_{\pi,1}(\boldsymbol{\eta}) = E_\boldsymbol{\eta} [n \log(2\pi\hat{\sigma}^2) + \log |\boldsymbol{\Sigma}| + \log |\mathbf{W} \mathbf{X}^T \boldsymbol{\Sigma}^{-1} \mathbf{X} + \mathbf{I}_p| + \tilde{\mathbf{y}}^T \mathbf{A} \tilde{\mathbf{y}} / \hat{\sigma}^2], \quad (6.7)$$

where  $\hat{\sigma}^2 = \mathbf{y}^T (\boldsymbol{\Sigma}^{-1} - \mathbf{P}) \mathbf{y} / n$  for  $\mathbf{P} = \boldsymbol{\Sigma}^{-1} \mathbf{X} (\mathbf{X}^T \boldsymbol{\Sigma}^{-1} \mathbf{X})^{-1} \mathbf{X}^T \boldsymbol{\Sigma}^{-1}$  and  $E_\boldsymbol{\eta}$  denotes the expectation with respect to the distribution of  $m(\tilde{\mathbf{y}}, \mathbf{y}|\boldsymbol{\beta}, \sigma^2) = m(\tilde{\mathbf{y}}|\boldsymbol{\beta}, \sigma^2) m(\mathbf{y}|\boldsymbol{\beta}, \sigma^2)$  for  $\boldsymbol{\eta} = (\boldsymbol{\beta}^T, \sigma^2)^T$ . Note that  $\boldsymbol{\beta}$  is the parameter of interest and  $\sigma^2$  is the nuisance parameter, which corresponds to  $\boldsymbol{\omega}$  in the previous section. Then we propose the information criterion.

**Proposition 6.1** *The information  $I_{\pi,1}(\boldsymbol{\eta})$  in (6.7) is unbiasedly estimated by the information criterion*

$$\text{IC}_{\pi,1} = -2 \log\{m_\pi(\mathbf{y}|\hat{\sigma}^2)\} + \frac{2n}{n-p-2}, \quad (6.8)$$

where

$$-2 \log\{m_\pi(\mathbf{y}|\hat{\sigma}^2)\} = n \log \hat{\sigma}^2 + \log |\boldsymbol{\Sigma}| + \log |\mathbf{W} \mathbf{X}^T \boldsymbol{\Sigma}^{-1} \mathbf{X} + \mathbf{I}_p| + \mathbf{y}^T \mathbf{A} \mathbf{y} / \hat{\sigma}^2 + (\text{const}),$$

namely,  $E_\boldsymbol{\eta}(\text{IC}_{\pi,1}) = I_{\pi,1}(\boldsymbol{\eta})$ .

If  $n^{-1}\mathbf{W}^{1/2}\mathbf{X}^T\boldsymbol{\Sigma}^{-1}\mathbf{X}\mathbf{W}^{1/2}$  converges to a  $p \times p$  positive definite matrix as  $n \rightarrow \infty$ ,  $\log|\mathbf{W}\mathbf{X}^T\boldsymbol{\Sigma}^{-1}\mathbf{X} + \mathbf{I}_p|$  can be approximated to  $p \log n$ . In that case,  $\text{IC}_{\pi,1}$  is approximately expressed as

$$\text{IC}_{\pi,1}^* = n \log \hat{\sigma}^2 + \log |\boldsymbol{\Sigma}| + p \log n + 2 + \mathbf{y}^T \mathbf{A} \mathbf{y} / \hat{\sigma}^2,$$

when  $n$  is large.

Alternatively, the KL risk  $r(\boldsymbol{\psi}; \hat{m})$  in (6.5) can be also used for evaluating the risk of the predictive density  $m_\pi(\tilde{\mathbf{y}}|\hat{\sigma}^2)$ , since the prior distribution is proper. The resulting criterion is

$$\text{IC}_{\pi,2} = n \log \hat{\sigma}^2 + \log |\boldsymbol{\Sigma}| + p \log n + p, \quad (6.9)$$

which is an asymptotically unbiased estimator of  $I_{\pi,2}(\sigma^2) = E_\pi[I_{\pi,1}(\boldsymbol{\eta})]$  up to constant where  $E_\pi$  denotes the expectation with respect to the prior distribution  $\pi(\boldsymbol{\beta}|\sigma^2)$ , namely  $E_\pi E_\eta(\text{IC}_{\pi,2}) \rightarrow I_{\pi,2}(\sigma^2) + (\text{const})$  as  $n \rightarrow \infty$ . It is interesting to point out that  $\text{IC}_{\pi,2}$  is analogous to the criterion proposed by Bozdogan (1987) known as the consistent AIC, who suggested to replace the penalty term  $2p$  in the AIC with  $p + p \log n$ .

**[ii] Uniform prior for  $\boldsymbol{\beta}$ .** We next assume the uniform prior for  $\boldsymbol{\beta}$ , namely  $\boldsymbol{\beta} \sim \text{uniform}(\mathbb{R}^p)$ . Though this is improper prior distribution, we can obtain the marginal likelihood function formally:

$$\begin{aligned} m_r(\mathbf{y}|\sigma^2) &= \int m(\mathbf{y}|\boldsymbol{\beta}, \sigma^2) d\boldsymbol{\beta} \\ &= (2\pi\sigma^2)^{-(n-p)/2} \cdot |\boldsymbol{\Sigma}|^{-1/2} \cdot |\mathbf{X}^T\boldsymbol{\Sigma}^{-1}\mathbf{X}|^{-1/2} \cdot \exp\{-\mathbf{y}^T(\boldsymbol{\Sigma}^{-1} - \mathbf{P})\mathbf{y}/(2\sigma^2)\}, \end{aligned}$$

which is known as the residual likelihood (Patterson and Thompson, 1971). Then we take the predictive density as  $\hat{m}(\tilde{\mathbf{y}}; \mathbf{y}) = m_r(\tilde{\mathbf{y}}|\hat{\sigma}^2)$  and the information (6.3) can be written as

$$I_r(\boldsymbol{\eta}) = E_\eta \left[ (n-p) \log(2\pi\tilde{\sigma}^2) + \log |\boldsymbol{\Sigma}| + \log |\mathbf{X}^T\boldsymbol{\Sigma}^{-1}\mathbf{X}| + \tilde{\mathbf{y}}^T(\boldsymbol{\Sigma}^{-1} - \mathbf{P})\tilde{\mathbf{y}}/\tilde{\sigma}^2 \right], \quad (6.10)$$

where  $\tilde{\sigma}^2 = \mathbf{y}^T(\boldsymbol{\Sigma}^{-1} - \mathbf{P})\mathbf{y}/(n-p)$ , which is the residual maximum likelihood (REML) estimator of  $\sigma^2$  based on the residual likelihood  $m_r(\mathbf{y}|\sigma^2)$ . Then we propose the information criterion.

**Proposition 6.2** *The information  $I_r(\boldsymbol{\eta})$  in (6.10) is unbiasedly estimated by the information criterion*

$$\text{IC}_r = -2 \log\{m_r(\mathbf{y}|\tilde{\sigma}^2)\} + \frac{2(n-p)}{n-p-2}, \quad (6.11)$$

where

$$-2 \log\{m_r(\mathbf{y}|\tilde{\sigma}^2)\} = (n-p) \log \tilde{\sigma}^2 + \log |\boldsymbol{\Sigma}| + \log |\mathbf{X}^T\boldsymbol{\Sigma}^{-1}\mathbf{X}| + \mathbf{y}^T(\boldsymbol{\Sigma}^{-1} - \mathbf{P})\mathbf{y}/\tilde{\sigma}^2 + (\text{const}),$$

namely,  $E_\eta(\text{IC}_r) = I_r(\boldsymbol{\eta})$ .

Note that  $\mathbf{y}^T(\boldsymbol{\Sigma}^{-1} - \mathbf{P})\mathbf{y}/\tilde{\sigma}^2 = n-p$ . If  $n^{-1}\mathbf{X}^T\boldsymbol{\Sigma}^{-1}\mathbf{X}$  converges to  $p \times p$  positive definite matrix as  $n \rightarrow \infty$ ,  $\log |\mathbf{X}^T\boldsymbol{\Sigma}^{-1}\mathbf{X}|$  can be approximated to  $p \log n$ . In that case, we can approximately write

$$\text{IC}_r^* = (n-p) \log \tilde{\sigma}^2 + \log |\boldsymbol{\Sigma}| + p \log n + \frac{(n-p)^2}{n-p-2}, \quad (6.12)$$

when  $n$  is large. It is important to note that  $\text{IC}_r^*$  is identical to the RIC proposed by Shi and Tsai (2002) up to constant. Since  $(n-p)^2/(n-p-2) = (n+2) + \{4/(n-p-2) - p\}$  and  $n+2$  is irrelevant to the model, we can subtract  $n+2$  from  $\text{IC}_r^*$  in (6.12), which results in the RIC exactly. Note the criterion based on  $m_r(\mathbf{y}|\sigma^2)$  and  $r(\boldsymbol{\psi}; m_r)$  cannot be constructed because the KL risk of it diverges to infinity.

### 6.2.2 Extension to the case of unknown covariance

In the derivation of the criteria, we have assumed that the scaled covariance matrix  $\Sigma$  of the error terms vector are known. However, it is often the case that  $\Sigma$  is unknown and is some function of the unknown parameter  $\phi$ , namely  $\Sigma = \Sigma(\phi)$ . In that case,  $\Sigma$  in each criterion is replaced with its plug-in estimator  $\Sigma(\hat{\phi})$ , where  $\hat{\phi}$  is some consistent estimator of  $\phi$ . This strategy is also used in many other studies, for example in Shi and Tsai (2002), who proposed the RIC. We suggest that the  $\phi$  is estimated based on the full model. The method to estimate the nuisance parameters by the full model is similar to the  $C_p$  criterion by Mallows (1973). The scaled covariance matrix  $\mathbf{W}$  of the prior distribution of  $\beta$  is also assumed to be known. In practice, its structure should be specified and we have to estimate the parameters  $\lambda$  involved in  $\mathbf{W}$  from the data. In the same manner as  $\Sigma$ ,  $\mathbf{W}$  in each criterion is replaced with  $\mathbf{W}(\hat{\lambda})$ . We propose that  $\lambda$  is estimated based on each candidate model under consideration because the structure of  $\mathbf{W}$  depends on the model.

We here give three examples for the regression model (6.6), a regression model with constant variance, a variance components model, and a regression model with ARMA errors, where the second and the third ones include the unknown parameter in the covariance matrix.

[1] **regression model with constant variance.** In the case where  $\Sigma = \mathbf{I}_n$ , (6.6) represents a multiple regression model with constant variance. In this model, the scaled covariance matrix  $\Sigma$  does not contain any unknown parameters.

[2] **variance components model.** Consider a variance components model (Henderson, 1950) described by

$$\mathbf{y} = \mathbf{X}\beta + \mathbf{Z}_2\mathbf{b}_2 + \cdots + \mathbf{Z}_r\mathbf{b}_r + \varepsilon, \quad (6.13)$$

where  $\mathbf{Z}_i$  is an  $n \times m_i$  matrix with  $\Sigma_i = \mathbf{Z}_i\mathbf{Z}_i^T$ ,  $\mathbf{b}_i$  is an  $m_i \times 1$  random vector having the distribution  $\mathcal{N}_{m_i}(\mathbf{0}, \theta_i\mathbf{I}_{m_i})$  for  $i \geq 2$ ,  $\varepsilon$  is an  $n \times 1$  random vector with  $\varepsilon \sim \mathcal{N}_n(\mathbf{0}, \Sigma_0 + \theta_1\Sigma_1)$  for known  $n \times n$  matrices  $\Sigma_0$  and  $\Sigma_1$ , and  $\mathbf{e}$ ,  $\mathbf{b}_2, \dots, \mathbf{b}_r$  are mutually independently distributed. The nested error regression model (NERM) is a special case of variance components model given by

$$y_{ik} = \mathbf{x}_{ik}^T\beta + b_i + \varepsilon_{ik}, \quad (i = 1, \dots, m; k = 1, \dots, n_i), \quad (6.14)$$

where  $b_i$ 's and  $\varepsilon_{ik}$ 's are mutually independently distributed as  $b_i \sim \mathcal{N}(0, \tau^2)$  and  $\varepsilon_{ik} \sim \mathcal{N}(0, \sigma^2)$  and  $n = \sum_{i=1}^m n_i$ . Note that the NERM in (6.14) is given by  $\theta_1 = \sigma^2$ ,  $\theta_2 = \tau^2$ ,  $\Sigma_1 = \mathbf{I}_n$  and  $\mathbf{Z}_2 = \text{diag}(\mathbf{1}_{n_1}, \dots, \mathbf{1}_{n_m})$ , where  $\mathbf{1}_l$  is the  $l$ -dimensional vector of ones, for variance components model (6.13). This model is often used for the clustered data and  $b_i$  can be seen as the random effect of the cluster (Battese et al., 1988). For such a model, when one is interested in the specific cluster or predicting the random effects, the conditional AIC proposed by Vaida and Blanchard (2005), which is based on the conditional likelihood given the random effects, is appropriate. However, when the aim of the analysis is focused on the population, the NERM can be seen as linear regression model and the random effects are involved in the error term, namely we can treat  $\mathbf{u} = \mathbf{Z}_2\mathbf{b}_2 + \varepsilon$ ,  $\Sigma = \Sigma(\phi) = \phi\Sigma_2 + \mathbf{I}_n$  for (6.6), where  $\phi = \tau^2/\sigma^2$  and  $\Sigma_2 = \mathbf{Z}_2\mathbf{Z}_2^T = \text{diag}(\mathbf{J}_{n_1}, \dots, \mathbf{J}_{n_m})$  for  $\mathbf{J}_l = \mathbf{1}_l\mathbf{1}_l^T$ . In that case, our proposed variable selection procedure is valid.

[3] **regression model with autoregressive moving average errors.** Consider the regression model (6.6), assuming the random errors are generated by an ARMA( $q, r$ ) process defined by

$$u_i - \phi_1 u_{i-1} - \cdots - \phi_q u_{i-q} = \varepsilon_i - \varphi_1 \varepsilon_{i-1} - \cdots - \varphi_r \varepsilon_{i-r},$$

where  $\{\varepsilon_i\}$  is a sequence of independent normal random variables having mean 0 and variance  $\tau^2$ . A special case of this model is the regression model with AR(1) errors satisfying  $u_1 \sim \mathcal{N}(0, \tau^2/(1 - \phi^2))$ ,  $u_i = \phi u_{i-1} + \varepsilon_i$ ,  $\varepsilon_i \sim \mathcal{N}(0, \tau^2)$  for  $i = 2, 3, \dots, n$ . When we define  $\sigma^2 = \tau^2/(1 - \phi^2)$ ,  $(i, j)$ -element of the scaled covariance matrix  $\Sigma$  in (6.6) is  $\phi^{|i-j|}$ .

### 6.3 Consistency of the criteria

In this section, we prove that the proposed criteria have the consistency property. Our asymptotic framework is that  $n$  goes to infinity and the true dimension of the regression coefficients  $p$  is fixed. Following Shi and Tsai (2002), we first show the criteria are consistent for the regression model with constant variance and the prespecified  $\mathbf{W}$ , and then extend the result to the regression model with general covariance matrix and the case where  $\mathbf{W}$  is estimated.

Let  $\hat{j}$  denote the model selected by some criterion. Following Shi and Tsai (2002), we make the assumptions.

$$(A1) \ E(u_1^4) < \infty.$$

$$(A2) \ 0 < \liminf_{n \rightarrow \infty} \min_{j \in \mathcal{J}} |\mathbf{X}(j)^T \mathbf{X}(j)/n| \text{ and } \limsup_{n \rightarrow \infty} \max_{j \in \mathcal{J}} |\mathbf{X}(j)^T \mathbf{X}(j)/n| < \infty.$$

$$(A3) \ \liminf_{n \rightarrow \infty} n^{-1} \inf_{j \in \mathcal{J}_-} \|\mathbf{X}(\omega)\beta_* - \mathbf{H}_j \mathbf{X}(\omega)\beta_*\|^2 > 0, \text{ where } \mathbf{H}_j = \mathbf{X}(j)(\mathbf{X}(j)^T \mathbf{X}(j))^{-1} \mathbf{X}(j)^T.$$

We can now obtain asymptotic properties of the criteria for the regression model with constant variance.

**Theorem 6.1** *If assumptions (A1)–(A3) are satisfied,  $\mathcal{J}_+$  is not empty, the  $u_i$ 's are independent and identically distributed (iid) and  $\mathbf{W}_j$  in the prior distribution of  $\beta_j$  is prespecified, then the criteria  $\text{IC}_{\pi,1}$ ,  $\text{IC}_{\pi,1}^*$ ,  $\text{IC}_{\pi,2}$ ,  $\text{IC}_r$  and  $\text{IC}_r^*$  are consistent, namely  $P(\hat{j} = j_*) \rightarrow 1$  as  $n \rightarrow \infty$ .*

The proof of Theorem 6.1 is given in Section 6.7.

We next consider the regression model with a general covariance structure and the case where  $\mathbf{W}_j$  is estimated by the data. In this case,  $\Sigma$  and  $\mathbf{W}_j$  are replaced with their plug-in estimators  $\Sigma(\hat{\phi})$  and  $\mathbf{W}_j(\hat{\lambda}_j)$ , respectively.

**Theorem 6.2** *Assume that  $\hat{\phi} - \phi_0$  and  $\hat{\lambda}_j - \lambda_{j,0}$  tend to 0 in probability as  $n \rightarrow \infty$  for all  $j \in \mathcal{J}$ . In addition, assume that the elements of  $\Sigma(\phi)$  and  $\mathbf{W}_j(\lambda_j)$  are continuous functions of  $\phi$  and  $\lambda_j$ , and  $\Sigma(\phi)$  and  $\mathbf{W}_j(\lambda_j)$  is positive definite in the neighborhood of  $\phi_0$  and  $\lambda_{j,0}$  for all  $j \in \mathcal{J}$ . If assumptions (A1)–(A3) are satisfied when  $\mathbf{X}(j)$  and  $\mathbf{u}$  are replaced with  $\Sigma^{-1/2} \mathbf{X}(j)$  and  $\mathbf{u}^* = \Sigma^{-1/2} \mathbf{u}$  respectively,  $\mathcal{J}_+$  is not empty and the  $u_i^*$ 's are iid, then the criteria  $\text{IC}_{\pi,1}$ ,  $\text{IC}_{\pi,1}^*$ ,  $\text{IC}_{\pi,2}$ ,  $\text{IC}_r$  and  $\text{IC}_r^*$  are consistent.*

For the proof of Theorem 6.2, we can use the same techniques as those for the proof of Theorem 6.1.

### 6.4 Simulations

In this section, we compare the numerical performance of the proposed criteria with some other conventional ones, which are AIC, BIC, the corrected AIC (AICC) by Sugiura (1978) and Hurvich and Tsai (1989). We shall consider the three regression models—regression model with constant variance, NERM, and regression model with AR(1) errors—which are taken as examples of (6.6) in Section 6.2.2. For the NERM, we consider the balanced sample case, namely

$n_1 = \dots = n_m (= n_0)$ . In each simulation, 1000 realizations are generated from (6.6) with  $\boldsymbol{\beta} = (1, 1, 1, 1, 0, 0, 0)^T$ , namely the full model is seven-dimensional and the true model is four-dimensional. All explanatory variables are randomly generated from the standard normal distribution. The signal-to-noise ratio ( $\text{SNR} = \{\text{var}(\mathbf{x}_i^T \boldsymbol{\beta}) / \text{var}(u_i)\}^{1/2}$ ) is controlled at 1, 3, and 5. In the NERM, three cases of variance ratio  $\phi = \tau^2 / \sigma^2$  are considered with  $\phi = 0.5, 1$  and 2. In the regression model with AR(1) errors, three correlation structures are considered with AR parameter  $\phi = 0.1, 0.5$  and 0.8.

When deriving the criteria  $\text{IC}_{\pi,1}$ ,  $\text{IC}_{\pi,1}^*$  and  $\text{IC}_{\pi,2}$ , we set the prior distribution of  $\boldsymbol{\beta}$  as  $\mathcal{N}_p(\mathbf{0}, \sigma^2 \lambda^{-1} \mathbf{I}_p)$ , namely  $\mathbf{W} = \lambda^{-1} \mathbf{I}_p$ . The hyperparameter  $\lambda$  is estimated by maximizing the marginal likelihood  $m_\pi(\mathbf{y} | \hat{\sigma}^2)$ , where the estimate  $\hat{\sigma}^2 = \mathbf{y}^T (\boldsymbol{\Sigma}^{-1} - \mathbf{P}) \mathbf{y} / n$  of  $\sigma^2$  is plugged in. The unknown parameter  $\phi$  involved in  $\boldsymbol{\Sigma}$  is estimated by some consistent estimator based on the full model. In the NERM,  $\phi = \tau^2 / \sigma^2$  is estimated by  $\hat{\tau}^{2\text{PR}} / \hat{\sigma}^{2\text{PR}}$ , where  $\hat{\tau}^{2\text{PR}}$  and  $\hat{\sigma}^{2\text{PR}}$  are unbiased estimators proposed by Prasad and Rao (1990). Let  $S_0 = \mathbf{y}^T \{\mathbf{I}_n - \mathbf{X}(\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T\} \mathbf{y}$  and  $S_1 = \mathbf{y}^T \{\mathbf{E} - \mathbf{E} \mathbf{X}(\mathbf{X}^T \mathbf{E} \mathbf{X})^{-1} \mathbf{X}^T \mathbf{E}\} \mathbf{y}$  where  $\mathbf{E} = \text{diag}(\mathbf{E}_1, \dots, \mathbf{E}_m)$ ,  $\mathbf{E}_i = \mathbf{I}_{n_0} - n_0^{-1} \mathbf{J}_{n_0}$  for  $i = 1, \dots, m$ . Then, the Prasad–Rao estimators of  $\sigma^2$  and  $\tau^2$  are

$$\hat{\sigma}^{2\text{PR}} = S_1 / (n - m - p), \quad \hat{\tau}^{2\text{PR}} = \{S_0 - (n - p) \hat{\sigma}^{2\text{PR}}\} / n^*$$

where  $n^* = n - \text{tr}[\mathbf{Z}_2^T \mathbf{X}(\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{Z}_2]$ . In the regression model with AR(1) errors, the AR parameter  $\phi$  is estimated by the maximum likelihood estimator based. Note that  $\phi$  is estimated based on the full model and that  $\sigma^2$  and  $\lambda$  is estimated based on each candidate model using the plug-in version of  $\boldsymbol{\Sigma}(\hat{\phi})$ .

The candidate models include all the subsets of the full model and select the model by the criteria. The performance of the criteria is measured by the number of selecting the true model and the prediction error of the selected model based on quadratic loss, namely  $\|\mathbf{X}(j) \hat{\boldsymbol{\beta}}_j - \mathbf{X}(\omega) \boldsymbol{\beta}_*\|^2 / n$ .

Tables 6.1–6.3 give the number of selecting the true model by the criteria and the average prediction error of the selected model by each criterion is shown in Tables 6.4–6.6 for each of the regression models. From these tables, we can see the following three facts. Firstly, the number of selecting the true model approaches 1000 for all the proposed criteria, that is the numerical evidence of the consistency of the criteria. Though the BIC is also consistent, the small sample performance is not as good as our criteria. Secondly, the proposed criteria are not only consistent but also have smaller prediction error even when the sample size is small. Especially,  $\text{IC}_{\pi,1}$  is the best for the most of the experiments except when both the sample size and SNR are small. AIC and AICC have good performance in that situation in terms of prediction error. Thirdly,  $\text{IC}_{\pi,1}$  and  $\text{IC}_r$  have better performance than their approximation  $\text{IC}_{\pi,1}^*$  and  $\text{IC}_r^*$ , respectively, but the difference gets smaller as  $n$  becomes larger.

Table 6.1: The number of selecting the true model by the criteria in 1000 realizations of the regression model with constant variance.

		SNR = 1	SNR = 3	SNR = 5
$n = 20$	AIC	130	428	428
	BIC	118	587	588
	AICC	89	749	755
	$IC_{\pi,1}$	115	843	905
	$IC_{\pi,1}^*$	73	732	738
	$IC_{\pi,2}$	73	731	737
	$IC_r$	143	797	882
	$IC_r^*$	147	828	898
$n = 40$	AIC	419	536	536
	BIC	424	800	800
	AICC	470	687	687
	$IC_{\pi,1}$	472	900	938
	$IC_{\pi,1}^*$	353	876	876
	$IC_{\pi,2}$	352	876	876
	$IC_r$	462	895	934
	$IC_r^*$	478	899	941
$n = 80$	AIC	546	553	553
	BIC	827	872	872
	AICC	604	613	613
	$IC_{\pi,1}$	750	934	968
	$IC_{\pi,1}^*$	839	928	928
	$IC_{\pi,2}$	838	928	928
	$IC_r$	722	937	968
	$IC_r^*$	739	936	969

Table 6.2: The number of selecting the true model by the criteria in 1000 realizations of the nested error regression model.

		$\phi = 0.5$			$\phi = 1$			$\phi = 2$		
SNR		1	3	5	1	3	5	1	3	5
$n_0 = 5$	AIC	75	382	385	104	387	393	137	393	403
$m = 4$	BIC	60	546	564	90	539	574	134	556	586
	AICC	57	694	736	86	693	742	140	691	748
	$IC_{\pi,1}$	71	821	902	119	829	915	181	835	941
	$IC_{\pi,1}^*$	35	646	702	66	656	715	115	662	721
	$IC_{\pi,2}$	34	642	701	70	653	715	116	657	720
	$IC_r$	244	789	888	310	818	900	432	838	924
	$IC_r^*$	78	723	912	106	723	922	149	698	941
$n_0 = 5$	AIC	220	458	458	235	465	465	259	473	473
$m = 8$	BIC	169	731	731	208	739	741	240	746	750
	AICC	219	607	607	251	612	612	284	625	627
	$IC_{\pi,1}$	319	891	936	369	913	943	436	928	953
	$IC_{\pi,1}^*$	107	838	839	151	836	841	199	843	853
	$IC_{\pi,2}$	112	838	839	158	836	841	202	841	853
	$IC_r$	436	890	934	512	903	943	593	925	953
	$IC_r^*$	209	901	944	230	901	949	247	905	962
$n_0 = 5$	AIC	418	522	522	417	528	528	418	545	545
$m = 16$	BIC	394	853	853	407	859	859	416	866	866
	AICC	452	594	594	447	603	603	443	616	616
	$IC_{\pi,1}$	622	926	955	618	941	959	624	951	968
	$IC_{\pi,1}^*$	299	911	910	332	913	913	354	915	915
	$IC_{\pi,2}$	300	910	910	331	913	913	358	915	915
	$IC_r$	691	925	954	695	942	960	709	953	968
	$IC_r^*$	443	932	959	438	946	964	421	956	972



Table 6.3: The number of selecting the true model by the criteria in 1000 realizations of the regression model with AR(1) errors.

		$\phi = 0.1$			$\phi = 0.5$			$\phi = 0.8$		
SNR		1	3	5	1	3	5	1	3	5
$n = 20$	AIC	110	347	346	125	373	372	193	377	414
	BIC	91	482	482	118	523	542	209	502	587
	AICC	84	642	646	117	652	688	228	584	715
	$IC_{\pi,1}$	101	741	834	138	774	862	274	742	901
	$IC_{\pi,1}^*$	64	620	625	83	625	667	194	554	688
	$IC_{\pi,2}$	63	618	624	88	621	667	197	553	685
	$IC_r$	123	698	801	224	769	846	562	808	901
	$IC_r^*$	122	702	790	144	715	826	241	646	825
$n = 40$	AIC	365	483	483	356	544	544	283	533	551
	BIC	369	756	756	334	791	793	286	737	797
	AICC	416	642	642	373	671	672	308	668	700
	$IC_{\pi,1}$	422	877	917	450	909	949	427	901	976
	$IC_{\pi,1}^*$	315	844	844	281	859	862	247	754	856
	$IC_{\pi,2}$	314	844	844	282	858	862	255	752	854
	$IC_r$	430	866	917	507	903	945	685	918	974
	$IC_r^*$	412	865	918	377	881	932	316	785	932
$n = 80$	AIC	516	521	521	483	552	552	333	553	553
	BIC	789	851	851	598	865	865	334	859	868
	AICC	586	593	593	525	614	614	367	637	638
	$IC_{\pi,1}$	738	926	949	691	936	962	550	961	979
	$IC_{\pi,1}^*$	795	912	912	560	905	905	289	899	919
	$IC_{\pi,2}$	792	912	912	559	905	905	296	889	919
	$IC_r$	713	921	951	724	938	961	692	966	980
	$IC_r^*$	714	920	951	600	911	952	382	908	958

Table 6.4: The prediction error of the best model selected by the criteria for the regression model with constant variance.

	SNR	1	3	5
$n = 20$	AIC	1.59	0.141	0.0504
	BIC	1.74	0.131	0.0468
	AICC	1.77	0.120	0.0421
	$IC_{\pi,1}$	1.70	0.111	0.0372
	$IC_{\pi,1}^*$	1.92	0.122	0.0429
	$IC_{\pi,2}$	1.92	0.122	0.0430
	$IC_r$	1.56	0.116	0.0383
	$IC_r^*$	1.57	0.114	0.0374
$n = 40$	AIC	0.708	0.0660	0.0238
	BIC	0.862	0.0568	0.0205
	AICC	0.732	0.0609	0.0219
	$IC_{\pi,1}$	0.754	0.0523	0.0180
	$IC_{\pi,1}^*$	1.05	0.0534	0.0192
	$IC_{\pi,2}$	1.05	0.0534	0.0192
	$IC_r$	0.716	0.0524	0.0181
	$IC_r^*$	0.718	0.0522	0.0179
$n = 80$	AIC	0.292	0.0321	0.0115
	BIC	0.265	0.0260	0.00936
	AICC	0.283	0.0310	0.0112
	$IC_{\pi,1}$	0.265	0.0244	0.00841
	$IC_{\pi,1}^*$	0.285	0.0245	0.00883
	$IC_{\pi,2}$	0.285	0.0245	0.00883
	$IC_r$	0.270	0.0243	0.00841
	$IC_r^*$	0.267	0.0243	0.00840

Table 6.5: The prediction error of the best model selected by the criteria for the nested error regression model.

		$\phi = 0.5$			$\phi = 1$			$\phi = 2$		
	SNR	1	3	5	1	3	5	1	3	5
$n_0 = 5$	AIC	1.80	0.150	0.0524	1.61	0.145	0.0494	1.44	0.140	0.0463
$m = 4$	BIC	1.96	0.159	0.0496	1.76	0.159	0.0473	1.50	0.158	0.0447
	AICC	2.00	0.172	0.0458	1.78	0.174	0.0443	1.53	0.179	0.0428
	$IC_{\pi,1}$	1.93	0.151	0.0417	1.72	0.156	0.0410	1.47	0.159	0.0400
	$IC_{\pi,1}^*$	2.14	0.190	0.0470	1.88	0.188	0.0452	1.60	0.186	0.0434
	$IC_{\pi,2}$	2.14	0.192	0.0470	1.87	0.190	0.0452	1.59	0.186	0.0434
	$IC_r$	1.48	0.135	0.0422	1.37	0.137	0.0413	1.21	0.139	0.0404
	$IC_r^*$	1.91	0.251	0.0413	1.72	0.271	0.0434	1.52	0.316	0.0471
$n_0 = 5$	AIC	0.983	0.0696	0.0251	0.907	0.0659	0.0237	0.824	0.0619	0.0223
$m = 8$	BIC	1.25	0.0622	0.0224	1.11	0.0615	0.0216	1.01	0.0613	0.0209
	AICC	1.10	0.0655	0.0236	0.989	0.0627	0.0226	0.899	0.0610	0.0215
	$IC_{\pi,1}$	0.989	0.0567	0.0197	0.888	0.0554	0.0196	0.804	0.0565	0.0194
	$IC_{\pi,1}^*$	1.41	0.0594	0.0211	1.21	0.0613	0.0207	1.07	0.0635	0.0202
	$IC_{\pi,2}$	1.40	0.0594	0.0211	1.20	0.0613	0.0207	1.07	0.0652	0.0202
	$IC_r$	0.743	0.0568	0.0197	0.666	0.0557	0.0196	0.602	0.0565	0.0194
	$IC_r^*$	1.16	0.0609	0.0196	1.07	0.0691	0.0195	1.01	0.0841	0.0193
$n_0 = 5$	AIC	0.451	0.0341	0.0123	0.440	0.0325	0.0117	0.434	0.0308	0.0111
$m = 16$	BIC	0.740	0.0294	0.0106	0.711	0.0289	0.0104	0.684	0.0284	0.0102
	AICC	0.489	0.0333	0.0120	0.483	0.0319	0.0115	0.481	0.0304	0.0109
	$IC_{\pi,1}$	0.433	0.0280	0.00980	0.435	0.0277	0.00983	0.438	0.0275	0.00980
	$IC_{\pi,1}^*$	0.864	0.0283	0.0102	0.812	0.0281	0.0101	0.770	0.0279	0.0100
	$IC_{\pi,2}$	0.862	0.0283	0.0102	0.811	0.0281	0.0101	0.766	0.0279	0.0100
	$IC_r$	0.348	0.0280	0.00981	0.346	0.0277	0.00982	0.344	0.0274	0.00980
	$IC_r^*$	0.658	0.0278	0.00977	0.664	0.0276	0.00979	0.683	0.0281	0.00978

Table 6.6: The prediction error of the best model selected by the criteria for the regression model with AR(1) errors.

	SNR	$\phi = 0.1$			$\phi = 0.5$			$\phi = 0.8$		
		1	3	5	1	3	5	1	3	5
$n = 20$	AIC	1.85	0.175	0.0628	1.71	0.173	0.0613	1.75	0.256	0.0759
	BIC	2.01	0.170	0.0595	1.82	0.192	0.0583	1.73	0.305	0.0802
	AICC	2.04	0.163	0.0549	1.85	0.203	0.0577	1.69	0.343	0.0909
	$IC_{\pi,1}$	1.96	0.153	0.0492	1.78	0.182	0.0538	1.66	0.312	0.0831
	$IC_{\pi,1}^*$	2.17	0.164	0.0559	1.95	0.211	0.0566	1.71	0.355	0.0962
	$IC_{\pi,2}$	2.17	0.164	0.0559	1.95	0.216	0.0566	1.72	0.354	0.0972
	$IC_r$	1.79	0.154	0.0505	1.59	0.157	0.0516	1.71	0.224	0.0720
	$IC_r^*$	1.84	0.159	0.0507	1.72	0.212	0.0557	1.77	0.361	0.114
$n = 40$	AIC	0.783	0.0733	0.0264	0.812	0.0685	0.0246	1.09	0.124	0.0365
	BIC	0.973	0.0639	0.0230	0.992	0.0640	0.0225	1.13	0.162	0.0408
	AICC	0.824	0.0681	0.0245	0.881	0.0662	0.0236	1.12	0.135	0.0360
	$IC_{\pi,1}$	0.839	0.0587	0.0204	0.832	0.0595	0.0207	1.07	0.136	0.0347
	$IC_{\pi,1}^*$	1.15	0.0601	0.0216	1.10	0.0638	0.0218	1.14	0.202	0.0432
	$IC_{\pi,2}$	1.15	0.0601	0.0216	1.10	0.0639	0.0218	1.15	0.202	0.0441
	$IC_r$	0.782	0.0592	0.0203	0.711	0.0598	0.0208	0.919	0.118	0.0347
	$IC_r^*$	0.813	0.0592	0.0203	0.857	0.0631	0.0209	1.12	0.195	0.0446
$n = 80$	AIC	0.311	0.0342	0.0123	0.365	0.0329	0.0118	0.666	0.0520	0.0187
	BIC	0.301	0.0282	0.0102	0.500	0.0290	0.0105	0.861	0.0613	0.0182
	AICC	0.303	0.0331	0.0119	0.377	0.0322	0.0116	0.693	0.0526	0.0186
	$IC_{\pi,1}$	0.286	0.0262	0.00920	0.365	0.0279	0.00983	0.643	0.0530	0.0179
	$IC_{\pi,1}^*$	0.331	0.0266	0.00958	0.566	0.0284	0.0102	0.914	0.0692	0.0181
	$IC_{\pi,2}$	0.333	0.0266	0.00958	0.566	0.0284	0.0102	0.910	0.0692	0.0181
	$IC_r$	0.290	0.0264	0.00918	0.330	0.0278	0.00984	0.514	0.0499	0.0179
	$IC_r^*$	0.292	0.0264	0.00919	0.414	0.0283	0.00991	0.778	0.0662	0.0180

## 6.5 Discussion

We have derived the variable selection criteria for linear regression model relative to the frequentist KL risk of the predictive density based on the Bayesian marginal likelihood. We have proved the consistency of the criteria and have showed that they perform well also in the sense of the prediction through simulations.

We gave some advantages of the approach based on frequentist's risk  $R(\boldsymbol{\eta}; \hat{m})$  in (6.1). We here explain them more clearly through comparison of the related Bayesian criteria. When the prior distribution  $\pi(\boldsymbol{\beta}|\boldsymbol{\lambda}, \boldsymbol{\omega})$  is proper, we can treat the Bayesian prediction risk

$$r(\boldsymbol{\psi}; \hat{m}) = \int R(\boldsymbol{\eta}; \hat{m})\pi(\boldsymbol{\beta}|\boldsymbol{\lambda}, \boldsymbol{\omega})d\boldsymbol{\beta}$$

in (6.5). When  $\boldsymbol{\lambda}$  and  $\boldsymbol{\omega}$  are known, the predictive density  $\hat{m}(\tilde{\mathbf{y}}; \mathbf{y})$  which minimizes  $r(\boldsymbol{\psi}; \hat{m})$  is the Bayesian predictive density (posterior predictive density)  $\hat{m}_\pi(\tilde{\mathbf{y}}|\mathbf{y}, \boldsymbol{\lambda}, \boldsymbol{\omega})$  given as

$$\int m(\tilde{\mathbf{y}}|\boldsymbol{\beta}, \boldsymbol{\omega})\pi(\boldsymbol{\beta}|\mathbf{y}, \boldsymbol{\lambda}, \boldsymbol{\omega})d\boldsymbol{\beta} = \frac{\int m(\tilde{\mathbf{y}}|\boldsymbol{\beta}, \boldsymbol{\omega})m(\mathbf{y}|\boldsymbol{\beta}, \boldsymbol{\omega})\pi(\boldsymbol{\beta}|\boldsymbol{\lambda}, \boldsymbol{\omega})d\boldsymbol{\beta}}{\int m(\mathbf{y}|\boldsymbol{\beta}, \boldsymbol{\omega})\pi(\boldsymbol{\beta}|\boldsymbol{\lambda}, \boldsymbol{\omega})d\boldsymbol{\beta}}.$$

When  $\boldsymbol{\lambda}$  and  $\boldsymbol{\omega}$  are unknown, we can consider the Bayesian risk of the plug-in predictive density  $\hat{m}_\pi(\tilde{\mathbf{y}}|\mathbf{y}, \hat{\boldsymbol{\lambda}}, \hat{\boldsymbol{\omega}})$ . Then the resulting criterion is known as the predictive likelihood (Akaike, 1980a) or the PIC (Kitagawa, 1997). The deviance information criterion (DIC) of Spiegelhalter et al. (2002) and the Bayesian predictive information criterion (BPIC) of Ando (2007) are related criteria based on the Bayesian prediction risk  $r(\boldsymbol{\psi}; \hat{m})$ .

The Akaike's Bayesian information criterion (ABIC) (Akaike, 1980b) is another information criterion based on the Bayesian marginal likelihood, given by

$$\text{ABIC} = -2 \log\{m_\pi(\mathbf{y}|\hat{\boldsymbol{\lambda}})\} + 2 \dim(\boldsymbol{\lambda}),$$

where the nuisance parameter  $\boldsymbol{\omega}$  is not considered. The ABIC measures the following KL risk:

$$\int \left[ \int \log \left\{ \frac{m_\pi(\tilde{\mathbf{y}}|\boldsymbol{\lambda})}{m_\pi(\tilde{\mathbf{y}}|\hat{\boldsymbol{\lambda}})} \right\} m_\pi(\tilde{\mathbf{y}}|\boldsymbol{\lambda})d\tilde{\mathbf{y}} \right] m_\pi(\mathbf{y}|\boldsymbol{\lambda})d\mathbf{y},$$

which is not the same as either  $R(\boldsymbol{\eta}; \hat{m})$  or  $r(\boldsymbol{\psi}; \hat{m})$ . The ABIC is the criterion for choosing the hyperparameter  $\boldsymbol{\lambda}$  in the same sense as the AIC. However, it is noted that the ABIC works as a model selection criterion for  $\boldsymbol{\beta}$  because it is based on the Bayesian marginal likelihood.

A drawback of such Bayesian criteria is that we cannot construct them for improper prior distributions  $\pi(\boldsymbol{\beta}|\boldsymbol{\lambda}, \boldsymbol{\omega})$ , since the corresponding Bayesian prediction risks do not exist. On the other hand, we can construct the corresponding criteria based on  $R(\boldsymbol{\eta}; \hat{m})$ , because the approach suggested in this paper measures the prediction risk in the framework of frequentists. In fact, putting the uniform improper prior on regression coefficients  $\boldsymbol{\beta}$  in the linear regression model, we get the RIC of Shi and Tsai (2002). Note that the criteria based on improper marginal likelihood works as variable selection only when the marginal likelihood itself does. For the case where the improper priors cannot be used for model selection, intrinsic prior was proposed in the literature (Berger and Pericchi, 1996; Casella and Moreno, 2006, and others), which is an objective and automatic procedure. As future work, it is worthwhile to consider such an automatic procedure in the framework of our proposed criteria.

## 6.6 Derivations of the criteria

### 6.6.1 Derivation of $IC_{\pi,1}$ in (6.8)

It is sufficient to show that the bias correction  $\Delta_{\pi,1} = I_{\pi,1}(\boldsymbol{\eta}) - E_{\boldsymbol{\eta}}[-2 \log\{m_{\pi}(\mathbf{y}|\hat{\sigma}^2)\}]$  is  $2n/(n-p-2)$ , where  $I_{\pi,1}(\boldsymbol{\eta})$  is given by (6.7). It follows that

$$\begin{aligned}\Delta_{\pi,1} &= E_{\boldsymbol{\eta}}(\tilde{\mathbf{y}}^T \mathbf{A} \tilde{\mathbf{y}} / \hat{\sigma}^2) - E_{\boldsymbol{\eta}}(\mathbf{y}^T \mathbf{A} \mathbf{y} / \hat{\sigma}^2) \\ &= E_{\boldsymbol{\eta}}(\tilde{\mathbf{y}}^T \mathbf{A} \tilde{\mathbf{y}}) \cdot E_{\boldsymbol{\eta}}(1/\hat{\sigma}^2) - E_{\boldsymbol{\eta}}(\mathbf{y}^T \mathbf{A} \mathbf{y} / \hat{\sigma}^2).\end{aligned}$$

Firstly,

$$\begin{aligned}E_{\boldsymbol{\eta}}(\tilde{\mathbf{y}}^T \mathbf{A} \tilde{\mathbf{y}}) &= E_{\boldsymbol{\eta}}[(\tilde{\mathbf{y}} - \mathbf{X}\boldsymbol{\beta} + \mathbf{X}\boldsymbol{\beta})^T \mathbf{A} (\tilde{\mathbf{y}} - \mathbf{X}\boldsymbol{\beta} + \mathbf{X}\boldsymbol{\beta})] \\ &= \sigma^2 \text{tr}(\mathbf{A}\mathbf{V}) + \boldsymbol{\beta}^T \mathbf{X}^T \mathbf{A} \mathbf{X} \boldsymbol{\beta}.\end{aligned}\tag{6.15}$$

Secondly, noting that  $n\hat{\sigma}^2 = \mathbf{y}^T(\boldsymbol{\Sigma}^{-1} - \mathbf{P})\mathbf{y} = \sigma^2 \mathbf{v}^T(\mathbf{I}_n - \mathbf{M})\mathbf{v}$  for

$$\begin{aligned}\mathbf{v} &= \boldsymbol{\Sigma}^{-1/2}(\mathbf{y} - \mathbf{X}\boldsymbol{\beta})/\sigma, \\ \mathbf{M} &= \mathbf{I}_n - \boldsymbol{\Sigma}^{-1/2} \mathbf{X}(\mathbf{X}^T \boldsymbol{\Sigma}^{-1} \mathbf{X})^{-1} \mathbf{X}^T \boldsymbol{\Sigma}^{-1/2},\end{aligned}\tag{6.16}$$

and that  $\mathbf{v}^T(\mathbf{I}_n - \mathbf{M})\mathbf{v} \sim \chi_{n-p}^2$ , we can obtain

$$\begin{aligned}E_{\boldsymbol{\eta}}(1/\hat{\sigma}^2) &= n E_{\boldsymbol{\eta}} \left[ \frac{1}{\mathbf{y}^T(\boldsymbol{\Sigma}^{-1} - \mathbf{P})\mathbf{y}} \right] = n E_{\boldsymbol{\eta}} \left[ \frac{1}{\sigma^2 \mathbf{v}^T(\mathbf{I}_n - \mathbf{M})\mathbf{v}} \right] \\ &= \frac{n}{\sigma^2(n-p-2)}.\end{aligned}\tag{6.17}$$

Finally,

$$\begin{aligned}E_{\boldsymbol{\eta}}(\mathbf{y}^T \mathbf{A} \mathbf{y} / \hat{\sigma}^2) &= n E_{\boldsymbol{\eta}} \left[ \frac{\mathbf{y}^T \mathbf{A} \mathbf{y}}{\mathbf{y}^T(\boldsymbol{\Sigma}^{-1} - \mathbf{P})\mathbf{y}} \right] = n E_{\boldsymbol{\eta}} \left[ \frac{\sigma^2 \mathbf{v}^T \boldsymbol{\Sigma}^{1/2} \mathbf{A} \boldsymbol{\Sigma}^{1/2} \mathbf{v} + \boldsymbol{\beta}^T \mathbf{X}^T \mathbf{A} \mathbf{X} \boldsymbol{\beta}}{\sigma^2 \mathbf{v}^T(\mathbf{I}_n - \mathbf{M})\mathbf{v}} \right] \\ &= n \times \left\{ \frac{\text{tr}(\mathbf{A}\boldsymbol{\Sigma})}{n-p-2} - \frac{2\text{tr}[\mathbf{A}\boldsymbol{\Sigma}(\boldsymbol{\Sigma}^{-1} - \mathbf{P})\boldsymbol{\Sigma}]}{(n-p)(n-p-2)} + \frac{\boldsymbol{\beta}^T \mathbf{X}^T \mathbf{A} \mathbf{X} \boldsymbol{\beta}}{\sigma^2(n-p-2)} \right\}.\end{aligned}\tag{6.18}$$

The last equation in the above can be derived by Lemma 4.6. Combining (6.15), (6.17) and (6.18), we get

$$\Delta_{\pi,1} = \frac{2n \cdot \text{tr}[\mathbf{A}\boldsymbol{\Sigma}(\boldsymbol{\Sigma}^{-1} - \mathbf{P})\boldsymbol{\Sigma}]}{(n-p)(n-p-2)}.$$

We can see that

$$\begin{aligned}\text{tr}[\mathbf{A}\boldsymbol{\Sigma}(\boldsymbol{\Sigma}^{-1} - \mathbf{P})\boldsymbol{\Sigma}] &= \text{tr}\{(\boldsymbol{\Sigma} + \mathbf{B})^{-1}(\boldsymbol{\Sigma} + \mathbf{B} - \mathbf{B})(\boldsymbol{\Sigma}^{-1} - \mathbf{P})\boldsymbol{\Sigma}\} \\ &= \text{tr}[(\boldsymbol{\Sigma}^{-1} - \mathbf{P})\boldsymbol{\Sigma}] - \text{tr}\{(\boldsymbol{\Sigma} + \mathbf{B})^{-1} \mathbf{B}(\boldsymbol{\Sigma}^{-1} - \mathbf{P})\boldsymbol{\Sigma}\} \\ &= \text{tr}(\mathbf{I}_n - \mathbf{M}) = n-p,\end{aligned}\tag{6.19}$$

since  $\mathbf{B}(\boldsymbol{\Sigma}^{-1} - \mathbf{P}) = \mathbf{X}\mathbf{W}\mathbf{X}^T(\boldsymbol{\Sigma}^{-1} - \mathbf{P}) = \mathbf{0}$ , then we obtain  $\Delta_{\pi,1} = 2n/(n-p-2)$ .  $\square$

### 6.6.2 Derivation of $IC_{\pi,2}$ in (6.9)

From the fact that  $E_{\boldsymbol{\eta}}(IC_{\pi,1}) = I_{\pi,1}(\boldsymbol{\eta})$  and that  $E_{\pi}E_{\boldsymbol{\eta}}(IC_{\pi,1}) = E_{\pi}[I_{\pi,1}(\boldsymbol{\eta})] = I_{\pi,2}(\sigma^2)$ , it suffices to show that  $E_{\pi}E_{\boldsymbol{\eta}}(IC_{\pi,1})$  is approximated to

$$\begin{aligned} E_{\pi}E_{\boldsymbol{\eta}}(IC_{\pi,1}) &\approx E_{\pi}E_{\boldsymbol{\eta}}[n \log \hat{\sigma}^2 + \log |\boldsymbol{\Sigma}| + p \log n + 2 + E_{\pi}E_{\boldsymbol{\eta}}(\mathbf{y}^T \mathbf{A} \mathbf{y} / \hat{\sigma}^2)] \\ &\approx E_{\pi}E_{\boldsymbol{\eta}}[n \log \hat{\sigma}^2 + \log |\boldsymbol{\Sigma}| + p \log n + p] + (n + 2) = E_{\pi}E_{\boldsymbol{\eta}}(IC_{\pi,2}) + (n + 2), \end{aligned}$$

when  $n$  is large. Note that  $n + 2$  is irrelevant to the model. It follows that

$$\begin{aligned} &E_{\boldsymbol{\eta}} \left( \frac{\mathbf{y}^T \mathbf{A} \mathbf{y}}{\hat{\sigma}^2} \right) \\ &= n \times E_{\boldsymbol{\eta}} \left[ \frac{\mathbf{y}^T \{ \boldsymbol{\Sigma}^{-1} - \boldsymbol{\Sigma}^{-1} \mathbf{X} (\mathbf{X}^T \boldsymbol{\Sigma}^{-1} \mathbf{X} + \mathbf{W}^{-1})^{-1} \mathbf{X}^T \boldsymbol{\Sigma}^{-1} \} \mathbf{y}}{\mathbf{y}^T \{ \boldsymbol{\Sigma}^{-1} - \boldsymbol{\Sigma}^{-1} \mathbf{X} (\mathbf{X}^T \boldsymbol{\Sigma}^{-1} \mathbf{X})^{-1} \mathbf{X}^T \boldsymbol{\Sigma}^{-1} \} \mathbf{y}} \right] \\ &= n + n \times E_{\boldsymbol{\eta}} \left[ \frac{\mathbf{y}^T \boldsymbol{\Sigma}^{-1} \mathbf{X} (\mathbf{X}^T \boldsymbol{\Sigma}^{-1} \mathbf{X} + \mathbf{W}^{-1})^{-1} \mathbf{W}^{-1} (\mathbf{X}^T \boldsymbol{\Sigma}^{-1} \mathbf{X})^{-1} \mathbf{X}^T \boldsymbol{\Sigma}^{-1} \mathbf{y}}{\mathbf{y}^T \{ \boldsymbol{\Sigma}^{-1} - \boldsymbol{\Sigma}^{-1} \mathbf{X} (\mathbf{X}^T \boldsymbol{\Sigma}^{-1} \mathbf{X})^{-1} \mathbf{X}^T \boldsymbol{\Sigma}^{-1} \} \mathbf{y}} \right] \\ &= n + \frac{n}{\sigma^2(n-p-2)} \times E_{\boldsymbol{\eta}} [\mathbf{y}^T \boldsymbol{\Sigma}^{-1} \mathbf{X} (\mathbf{X}^T \boldsymbol{\Sigma}^{-1} \mathbf{X} + \mathbf{W}^{-1})^{-1} \mathbf{W}^{-1} (\mathbf{X}^T \boldsymbol{\Sigma}^{-1} \mathbf{X})^{-1} \mathbf{X}^T \boldsymbol{\Sigma}^{-1} \mathbf{y}] \\ &= n + \frac{n}{\sigma^2(n-p-2)} \times \left[ \sigma^2 \cdot \text{tr} \{ (\mathbf{X}^T \boldsymbol{\Sigma}^{-1} \mathbf{X} + \mathbf{W}^{-1})^{-1} \mathbf{W}^{-1} \} \right. \\ &\quad \left. + \boldsymbol{\beta}^T \mathbf{X}^T \boldsymbol{\Sigma}^{-1} \mathbf{X} (\mathbf{X}^T \boldsymbol{\Sigma}^{-1} \mathbf{X} + \mathbf{W}^{-1})^{-1} \mathbf{W}^{-1} \boldsymbol{\beta} \right], \end{aligned}$$

and that

$$E_{\pi}[\boldsymbol{\beta}^T \mathbf{X}^T \boldsymbol{\Sigma}^{-1} \mathbf{X} (\mathbf{X}^T \boldsymbol{\Sigma}^{-1} \mathbf{X} + \mathbf{W}^{-1})^{-1} \mathbf{W}^{-1} \boldsymbol{\beta}] = \sigma^2 \cdot \text{tr} [\mathbf{X}^T \boldsymbol{\Sigma}^{-1} \mathbf{X} (\mathbf{X}^T \boldsymbol{\Sigma}^{-1} \mathbf{X} + \mathbf{W}^{-1})^{-1}].$$

If  $n^{-1} \mathbf{X}^T \boldsymbol{\Sigma}^{-1} \mathbf{X}$  converges to  $p \times p$  positive definite matrix as  $n \rightarrow \infty$ ,  $\text{tr} [(\mathbf{X}^T \boldsymbol{\Sigma}^{-1} \mathbf{X} + \mathbf{W}^{-1})^{-1} \mathbf{W}^{-1}] \rightarrow 0$  and  $\text{tr} [\mathbf{X}^T \boldsymbol{\Sigma}^{-1} \mathbf{X} (\mathbf{X}^T \boldsymbol{\Sigma}^{-1} \mathbf{X} + \mathbf{W}^{-1})^{-1}] \rightarrow p$ . Then we can obtain  $E_{\pi}E_{\boldsymbol{\eta}}(\mathbf{y}^T \mathbf{A} \mathbf{y} / \hat{\sigma}^2 - n) \rightarrow p$ , which we want to show.

### 6.6.3 Derivation of $IC_r$ in (6.11)

We shall show that the bias correction  $\Delta_r = I_r(\boldsymbol{\eta}) - E_{\boldsymbol{\eta}}[-2 \log \{m_r(\mathbf{y} | \tilde{\sigma}^2)\}]$  is  $2(n-p)/(n-p-2)$ , where  $I_r(\boldsymbol{\eta})$  is given by (6.10). Then,

$$\begin{aligned} \Delta_r &= E_{\boldsymbol{\eta}}[\tilde{\mathbf{y}}^T (\boldsymbol{\Sigma}^{-1} - \mathbf{P}) \tilde{\mathbf{y}} / \tilde{\sigma}^2] - E_{\boldsymbol{\eta}}[\mathbf{y}^T (\boldsymbol{\Sigma}^{-1} - \mathbf{P}) \mathbf{y} / \tilde{\sigma}^2] \\ &= E_{\boldsymbol{\eta}}[\tilde{\mathbf{y}}^T (\boldsymbol{\Sigma}^{-1} - \mathbf{P}) \tilde{\mathbf{y}}] \cdot E_{\boldsymbol{\eta}}(1/\tilde{\sigma}^2) - (n-p). \end{aligned}$$

Since  $E_{\boldsymbol{\eta}}[\tilde{\mathbf{y}}^T (\boldsymbol{\Sigma}^{-1} - \mathbf{P}) \tilde{\mathbf{y}}] = (n-p)\sigma^2$  and  $E_{\boldsymbol{\eta}}(1/\tilde{\sigma}^2) = (n-p)/\{\sigma^2(n-p-2)\}$ , we get  $\Delta_r = 2(n-p)/(n-p-2)$ .  $\square$

## 6.7 Proof of Theorem 6.1

We only prove the consistency of  $IC_{\pi,1}$ . The proof of the consistency of the other criteria can be done in the same manner. Because we see that

$$P(\hat{j} = j) \leq P\{IC_{\pi,1}(j) < IC_{\pi,1}(j_*)\}$$

for any  $j \in \mathcal{J} \setminus \{j_*\}$ , it suffices to show that  $P\{\text{IC}_{\pi,1}(j) < \text{IC}_{\pi,1}(j_*)\} \rightarrow 0$ , or equivalently  $P\{\text{IC}_{\pi,1}(j) - \text{IC}_{\pi,1}(j_*) > 0\} \rightarrow 1$  as  $n \rightarrow \infty$ . When  $\boldsymbol{\Sigma} = \mathbf{I}_n$ , we obtain

$$\text{IC}_{\pi,1}(j) - \text{IC}_{\pi,1}(j_*) = I_1 + I_2 + I_3,$$

where

$$\begin{aligned} I_1 &= n \log(\hat{\sigma}_j^2 / \hat{\sigma}_*^2) + \mathbf{y}^\top \mathbf{A}_j \mathbf{y} / \hat{\sigma}_j^2 - \mathbf{y}^\top \mathbf{A}_* \mathbf{y} / \hat{\sigma}_*^2, \\ I_2 &= \log |\mathbf{X}(j)^\top \mathbf{X}(j) + \mathbf{W}_j^{-1}| - \log |\mathbf{X}(j_*)^\top \mathbf{X}(j_*) + \mathbf{W}_*^{-1}|, \\ I_3 &= \log\{|\mathbf{W}_j|/|\mathbf{W}_*|\} + \frac{2n}{n - p_j - 2} - \frac{2n}{n - p_* - 2}, \end{aligned}$$

for  $\hat{\sigma}_j^2 = \mathbf{y}^\top (\mathbf{I}_n - \mathbf{H}_j) \mathbf{y} / n$ ,  $\hat{\sigma}_*^2 = \hat{\sigma}_{j_*}^2$ ,  $\mathbf{A}_j = \mathbf{I}_n - \mathbf{X}(j)(\mathbf{X}(j)^\top \mathbf{X}(j) + \mathbf{W}_j^{-1})^{-1} \mathbf{X}(j)^\top$ ,  $\mathbf{A}_* = \mathbf{A}_{j_*}$  and  $\mathbf{W}_* = \mathbf{W}_{j_*}$ . We evaluate asymptotic behaviors of  $I_1$ ,  $I_2$  and  $I_3$  for  $j \in \mathcal{J}_-$  and  $j \in \mathcal{J}_+ \setminus \{j_*\}$ , separately.

[Case of  $j \in \mathcal{J}_-$ . Firstly, we evaluate  $I_1$ . We decompose  $I_1 = I_{11} + I_{12}$ , where  $I_{11} = n \log(\hat{\sigma}_j^2 / \hat{\sigma}_*^2)$  and  $I_{12} = \mathbf{y}^\top \mathbf{A}_j \mathbf{y} / \hat{\sigma}_j^2 - \mathbf{y}^\top \mathbf{A}_* \mathbf{y} / \hat{\sigma}_*^2$ . It follows that

$$\begin{aligned} \hat{\sigma}_j^2 - \hat{\sigma}_*^2 &= (\mathbf{X}(\omega) \boldsymbol{\beta}_* + \mathbf{u})^\top (\mathbf{I}_n - \mathbf{H}_j) (\mathbf{X}(\omega) \boldsymbol{\beta}_* + \mathbf{u}) / n - \mathbf{u}^\top (\mathbf{I}_n - \mathbf{H}_*) \mathbf{u} / n \\ &= \|\mathbf{X}(\omega) \boldsymbol{\beta}_* - \mathbf{H}_j \mathbf{X}(\omega) \boldsymbol{\beta}_*\|^2 / n + o_p(1), \end{aligned}$$

where  $\mathbf{H}_* = \mathbf{H}_{j_*}$ . Then we can see that

$$n^{-1} I_{11} = \log \left( 1 + \frac{\hat{\sigma}_j^2 - \hat{\sigma}_*^2}{\hat{\sigma}_*^2} \right) = \log \left\{ 1 + \frac{\|\mathbf{X}(\omega) \boldsymbol{\beta}_* - \mathbf{H}_j \mathbf{X}(\omega) \boldsymbol{\beta}_*\|^2}{n \sigma^2} \right\} + o_p(1), \quad (6.20)$$

and it follows from the assumption (A3) that

$$\liminf_{n \rightarrow \infty} \log \left\{ 1 + \frac{\|\mathbf{X}(\omega) \boldsymbol{\beta}_* - \mathbf{H}_j \mathbf{X}(\omega) \boldsymbol{\beta}_*\|^2}{n \sigma^2} \right\} > 0. \quad (6.21)$$

Because  $\mathbf{y}^\top \mathbf{A}_j \mathbf{y} / (n \hat{\sigma}_j^2) = 1 + o_p(1)$  and  $\mathbf{y}^\top \mathbf{A}_* \mathbf{y} / (n \hat{\sigma}_*^2) = 1 + o_p(1)$ , we obtain

$$n^{-1} I_{12} = o_p(1). \quad (6.22)$$

Secondly, we evaluate  $I_2$ . It follows that

$$\log |\mathbf{X}(j)^\top \mathbf{X}(j) + \mathbf{W}_j^{-1}| = p_j \log n + \log |\mathbf{X}(j)^\top \mathbf{X}(j) / n + \mathbf{W}_j^{-1} / n| = p_j \log n + O(1).$$

It can be also seen that  $\log |\mathbf{X}(j_*)^\top \mathbf{X}(j_*) + \mathbf{W}_*^{-1}| = p_* \log n + O(1)$ . Then,

$$n^{-1} I_2 = (p_j - p_*) n^{-1} \log n + o(1) = o(1). \quad (6.23)$$

Lastly, it is easy to see that

$$n^{-1} I_3 = o(1). \quad (6.24)$$

From (6.20)–(6.24), it follows that

$$P\{\text{IC}_{\pi,1}(j) - \text{IC}_{\pi,1}(j_*) > 0\} \rightarrow 1, \quad (6.25)$$

for all  $j \in \mathcal{J}_-$ .



[Case of  $j \in \mathcal{J}_+ \setminus \{j_*\}$ ]. Firstly, we evaluate  $I_1$ . From the fact that

$$\hat{\sigma}_*^2 - \hat{\sigma}_j^2 = \mathbf{u}^\top (\mathbf{H}_j - \mathbf{H}_*) \mathbf{u} / n = O_p(n^{-1}), \quad (6.26)$$

it follows that

$$\begin{aligned} (\log n)^{-1} I_{11} &= (\log n)^{-1} \cdot n \log \left\{ \frac{\hat{\sigma}_*^2 - (\hat{\sigma}_*^2 - \hat{\sigma}_j^2)}{\hat{\sigma}_*^2} \right\} \\ &= (\log n)^{-1} \cdot n \cdot \log\{1 + O_p(n^{-1})\} = o_p(1). \end{aligned} \quad (6.27)$$

As for  $I_{12}$ , from (6.26) and  $\mathbf{y}^\top \mathbf{A}_j \mathbf{y} - \mathbf{y}^\top \mathbf{A}_0 \mathbf{y} = O_p(1)$ , we can obtain

$$\begin{aligned} I_{12} &= \mathbf{y}^\top \mathbf{A}_j \mathbf{y} / \hat{\sigma}_j^2 - \mathbf{y}^\top \mathbf{A}_* \mathbf{y} / \hat{\sigma}_*^2 \\ &= (\mathbf{y}^\top \mathbf{A}_j \mathbf{y} - \mathbf{y}^\top \mathbf{A}_* \mathbf{y}) / \hat{\sigma}_*^2 + O_p(1) = O_p(1). \end{aligned}$$

Then,

$$(\log n)^{-1} I_{12} = o_p(1). \quad (6.28)$$

Secondly, we evaluate  $I_2$ . Since  $p_j > p_*$  for all  $j \in \mathcal{J}_+ \setminus \{j_*\}$ ,

$$\liminf_{n \rightarrow \infty} (\log n)^{-1} I_2 = p_j - p_* > 0. \quad (6.29)$$

Finally, it is easy to see that

$$(\log n)^{-1} I_3 = o(1). \quad (6.30)$$

From (6.27)–(6.30), it follows that

$$P\{\text{IC}_{\pi,1}(j) - \text{IC}_{\pi,2}(j_*) > 0\} \rightarrow 1, \quad (6.31)$$

for all  $j \in \mathcal{J}_+ \setminus \{j_*\}$ .

Combining (6.25) and (6.31), we obtain

$$P\{\text{IC}_{\pi,1}(j) - \text{IC}_{\pi,1}(j_*) > 0\} \rightarrow 1,$$

for all  $j \in \mathcal{J} \setminus \{j_*\}$ , which shows that  $\text{IC}_{\pi,1}$  is consistent.  $\square$



## Chapter 7

# Variants of conditional AIC in linear mixed models

In this chapter, we consider information criteria which measure prediction risks of several predictive densities in terms of expected Kullback–Leibler divergence based on conditional likelihood given random effects for variable selection problem in linear mixed model. When the predictive density considered is plug-in predictive density, the resulting criterion is the conditional AIC proposed by Vaida and Blanchard (2005). We consider two types of predictive densities, both of which are superior to plug-in predictive density in some sense. The first one is the Bayesian predictive density, which is the best predictive density in the sense of minimizing the expected Kullback–Leibler divergence. The second one is the predictive density based on the Bayesian marginal likelihood. The resulting criterion is related to the ones introduced in the last chapter.

### 7.1 Motivation

For variable selection problem in linear mixed model, Vaida and Blanchard (2005) introduced the conditional Akaike information (cAI), which is relevant to expected Kullback–Leibler (KL) divergence based on the conditional likelihood given random effects. This risk function is appropriate when one is interested in predicting the random effects. The cAI or its estimator conditional AIC (cAIC) measures the risk of the plug-in predictive density. However, there is no reason to restrict the predictive density to plug-in predictive density. Then, in this chapter, we consider two types of predictive densities.

The first one is the Bayesian (posterior) predictive density, which is the best predictive density among any predictive density in the sense of minimizing the expected KL divergence. Information criterion based on the Bayesian predictive density is known as the predictive likelihood (Akaike, 1980a) or the predictive information criterion (PIC) (Kitagawa, 1997). We construct PIC's in linear mixed model for two situations. In one situation, we assume that the vector of regression coefficients is unknown parameter. In the other situation, we consider the prior distribution of regression coefficients.

The second one is a predictive density based on Bayesian marginal likelihood assuming prior distribution on the regression coefficients of fixed effects. The resulting criterion is relevant to the AIC variant using Bayesian marginal likelihood, which was introduced in the last chapter. Like the AIC variant, there are three advantages of this method. Firstly, this criterion is less influenced by prior misspecification because risk function does not take expectation with respect to the prior distribution of regression coefficients. Secondly, non-informative improper prior can

be also used for constructing the criterion. When the uniform prior is assumed on the regression coefficients, the Bayesian marginal likelihood is identical to residual likelihood. Then we call the resulting criterion the conditional residual information criterion (cRIC). Lastly, the criteria have the consistency property for selecting the true model.

The rest of this chapter is organized as follows. In Section 7.2, we explain about the setup of variable selection and about the concept of cAIC variants in detail. In Section 7.3, we propose two PIC's in linear mixed model. In Section 7.4, we propose the cAIC variants based on Bayesian marginal likelihood. The numerical performance of the proposed criteria is investigated by simulations in Section 7.5.

## 7.2 Variants of conditional Akaike information

### 7.2.1 Setup of variable selection

The candidate model  $j$  is the linear mixed model

$$\mathbf{y} = \mathbf{X}(j)\boldsymbol{\beta}_j + \mathbf{Z}\mathbf{b}_j + \boldsymbol{\varepsilon}_j, \quad (7.1)$$

where  $\mathbf{y}$  is an  $n \times 1$  observation vector of response variables,  $\mathbf{X}(j)$  and  $\mathbf{Z}$  are  $n \times p_j$  and  $n \times q$  matrices of covariates, respectively,  $\boldsymbol{\beta}_j$  is a  $p_j \times 1$  vector of coefficients,  $\mathbf{b}_j$  is a  $q \times 1$  vector of random effects, and  $\boldsymbol{\varepsilon}_j$  is an  $n \times 1$  vector of random errors. Let  $\mathbf{b}_j$  and  $\boldsymbol{\varepsilon}_j$  be mutually independent and  $\mathbf{b}_j \sim \mathcal{N}_q(\mathbf{0}, \sigma_j^2 \mathbf{G})$ ,  $\boldsymbol{\varepsilon}_j \sim \mathcal{N}_n(\mathbf{0}, \sigma_j^2 \mathbf{I}_n)$ , where  $\mathbf{G}$  is a  $q \times q$  positive definite matrix and  $\sigma_j^2$  is a scalar. We assume that  $\mathbf{G}$  is known and  $\sigma_j^2$  is unknown.

To derive the criterion, we assume that the true model  $j_*$  is included by each candidate model  $j$ , namely each candidate model  $j$  is overspecified. In that situation, the true model  $j_*$  can be expressed by

$$\mathbf{y} = \mathbf{X}(j)\boldsymbol{\beta}_j^* + \mathbf{Z}\mathbf{b}_* + \boldsymbol{\varepsilon}_*, \quad (7.2)$$

where  $\mathbf{b}_* \sim \mathcal{N}_q(\mathbf{0}, \sigma_*^2 \mathbf{G})$ ,  $\boldsymbol{\varepsilon}_* \sim \mathcal{N}_n(\mathbf{0}, \sigma_*^2 \mathbf{I}_n)$  and  $\boldsymbol{\beta}_j^*$  is a  $p_j \times 1$  vector of regression coefficients, whose  $p_j - p_*$  components are exactly 0 and the rest of components are not 0, for  $p_* = p_{j_*}$ . Henceforth, we abbreviate the model index  $j$  for notational convenience when there is no confusion. We also abbreviate  $\boldsymbol{\beta}_j^*$  as  $\boldsymbol{\beta}$ ,  $\mathbf{b}_*$  as  $\mathbf{b}$  and  $\sigma_*^2$  as  $\sigma^2$ . Let  $f(\mathbf{y}|\mathbf{b}, \boldsymbol{\beta}, \sigma^2)$  and  $p(\mathbf{b}|\sigma^2)$  denote the conditional density function of  $\mathbf{y}$  given  $\mathbf{b}$  and the density function of  $\mathbf{b}$ , respectively.

### 7.2.2 Conditional Kullback–Leibler risk

A conventional method of selecting the explanatory variable in linear mixed model is to use the marginal AIC (mAIC). The mAIC is based on the marginal likelihood integrating out the random effects  $\mathbf{b}$ , which is given by

$$\text{mAIC} = -2 \log\{m(\mathbf{y}|\hat{\boldsymbol{\beta}}_j, \hat{\sigma}_j^2)\} + 2p_j,$$

where  $m(\mathbf{y}|\boldsymbol{\beta}, \sigma^2) = \int f(\mathbf{y}|\mathbf{b}, \boldsymbol{\beta}, \sigma^2)p(\mathbf{b}|\sigma^2)d\mathbf{b}$  and  $\hat{\boldsymbol{\beta}}_j$  is the maximum likelihood estimator of  $\boldsymbol{\beta}_j$ . However, the mAIC is not appropriate for the focus on the prediction of specific clusters or random effects. Then, Vaida and Blanchard (2005) considered the expected Kullback–Leibler (KL) divergence based on the conditional density, which is given by

$$R_c(\boldsymbol{\eta}; \hat{f}_j) = \iint \left[ \int \log \left\{ \frac{f(\tilde{\mathbf{y}}|\mathbf{b}, \boldsymbol{\beta}, \sigma^2)}{\hat{f}_j(\tilde{\mathbf{y}}; \mathbf{y})} \right\} f(\tilde{\mathbf{y}}|\mathbf{b}, \boldsymbol{\beta}, \sigma^2)d\tilde{\mathbf{y}} \right] f(\mathbf{y}|\mathbf{b}, \boldsymbol{\beta}, \sigma^2)p(\mathbf{b}|\sigma^2)d\mathbf{y}d\mathbf{b}, \quad (7.3)$$

where  $\tilde{\mathbf{y}}$  is an independent replication of  $\mathbf{y}$  given  $\mathbf{b}$ ,  $\hat{f}_j(\tilde{\mathbf{y}}; \mathbf{y})$  is some predictive density of  $f(\tilde{\mathbf{y}}|\mathbf{b}, \boldsymbol{\beta}, \sigma^2)$  based on the candidate model  $j$ , and  $\boldsymbol{\eta} = (\boldsymbol{\beta}^\top, \sigma^2)^\top$ . We call  $R_c(\boldsymbol{\eta}; \hat{f}_j)$  in (7.3) the conditional KL risk. We can provide an information criterion as an (asymptotically) unbiased estimator of the information

$$I_c(\boldsymbol{\eta}; \hat{f}_j) = \iiint -2 \log \{ \hat{f}_j(\tilde{\mathbf{y}}; \mathbf{y}) \} f(\tilde{\mathbf{y}}|\mathbf{b}, \boldsymbol{\beta}, \sigma^2) f(\mathbf{y}|\mathbf{b}, \boldsymbol{\beta}, \sigma^2) p(\mathbf{b}|\sigma^2) d\tilde{\mathbf{y}} d\mathbf{y} d\mathbf{b},$$

which is a part of (7.3) (multiplied by 2). Vaida and Blanchard (2005) proposed the conditional AIC (cAIC) as an unbiased estimator of  $I_c(\boldsymbol{\eta}; \hat{f}_j)$  for  $\hat{f}_j = f(\tilde{\mathbf{y}}|\hat{\mathbf{b}}_j, \hat{\boldsymbol{\beta}}_j, \hat{\sigma}_j^2)$ , which is called plug-in predictive density, where  $\hat{\mathbf{b}}_j$ ,  $\hat{\boldsymbol{\beta}}_j$  and  $\hat{\sigma}_j^2$  are some predictor of  $\mathbf{b}_j$  and estimators of  $\boldsymbol{\beta}_j$  and  $\sigma_j^2$ , respectively. The cAIC is of the form

$$\text{cAIC} = -2 \log \{ f(\mathbf{y}|\hat{\mathbf{b}}_j, \hat{\boldsymbol{\beta}}_j, \hat{\sigma}_j^2) \} + \Delta_{\text{cAI}},$$

where  $\Delta_{\text{cAI}} = I_c(\boldsymbol{\eta}; \hat{f}_j) - E[-2 \log \{ f(\mathbf{y}|\hat{\mathbf{b}}_j, \hat{\boldsymbol{\beta}}_j, \hat{\sigma}_j^2) \}]$ .

Although the conditional KL risk  $R_c(\boldsymbol{\eta}; \hat{f}_j)$  is appropriate for variable selection in linear mixed model, there is no reason to restrict the predictive density  $\hat{f}_j$  to plug-in predictive density. Thus we consider two types of predictive densities: the first one is the Bayesian (posterior) predictive density and the second one is a predictive density based on Bayesian marginal likelihood, both of which are superior to the plug-in predictive density in some sense. In the following sections, we derive information criteria based on these predictive densities and discuss their properties.

## 7.3 Predictive information criterion

### 7.3.1 Bayesian predictive density in linear mixed model

Aitchison (1975) showed that the best predictive density  $\hat{f}_j$  which minimizes the risk  $R_c(\boldsymbol{\eta}; \hat{f}_j)$  is

$$\hat{f}_j^{\text{BP}}(\tilde{\mathbf{y}}|\mathbf{y}, \boldsymbol{\eta}) = \frac{\int f(\tilde{\mathbf{y}}|\mathbf{b}, \boldsymbol{\beta}, \sigma^2) f(\mathbf{y}|\mathbf{b}, \boldsymbol{\beta}, \sigma^2) p(\mathbf{b}|\sigma^2) d\mathbf{b}}{\int f(\mathbf{y}|\mathbf{b}, \boldsymbol{\beta}, \sigma^2) p(\mathbf{b}|\sigma^2) d\mathbf{b}}, \quad (7.4)$$

when  $\boldsymbol{\eta} = (\boldsymbol{\beta}^\top, \sigma^2)^\top$  is known. The predictive density  $\hat{f}_j^{\text{BP}}$  is known as the Bayesian (or posterior) predictive density. Note that the linear mixed model can be seen as a Bayesian model where the random effects  $\mathbf{b}$  has prior distribution  $p(\mathbf{b}|\sigma^2)$ . The cAIC measures the risk  $R_c(\boldsymbol{\eta}; \hat{f}_j)$  for  $\hat{f}_j = f(\tilde{\mathbf{y}}|\hat{\mathbf{b}}_j, \hat{\boldsymbol{\beta}}_j, \hat{\sigma}_j^2)$ , which is plug-in predictive density, and it follows that

$$R_c(\boldsymbol{\eta}; f(\tilde{\mathbf{y}}|\hat{\mathbf{b}}_j, \boldsymbol{\beta}, \sigma^2)) \geq R_c(\boldsymbol{\eta}; \hat{f}_j^{\text{BP}}(\tilde{\mathbf{y}}|\mathbf{y}, \boldsymbol{\eta})),$$

when  $\boldsymbol{\eta}$  is known. Thus it is natural to use the Bayesian predictive density to measure the prediction risk of the model. Information criterion based on the Bayesian predictive density is known as the predictive likelihood (Akaike, 1980a) or the predictive information criterion (PIC) (Kitagawa, 1997).

The next proposition shows the Bayesian predictive distribution in linear mixed model (7.2).

**Proposition 7.1** *In the linear mixed model (7.2), the Bayesian predictive distribution whose density function is given by (7.4) is*

$$\hat{f}_j^{\text{BP}}(\tilde{\mathbf{y}}|\mathbf{y}, \boldsymbol{\eta}) = \mathcal{N}(\mathbf{X}\boldsymbol{\beta} + \mathbf{Z}\hat{\mathbf{b}}^{\text{B}}, \sigma^2\mathbf{V}), \quad (7.5)$$

where  $\hat{\mathbf{b}}^{\text{B}} = \mathbf{G}\mathbf{Z}^\top \boldsymbol{\Sigma}^{-1}(\mathbf{y} - \mathbf{X}\boldsymbol{\beta})$  and  $\mathbf{V} = 2\mathbf{I}_n - \boldsymbol{\Sigma}^{-1}$ .

### 7.3.2 Derivation of PIC in linear mixed model

In this subsection, we derive the PIC in linear mixed model based on the Bayesian predictive distribution (7.5). Because  $\boldsymbol{\eta} = (\boldsymbol{\beta}^\top, \sigma^2)^\top$  is unknown, we estimate it by maximum likelihood based on each candidate model  $j$  as follows:

$$\begin{aligned}\widehat{\boldsymbol{\beta}}_j &= (\mathbf{X}(j)^\top \boldsymbol{\Sigma}^{-1} \mathbf{X}(j))^{-1} \mathbf{X}(j)^\top \boldsymbol{\Sigma}^{-1} \mathbf{y}, \\ \widehat{\sigma}_j^2 &= (\mathbf{y} - \mathbf{X}(j) \widehat{\boldsymbol{\beta}}_j)^\top \boldsymbol{\Sigma}^{-1} (\mathbf{y} - \mathbf{X}(j) \widehat{\boldsymbol{\beta}}_j) / n.\end{aligned}$$

Then we consider the following information:

$$\begin{aligned}I_c(\boldsymbol{\eta}; \hat{f}_j^{\text{BP}}) &= \iiint -2 \log \{ \hat{f}^{\text{BP}}(\tilde{\mathbf{y}}|\mathbf{y}, \widehat{\boldsymbol{\beta}}_j, \widehat{\sigma}_j^2) \} f(\tilde{\mathbf{y}}|\mathbf{b}, \boldsymbol{\beta}, \sigma^2) f(\mathbf{y}|\mathbf{b}, \boldsymbol{\beta}, \sigma^2) p(\mathbf{b}|\sigma^2) d\tilde{\mathbf{y}} d\mathbf{y} d\mathbf{b}, \\ &= \int \left[ \int -2 \log \{ \hat{f}^{\text{BP}}(\tilde{\mathbf{y}}|\mathbf{y}, \widehat{\boldsymbol{\beta}}_j, \widehat{\sigma}_j^2) \} \hat{f}^{\text{BP}}(\tilde{\mathbf{y}}|\mathbf{y}, \boldsymbol{\eta}) d\tilde{\mathbf{y}} \right] m(\mathbf{y}|\boldsymbol{\beta}, \sigma^2) d\mathbf{y} \\ &= \text{PI} \quad (\text{say}),\end{aligned}\tag{7.6}$$

We want to construct an information criterion as an unbiased estimator of PI, which is of the form

$$\text{PIC} = -2 \log \{ \hat{f}^{\text{BP}}(\mathbf{y}|\mathbf{y}, \widehat{\boldsymbol{\beta}}_j, \widehat{\sigma}_j^2) \} + \Delta_{\text{PI}},$$

where

$$\begin{aligned}-2 \log \{ \hat{f}^{\text{BP}}(\mathbf{y}|\mathbf{y}, \widehat{\boldsymbol{\beta}}_j, \widehat{\sigma}_j^2) \} &= n \log(2\pi \widehat{\sigma}_j^2) + \log |\mathbf{V}| + (\mathbf{y} - \mathbf{X} \widehat{\boldsymbol{\beta}}_j - \mathbf{Z} \hat{\mathbf{b}}_j)^\top \mathbf{V}^{-1} (\mathbf{y} - \mathbf{X} \widehat{\boldsymbol{\beta}}_j - \mathbf{Z} \hat{\mathbf{b}}_j) / \widehat{\sigma}_j^2, \\ \Delta_{\text{PI}} &= \text{PI} - E[-2 \log \{ \hat{f}^{\text{BP}}(\mathbf{y}|\mathbf{y}, \widehat{\boldsymbol{\beta}}_j, \widehat{\sigma}_j^2) \}]\end{aligned}$$

for  $\hat{\mathbf{b}}_j = \mathbf{GZ}^\top \boldsymbol{\Sigma}^{-1} (\mathbf{y} - \mathbf{X} \widehat{\boldsymbol{\beta}}_j)$ . Then we have to evaluate the bias correction  $\Delta_{\text{PI}}$ .

Firstly, taking expectation of

$$-2 \log \{ \hat{f}^{\text{BP}}(\tilde{\mathbf{y}}|\mathbf{y}, \widehat{\boldsymbol{\beta}}_j, \widehat{\sigma}_j^2) \} = n \log(2\pi \widehat{\sigma}_j^2) + \log |\mathbf{V}| + (\tilde{\mathbf{y}} - \mathbf{X} \widehat{\boldsymbol{\beta}}_j - \mathbf{Z} \hat{\mathbf{b}}_j)^\top \mathbf{V}^{-1} (\tilde{\mathbf{y}} - \mathbf{X} \widehat{\boldsymbol{\beta}}_j - \mathbf{Z} \hat{\mathbf{b}}_j) / \widehat{\sigma}_j^2,$$

with respect to the distribution of  $\tilde{\mathbf{y}}|\mathbf{y} \sim \hat{f}^{\text{BP}}(\tilde{\mathbf{y}}|\mathbf{y}, \boldsymbol{\eta}) = \mathcal{N}(\mathbf{X}\boldsymbol{\beta} + \mathbf{Z}\hat{\mathbf{b}}^{\text{B}}, \sigma^2 \mathbf{V})$ , we can rewrite the PI in (7.6) as

$$\begin{aligned}\text{PI} &= E \left[ n \log(2\pi \widehat{\sigma}_j^2) + \log |\mathbf{V}| + n\sigma^2 / \widehat{\sigma}_j^2 \right. \\ &\quad \left. + \{ \mathbf{X}(\widehat{\boldsymbol{\beta}}_j - \boldsymbol{\beta}) + \mathbf{Z}(\hat{\mathbf{b}}_j - \hat{\mathbf{b}}^{\text{B}}) \}^\top \mathbf{V}^{-1} \{ \mathbf{X}(\widehat{\boldsymbol{\beta}}_j - \boldsymbol{\beta}) + \mathbf{Z}(\hat{\mathbf{b}}_j - \hat{\mathbf{b}}^{\text{B}}) \} \right] \\ &= E \left[ n \log(2\pi \widehat{\sigma}_j^2) + \log |\mathbf{V}| + n\sigma^2 / \widehat{\sigma}_j^2 \right. \\ &\quad \left. + (\widehat{\boldsymbol{\beta}}_j - \boldsymbol{\beta})^\top \mathbf{X}^\top \boldsymbol{\Sigma}^{-1} \mathbf{V}^{-1} \boldsymbol{\Sigma}^{-1} \mathbf{X}(\widehat{\boldsymbol{\beta}}_j - \boldsymbol{\beta}) / \widehat{\sigma}_j^2 \right],\end{aligned}$$

noting that  $\mathbf{X}(\widehat{\boldsymbol{\beta}}_j - \boldsymbol{\beta}) + \mathbf{Z}(\hat{\mathbf{b}}_j - \hat{\mathbf{b}}^{\text{B}}) = \boldsymbol{\Sigma}^{-1} \mathbf{X}(\widehat{\boldsymbol{\beta}}_j - \boldsymbol{\beta})$ . Next,  $-2 \log \{ \hat{f}^{\text{BP}}(\mathbf{y}|\mathbf{y}, \widehat{\boldsymbol{\beta}}_j, \widehat{\sigma}_j^2) \}$  is rewritten as

$$\begin{aligned}&-2 \log \{ \hat{f}^{\text{BP}}(\mathbf{y}|\mathbf{y}, \widehat{\boldsymbol{\beta}}_j, \widehat{\sigma}_j^2) \} \\ &= n \log(2\pi \widehat{\sigma}_j^2) + \log |\mathbf{V}| + \{ \mathbf{u} - \mathbf{X}(\widehat{\boldsymbol{\beta}}_j - \boldsymbol{\beta}) \}^\top \boldsymbol{\Sigma}^{-1} \mathbf{V}^{-1} \boldsymbol{\Sigma}^{-1} \{ \mathbf{u} - \mathbf{X}(\widehat{\boldsymbol{\beta}}_j - \boldsymbol{\beta}) \} / \widehat{\sigma}_j^2,\end{aligned}$$

where  $\mathbf{u} = \mathbf{y} - \mathbf{X}\boldsymbol{\beta}$ . Then the bias correction  $\Delta_{\text{PI}}$  is

$$\begin{aligned}\Delta_{\text{PI}} &= E \left[ n\sigma^2 / \widehat{\sigma}_j^2 - \mathbf{u}^\top \boldsymbol{\Sigma}^{-1} \mathbf{V}^{-1} \boldsymbol{\Sigma}^{-1} \mathbf{u} / \widehat{\sigma}_j^2 + 2\mathbf{u}^\top \boldsymbol{\Sigma}^{-1} \mathbf{V}^{-1} \boldsymbol{\Sigma}^{-1} \mathbf{X}(\widehat{\boldsymbol{\beta}}_j - \boldsymbol{\beta}) / \widehat{\sigma}_j^2 \right] \\ &= I_1 - I_2 + 2I_3 \quad (\text{say}).\end{aligned}$$

Thus, it suffices to evaluate  $I_1$ ,  $I_2$  and  $I_3$ .

Noting that  $n\hat{\sigma}_j^2/\sigma^2 \sim \chi_{n-p_j}^2$ , we can evaluate  $I_1$  as

$$I_1 = \frac{n^2}{n - p_j - 2}.$$

Next,  $I_2$  is rewritten as

$$I_2 = n \cdot E \left[ \frac{\mathbf{v}^\top \boldsymbol{\Sigma}^{-1/2} \mathbf{V}^{-1} \boldsymbol{\Sigma}^{-1/2} \mathbf{v}}{\mathbf{v}^\top (\mathbf{I}_n - \mathbf{M}) \mathbf{v}} \right],$$

where  $\mathbf{v} = \boldsymbol{\Sigma}^{-1/2} \mathbf{u} / \sigma$  and  $\mathbf{M} = \boldsymbol{\Sigma}^{-1/2} \mathbf{X} (\mathbf{X}^\top \boldsymbol{\Sigma}^{-1} \mathbf{X})^{-1} \mathbf{X}^\top \boldsymbol{\Sigma}^{-1/2}$ . Then, it follows from Lemma 4.6 that

$$\begin{aligned} I_2 &= n \times \left\{ \frac{\text{tr}(\mathbf{V}^{-1} \boldsymbol{\Sigma}^{-1})}{n - p_j - 2} - \frac{2 \text{tr}[\boldsymbol{\Sigma}^{-1/2} \mathbf{V}^{-1} \boldsymbol{\Sigma}^{-1/2} (\mathbf{I}_n - \mathbf{M})]}{(n - p_j)(n - p_j - 2)} \right\} \\ &= \frac{n \cdot \text{tr}(\mathbf{V}^{-1} \boldsymbol{\Sigma}^{-1})}{n - p_j} + \frac{2n \cdot \text{tr}(\mathbf{V}^{-1} \mathbf{P})}{(n - p_j)(n - p_j - 2)}, \end{aligned}$$

where  $\mathbf{P} = \boldsymbol{\Sigma}^{-1} \mathbf{X} (\mathbf{X}^\top \boldsymbol{\Sigma}^{-1} \mathbf{X})^{-1} \mathbf{X}^\top \boldsymbol{\Sigma}^{-1}$ . Lastly,  $I_3$  can be rewritten as

$$\begin{aligned} I_3 &= n \cdot E \left[ \frac{\mathbf{v}^\top \boldsymbol{\Sigma}^{-1/2} \mathbf{V}^{-1} \boldsymbol{\Sigma}^{-1/2} \mathbf{M} \mathbf{v}}{\mathbf{v}^\top (\mathbf{I}_n - \mathbf{M}) \mathbf{v}} \right] \\ &= n \cdot E \left[ \frac{\mathbf{v}^\top \mathbf{M} \boldsymbol{\Sigma}^{-1/2} \mathbf{V}^{-1} \boldsymbol{\Sigma}^{-1/2} \mathbf{M} \mathbf{v}}{\mathbf{v}^\top (\mathbf{I}_n - \mathbf{M}) \mathbf{v}} \right] + E \left[ \frac{\mathbf{v}^\top (\mathbf{I}_n - \mathbf{M}) \boldsymbol{\Sigma}^{-1/2} \mathbf{V}^{-1} \boldsymbol{\Sigma}^{-1/2} \mathbf{M} \mathbf{v}}{\mathbf{v}^\top (\mathbf{I}_n - \mathbf{M}) \mathbf{v}} \right]. \end{aligned}$$

Because  $\mathbf{M} \mathbf{v}$  and  $(\mathbf{I}_n - \mathbf{M}) \mathbf{v}$  are independent and  $E(\mathbf{M} \mathbf{v}) = \mathbf{0}$ , the second term of the equation above is 0. Then,

$$I_3 = \frac{n \cdot \text{tr}(\mathbf{V}^{-1} \mathbf{P})}{n - p_j - 2}.$$

Thus we can obtain

$$\Delta_{\text{PI}} = \frac{n^2}{n - p_j - 2} - \frac{n \cdot \text{tr}(\mathbf{V}^{-1} \boldsymbol{\Sigma}^{-1})}{n - p_j} + \frac{2n(n - p_j - 1) \text{tr}(\mathbf{V}^{-1} \mathbf{P})}{(n - p_j)(n - p_j - 2)},$$

and propose the following PIC:

$$\text{PIC} = -2 \log \{ \hat{f}^{\text{BP}}(\mathbf{y} | \mathbf{y}, \hat{\boldsymbol{\beta}}_j, \hat{\sigma}_j^2) \} + \Delta_{\text{PI}}. \quad (7.7)$$

**Theorem 7.1** *The PIC in (7.7) is an unbiased estimator of PI in (7.6), namely  $E(\text{PIC}) = \text{PI}$ .*

### 7.3.3 Another PIC putting prior on regression coefficients

Although the PIC in (7.7) works as a variable selection criterion, the predictive density function  $\hat{f}^{\text{BP}}(\tilde{\mathbf{y}} | \mathbf{y}, \hat{\boldsymbol{\beta}}_j, \hat{\sigma}_j^2)$  seems not to be very appropriate for evaluation of the risk of the candidate model  $j$ . This is because the Bayesian predictive density  $\hat{f}^{\text{BP}}(\tilde{\mathbf{y}} | \mathbf{y}, \boldsymbol{\eta})$ , which minimizes the conditional KL risk, is derived under the condition that  $\boldsymbol{\eta} = (\boldsymbol{\beta}^\top, \sigma^2)^\top$  is known, although  $\boldsymbol{\beta}$  varies with the model  $j$ . Then, we consider in this subsection another approach to construct the PIC.

To this end, we put the prior on the regression coefficients  $\boldsymbol{\beta}$  as well as on the random effects  $\mathbf{b}$ . We assume that the prior distribution of  $\boldsymbol{\beta}$  is  $\boldsymbol{\beta} \sim \pi(\boldsymbol{\beta}|\sigma^2) = \mathcal{N}_p(\mathbf{0}, \sigma^2 \mathbf{W})$ , where  $\mathbf{W}$  is a known positive definite matrix, and that  $\boldsymbol{\beta}$  is independent of  $\mathbf{b}$ . Let  $\boldsymbol{\theta} = (\boldsymbol{\beta}^\top, \mathbf{b}^\top)^\top$  and  $\psi(\boldsymbol{\theta}|\sigma^2)$  denote the prior density function of  $\boldsymbol{\theta}$ . Then, we can see the linear mixed model as the Bayesian model

$$\begin{aligned} \mathbf{y}|\boldsymbol{\theta} &\sim g(\mathbf{y}|\boldsymbol{\theta}, \sigma^2) = \mathcal{N}(\mathbf{X}\boldsymbol{\beta} + \mathbf{Z}\mathbf{b}, \sigma^2 \mathbf{I}_n), \\ \boldsymbol{\theta} &\sim \psi(\boldsymbol{\theta}|\sigma^2) = p(\mathbf{b}|\sigma^2)\pi(\boldsymbol{\beta}|\sigma^2) = \mathcal{N}(\mathbf{0}, \sigma^2 \mathbf{G}) \cdot \mathcal{N}(\mathbf{0}, \sigma^2 \mathbf{W}), \end{aligned} \quad (7.8)$$

In this model, the Bayesian predictive density is

$$\hat{g}^{\text{BP}}(\mathbf{z}|\mathbf{y}, \sigma^2) = \frac{\int g(\mathbf{z}|\boldsymbol{\theta}, \sigma^2)g(\mathbf{y}|\boldsymbol{\theta}, \sigma^2)\psi(\boldsymbol{\theta}|\sigma^2)d\boldsymbol{\theta}}{\int g(\mathbf{y}|\boldsymbol{\theta}, \sigma^2)\psi(\boldsymbol{\theta}|\sigma^2)d\boldsymbol{\theta}}, \quad (7.9)$$

where  $\mathbf{z}$  is independent replication of  $\mathbf{y}$  given  $\boldsymbol{\theta}$ . Note that  $\hat{g}^{\text{BP}}(\mathbf{z}|\mathbf{y}, \sigma^2)$  is the predictive density which minimizes the following conditional KL risk:

$$\bar{R}_c(\sigma^2; \hat{g}) = \iiint \left[ \int \log \left\{ \frac{g(\mathbf{z}|\boldsymbol{\theta}, \sigma^2)}{\hat{g}(\mathbf{z}; \mathbf{y})} \right\} g(\mathbf{z}|\boldsymbol{\theta}, \sigma^2) d\mathbf{z} \right] g(\mathbf{y}|\boldsymbol{\theta}, \sigma^2) \psi(\boldsymbol{\theta}|\sigma^2) d\mathbf{y} d\boldsymbol{\theta},$$

when  $\sigma^2$  is known, namely for any predictive density  $\hat{g}(\mathbf{z}; \mathbf{y})$  the following inequality holds:

$$\bar{R}_c(\sigma^2; \hat{g}(\mathbf{z}; \mathbf{y})) \geq \bar{R}_c(\sigma^2; \hat{g}^{\text{BP}}(\mathbf{z}|\mathbf{y}, \sigma^2))$$

It is also important to note that  $\hat{g}^{\text{BP}}(\mathbf{z}|\mathbf{y}, \sigma^2)$  implicitly depends on the candidate model  $j$ . Thus an information criterion based on  $\hat{g}^{\text{BP}}(\mathbf{z}|\mathbf{y}, \sigma^2)$  is appropriate for variable selection.

The next proposition shows the Bayesian predictive distribution in the Bayesian model (7.8).

**Proposition 7.2** *In the Bayesian model (7.8), the Bayesian predictive distribution whose density function is given by (7.9) is*

$$\hat{g}^{\text{BP}}(\mathbf{z}|\mathbf{y}, \sigma^2) = \mathcal{N}(\{\mathbf{I}_n - (\mathbf{B} + \boldsymbol{\Sigma})^{-1}\}\mathbf{y}, \sigma^2 \mathbf{V}_2), \quad (7.10)$$

where  $\mathbf{B} = \mathbf{X}\mathbf{W}\mathbf{X}^\top$  and  $\mathbf{V}_2 = 2\mathbf{I}_n - (\mathbf{B} + \boldsymbol{\Sigma})^{-1}$ .

We consider the following information:

$$\begin{aligned} \text{PI}_2 &= \iiint -2 \log \{\hat{g}^{\text{BP}}(\mathbf{z}|\mathbf{y}, \hat{\sigma}_j^2)\} g(\mathbf{z}|\boldsymbol{\theta}, \sigma^2) g(\mathbf{y}|\boldsymbol{\theta}, \sigma^2) \psi(\boldsymbol{\theta}|\sigma^2) d\mathbf{z} d\mathbf{y} d\boldsymbol{\theta} \\ &= \int \left[ \int -2 \log \{\hat{g}^{\text{BP}}(\mathbf{z}|\mathbf{y}, \hat{\sigma}_j^2)\} \hat{g}^{\text{BP}}(\mathbf{z}|\mathbf{y}, \sigma^2) d\mathbf{z} \right] g_\psi(\mathbf{y}|\sigma^2) d\mathbf{y}, \end{aligned} \quad (7.11)$$

where  $g_\psi(\mathbf{y}|\sigma^2) = \int g(\mathbf{y}|\boldsymbol{\theta}, \sigma^2) \psi(\boldsymbol{\theta}|\sigma^2) d\boldsymbol{\theta}$ . We want to construct an information criterion as an unbiased estimator of  $\text{PI}_2$ , which is of the form

$$\text{PIC}_2 = -2 \log \{\hat{g}^{\text{BP}}(\mathbf{y}|\mathbf{y}, \hat{\sigma}_j^2)\} + \Delta_{\text{PI}_2},$$

where

$$\begin{aligned} -2 \log \{\hat{g}^{\text{BP}}(\mathbf{y}|\mathbf{y}, \hat{\sigma}_j^2)\} &= n \log(2\pi \hat{\sigma}_j^2) + \log |\mathbf{V}_2| + \mathbf{y}^\top (\mathbf{B} + \boldsymbol{\Sigma})^{-1} \mathbf{V}_2^{-1} (\mathbf{B} + \boldsymbol{\Sigma})^{-1} \mathbf{y} / \hat{\sigma}_j^2, \\ \Delta_{\text{PI}_2} &= \text{PI}_2 - E_{g_\psi}[-2 \log \{\hat{g}^{\text{BP}}(\mathbf{y}|\mathbf{y}, \hat{\sigma}_j^2)\}], \end{aligned}$$



and  $E_{g_\psi}$  denotes the expectation with respect to the distribution of  $g_\psi(\mathbf{y}|\sigma^2)$ . Then we have to evaluate the bias correction  $\Delta_{\text{PI}_2}$ .

Firstly, taking expectation of

$$\begin{aligned} & -2 \log \{ \hat{g}^{\text{BP}}(\mathbf{z}|\mathbf{y}, \hat{\sigma}_j^2) \} \\ & = n \log(2\pi\hat{\sigma}_j^2) + \log |\mathbf{V}_2| + [\mathbf{z} - \{\mathbf{I}_n - (\mathbf{B} + \boldsymbol{\Sigma})^{-1}\}\mathbf{y}]^T \mathbf{V}_2^{-1} [\mathbf{z} - \{\mathbf{I}_n - (\mathbf{B} + \boldsymbol{\Sigma})^{-1}\}\mathbf{y}] / \hat{\sigma}_j^2, \end{aligned}$$

with respect to the distribution of  $\hat{g}^{\text{BP}}(\mathbf{z}|\mathbf{y}, \sigma^2) = \mathcal{N}(\{\mathbf{I}_n - (\mathbf{B} + \boldsymbol{\Sigma})^{-1}\}\mathbf{y}, \sigma^2 \mathbf{V}_2)$ , we can rewrite  $\text{PI}_2$  in (7.11) as

$$\text{PI}_2 = E_{g_\psi} [n \log(2\pi\hat{\sigma}_j^2) + \log |\mathbf{V}_2| + n\sigma^2/\hat{\sigma}_j^2].$$

Then, the bias correction  $\Delta_{\text{PI}_2}$  can be expressed as

$$\begin{aligned} \Delta_{\text{PI}_2} & = E_{g_\psi} [n\sigma^2/\hat{\sigma}_j^2 - \mathbf{y}^T (\mathbf{B} + \boldsymbol{\Sigma})^{-1} \mathbf{V}_2^{-1} (\mathbf{B} + \boldsymbol{\Sigma})^{-1} \mathbf{y} / \hat{\sigma}_j^2] \\ & = I_4 - I_5 \quad (\text{say}). \end{aligned}$$

Thus, it suffices to evaluate  $I_4$  and  $I_5$ .

Noting the fact that

$$\begin{aligned} g_\psi(\mathbf{y}|\sigma^2) & = \int g(\mathbf{y}|\boldsymbol{\theta}, \sigma^2) \psi(\boldsymbol{\theta}|\sigma^2) d\boldsymbol{\theta} = \iint f(\mathbf{y}|\boldsymbol{\beta}, \mathbf{b}, \sigma^2) p(\mathbf{b}|\sigma^2) \pi(\boldsymbol{\beta}|\sigma^2) d\mathbf{b} d\boldsymbol{\beta} \\ & = \int m(\mathbf{y}|\boldsymbol{\beta}, \sigma^2) \pi(\boldsymbol{\beta}|\sigma^2) d\boldsymbol{\beta}, \end{aligned}$$

$I_4$  can be rewritten as

$$I_4 = E_{g_\psi} [n\sigma^2/\hat{\sigma}_j^2] = n^2 E_{\boldsymbol{\beta}} E_{\mathbf{y}|\boldsymbol{\beta}} [n\sigma^2/\hat{\sigma}_j^2],$$

where  $E_{\mathbf{y}|\boldsymbol{\beta}}$  and  $E_{\boldsymbol{\beta}}$  denote the expectation with respect to the distribution of  $\mathbf{y}|\boldsymbol{\beta} \sim m(\mathbf{y}|\boldsymbol{\beta}, \sigma^2)$  and  $\boldsymbol{\beta} \sim \pi(\boldsymbol{\beta}|\sigma^2)$ , respectively. Then,  $I_4$  is evaluated as

$$I_4 = \frac{n^2}{n - p_j - 2}.$$

Next,  $I_5$  can be decomposed as

$$\begin{aligned} I_5 & = n E_{g_\psi} \left[ \frac{\mathbf{v}^T \boldsymbol{\Sigma}^{1/2} (\mathbf{B} + \boldsymbol{\Sigma})^{-1} \mathbf{V}_2^{-1} (\mathbf{B} + \boldsymbol{\Sigma})^{-1} \boldsymbol{\Sigma}^{1/2} \mathbf{v}}{\mathbf{v}^T (\mathbf{I}_n - \mathbf{M}) \mathbf{v}} \right] + n E_{g_\psi} \left[ \frac{\boldsymbol{\beta}^T \mathbf{X}^T (\mathbf{B} + \boldsymbol{\Sigma})^{-1} \mathbf{V}_2^{-1} (\mathbf{B} + \boldsymbol{\Sigma})^{-1} \mathbf{X} \boldsymbol{\beta}}{\sigma^2 \cdot \mathbf{v}^T (\mathbf{I}_n - \mathbf{M}) \mathbf{v}} \right] \\ & = I_{51} + I_{52} \quad (\text{say}). \end{aligned}$$

Using Lemma 4.6, we can evaluate  $I_{51}$  as

$$\begin{aligned} I_{51} & = \frac{n \cdot \text{tr} [\mathbf{V}_2^{-1} (\mathbf{B} + \boldsymbol{\Sigma})^{-1} \boldsymbol{\Sigma} (\mathbf{B} + \boldsymbol{\Sigma})^{-1}]}{n - p_j - 2} - \frac{2n \cdot \text{tr} [\boldsymbol{\Sigma}^{1/2} (\mathbf{B} + \boldsymbol{\Sigma})^{-1} \mathbf{V}_2^{-1} (\mathbf{B} + \boldsymbol{\Sigma})^{-1} \boldsymbol{\Sigma}^{1/2} (\mathbf{I}_n - \mathbf{M})]}{(n - p_j)(n - p_j - 2)} \\ & = \frac{n \cdot \text{tr} [\mathbf{V}_2^{-1} (\mathbf{B} + \boldsymbol{\Sigma})^{-1} \boldsymbol{\Sigma} (\mathbf{B} + \boldsymbol{\Sigma})^{-1}]}{n - p_j - 2} - \frac{2n \cdot \text{tr} [\mathbf{V}_2^{-1} (\boldsymbol{\Sigma}^{-1} \mathbf{B} + \mathbf{I}_n)^{-1} (\boldsymbol{\Sigma}^{-1} - \mathbf{P}) (\mathbf{B} \boldsymbol{\Sigma}^{-1} + \mathbf{I}_n)^{-1}]}{(n - p_j)(n - p_j - 2)}. \end{aligned}$$

$I_{52}$  can be evaluated as

$$\begin{aligned}
I_{52} &= nE_{\beta}E_{\mathbf{y}|\beta} \left[ \frac{\boldsymbol{\beta}^T \mathbf{X}^T (\mathbf{B} + \boldsymbol{\Sigma})^{-1} \mathbf{V}_2^{-1} (\mathbf{B} + \boldsymbol{\Sigma})^{-1} \mathbf{X} \boldsymbol{\beta}}{\sigma^2 \cdot \mathbf{v}^T (\mathbf{I}_n - \mathbf{M}) \mathbf{v}} \right] \\
&= nE_{\beta} \left[ \frac{\boldsymbol{\beta}^T \mathbf{X}^T (\mathbf{B} + \boldsymbol{\Sigma})^{-1} \mathbf{V}_2^{-1} (\mathbf{B} + \boldsymbol{\Sigma})^{-1} \mathbf{X} \boldsymbol{\beta}}{\sigma^2 (n - p_j - 2)} \right] \\
&= \frac{n \cdot \text{tr} [\mathbf{X}^T (\mathbf{B} + \boldsymbol{\Sigma})^{-1} \mathbf{V}_2^{-1} (\mathbf{B} + \boldsymbol{\Sigma})^{-1} \mathbf{X} \mathbf{W}]}{n - p_j - 2} \\
&= \frac{n \cdot \text{tr} [\mathbf{V}_2^{-1} (\mathbf{B} + \boldsymbol{\Sigma})^{-1} \mathbf{B} (\mathbf{B} + \boldsymbol{\Sigma})^{-1}]}{n - p_j - 2}.
\end{aligned}$$

Then, it follows that

$$I_5 = \frac{n \cdot \text{tr} [\mathbf{V}_2^{-1} (\mathbf{B} + \boldsymbol{\Sigma})^{-1}]}{n - p_j - 2} - \frac{2n \cdot \text{tr} [\mathbf{V}_2^{-1} (\boldsymbol{\Sigma}^{-1} \mathbf{B} + \mathbf{I}_n)^{-1} (\boldsymbol{\Sigma}^{-1} - \mathbf{P}) (\mathbf{B} \boldsymbol{\Sigma}^{-1} + \mathbf{I}_n)^{-1}]}{(n - p_j)(n - p_j - 2)}$$

Thus we can obtain

$$\Delta_{\text{PI}_2} = \frac{n^2}{n - p_j - 2} - \frac{n \cdot \text{tr} [\mathbf{V}_2^{-1} (\mathbf{B} + \boldsymbol{\Sigma})^{-1}]}{n - p_j - 2} + \frac{2n \cdot \text{tr} [\mathbf{V}_2^{-1} (\boldsymbol{\Sigma}^{-1} \mathbf{B} + \mathbf{I}_n)^{-1} (\boldsymbol{\Sigma}^{-1} - \mathbf{P}) (\mathbf{B} \boldsymbol{\Sigma}^{-1} + \mathbf{I}_n)^{-1}]}{(n - p_j)(n - p_j - 2)},$$

and propose the following information criterion:

$$\text{PIC}_2 = -2 \log \{ \hat{g}^{\text{BP}}(\mathbf{y}|\mathbf{y}, \hat{\sigma}_j^2) \} + \Delta_{\text{PI}_2}. \quad (7.12)$$

**Theorem 7.2** *The  $\text{PIC}_2$  in (7.12) is an unbiased estimator of  $\text{PI}_2$  in (7.11), namely  $E_{g_{\psi}}(\text{PIC}_2) = \text{PI}_2$ .*

## 7.4 Conditional AIC variant based on Bayesian marginal likelihood

### 7.4.1 Conditional KL risk of predictive density based on Bayesian marginal likelihood

In Chapter 6, we considered the KL risk of predictive density based on Bayesian marginal likelihood. We derived criteria for variable selection in linear regression model by putting the prior on regression coefficients. The key point of the criteria is that the risk of the predictive density is measured by KL divergence not considering the prior distribution of the regression coefficients. In other words, we considered the risk of the Bayesian model (marginal likelihood) by frequentist point of view.

In this section, we consider the conditional KL risk (7.3) for the predictive density  $\hat{f}_j$  which is based on the Bayesian marginal likelihood putting the prior on the unknown parameter of interest. Especially, we assume the prior distribution of the vector of the regression coefficients and derive variable selection criteria in linear mixed model.

### 7.4.2 Case of normal prior

In this subsection, we consider normal prior on regression coefficients. We assume the prior distribution of  $\boldsymbol{\beta}$ ,

$$\pi(\boldsymbol{\beta}|\sigma^2) \sim \mathcal{N}(\mathbf{0}, \sigma^2 \mathbf{W}),$$

where  $\mathbf{W}$  is a  $p \times p$  matrix suitably chosen with full rank. Then the marginal likelihood of  $\mathbf{y}$  is

$$\begin{aligned} m_\pi(\mathbf{y}|\sigma^2) &= \int \int f(\mathbf{y}|\mathbf{b}, \boldsymbol{\beta}, \sigma^2) p(\mathbf{b}|\sigma^2) \pi(\boldsymbol{\beta}|\sigma^2) d\mathbf{b} d\boldsymbol{\beta} = \int m(\mathbf{y}|\boldsymbol{\beta}, \sigma^2) \pi(\boldsymbol{\beta}|\sigma^2) d\boldsymbol{\beta} \\ &= (2\pi\sigma^2)^{-n/2} \cdot |\boldsymbol{\Sigma}|^{-1/2} \cdot |\mathbf{W}|^{-1/2} \cdot |\mathbf{X}^\top \boldsymbol{\Sigma}^{-1} \mathbf{X} + \mathbf{W}^{-1}|^{-1/2} \cdot \exp\{-\mathbf{y}^\top \mathbf{A} \mathbf{y} / (2\sigma^2)\}, \end{aligned}$$

where  $\mathbf{A} = \boldsymbol{\Sigma}^{-1} - \boldsymbol{\Sigma}^{-1} \mathbf{X} (\mathbf{X}^\top \boldsymbol{\Sigma}^{-1} \mathbf{X} + \mathbf{W}^{-1})^{-1} \mathbf{X}^\top \boldsymbol{\Sigma}^{-1}$ . Note that  $\mathbf{A} = (\boldsymbol{\Sigma} + \mathbf{B})^{-1}$  for  $\mathbf{B} = \mathbf{X} \mathbf{W} \mathbf{X}^\top$ , namely  $m_\pi(\mathbf{y}|\sigma^2) \sim \mathcal{N}(\mathbf{0}, \sigma^2(\boldsymbol{\Sigma} + \mathbf{B}))$ . Then we consider the following information:

$$\begin{aligned} I_c(\boldsymbol{\eta}; m_\pi) &= \iiint -2 \log \{m_\pi(\tilde{\mathbf{y}}|\hat{\sigma}_j^2)\} f(\tilde{\mathbf{y}}|\mathbf{b}, \boldsymbol{\beta}, \sigma^2) f(\mathbf{y}|\mathbf{b}, \boldsymbol{\beta}, \sigma^2) p(\mathbf{b}|\sigma^2) d\tilde{\mathbf{y}} d\mathbf{y} d\mathbf{b} \\ &= \int \left[ \int -2 \log \{m_\pi(\tilde{\mathbf{y}}|\hat{\sigma}_j^2)\} \hat{f}^{\text{BP}}(\tilde{\mathbf{y}}|\mathbf{y}, \boldsymbol{\eta}) d\tilde{\mathbf{y}} \right] m(\mathbf{y}|\boldsymbol{\beta}, \sigma^2) d\mathbf{y}. \end{aligned} \quad (7.13)$$

We want to construct an information criterion as an unbiased estimator  $I_c(\boldsymbol{\eta}; m_\pi)$ , which is of the form

$$\text{IC}_{c,\pi} = -2 \log \{m_\pi(\mathbf{y}|\hat{\sigma}_j^2)\} + \Delta_{c,\pi},$$

where

$$\begin{aligned} -2 \log \{m_\pi(\mathbf{y}|\hat{\sigma}_j^2)\} &= n \log(2\pi\hat{\sigma}_j^2) + \log |\boldsymbol{\Sigma}| + \log |\mathbf{W} \mathbf{X}^\top \boldsymbol{\Sigma}^{-1} \mathbf{X} + \mathbf{I}_p| + \mathbf{y}^\top \mathbf{A} \mathbf{y} / \hat{\sigma}_j^2, \\ \Delta_{c,\pi} &= I_c(\boldsymbol{\eta}; m_\pi) - E[-2 \log \{m_\pi(\mathbf{y}|\hat{\sigma}_j^2)\}]. \end{aligned}$$

Note that the expectation in the equation above is the one with respect to the distribution of  $\mathbf{y} \sim m(\mathbf{y}|\boldsymbol{\beta}, \sigma^2)$ . Then we have to evaluate the bias correction  $\Delta_{c,\pi}$ .

Firstly, taking expectation of

$$-2 \log \{m_\pi(\tilde{\mathbf{y}}|\hat{\sigma}_j^2)\} = n \log(2\pi\hat{\sigma}_j^2) + \log |\boldsymbol{\Sigma}| + \log |\mathbf{W} \mathbf{X}^\top \boldsymbol{\Sigma}^{-1} \mathbf{X} + \mathbf{I}_p| + \tilde{\mathbf{y}}^\top \mathbf{A} \tilde{\mathbf{y}} / \hat{\sigma}_j^2$$

with respect to the distribution of  $\tilde{\mathbf{y}}|\mathbf{y} \sim \hat{f}^{\text{BP}}(\tilde{\mathbf{y}}|\mathbf{y}, \boldsymbol{\eta}) = \mathcal{N}(\mathbf{X}\boldsymbol{\beta} + \mathbf{Z}\hat{\mathbf{b}}^{\text{B}}, \sigma^2 \mathbf{V})$ , we can rewrite  $I_c(\boldsymbol{\eta}; m_\pi)$  in (7.13) as

$$\begin{aligned} I_c(\boldsymbol{\eta}; m_\pi) &= E \left[ n \log(2\pi\hat{\sigma}_j^2) + \log |\boldsymbol{\Sigma}| + \log |\mathbf{W} \mathbf{X}^\top \boldsymbol{\Sigma}^{-1} \mathbf{X} + \mathbf{I}_p| + \text{tr}(\mathbf{A} \mathbf{V}) \cdot \sigma^2 / \hat{\sigma}_j^2 \right. \\ &\quad \left. + (\mathbf{X}\boldsymbol{\beta} + \mathbf{Z}\hat{\mathbf{b}}^{\text{B}})^\top \mathbf{A} (\mathbf{X}\boldsymbol{\beta} + \mathbf{Z}\hat{\mathbf{b}}^{\text{B}}) / \hat{\sigma}_j^2 \right] \\ &= E \left[ n \log(2\pi\hat{\sigma}_j^2) + \log |\boldsymbol{\Sigma}| + \log |\mathbf{W} \mathbf{X}^\top \boldsymbol{\Sigma}^{-1} \mathbf{X} + \mathbf{I}_p| + \text{tr}(\mathbf{A} \mathbf{V}) \cdot \sigma^2 / \hat{\sigma}_j^2 \right. \\ &\quad \left. + \{\mathbf{X}\boldsymbol{\beta} + (\mathbf{I}_n - \boldsymbol{\Sigma}^{-1})\mathbf{u}\}^\top \mathbf{A} \{\mathbf{X}\boldsymbol{\beta} + (\mathbf{I}_n - \boldsymbol{\Sigma}^{-1})\mathbf{u}\} / \hat{\sigma}_j^2 \right], \end{aligned}$$

noting that  $\mathbf{Z}\hat{\mathbf{b}}^{\text{B}} = \mathbf{Z} \mathbf{G} \mathbf{Z}^\top \boldsymbol{\Sigma}^{-1} (\mathbf{y} - \mathbf{X}\boldsymbol{\beta}) = (\mathbf{I}_n - \boldsymbol{\Sigma}^{-1})\mathbf{u}$ . Next,  $-2 \log \{m_\pi(\mathbf{y}|\hat{\sigma}_j^2)\}$  is rewritten as

$$\begin{aligned} &-2 \log \{m_\pi(\mathbf{y}|\hat{\sigma}_j^2)\} \\ &= n \log(2\pi\hat{\sigma}_j^2) + \log |\boldsymbol{\Sigma}| + \log |\mathbf{W} \mathbf{X}^\top \boldsymbol{\Sigma}^{-1} \mathbf{X} + \mathbf{I}_p| + (\mathbf{u}^\top \mathbf{A} \mathbf{u} + \boldsymbol{\beta}^\top \mathbf{X}^\top \mathbf{A} \mathbf{X} \boldsymbol{\beta} + 2\mathbf{u}^\top \mathbf{A} \mathbf{X} \boldsymbol{\beta}) / \hat{\sigma}_j^2 \end{aligned}$$

Then the bias correction  $\Delta_{c,\pi}$  is

$$\begin{aligned} \Delta_{c,\pi} &= E \left[ \text{tr}(\mathbf{A} \mathbf{V}) \cdot \sigma^2 / \hat{\sigma}_j^2 + \mathbf{u}^\top (\boldsymbol{\Sigma}^{-1} \mathbf{A} \boldsymbol{\Sigma}^{-1} - 2\mathbf{A} \boldsymbol{\Sigma}^{-1}) \mathbf{u} / \hat{\sigma}_j^2 \right] \\ &= J_1 + J_2 \quad (\text{say}). \end{aligned}$$

Thus, it suffices to evaluate  $J_1$  and  $J_2$ .

Noting that  $n\hat{\sigma}_j^2/\sigma^2 \sim \chi_{n-p_j}^2$ , we can evaluate  $J_1$  as

$$J_1 = \frac{n \cdot \text{tr}(\mathbf{A}\mathbf{V})}{n - p_j - 2} = -\frac{n \times \{\text{tr}(\mathbf{A}\boldsymbol{\Sigma}^{-1}) - 2\text{tr}(\mathbf{A})\}}{n - p_j - 2}$$

Next,  $J_2$  is rewritten as

$$J_2 = n \cdot E \left[ \frac{\mathbf{v}^\top (\boldsymbol{\Sigma}^{-1/2} \mathbf{A} \boldsymbol{\Sigma}^{-1/2} - \boldsymbol{\Sigma}^{1/2} \mathbf{A} \boldsymbol{\Sigma}^{-1/2} - \boldsymbol{\Sigma}^{-1/2} \mathbf{A} \boldsymbol{\Sigma}^{1/2}) \mathbf{v}}{\mathbf{v}^\top (\mathbf{I}_n - \mathbf{M}) \mathbf{v}} \right],$$

which can be evaluated by Lemma 4.6 as

$$\begin{aligned} J_2 &= n \times \left\{ \frac{\text{tr}(\mathbf{A}\boldsymbol{\Sigma}^{-1}) - 2\text{tr}(\mathbf{A})}{n - p_j - 2} - \frac{2\text{tr}[\boldsymbol{\Sigma}^{-1/2} \mathbf{A} \boldsymbol{\Sigma}^{-1/2} (\mathbf{I}_n - \mathbf{M}) - 2\boldsymbol{\Sigma}^{1/2} \mathbf{A} \boldsymbol{\Sigma}^{-1/2} (\mathbf{I}_n - \mathbf{M})]}{(n - p_j)(n - p_j - 2)} \right\} \\ &= n \times \left\{ \frac{\text{tr}(\mathbf{A}\boldsymbol{\Sigma}^{-1}) - 2\text{tr}(\mathbf{A})}{n - p_j - 2} - \frac{2\text{tr}[\mathbf{A}(\boldsymbol{\Sigma}^{-1} - \mathbf{P})] - 4\text{tr}[\mathbf{A}(\boldsymbol{\Sigma}^{-1} - \mathbf{P})\boldsymbol{\Sigma}]}{(n - p_j)(n - p_j - 2)} \right\} \\ &= n \times \left\{ \frac{\text{tr}(\mathbf{A}\boldsymbol{\Sigma}^{-1}) - 2\text{tr}(\mathbf{A})}{n - p_j} + \frac{2\text{tr}(\mathbf{A}\mathbf{P}) - 4\text{tr}(\mathbf{A}\mathbf{P}\boldsymbol{\Sigma})}{(n - p_j)(n - p_j - 2)} \right\} \end{aligned}$$

Thus we can obtain

$$\Delta_{c,\pi} = n \times \left\{ \frac{-2\text{tr}[\mathbf{A}(\boldsymbol{\Sigma}^{-1} - \mathbf{P})] + 4\text{tr}[\mathbf{A}(\mathbf{I}_n - \mathbf{P}\boldsymbol{\Sigma})]}{(n - p_j)(n - p_j - 2)} \right\},$$

and propose the following information criterion:

$$\text{IC}_{c,\pi} = -2 \log\{m_\pi(\mathbf{y}|\hat{\sigma}_j^2)\} + \Delta_{c,\pi}. \quad (7.14)$$

**Theorem 7.3** *The information criterion  $\text{IC}_{c,\pi}$  in (7.14) is an unbiased estimator of  $I_c(\boldsymbol{\eta}; m_\pi)$  in (7.13), namely  $E(\text{IC}_{c,\pi}) = I_c(\boldsymbol{\eta}; m_\pi)$ .*

### 7.4.3 Conditional RIC

In this subsection, we assume the uniform prior for  $\boldsymbol{\beta}$ , namely  $\boldsymbol{\beta} \sim \text{uniform}(\mathbb{R}^p)$ . Although this is improper prior distribution, we can obtain the marginal likelihood function formally as follows:

$$\begin{aligned} m_r(\mathbf{y}|\sigma^2) &= \int m(\mathbf{y}|\boldsymbol{\beta}, \sigma^2) d\boldsymbol{\beta} \\ &= (2\pi\sigma^2)^{-(n-p)/2} \cdot |\boldsymbol{\Sigma}|^{-1/2} \cdot |\mathbf{X}^\top \boldsymbol{\Sigma}^{-1} \mathbf{X}|^{-1/2} \cdot \exp\{\mathbf{y}^\top (\boldsymbol{\Sigma}^{-1} - \mathbf{P}) \mathbf{y} / (2\sigma^2)\}, \end{aligned}$$

which is the same as the residual likelihood (Patterson and Thompson, 1971). In the last chapter, we measured the risk of predictive density  $m_r(\tilde{\mathbf{y}}|\tilde{\sigma}_j^2)$  in terms of the following KL divergence:

$$R(\boldsymbol{\eta}; m_r) = \int \left[ \int \log \left\{ \frac{m(\tilde{\mathbf{y}}|\boldsymbol{\beta}, \sigma^2)}{m_r(\tilde{\mathbf{y}}|\tilde{\sigma}_j^2)} \right\} m(\tilde{\mathbf{y}}|\boldsymbol{\beta}, \sigma^2) d\tilde{\mathbf{y}} \right] m(\mathbf{y}|\boldsymbol{\beta}, \sigma^2) d\mathbf{y},$$

and showed that the resulting criterion is identical to residual information criterion (RIC) proposed by Shi and Tsai (2002). Azari et al. (2006) applied the RIC to variable selection in linear mixed model for longitudinal data analysis. However, the KL risk  $R(\boldsymbol{\eta}; m_r)$  is not appropriate

when one is interested in predicting random effects, which was pointed out by Vaida and Blanchard (2005). Then we propose to measure the prediction risk of  $m_r(\tilde{\mathbf{y}}|\tilde{\sigma}_j^2)$  by the conditional KL risk  $R_c(\boldsymbol{\eta}; \hat{f}_j)$  in (7.3) for  $\hat{f}_j = m_r(\mathbf{y}|\tilde{\sigma}_j)$  and call the resulting criterion the conditional RIC (cRIC).

Then we consider the following information:

$$\begin{aligned} I_c(\boldsymbol{\eta}; m_r) &= \iiint -2 \log\{m_r(\tilde{\mathbf{y}}|\tilde{\sigma}_j^2)\} f(\tilde{\mathbf{y}}|\mathbf{b}, \boldsymbol{\beta}, \sigma^2) f(\mathbf{y}|\mathbf{b}, \boldsymbol{\beta}, \sigma^2) p(\mathbf{b}|\sigma^2) d\tilde{\mathbf{y}} d\mathbf{y} d\mathbf{b} \\ &= \int \left[ \int -2 \log\{m_r(\tilde{\mathbf{y}}|\tilde{\sigma}_j^2)\} \hat{f}^{\text{BP}}(\tilde{\mathbf{y}}|\mathbf{y}, \boldsymbol{\eta}) d\tilde{\mathbf{y}} \right] m(\mathbf{y}|\boldsymbol{\beta}, \sigma^2) d\mathbf{y} \\ &= \text{cRI} \quad (\text{say}), \end{aligned} \tag{7.15}$$

where  $\tilde{\sigma}_j^2 = \mathbf{y}^T(\boldsymbol{\Sigma}^{-1} - \mathbf{P})\mathbf{y}/(n - p_j)$  is the residual maximum likelihood (REML) estimator of  $\sigma_j^2$ . We want to construct an information criterion as an unbiased estimator of  $I_c(\boldsymbol{\eta}; m_r)$ , which is of the form

$$\text{cRIC} = -2 \log\{m_r(\mathbf{y}|\tilde{\sigma}_j^2)\} + \Delta_{\text{cRI}},$$

where

$$\begin{aligned} -2 \log\{m_r(\mathbf{y}|\tilde{\sigma}_j^2)\} &= (n - p_j) \log(2\pi\tilde{\sigma}_j^2) + \log|\boldsymbol{\Sigma}| + \log|\mathbf{X}^T\boldsymbol{\Sigma}^{-1}\mathbf{X}| + \mathbf{y}^T(\boldsymbol{\Sigma}^{-1} - \mathbf{P})\mathbf{y}/\tilde{\sigma}_j^2, \\ \Delta_{\text{cRI}} &= \text{cRI} - E[-2 \log\{m_r(\mathbf{y}|\tilde{\sigma}_j^2)\}]. \end{aligned}$$

Note that the expectation in the equation above is the one with respect to the distribution of  $\mathbf{y} \sim m(\mathbf{y}|\boldsymbol{\beta}, \sigma^2)$ . Then we have to evaluate the bias correction  $\Delta_{\text{cRI}}$ .

Firstly, taking expectation of

$$\begin{aligned} &-2 \log\{m_r(\tilde{\mathbf{y}}|\tilde{\sigma}_j^2)\} \\ &= (n - p_j) \log(2\pi\tilde{\sigma}_j^2) + \log|\boldsymbol{\Sigma}| + \log|\mathbf{X}^T\boldsymbol{\Sigma}^{-1}\mathbf{X}| + \tilde{\mathbf{y}}^T(\boldsymbol{\Sigma}^{-1} - \mathbf{P})\tilde{\mathbf{y}}/\tilde{\sigma}_j^2, \end{aligned}$$

with respect to the distribution of  $\hat{f}^{\text{BP}}(\tilde{\mathbf{y}}|\mathbf{y}, \boldsymbol{\eta}) \sim \mathcal{N}(\mathbf{X}\boldsymbol{\beta} + \mathbf{Z}\hat{\mathbf{b}}^{\text{B}}, \sigma^2\mathbf{V})$ , we can rewrite the cRI in (7.15) as

$$\begin{aligned} \text{cRI} &= E \left[ (n - p_j) \log(2\pi\tilde{\sigma}_j^2) + \log|\boldsymbol{\Sigma}| + \log|\mathbf{X}^T\boldsymbol{\Sigma}^{-1}\mathbf{X}| + \text{tr}[\mathbf{V}(\boldsymbol{\Sigma}^{-1} - \mathbf{P})] \cdot \sigma^2/\tilde{\sigma}_j^2 \right. \\ &\quad \left. + \{\mathbf{X}\boldsymbol{\beta} + (\mathbf{I}_n - \boldsymbol{\Sigma}^{-1})\mathbf{u}\}^T(\boldsymbol{\Sigma}^{-1} - \mathbf{P})\{\mathbf{X}\boldsymbol{\beta} + (\mathbf{I}_n - \boldsymbol{\Sigma}^{-1})\mathbf{u}\}/\tilde{\sigma}_j^2 \right] \\ &= E \left[ (n - p_j) \log(2\pi\tilde{\sigma}_j^2) + \log|\boldsymbol{\Sigma}| + \log|\mathbf{X}^T\boldsymbol{\Sigma}^{-1}\mathbf{X}| + \text{tr}[\mathbf{V}(\boldsymbol{\Sigma}^{-1} - \mathbf{P})] \cdot \sigma^2/\tilde{\sigma}_j^2 \right. \\ &\quad \left. + \mathbf{u}^T(\mathbf{I}_n - \boldsymbol{\Sigma}^{-1})(\boldsymbol{\Sigma}^{-1} - \mathbf{P})(\mathbf{I}_n - \boldsymbol{\Sigma}^{-1})\mathbf{u}/\tilde{\sigma}_j^2 \right], \end{aligned}$$

noting that  $(\boldsymbol{\Sigma}^{-1} - \mathbf{P})\mathbf{X} = \mathbf{0}$ . From the fact that

$$\begin{aligned} &-2 \log\{m_r(\mathbf{y}|\tilde{\sigma}_j^2)\} \\ &= (n - p_j) \log(2\pi\tilde{\sigma}_j^2) + \log|\boldsymbol{\Sigma}| + \log|\mathbf{X}^T\boldsymbol{\Sigma}^{-1}\mathbf{X}| + (n - p_j), \end{aligned}$$

the bias correction can be reduced to

$$\begin{aligned} \Delta_{\text{cRI}} &= E \left[ \text{tr}[\mathbf{V}(\boldsymbol{\Sigma}^{-1} - \mathbf{P})] \cdot \sigma^2/\tilde{\sigma}_j^2 + \mathbf{u}^T(\mathbf{I}_n - \boldsymbol{\Sigma}^{-1})(\boldsymbol{\Sigma}^{-1} - \mathbf{P})(\mathbf{I}_n - \boldsymbol{\Sigma}^{-1})\mathbf{u}/\tilde{\sigma}_j^2 \right] - (n - p_j) \\ &= J_3 + J_4 - (n - p_j), \quad (\text{say}). \end{aligned}$$

Thus, it suffices to evaluate  $J_3$  and  $J_4$ .

It is easy to see that

$$J_3 = \frac{(n - p_j) \cdot \text{tr}[\mathbf{V}(\boldsymbol{\Sigma}^{-1} - \mathbf{P})]}{n - p_j - 2} = \frac{(n - p_j) \cdot \text{tr}[(2\mathbf{I}_n - \boldsymbol{\Sigma}^{-1})(\boldsymbol{\Sigma}^{-1} - \mathbf{P})]}{n - p_j - 2},$$

because  $(n - p_j)\tilde{\sigma}_j^2/\sigma^2 \sim \chi_{n-p_j}^2$ . Next,  $J_4$  is rewritten as

$$J_4 = (n - p_j) \cdot E \left[ \frac{\mathbf{v}^\top \boldsymbol{\Sigma}^{1/2} (\mathbf{I}_n - \boldsymbol{\Sigma}^{-1}) (\boldsymbol{\Sigma}^{-1} - \mathbf{P}) (\mathbf{I}_n - \boldsymbol{\Sigma}^{-1}) \boldsymbol{\Sigma}^{1/2} \mathbf{v}}{\mathbf{v}^\top (\mathbf{I}_n - \mathbf{M}) \mathbf{v}} \right],$$

which can be evaluated by Lemma 4.6 as

$$\begin{aligned} J_4 &= (n - p_j) \\ &\times \left\{ \frac{\text{tr}[(\boldsymbol{\Sigma}^{-1} - \mathbf{P})(\mathbf{I}_n - \boldsymbol{\Sigma}^{-1})\boldsymbol{\Sigma}(\mathbf{I}_n - \boldsymbol{\Sigma}^{-1})]}{n - p_j - 2} - \frac{2\text{tr}[(\mathbf{I}_n - \boldsymbol{\Sigma}^{-1})(\mathbf{I}_n - \mathbf{M})(\mathbf{I}_n - \boldsymbol{\Sigma}^{-1})(\mathbf{I}_n - \mathbf{M})]}{(n - p_j)(n - p_j - 2)} \right\} \\ &= (n - p_j) \\ &\times \left\{ \frac{\text{tr}[(\boldsymbol{\Sigma}^{-1} - \mathbf{P})(\boldsymbol{\Sigma} - 2\mathbf{I}_n + \boldsymbol{\Sigma}^{-1})]}{n - p_j - 2} - \frac{2\text{tr}[(\boldsymbol{\Sigma} - \mathbf{I}_n)(\boldsymbol{\Sigma}^{-1} - \mathbf{P})(\boldsymbol{\Sigma} - \mathbf{I}_n)(\boldsymbol{\Sigma}^{-1} - \mathbf{P})]}{(n - p_j)(n - p_j - 2)} \right\}. \end{aligned}$$

Thus we can obtain

$$\Delta_{\text{cRI}} = \frac{2(n - p_j)}{n - p_j - 2} - \frac{2}{n - p_j - 2} \text{tr}[(\boldsymbol{\Sigma} - \mathbf{I}_n)(\boldsymbol{\Sigma}^{-1} - \mathbf{P})(\boldsymbol{\Sigma} - \mathbf{I}_n)(\boldsymbol{\Sigma}^{-1} - \mathbf{P})]$$

and propose the following cRIC:

$$\text{cRIC} = -2 \log\{m_r(\mathbf{y}|\tilde{\sigma}_j^2)\} + \Delta_{\text{cRI}}. \quad (7.16)$$

**Theorem 7.4** *The cRIC in (7.16) is an unbiased estimator of cRI in (7.15), namely  $E(\text{cRIC}) = \text{cRI}$ .*

## 7.5 Simulations

In this section, we compare the numerical performance of the proposed criteria, PIC, PIC<sub>2</sub>, IC<sub>c,π</sub> and cRIC with the conventional cAIC of Vaida and Blanchard (2005). We handle the nested error regression model (NERM) and consider the same setting as that of Section 6.4. When we derive the criteria PIC<sub>2</sub> and IC<sub>c,π</sub>, we set the prior distribution of  $\boldsymbol{\beta}$  as  $\mathcal{N}_p(\mathbf{0}, \sigma^2 \lambda^{-1} \mathbf{I}_p)$ , namely  $\mathbf{W} = \lambda^{-1} \mathbf{I}_p$ . The hyperparameter  $\lambda$  is estimated by maximizing the marginal likelihood  $m_\pi(\mathbf{y}|\hat{\sigma}_j^2)$ , where the estimate  $\hat{\sigma}_j^2 = \mathbf{y}^\top (\boldsymbol{\Sigma}^{-1} - \mathbf{P}) \mathbf{y} / n$  of  $\sigma^2$  is plugged in. The unknown parameter  $\phi = \tau^2 / \sigma^2$  included by  $\boldsymbol{\Sigma}$  is estimated by consistent estimator based on the full model in the same way as Section 6.4. The class of the candidate models includes all the subsets of the full model and select the model by the criteria. The performance of the criteria is measured by the number of selecting the true model and the prediction error of the selected model based on quadratic loss, namely  $\|\mathbf{X}(\hat{j})\hat{\boldsymbol{\beta}}_{\hat{j}} - \mathbf{X}(\omega)\boldsymbol{\beta}_*\|^2/n$ .

Table 7.1 and 7.2 show the number of selecting the true model by the criteria and the average prediction error of the selected model by each criterion, respectively. From the table we can see the following facts. Firstly, the number of selecting the true model approaches 1000 for IC<sub>c,π</sub> and cRIC, which is the numerical evidence of the consistency of the criteria. Especially, IC<sub>c,π</sub> performs well for almost all situations in terms of prediction error as well as selecting the true model. However, for the case of small sample size and noisy data, namely signal-to-noise ratio (SNR) is small, PIC<sub>2</sub> performs the best among the criteria.

Table 7.1: The number of selecting the true model by the criteria in 1000 realizations

		$\phi = 0.5$			$\phi = 1$			$\phi = 2$		
SNR		1	3	5	1	3	5	1	3	5
$n_0 = 5$	cAIC	121	718	732	140	694	729	175	677	740
$m = 4$	PIC	124	706	715	157	683	703	194	673	713
	PIC <sub>2</sub>	155	555	558	171	550	557	201	539	559
	IC <sub>c,<math>\pi</math></sub>	146	850	905	176	849	913	237	849	937
	cRIC	20	531	730	41	590	758	96	679	804
$n_0 = 5$	cAIC	355	631	631	357	645	645	337	652	653
$m = 8$	PIC	364	602	602	359	593	593	359	601	601
	PIC <sub>2</sub>	330	479	479	322	487	489	330	495	496
	IC <sub>c,<math>\pi</math></sub>	427	894	938	448	918	952	489	934	960
	cRIC	153	702	850	234	754	863	358	802	895
$n_0 = 5$	cAIC	576	625	625	542	630	630	492	629	629
$m = 16$	PIC	561	596	596	535	595	595	507	596	596
	PIC <sub>2</sub>	475	509	511	452	505	507	438	514	515
	IC <sub>c,<math>\pi</math></sub>	736	938	958	720	947	963	699	955	971
	cRIC	436	832	900	498	856	921	578	893	935

Table 7.2: The prediction error of the best model selected by the criteria

		$\phi = 0.5$			$\phi = 1$			$\phi = 2$		
SNR		1	3	5	1	3	5	1	3	5
$n_0 = 5$	cAIC	1.44	0.120	0.0430	1.11	0.0962	0.0343	0.708	0.0652	0.0233
$m = 4$	PIC	1.43	0.121	0.0433	1.10	0.0968	0.0346	0.698	0.0658	0.0235
	PIC <sub>2</sub>	1.35	0.127	0.0455	1.05	0.101	0.0362	0.675	0.0689	0.0247
	IC <sub>c,<math>\pi</math></sub>	1.41	0.115	0.0400	1.09	0.0915	0.0320	0.705	0.0615	0.0214
	cRIC	1.28	0.127	0.0431	1.01	0.100	0.0341	0.679	0.0664	0.0229
$n_0 = 5$	cAIC	0.846	0.0854	0.0307	0.638	0.0672	0.0242	0.418	0.0455	0.0164
$m = 8$	PIC	0.845	0.0859	0.0309	0.638	0.0679	0.0244	0.422	0.0460	0.0166
	PIC <sub>2</sub>	0.833	0.0880	0.0317	0.644	0.0695	0.0250	0.429	0.0472	0.0170
	IC <sub>c,<math>\pi</math></sub>	0.842	0.0793	0.0280	0.628	0.0621	0.0220	0.405	0.0418	0.0149
	cRIC	0.840	0.0840	0.0289	0.655	0.0656	0.0228	0.437	0.0439	0.0153
$n_0 = 5$	cAIC	0.564	0.0622	0.0224	0.459	0.0508	0.0183	0.322	0.0356	0.0128
$m = 16$	PIC	0.566	0.0625	0.0225	0.461	0.0510	0.0184	0.323	0.0358	0.0129
	PIC <sub>2</sub>	0.572	0.0633	0.0228	0.466	0.0516	0.0186	0.327	0.0362	0.0130
	IC <sub>c,<math>\pi</math></sub>	0.547	0.0582	0.0208	0.443	0.0476	0.0170	0.308	0.0334	0.0120
	cRIC	0.577	0.0599	0.0212	0.466	0.0488	0.0173	0.323	0.0340	0.0121

## 7.6 Proofs

### 7.6.1 Proof of Proposition 7.1

The marginal joint distribution of  $\tilde{\mathbf{y}}$  and  $\mathbf{y}$  is

$$\begin{pmatrix} \tilde{\mathbf{y}} \\ \mathbf{y} \end{pmatrix} \sim \mathcal{N} \left( \begin{pmatrix} \mathbf{X}\boldsymbol{\beta} \\ \mathbf{X}\boldsymbol{\beta} \end{pmatrix}, \begin{pmatrix} \boldsymbol{\Lambda}_{11} & \boldsymbol{\Lambda}_{12} \\ \boldsymbol{\Lambda}_{21} & \boldsymbol{\Lambda}_{22} \end{pmatrix} \right),$$

where  $\boldsymbol{\Lambda}_{11} = \boldsymbol{\Lambda}_{22} = \sigma^2 \boldsymbol{\Sigma}$ ,

$$\begin{aligned} \boldsymbol{\Lambda}_{12} &= \boldsymbol{\Lambda}_{21} \\ &= E[(\mathbf{y} - \mathbf{X}\boldsymbol{\beta})(\tilde{\mathbf{y}} - \mathbf{X}\boldsymbol{\beta})^\top] \\ &= E[(\mathbf{Z}\mathbf{b} + \boldsymbol{\varepsilon})(\mathbf{Z}\mathbf{b} + \tilde{\boldsymbol{\varepsilon}})^\top] \\ &= \sigma^2 \mathbf{Z}\mathbf{G}\mathbf{Z}^\top, \end{aligned}$$

and  $\tilde{\boldsymbol{\varepsilon}}$  is independent replication of  $\boldsymbol{\varepsilon}$ , which is also independent of  $\mathbf{b}$ . Then, from the property of multivariate normal distribution, it follows that

$$\hat{f}^{\text{BP}}(\tilde{\mathbf{y}}|\mathbf{y}, \boldsymbol{\eta}) = \mathcal{N}(\mathbf{X}\boldsymbol{\beta} + \boldsymbol{\Lambda}_{12}\boldsymbol{\Lambda}_{22}^{-1}(\mathbf{y} - \mathbf{X}\boldsymbol{\beta}), \boldsymbol{\Lambda}_{11} - \boldsymbol{\Lambda}_{12}\boldsymbol{\Lambda}_{22}^{-1}\boldsymbol{\Lambda}_{21}),$$

where  $\boldsymbol{\Lambda}_{12}\boldsymbol{\Lambda}_{22}^{-1}(\mathbf{y} - \mathbf{X}\boldsymbol{\beta}) = \mathbf{Z}\mathbf{G}\mathbf{Z}^\top\boldsymbol{\Sigma}^{-1}(\mathbf{y} - \mathbf{X}\boldsymbol{\beta}) = \mathbf{Z}\hat{\mathbf{b}}^{\text{B}}$  and

$$\begin{aligned} \boldsymbol{\Lambda}_{11} - \boldsymbol{\Lambda}_{12}\boldsymbol{\Lambda}_{22}^{-1}\boldsymbol{\Lambda}_{21} &= \sigma^2\boldsymbol{\Sigma} - \sigma^2\mathbf{Z}\mathbf{G}\mathbf{Z}^\top\boldsymbol{\Sigma}^{-1}\mathbf{Z}\mathbf{G}\mathbf{Z}^\top \\ &= \sigma^2\boldsymbol{\Sigma} - \sigma^2(\boldsymbol{\Sigma} - \mathbf{I}_n)\boldsymbol{\Sigma}^{-1}(\boldsymbol{\Sigma} - \mathbf{I}_n) \\ &= \sigma^2(2\mathbf{I}_n - \boldsymbol{\Sigma}^{-1}) = \sigma^2\mathbf{V}, \end{aligned}$$

which shows Proposition 7.1. □

### 7.6.2 Proof of Proposition 7.2

The marginal joint distribution of  $\mathbf{z}$  and  $\mathbf{y}$  is

$$\begin{pmatrix} \mathbf{z} \\ \mathbf{y} \end{pmatrix} \sim \mathcal{N} \left( \begin{pmatrix} \mathbf{0} \\ \mathbf{0} \end{pmatrix}, \begin{pmatrix} \boldsymbol{\Gamma}_{11} & \boldsymbol{\Gamma}_{12} \\ \boldsymbol{\Gamma}_{21} & \boldsymbol{\Gamma}_{22} \end{pmatrix} \right),$$

where  $\boldsymbol{\Gamma}_{11} = \boldsymbol{\Gamma}_{22} = \sigma^2(\mathbf{B} + \mathbf{Z}\mathbf{G}\mathbf{Z}^\top + \mathbf{I}_n) = \sigma^2(\mathbf{C} + \mathbf{I}_n)$ ,

$$\begin{aligned} \boldsymbol{\Gamma}_{12} &= \boldsymbol{\Gamma}_{21} \\ &= E[\mathbf{y}\mathbf{z}^\top] \\ &= \sigma^2\mathbf{C}, \end{aligned}$$

and  $\mathbf{C} = \mathbf{B} + \mathbf{Z}\mathbf{G}\mathbf{Z}^\top$ . Then, from the property of multivariate normal distribution, it follows that

$$\hat{g}^{\text{BP}} = \mathcal{N}(\boldsymbol{\Gamma}_{12}\boldsymbol{\Gamma}_{22}^{-1}\mathbf{y}, \boldsymbol{\Gamma}_{11} - \boldsymbol{\Gamma}_{12}\boldsymbol{\Gamma}_{22}^{-1}\boldsymbol{\Gamma}_{21}),$$

where

$$\begin{aligned} \boldsymbol{\Gamma}_{12}\boldsymbol{\Gamma}_{22}^{-1} &= \mathbf{C}(\mathbf{C} + \mathbf{I}_n)^{-1} \\ &= (\mathbf{C} + \mathbf{I}_n - \mathbf{I}_n)(\mathbf{C} + \mathbf{I}_n)^{-1} \\ &= \mathbf{I}_n - (\mathbf{C} + \mathbf{I}_n)^{-1} \\ &= \mathbf{I}_n - (\mathbf{B} + \boldsymbol{\Sigma})^{-1}, \end{aligned}$$



and

$$\begin{aligned}
\mathbf{\Gamma}_{11} - \mathbf{\Gamma}_{12}\mathbf{\Gamma}_{22}^{-1}\mathbf{\Gamma}_{21} &= \sigma^2(\mathbf{C} + \mathbf{I}_n) - \sigma^2\mathbf{C}(\mathbf{C} + \mathbf{I}_n)^{-1}\mathbf{C} \\
&= \sigma^2(\mathbf{C} + \mathbf{I}_n) - \sigma^2(\mathbf{C} + \mathbf{I}_n - \mathbf{I}_n)(\mathbf{C} + \mathbf{I}_n)^{-1}(\mathbf{C} + \mathbf{I}_n - \mathbf{I}_n) \\
&= \sigma^2\{2\mathbf{I}_n - (\mathbf{C} + \mathbf{I}_n)^{-1}\} \\
&= \sigma^2\{2\mathbf{I}_n - (\mathbf{B} + \mathbf{\Sigma})^{-1}\} = \sigma^2\mathbf{V}_2,
\end{aligned}$$

which shows Proposition 7.2.



# Acknowledgment

I thank my supervisor Prof. Tatsuya Kubokawa for all of his help and lots of encouragement to me during the master's and doctoral courses. I also thank Prof. Malay Ghosh, Prof. Muni Srivastava and Prof. J.N.K. Rao for their valuable comments and suggestions. Prof. Naoto Kunitomo, Prof. Yoshihiro Yajima, Prof. Yasuhiro Omori, Prof. Katsumi Shimotsu, Prof. Akimichi Takemura, Prof. Fumiyasu Komaki, Dr. Yuzo Maruyama and Dr. Kengo Kato gave me many helpful comments about my research. Mr. Shinichiro Shirota, Mr. Masaaki Imaizumi and Mr. Shonosuke Sugasawa made my student life very meaningful and enjoyable in various forms. I am also grateful for the discussion with Dr. Genya Kobayashi, Dr. Kota Ogawasara and all of my academic colleagues. I was financially supported by JSPS. Lastly, I thank my father, mother, sister, grand parents and their families for their support.



# Bibliography

- Aitchison, J. (1975). Goodness of prediction fit. *Biometrika*, **62**, 547–554.
- Akaike, H. (1973). Information theory and an extension of the maximum likelihood principle. In *2nd International Symposium on Information Theory*, (B.N. Petrov and Csaki, F, eds.), 267–281, Akademia Kiado, Budapest.
- Akaike, H. (1974). A new look at the statistical model identification. System identification and time-series analysis. *IEEE Transactions on Automatic Control*, **AC-19**, 716–723.
- Akaike, H. (1980a). On the use of predictive likelihood of a Gaussian model. *Annals of the Institute of Statistical Mathematics*, **32**, 311–324.
- Akaike, H. (1980b). Likelihood and the Bayes procedure. In *Bayesian Statistics*, (N.J. Bernard, M.H. Degroot, D.V. Lindaley and A.F.M. Simith, eds.), Valencia, Spain, University Press, 141–166.
- Ando, T. (2007). Bayesian predictive information criterion for the evaluation of hierarchical Bayesian and empirical Bayes models. *Biometrika*, **94**, 443–458.
- Azari, R., Li, L. and Tsai, C.-L. (2006). Longitudinal data model selection. *Computational Statistics and Data Analysis*, **50**, 3053–3066.
- Battese, G.E., Harter, R.M. and Fuller, W.A. (1988). An error-components model for prediction of county crop areas using survey and satellite data. *Journal of the American Statistical Association*, **83**, 28–36.
- Berger, J.O. and Pericchi, L.R. (1996). The intrinsic Bayes factor for model selection and prediction. *Journal of the American Statistical Association*, **91**, 109–122.
- Bozdogan, H. (1987). Model selection and Akaike’s information criterion (AIC): The general theory and its analytical extensions. *Psychometrika*, **52**, 345–370.
- Burnham, K. P. and Anderson, D. R. (2002). *Model Selection and Multimodel Inference: A Practical Information Theoretic Approach*, 2nd ed. New York: Springer.
- Casella, G. and Moreno, E. (2006). Objective Bayesian variable selection. *Journal of the American Statistical Association*, **101**, 157–167.
- Datta, G. and Ghosh, M. (2012). Small area shrinkage estimation. *Statistical Science*, **27**, 95–114.
- Donohue, M. C., Overholser, R., Xu, R., and Vaida, F. (2011). Conditional Akaike information under generalized linear and proportional hazards mixed models. *Biometrika*, **98**, 685–700.

- Fay, R.E. and Herriot, R.A. (1979). Estimates of income for small places: An application of James–Stein procedures to census data. *Journal of the American Statistical Association*, **74**, 269–277.
- Fujikoshi, Y. and Satoh, K. (1997). Modified AIC and  $C_p$  in multivariate linear regression. *Biometrika*, **84**, 707–716.
- Ghosh, M. and Maiti, T. (2004). Small-area estimation based on natural exponential family quadratic variance function models and survey weights. *Biometrika*, **91**, 95–112.
- Ghosh, M. and Maiti, T. (2008). Empirical Bayes confidence intervals for means of natural exponential family-quadratic variance function distributions with application to small area estimation. *Scandinavian Journal of Statistics*, **35**, 484–495.
- Godambe, V.P. and Thompson, M.E. (1989). An extension of quasi-likelihood estimation (with Discussion). *Journal of Statistical Planning and Inference*, **22**, 137–152.
- Greven, S., and Kneib, T. (2010). On the behaviour of marginal and conditional AIC in linear mixed models. *Biometrika*, **97**, 773–789.
- Henderson, C.R. (1950). Estimation of genetic parameters. *The Annals of Mathematical Statistics*, **21**, 309–310.
- Hodges, J.S. and Sargent, D.J. (2001). Counting degrees of freedom in hierarchical and the richly-parameterised models. *Biometrika*, **88**, 367–379.
- Jiang, J., Rao, J.S., Gu, Z. and Nguyen, T. (2008). Fence methods for mixed model selection. *The Annals of Statistics*, **36**, 1669–1692.
- Hansen, B.E. (2007). Least squares model averaging. *Econometrica*, **75**, 1175–1189.
- Hurvich, C.M. and Tsai, C.-L. (1989). Regression and time series model selection in small samples. *Biometrika*, **76**, 297–307.
- Kawakubo, Y. and Kubokawa, T. (2014). Modified conditional AIC in linear mixed models. *Journal of Multivariate Analysis*, **129**, 44–56.
- Kitagawa, G. (1997). Information criteria for the predictive evaluation of Bayesian models. *Communications in Statistics — Theory and Methods*, **26**, 2223–2246.
- Kubokawa, T. (2011). Conditional and unconditional methods for selecting variables in linear mixed model. *Journal of Multivariate Analysis*, **102**, 641–660.
- Kubokawa, T., Hasukawa, M. and Takahashi, K. (2014). On measuring uncertainty of benchmarked predictors with application to disease risk estimate. *Scandinavian Journal of Statistics*, **41**, 394–413.
- Kubokawa, T. and Nagashima, B. (2012). Parametric Bootstrap methods for bias correction in linear mixed models. *Journal of Multivariate Analysis*, **106**, 1–16.
- Liang, H., Wu, H. and Zou, G. (2008). A note on conditional AIC for linear mixed-effects models. *Biometrika*, **95**, 773–778.
- Lohr, S.L. and Rao, J.N.K. (2009). Jackknife estimation of mean squared error of small area predictors in nonlinear mixed models. *Biometrika*, **96**, 457–468.

- Mallows, C.L. (1973). Some comments on  $C_p$ . *Technometrics*, **15**, 661–675.
- Morris, C. (1982). Natural exponential families with quadratic variance functions. *The Annals of Statistics*, **10**, 65–80.
- Morris, C. (1983). Natural exponential families with quadratic variance functions: statistical theory. *The Annals of Statistics*, **11**, 515–529.
- Müller, S., Scealy, J.L. and Welsh, A.H. (2013). Model selection in linear mixed models. *Statistical Science*, **28**, 135–167.
- Nishii, R. (1984). Asymptotic properties of criteria for selection of variables in multiple regression. *The Annals of Statistics*, **12**, 758–765.
- Patterson, H.D. and Thompson, R. (1971). Recovery of inter-block information when block sizes are unequal. *Biometrika*, **58**, 545–554.
- Pfeffermann, D. (2013). New important developments in small area estimation. *Statistical Science*, **28**, 40–68.
- Prasad, N.G.N. and Rao, J.N.K. (1990). The estimation of the mean squared error of small-area estimators. *Journal of the American Statistical Association*, **85**, 163–171.
- Rao, J.N.K. and Molina, I. (2015). *Small Area Estimation*, Wiley.
- Saefken, B., Kneib, T., van Waveren, C.S. and Greven, S. (2014). A unifying approach to the estimation of the conditional Akaike information in generalized linear mixed models. *Electronic Journal of Statistics*, **8**, 201–225.
- Satoh, K. (1997). AIC-type model selection criterion for multivariate linear regression with a future experiment. *Journal of the Japan Statistical Society*, **27**, 135–140.
- Schwarz, G. (1978). Estimating the dimension of a model. *The Annals of Statistics*, **6**, 461–464.
- Shao, J. (1997). An asymptotic theory for linear model selection. *Statistica Sinica*, **7**, 221–264.
- Shi, P. and Tsai, C.-L. (2002). Regression model selection—a residual likelihood approach. *Journal of the Royal Statistical Society series B*, **64**, 237–252.
- Shibata, R. (1981). An optimal selection of regression variables. *Biometrika*, **68**, 45–54.
- Shimodaira, H. (2000). Improving predictive inference under covariate shift by weighting the log-likelihood function. *Journal of Statistical Planning and Inference*, **90**, 227–244.
- Spiegelhalter, D.J., Best, N.G., Carlin, B.P. and van der Linde, A. (2002). Bayesian measures of model complexity and fit. *Journal of the Royal Statistical Society series B*, **64**, 583–639.
- Srivastava, M.S. and Kubokawa, T. (2010). Conditional information criteria for selecting variables in linear mixed models. *Journal of Multivariate Analysis*, **101**, 1970–1980.
- Sugiura, N. (1978). Further analysis of the data by Akaike’s information criterion and the finite corrections. *Communications in Statistics — Theory and Methods*, **7**, 13–26.
- Vaida, F. and Blanchard, S. (2005). Conditional Akaike information for mixed-effects models. *Biometrika*, **92**, 351–370.

- Wedderburn, R.W.M. (1974). Quasi-likelihood functions, generalized linear models, and the Gauss–Newton method. *Biometrika*, **61**, 439–447.
- Ye, J. (1998). On measuring and correcting the effects of data mining and model selection. *Journal of the American Statistical Association*, **93**, 120–131.
- Yu, D. and Yau, K.K.W. (2012). Conditional Akaike information criterion of generalized linear mixed models. *Computational Statistics and Data Analysis*, **56**, 629–644.
- Yu, D., Zhang, X. and Yau, K.K.W. (2013). Information based model selection criteria for generalized linear mixed models with unknown variance component parameters. *Journal of Multivariate Analysis*, **116**, 245–262.
- Zellner, A. (1986). On assessing prior distributions and Bayesian regression analysis with g-prior distributions. In *Bayesian Inference and Decision Techniques: Essays in Honor of Bruno de Finetti*, (P.K. Goel and A. Zellner, eds.), pp. 233–243, Amsterdam: North-Holland/Elsevier.
- Zhang, X., Zou, G. and Liang, H. (2014). Model averaging and weight choice in linear mixed-effects models. *Biometrika*, **101**, 205–218.