

Some Asymptotics for Geostatistical Model Selection

Zhi-Hao Zhang¹, Hsin-Cheng Huang^{2,*} and Ching-Kang Ing²

¹ 1001 Ta Hsueh Road, Institute of Statistics, Hsinchu 300, Taiwan; jhow@stat.sinica.edu.tw

² 128 Section 2 Academia Road, Institute of Statistical Science, Academia Sinica, Taipei 115, Taiwan; hchuang@stat.sinica.edu.tw, cking@stat.sinica.edu.tw

*Corresponding author

Abstract. Information criteria such as AIC and BIC are often used for model selection. However, their asymptotic behaviors in geostatistical model selection have not been well studied. In this article, we provide some asymptotic results for the generalized information criterion, including both AIC and BIC.

Keywords. Akaike information criterion; Asymptotic efficiency; Bayesian information criterion; Consistency; Variable selection.

1 Geostatistical Models

Consider a spatial process $\{S(\mathbf{s}) : \mathbf{s} \in D\}$ of interest defined over a d -dimensional region $D \subset \mathbb{R}^d$. Suppose that we observe data $\mathbf{Z} = (Z(\mathbf{s}_1), \dots, Z(\mathbf{s}_n))'$ at locations $\mathbf{s}_1, \dots, \mathbf{s}_n \in D$, according to the measurement equation:

$$Z(\mathbf{s}_i) = S(\mathbf{s}_i) + \varepsilon(\mathbf{s}_i); \quad i = 1, \dots, n,$$

where $\varepsilon(\mathbf{s}_1), \dots, \varepsilon(\mathbf{s}_n) \sim N(0, \sigma_\varepsilon^2)$ are white noise variables representing measurement errors, and are independent of the process $S(\cdot)$. In addition, we observe p explanatory variables, $\mathbf{x}_i = (x_1(\mathbf{s}_i), \dots, x_p(\mathbf{s}_i))'$ for $i = 1, \dots, n$. We model $Z(\mathbf{s}_i)$ in terms of a linear combination of \mathbf{x}_i by considering the following geostatistical regression model:

$$\begin{aligned} Z(\mathbf{s}_i) &= \mu(\mathbf{s}_i) + \eta(\mathbf{s}_i) + \varepsilon(\mathbf{s}_i) \\ &= \beta_0 + \sum_{j=1}^p \beta_j x_j(\mathbf{s}_i) + \eta(\mathbf{s}_i) + \varepsilon(\mathbf{s}_i); \quad i = 1, \dots, n, \end{aligned} \quad (1)$$

where $\mu(\cdot) = \beta_0 + \sum_{j=1}^p \beta_j x_j(\cdot)$ is the deterministic mean process, β_j 's are unknown regression coefficients, and $\eta(\cdot)$ is a zero-mean, L_2 -continuous spatially dependent Gaussian process.

Instead of using all p variables, it is sometimes preferable to select only a subset of important variables. By doing this, we are able to trade off some bias for smaller variance. We consider a class of candidate models indexed by $\alpha \in \mathcal{A} \subset 2^{\{1, \dots, p\}}$ with each α corresponding to a subset of p variables. Then the geostatistical regression model corresponding to α can be written as:

$$\mathbf{Z} = \mathbf{X}_\alpha \beta_\alpha + \boldsymbol{\eta} + \boldsymbol{\varepsilon},$$

where β_α is the parameter vector consisting of β_0 and $\{\beta_j : j \in \alpha\}$, \mathbf{X}_α is the $n \times (p_\alpha + 1)$ design matrix corresponding to α with p_α being the number of elements in α , $\boldsymbol{\eta} = (\eta(\mathbf{s}_1), \dots, \eta(\mathbf{s}_n))'$, and $\boldsymbol{\varepsilon} = (\varepsilon(\mathbf{s}_1), \dots, \varepsilon(\mathbf{s}_n))'$. Let $\alpha^c = \{j : |\beta_j| > 0\}$, which is the smallest correct model, and let $\mathcal{A}^c = \{\alpha \in \mathcal{A} : \alpha^c \subset \alpha\}$ be the set of all correct models.

2 Generalized Information Criterion

Suppose that $\boldsymbol{\Sigma} = \text{var}(\mathbf{Z})$ is known. Then the log-likelihood function of \mathbf{Z} is

$$l(\beta_\alpha; \mathbf{Z}) = \text{constant} - \frac{1}{2} (\mathbf{Z} - \mathbf{X}_\alpha \beta_\alpha)' \boldsymbol{\Sigma}^{-1} (\mathbf{Z} - \mathbf{X}_\alpha \beta_\alpha).$$

The maximum likelihood estimate of β_α can be written as $\hat{\beta}_\alpha = (\mathbf{X}_\alpha' \boldsymbol{\Sigma}^{-1} \mathbf{X}_\alpha)^{-1} \mathbf{X}_\alpha' \boldsymbol{\Sigma}^{-1} \mathbf{Z}$. A commonly used loss function for the parameters fitted by model α is the Kullback-Leibler loss, which satisfies

$$L(\alpha) = \frac{1}{2} (\hat{\boldsymbol{\mu}}_\alpha - \boldsymbol{\mu})' \boldsymbol{\Sigma}^{-1} (\hat{\boldsymbol{\mu}}_\alpha - \boldsymbol{\mu}),$$

where $\boldsymbol{\mu} = E(\mathbf{Z})$ and $\hat{\boldsymbol{\mu}}_\alpha = \mathbf{X}_\alpha \hat{\beta}_\alpha$. We consider the generalized information criterion (GIC) [2]:

$$\text{GIC}_\lambda(\alpha) = \text{constant} + (\mathbf{Z} - \hat{\boldsymbol{\mu}}_\alpha)' \boldsymbol{\Sigma}^{-1} (\mathbf{Z} - \hat{\boldsymbol{\mu}}_\alpha) + \lambda p_\alpha,$$

where λ is a penalty parameter with a smaller λ corresponding to a larger model, and vice versa. The criterion includes the commonly used Akaike information criterion (AIC) [1] corresponding to $\lambda = 2$, and the Bayesian information criterion (BIC) [3] corresponding to $\lambda = \log(n)$ as special cases.

Let $\hat{\alpha}_2 = \arg \min_{\alpha \in \mathcal{A}} \text{GIC}_2(\alpha)$ be the model selected by AIC. The following theorem provides sufficient conditions under which AIC is asymptotically efficient.

Theorem 1 Suppose that $\min_{\alpha \in \mathcal{A} \setminus \mathcal{A}^c} E(L(\alpha)) \rightarrow \infty$, as $n \rightarrow \infty$.

(i) If $\mathcal{A}^c = \emptyset$, then $L(\hat{\alpha}_2) / \min_{\alpha \in \mathcal{A}} L(\alpha) \xrightarrow{P} 1$, as $n \rightarrow \infty$.

(ii) If $\mathcal{A}^c = \{\alpha^c\}$, then $P(\hat{\alpha}_2 = \alpha^c) \rightarrow 1$, as $n \rightarrow \infty$.

However, AIC is not able to distinguish among \mathcal{A}^c when $|\mathcal{A}^c| \geq 2$, and hence may select an over-fitted model. To avoid over-fitting, it is natural to consider a larger penalty. Let $\hat{\alpha}_\lambda = \arg \min_{\alpha \in \mathcal{A}} \text{GIC}_\lambda(\alpha)$ be the model selected by GIC_λ . The following theorem provides some results when λ is large.

Theorem 2 Suppose that $\lambda \rightarrow \infty$ and $\min_{\alpha \in \mathcal{A} \setminus \mathcal{A}^c} \lambda^{-1} E(L(\alpha)) \rightarrow \infty$, as $n \rightarrow \infty$. Then

$$L(\hat{\alpha}_\lambda) / \min_{\alpha \in \mathcal{A}} L(\alpha) \xrightarrow{P} 1, \quad \text{as } n \rightarrow \infty.$$

In addition, if $\mathcal{A}^c \neq \emptyset$, then $P(\hat{\alpha}_\lambda = \alpha^c) \rightarrow 1$, as $n \rightarrow \infty$.

In what follows, we provide an example in the one-dimensional space. Suppose that the data $\{\mathbf{x}_i : i = 1, \dots, n\}$ and \mathbf{Z} are sampled at $s_i = in^{\delta-1} \in D = [0, n^\delta]$; $i = 1, \dots, n$, for some $0 \leq \delta \leq 1$, according to (1) with p fixed, where $\text{cov}(\eta(s), \eta(s+h)) = \sigma_\eta^2 \exp(-\kappa_\eta |h|)$. Note that different δ values correspond to different asymptotic frameworks. When $\delta = 0$, increasingly dense observations are sampled in a bounded fixed region $D = [0, 1]$, corresponding to the fixed-domain asymptotic framework. On the other hand, when $0 < \delta \leq 1$, the region D increases with the sample size, corresponding to the increasing-domain asymptotic framework. Suppose that $x_j(\cdot)$'s are independently generated from zero-mean Gaussian processes with $\text{cov}(\eta(s), \eta(s+h)) = \sigma_j^2 \exp(-\kappa_j |h|)$, for $j = 1, \dots, p$, where $\sigma_j^2 > 0$ and $\kappa_j > 0$. We have the following results showing whether selection consistency is satisfied may depend on which asymptotic framework is chosen.

Theorem 3 Under the setup above, suppose that $\kappa_\eta > 0$, $\sigma_\eta^2 > 0$, and $\sigma_\varepsilon^2 > 0$ are known. If $\mathcal{A}^c \neq \emptyset$, and $\lambda \rightarrow \infty$ and $\lambda n^{-(1+\delta)/2} \rightarrow 0$, as $n \rightarrow \infty$, then $P(\hat{\alpha}_\lambda = \alpha^c) \rightarrow 1$, as $n \rightarrow \infty$.

Corollary 1 Under the conditions of Theorem 3, consider two sampling schemes with δ given by δ_1 and δ_2 , where $\delta_1 < \delta_2$. Suppose that $\lambda = n^{(2+\delta_1+\delta_2)/4}$. Then $\limsup_{n \rightarrow \infty} P(\hat{\alpha}_\lambda = \alpha^c) < 1$ for $\delta = \delta_1$, and $\lim_{n \rightarrow \infty} P(\hat{\alpha}_\lambda = \alpha^c) = 1$ for $\delta = \delta_2$.

References

- [1] Akaike, H. (1973). Information theory and the maximum likelihood principle. *International Symposium on Information Theory* (V. Petrov and F. Csáki eds.), Akademiai Kiádo, Budapest, 267–281.
- [2] Nishii, R. (1984). Asymptotic properties of criteria for selection of variables in multiple regression. *The Annals of Statistics* **12**, 758–765.
- [3] Schwartz, G. (1978). Estimating the dimensions of a model. *The Annals of Statistics* **6**, 461–464.