

Supplement to “Asymptotic Behavior of Cox’s Partial Likelihood and its Application to Variable Selection”

Runze Li¹, Jian-Jian Ren², Guangren Yang³ and Ye Yu⁴

¹*Pennsylvania State University*, ²*University of Maryland*,

³*Jinan University* and ⁴*Wells Fargo Bank*

This supplement consists of the proof of Theorem 2 in the main text.

Proof of Theorem 2. To prove Theorem 2, we show the following two lemmas. Theorem 2(A) and 2(B) follow Lemma 1 and 2 respectively.

Lemma 1. *Suppose that the partial likelihood function of the Cox model satisfies Conditions (A)-(D) in Fan and Li (2002). Assume that there exists a positive constant M such that $\kappa_n < M$. Then under Condition (E4), we have*

$$P\left\{\inf_{\lambda \in \Omega_-} GIC_{\kappa_n}(\widehat{\boldsymbol{\beta}}_\lambda) > GIC_{\kappa_n}(\widehat{\boldsymbol{\beta}}_{\bar{\alpha}}^*)\right\} \rightarrow 1 \quad \text{as } n \rightarrow \infty, \quad (\text{S.1})$$

$$\liminf_{n \rightarrow \infty} P\left\{\inf_{\lambda \in \Omega_0} GIC_{\kappa_n}(\widehat{\boldsymbol{\beta}}_\lambda) > GIC_{\kappa_n}(\widehat{\boldsymbol{\beta}}_{\bar{\alpha}}^*)\right\} \geq \pi. \quad (\text{S.2})$$

Proof. Recall that for any given λ , we can obtain a selected model α_λ by penalized variable selection. And based on this selected model α_λ , we are able to obtain its corresponding non-penalized estimates $\widehat{\boldsymbol{\beta}}_{\alpha_\lambda}^*$ by maximizing the corresponding partial likelihood. Then

$$\ell_c(\widehat{\boldsymbol{\beta}}_{\alpha_\lambda}^*) \geq \ell_c(\widehat{\boldsymbol{\beta}}_\lambda), \quad (\text{S.3})$$

and $-2\ell_c(\widehat{\boldsymbol{\beta}}_\lambda) + \kappa_n \text{df}_\lambda > -2\ell_c(\widehat{\boldsymbol{\beta}}_{\alpha_\lambda}^*)$ Thus,

$$GIC_{\kappa_n}(\widehat{\boldsymbol{\beta}}_\lambda) > -2\ell_c(\widehat{\boldsymbol{\beta}}_{\alpha_\lambda}^*). \quad (\text{S.4})$$

Subtract $\text{GIC}_{\kappa_n}(\widehat{\boldsymbol{\beta}}_{\bar{\alpha}}^*)$ from both size of (S.4), we can obtain that

$$\text{GIC}_{\kappa_n}(\widehat{\boldsymbol{\beta}}_{\lambda}) - \text{GIC}_{\kappa_n}(\widehat{\boldsymbol{\beta}}_{\bar{\alpha}}^*) > -2\ell_c(\widehat{\boldsymbol{\beta}}_{\alpha_\lambda}^*) - \{-2\ell_c(\widehat{\boldsymbol{\beta}}_{\bar{\alpha}}^*) + \kappa_n \text{df}_{\bar{\alpha}}\}.$$

For any $\lambda \in \Omega_- = \{\lambda : \alpha \not\geq \alpha_0\}$, we can take $\inf_{\lambda \in \Omega_-}$ over $\text{GIC}_{\kappa_n}(\widehat{\boldsymbol{\beta}}_{\lambda})$. Under Condition (E4) and $\kappa_n < M$, for any $\lambda \in \Omega_-$, we have

$$\begin{aligned} & P\left\{\inf_{\lambda \in \Omega_-} \text{GIC}_{\kappa_n}(\widehat{\boldsymbol{\beta}}_{\lambda}) - \text{GIC}_{\kappa_n}(\widehat{\boldsymbol{\beta}}_{\bar{\alpha}}^*) > 0\right\} \\ & \geq P\left\{\inf_{\lambda \in \Omega_-} \frac{-2\ell_c(\widehat{\boldsymbol{\beta}}_{\alpha_\lambda}^*)}{n} - \frac{-2\ell_c(\widehat{\boldsymbol{\beta}}_{\bar{\alpha}}^*)}{n} - \frac{\kappa_n \text{df}_{\bar{\alpha}}}{n} > 0\right\} \\ & = P\left\{\min_{\alpha \not\geq \alpha_0} \left[\frac{-2\ell_c(\widehat{\boldsymbol{\beta}}_{\alpha}^*)}{n} - \log(n)\rho_1\right] - \left[\frac{-2\ell_c(\widehat{\boldsymbol{\beta}}_{\bar{\alpha}}^*)}{n} - \log(n)\rho_1\right] - \frac{\kappa_n \text{df}_{\bar{\alpha}}}{n} > 0\right\} \quad (\text{S.5}) \end{aligned}$$

$$= P\left\{\min_{\alpha \not\geq \alpha_0} c_{\alpha} - c_{\bar{\alpha}} + o_P(1) > 0\right\} \rightarrow 1, \quad (\text{S.6})$$

as $n \rightarrow \infty$. (S.5) is due to the finiteness of \mathcal{A} , and (S.6) uses both (E4) and the fact that deviance tends to be smaller as covariate dimension increases. (S.1) follows from the above equations.

For any $\lambda \in \Omega_0, \alpha_\lambda = \alpha_0$, it follows by (2.5)

$$\begin{aligned} & P\left\{\inf_{\lambda \in \Omega_0} \text{GIC}_{\kappa_n}(\widehat{\boldsymbol{\beta}}_{\lambda}) - \text{GIC}_{\kappa_n}(\widehat{\boldsymbol{\beta}}_{\bar{\alpha}}^*) > 0\right\} \\ & \geq P\left\{\inf_{\lambda \in \Omega_0} -2\ell_c(\widehat{\boldsymbol{\beta}}_{\alpha_\lambda}^*) - [-2\ell_c(\widehat{\boldsymbol{\beta}}_{\bar{\alpha}}^*)] - \kappa_n \text{df}_{\bar{\alpha}} > 0\right\} \\ & = P\left\{-2[\ell_c(\widehat{\boldsymbol{\beta}}_{\alpha_0}^*) - \ell_c(\widehat{\boldsymbol{\beta}}_{\bar{\alpha}}^*)] - \kappa_n \text{df}_{\bar{\alpha}} > 0\right\} \\ & \geq P\left\{-2[\ell_c(\widehat{\boldsymbol{\beta}}_{\alpha_0}^*) - \ell_c(\widehat{\boldsymbol{\beta}}_{\bar{\alpha}}^*)] > M \text{df}_{\bar{\alpha}}\right\} \quad (\text{S.7}) \end{aligned}$$

$$\rightarrow P\left\{\chi_{\text{df}_{\bar{\alpha}} - \text{df}_{\alpha_0}}^2 \geq M \text{df}_{\bar{\alpha}}\right\} > 0. \quad (\text{S.8})$$

(S.7) is due to $\kappa_n < M$, and (S.8) uses the fact that $\widehat{\boldsymbol{\beta}}_{\alpha_0}^*$ and $\widehat{\boldsymbol{\beta}}_{\bar{\alpha}}^*$ are asymptotically normal under regular condition (A)-(D) in Fan and Li (2002). Hence, the likelihood ratio test statistics $-2[\ell_c(\widehat{\boldsymbol{\beta}}_{\alpha_0}^*) - \ell_c(\widehat{\boldsymbol{\beta}}_{\bar{\alpha}}^*)] \xrightarrow{\mathcal{L}} \chi_{\text{df}_{\bar{\alpha}} - \text{df}_{\alpha_0}}^2$. (S.2) follows by taking $\pi = P\{\chi_{\text{df}_{\bar{\alpha}} - \text{df}_{\alpha_0}}^2 \geq$

$Mdf_{\hat{\alpha}}\}$. This completes the proof of Lemma 1.

Lemma 2. *Suppose that the partial likelihood function of the Cox model satisfies Conditions (A)-(D) in Fan and Li (2002). Then under Condition (E1)-(E4), and let $\lambda_n = \kappa_n/\sqrt{n}$. If κ_n satisfies $\kappa_n \rightarrow \infty$ and $\lambda_n \rightarrow 0$ as $n \rightarrow \infty$, we have*

$$P\{GIC_{\kappa_n}(\hat{\beta}_{\lambda_n}) = GIC_{\kappa_n}(\hat{\beta}_{\alpha_0}^*)\} \rightarrow 1, \quad (\text{S.9})$$

$$P\left\{\inf_{\lambda \in (\Omega_- \cup \Omega_+)} GIC_{\kappa_n}(\hat{\beta}_{\lambda}) > GIC_{\kappa_n}(\hat{\beta}_{\lambda_n})\right\} \rightarrow 1. \quad (\text{S.10})$$

Proof. With loss of generality, assume that the first d_{α_0} component of β_0 are nonzero for the true model while the rest are zeros. By Conditions (A)-(D) in Fan and Li (2002) together with Condition (E3), Fan and Li (2002) showed that

$$\begin{aligned} \hat{\beta}_{\lambda_n j} &\xrightarrow{p} 0 \text{ for } j = d_{\alpha_0} + 1, \dots, d, \\ \frac{\partial}{\partial \beta_j} \ell_c(\hat{\beta}_{\lambda_n j}) - p'_{\lambda_n}(|\hat{\beta}_{\lambda_n j}|) \text{sgn}(\hat{\beta}_{\lambda_n j}) &\xrightarrow{p} 0 \text{ for } j = 1, \dots, d_{\alpha_0}, \end{aligned} \quad (\text{S.11})$$

where $\hat{\beta}_{\lambda_n j}$ is the j th component of $\hat{\beta}_{\lambda_n}$. Under Condition (E1) and (E2), for $j = 1, \dots, d_{\alpha_0}$, there exists an m such that

$$p'_{\lambda_n}(|\hat{\beta}_{\lambda_n j}|) = 0 \text{ for } |\hat{\beta}_{\lambda_n j}| \geq \min\{|\beta_{\lambda_n j}|\} \geq m\lambda_n.$$

By (S.11), with probability tending to 1, we have,

$$\frac{\partial}{\partial \beta_j} \ell_c(\hat{\beta}_{\lambda_n j}) = 0, \text{ for } j = 1, \dots, d_{\alpha_0},$$

This is the score equation for the unpenalized partial likelihood under the true model α_0 .

Therefore, with probability tending to 1, we have

$$\begin{aligned}\widehat{\boldsymbol{\beta}}_{\lambda_n} &= \widehat{\boldsymbol{\beta}}_{\alpha_0}^*, \\ \ell_c(\widehat{\boldsymbol{\beta}}_{\lambda_n}) &= \ell_c(\widehat{\boldsymbol{\beta}}_{\alpha_0}^*).\end{aligned}$$

Thus $\text{df}_{\alpha_{\lambda_n}} = \text{df}_{\alpha_0}$ with probability tending to 1. Hence it follows that,

$$\begin{aligned}P\{\text{GIC}_{\kappa_n}(\widehat{\boldsymbol{\beta}}_{\lambda_n}) &= \text{GIC}_{\kappa_n}(\widehat{\boldsymbol{\beta}}_{\alpha_0}^*)\} \\ &= P\{-2\ell_c(\widehat{\boldsymbol{\beta}}_{\lambda_n}) + \kappa_n \text{df}_{\alpha_{\lambda_n}} + 2\ell_c(\widehat{\boldsymbol{\beta}}_{\alpha_0}^*) - \kappa_n \text{df}_{\alpha_0} = 0\} \\ &= P\{-2[\ell_c(\widehat{\boldsymbol{\beta}}_{\lambda_n}) - \ell_c(\widehat{\boldsymbol{\beta}}_{\alpha_0}^*)] + \kappa_n(\text{df}_{\alpha_{\lambda_n}} - \text{df}_{\alpha_0}) = 0\} \\ &\rightarrow 1.\end{aligned}$$

This validates (S.9).

Next, we want to show that $\text{GIC}_{\kappa_n}(\widehat{\boldsymbol{\beta}}_{\lambda}) > \text{GIC}_{\kappa_n}(\widehat{\boldsymbol{\beta}}_{\lambda_n})$ for any λ that cannot result in the true model. First, we consider λ that could result in underfitting models, namely, $\lambda \in \Omega_- = \{\lambda : \alpha_{\lambda} \not\supseteq \alpha_0\}$. By (S.4) and (S.9), with probability tending to 1, it follows that

$$\text{GIC}_{\kappa_n}(\widehat{\boldsymbol{\beta}}_{\lambda}) - \text{GIC}_{\kappa_n}(\widehat{\boldsymbol{\beta}}_{\lambda_n}) > -2\ell_c(\widehat{\boldsymbol{\beta}}_{\alpha_{\lambda}}^*) - [-2\ell_c(\widehat{\boldsymbol{\beta}}_{\alpha_0}^*)] - \kappa_n \text{df}_{\alpha_0}.$$

For any $\lambda \in \Omega_- = \{\lambda : \alpha_{\lambda} \not\supseteq \alpha_0\}$, we can take $\inf_{\lambda \in \Omega_-}$ over $\text{GIC}_{\kappa_n}(\widehat{\boldsymbol{\beta}}_{\lambda})$. Under Condition (E4)

and $\kappa_n/\sqrt{n} \rightarrow 0$, for any $\lambda \in \Omega_-$, we have

$$\begin{aligned}
& P\left\{\inf_{\lambda \in \Omega_-} \text{GIC}_{\kappa_n}(\widehat{\boldsymbol{\beta}}_\lambda) - \text{GIC}_{\kappa_n}(\widehat{\boldsymbol{\beta}}_{\lambda_n}) > 0\right\} \\
& \geq P\left\{\inf_{\lambda \in \Omega_-} \frac{-2\ell_c(\widehat{\boldsymbol{\beta}}_{\alpha_\lambda}^*)}{n} - \frac{-2\ell_c(\widehat{\boldsymbol{\beta}}_{\alpha_0}^*)}{n} - \frac{\kappa_n \text{df}_{\alpha_0}}{n} > 0\right\} \\
& = P\left\{\min_{\alpha \not\supset \alpha_0} \left[\frac{-2\ell_c(\widehat{\boldsymbol{\beta}}_\alpha^*)}{n} - \log(n)\rho_1\right] - \left[\frac{-2\ell_c(\widehat{\boldsymbol{\beta}}_{\alpha_0}^*)}{n} - \log(n)\rho_1\right] - \frac{\kappa_n \text{df}_{\alpha_0}}{n} > 0\right\} \\
& = P\left\{\min_{\alpha \not\supset \alpha_0} c_\alpha - c_{\alpha_0} + o_P(1) > 0\right\} \rightarrow 1, \tag{S.12}
\end{aligned}$$

as $n \rightarrow \infty$. (S.12) is due to Condition (E4). This implies that

$$P\left\{\inf_{\lambda \in \Omega_-} \text{GIC}_{\kappa_n}(\widehat{\boldsymbol{\beta}}_\lambda) > \text{GIC}_{\kappa_n}(\widehat{\boldsymbol{\beta}}_{\lambda_n})\right\} \rightarrow 1. \tag{S.13}$$

For any $\lambda \in \Omega_+ = \{\lambda : \alpha_\lambda \supset \alpha_0\}$, we have

$$\begin{aligned}
& \text{GIC}_{\kappa_n}(\widehat{\boldsymbol{\beta}}_\lambda) - \text{GIC}_{\kappa_n}(\widehat{\boldsymbol{\beta}}_{\lambda_n}) \\
& = -2\ell_c(\widehat{\boldsymbol{\beta}}_\lambda) - [-2\ell_c(\widehat{\boldsymbol{\beta}}_{\lambda_n})] + \kappa_n(\text{df}_{\alpha_\lambda} - \text{df}_{\alpha_{\lambda_n}}) \\
& \geq -2\ell_c(\widehat{\boldsymbol{\beta}}_{\alpha_\lambda}^*) - [-2\ell_c(\widehat{\boldsymbol{\beta}}_{\alpha_0}^*)] + \kappa_n \tau_n, \tag{S.14}
\end{aligned}$$

where $\tau_n > 0$ due to the fact that $\text{df}_{\alpha_\lambda} - \text{df}_{\alpha_{\lambda_n}} = \tau_n > 0$ when n is large. And (S.14) follows (S.3). We then take $\inf_{\lambda \in \Omega_+}$ over $\text{GIC}_{\kappa_n}(\widehat{\boldsymbol{\beta}}_\lambda)$. Under Condition (E4) and $\kappa_n/\sqrt{n} \rightarrow 0$, for any $\lambda \in \Omega_+$, we have

$$\begin{aligned}
& \inf_{\lambda \in \Omega_+} \text{GIC}_{\kappa_n}(\widehat{\boldsymbol{\beta}}_\lambda) - \text{GIC}_{\kappa_n}(\widehat{\boldsymbol{\beta}}_{\lambda_n}) \\
& \geq \min_{\alpha \not\supset \alpha_0} -2[\ell_c(\widehat{\boldsymbol{\beta}}_\alpha^*) - \ell_c(\widehat{\boldsymbol{\beta}}_{\alpha_0}^*)] + \kappa_n \tau_n \tag{S.15}
\end{aligned}$$

$$= \kappa_n \tau_n \{1 + o_p(1)\}. \tag{S.16}$$

(S.16) uses the fact that $2[\ell_c(\widehat{\boldsymbol{\beta}}_\alpha^*) - \ell_c(\widehat{\boldsymbol{\beta}}_{\alpha_0}^*)] \rightarrow \chi_{\text{df}_\alpha - \text{df}_{\alpha_0}}^2$ for $\alpha \supset \alpha_0$ together with that $\kappa_n \rightarrow \infty$. Therefore, (S.15) is positive as $n \rightarrow \infty$. Hence, we have,

$$P \left\{ \inf_{\lambda \in \Omega_+} \text{GIC}_{\kappa_n}(\widehat{\boldsymbol{\beta}}_\lambda) > \text{GIC}_{\kappa_n}(\widehat{\boldsymbol{\beta}}_{\lambda_n}) \right\} \rightarrow 1. \quad (\text{S.17})$$

Based on (S.13) and (S.17) together, we prove (S.10). Consequently, this completes the proof of Lemma 2.

Proofs of Theorem 2. Lemma 1 implies that for any λ producing the underfitted model, its associated $\text{GIC}_{\kappa_n}(\widehat{\boldsymbol{\beta}}_\lambda)$ is consistently larger than $\text{GIC}_{\kappa_n}(\widehat{\boldsymbol{\beta}}_{\alpha}^*)$. Thus, the optimal model selected by minimizing the $\text{GIC}_{\kappa_n}(\boldsymbol{\beta})$ must be either the true model or overfitted models with probability tending to one. In addition, Lemma 1 indicates that there is a nonzero probability that the smallest value of $\text{GIC}_{\kappa_n}(\widehat{\boldsymbol{\beta}}_\lambda)$ associated with the true model is larger than that of the full model. As a result, there is a positive probability that any λ associated with the true model cannot be selected by $\text{GIC}_{\kappa_n}(\boldsymbol{\beta})$ as the regularization parameter. Theorem 2(A) follows.

Lemma 2 indicates that the model identified by λ_n converges to the true model as the sample size gets large. In addition, it shows that those λ 's, which fail to identify the true model, cannot be selected by $\text{GIC}_{\kappa_n}(\boldsymbol{\beta})$ asymptotically. Theorem 2(B) follows.

We next show Theorem 2(C). Note that $(1 - \text{df}_\lambda/n)^2 = 1 + 2\text{df}_\lambda/n + O(\{\text{df}_\lambda/n\}^2)$. By the definition of the GCV, it follows that

$$2n\text{GCV}(\lambda) = -2\ell_c(\widehat{\boldsymbol{\beta}}_\lambda) + 4(-\ell_c(\widehat{\boldsymbol{\beta}}_\lambda)/n)\text{df}_\lambda + O_p(\{\text{df}_\lambda/n\}^2\ell_c(\widehat{\boldsymbol{\beta}}_\lambda))$$

Theorem 1 implies $-\ell_c(\widehat{\boldsymbol{\beta}}_\lambda)/(n \log(n)) \rightarrow \rho_1 > 0$ as $n \rightarrow \infty$, then

$$\begin{aligned} 2n\text{GCV}(\lambda) &= -2\ell_c(\widehat{\boldsymbol{\beta}}_\lambda) + 4\rho_1 \log(n)\text{df}_\lambda \{1 + o_p(1)\} + o_p(1) \\ &= -2\ell_c(\widehat{\boldsymbol{\beta}}_\lambda) + \kappa_{gcv}\text{df}_{\alpha_\lambda} \{1 + o_p(1)\}, \end{aligned}$$

where $\kappa_{gcv} = 4\rho_1 \log(n)$. Theorem 2(C) follows by using the following the proof of Lemma 2.