

# Sieve maximum likelihood estimation for a general class of accelerated hazards models with bundled parameters

XINGQIU ZHAO<sup>1</sup>, YUANSHAN WU<sup>2</sup> and GUOSHENG YIN<sup>3</sup>

<sup>1</sup>*Department of Applied Mathematics, The Hong Kong Polytechnic University, Hong Kong.  
E-mail: xingqiu.zhao@polyu.edu.hk*

<sup>2</sup>*School of Mathematics and Statistics, Wuhan University, Wuhan, Hubei, 430072, China.  
E-mail: wuyuanshan@gmail.com*

<sup>3</sup>*Department of Statistics and Actuarial Science, The University of Hong Kong, Hong Kong.  
E-mail: gyin@hku.hk*

In semiparametric hazard regression, nonparametric components may involve unknown regression parameters. Such intertwining effects make model estimation and inference much more difficult than the case in which the parametric and nonparametric components can be separated out. We study the sieve maximum likelihood estimation for a general class of hazard regression models, which include the proportional hazards model, the accelerated failure time model, and the accelerated hazards model. Coupled with the cubic B-spline, we propose semiparametric efficient estimators for the parameters that are bundled inside the nonparametric component. We overcome the challenges due to intertwining effects of the bundled parameters, and establish the consistency and asymptotic normality properties of the estimators. We carry out simulation studies to examine the finite-sample properties of the proposed method, and demonstrate its efficiency gain over the conventional estimating equation approach. For illustration, we apply our proposed method to a study of bone marrow transplantation for patients with acute leukemia.

*Keywords:* accelerated failure time model; B-spline; proportional hazards model; semiparametric efficiency bound; sieve maximum likelihood estimator; survival data

## 1. Introduction

The Cox [6] proportional hazards model has been routinely used in survival analysis. Under the proportional hazards assumption, the Cox model takes the form of

$$\lambda(t|\mathbf{Z}) = \lambda_0(t)e^{\beta_0^T \mathbf{Z}}, \quad (1)$$

where  $\lambda_0(\cdot)$  is the unknown baseline hazard function,  $\mathbf{Z}$  is the covariate vector, and  $\beta_0$  is the regression parameter of interest. Nevertheless, such constant proportionality between hazard functions may not hold in practice. As a result, alternative modeling structures, such as the accelerated failure time (AFT) model, have been proposed, which directly model the logarithm of the failure time in a linear regression form,

$$\log(T) = -\beta_0^T \mathbf{Z} + \epsilon, \quad (2)$$

where  $T$  is the failure time, and the distribution of the error  $\epsilon$  is unspecified. In fact, model (2) can be rewritten as

$$S(t|\mathbf{Z}) = S_0(te^{\beta_0^T \mathbf{Z}}), \quad (3)$$

where  $S(t|\mathbf{Z})$  is the conditional survival function given covariate  $\mathbf{Z}$ , and  $S_0(t)$  is the baseline survival function. The inference procedure for model (2) or (3) can typically be carried out using the least squares or rank methods (Prentice [14], Buckley and James [1], Ritov [15], Tsiatis [20], Wei, Ying and Lin [22], Lai and Ying [13], and Jin *et al.* [10]), and the corresponding variance are often estimated by resampling algorithm, such as bootstrap. Clearly, the nonparametric function  $S_0(\cdot)$  involves the parametric component  $e^{\beta_0^T \mathbf{Z}}$ , which makes it difficult to derive the nonparametric maximum likelihood estimator (NPMLE). Whereas, Zeng and Lin [24] developed an efficient estimator for the AFT model (3) by maximizing a kernel-smoothed profile likelihood function for the regression parameter. However, their approach is restricted to the log-transformed linear model (2). Recently, when the failure time  $T$  is subject to any completely known and strictly increasing transformation, Ding and Nan [7] proposed a sieve maximum likelihood estimator (MLE) for the censored linear regression model where the bundled parameter problem is involved. Owing to its flexibility, the sieve MLE method has been widely adopted in various semiparametric models, such as the partly linear Cox model (Huang [8]), transformed hazard models (Zeng, Yin and Ibrahim [25]), and the proportional odds model for survival data under various interval censoring mechanisms (Rossini and Tsiatis [16], Huang and Rossini [9], Shen [18]). Chen [2] provided a comprehensive review on the sieve method in the semiparametric models.

Despite the popularity of the Cox model, it assumes the treatment effect to take place immediately after patients are randomized to different treatment groups; that is, the hazards for different groups are different from time  $t = 0$ . However, in a randomized clinical trial, the treatment groups are essentially identical at  $t = 0$  due to randomization. Randomization makes different groups alike except for treatments. Particularly in oncology, it often takes some time to observe efficacy effects of the treatment, for example, tumor shrinkage. In other words, it may take a certain period of lag time for the treatment to fully exert the therapeutic effect instead of being immediately effective. Along this direction, Chen and Wang [4] proposed the accelerated hazards model by replacing the survival functions in (3) with the corresponding hazard functions, and thus the conditional hazard function of failure time  $T$  given covariate  $\mathbf{Z}$  takes the form of

$$\lambda(t|\mathbf{Z}) = \lambda_0(te^{\beta_0^T \mathbf{Z}}). \quad (4)$$

This model is intuitive in the sense that the hazard functions for different values of  $\mathbf{Z}$  in (4) are the same at time  $t = 0$ . As time goes by, the hazards in different groups would gradually change due to different treatment effects. In a more general framework, Chen and Jewell [3] proposed a class of hazards regression model,

$$\lambda(t|\mathbf{Z}) = \lambda_0(te^{\beta_0^T \mathbf{Z}})e^{\gamma_0^T \mathbf{Z}}, \quad (5)$$

where  $\beta_0$  and  $\gamma_0$  are vectors of regression parameters. Based on different parametrization, model (5) includes the proportional hazards model ( $\beta_0 = 0$ ), the AFT model ( $\beta_0 = \gamma_0$ ), and the acceler-

ated hazards model ( $\boldsymbol{\gamma}_0 = 0$ ) as special cases. Chen and Jewell [3] developed martingale estimating equations for parameter estimation and inference, which may not be semiparametric efficient. Due to the discontinuity of estimating equations with respect to the regression parameters, the estimation procedure may suffer from potential multiple roots. Furthermore, the variance estimation depends on the derivation of the baseline function, which makes it difficult to calculate in practice.

To enhance the estimation efficiency and modeling flexibility, we study the sieve maximum likelihood estimation for a general class of accelerated hazards regression models in the form of

$$\Lambda(t|\mathbf{Z}, \mathbf{X}) = \Lambda_0(te^{\boldsymbol{\beta}_0^T \mathbf{Z}})e^{\boldsymbol{\gamma}_0^T \mathbf{X}}, \quad (6)$$

where  $\Lambda_0(\cdot)$  is an unknown baseline cumulative hazards function, and  $\boldsymbol{\beta}_0$  and  $\boldsymbol{\gamma}_0$  are unknown vectors of regression parameters. Covariates  $\mathbf{Z}$  and  $\mathbf{X}$  are allowed to share some common components. It is easy to see that model (6) reduces to the proportional hazards model when  $\boldsymbol{\beta}_0 = 0$  and to the AFT model when  $\boldsymbol{\gamma}_0 = 0$ . In the case where  $\mathbf{Z}$  is the same as  $\mathbf{X}$ , model (6) reduces to the accelerated hazards model when  $\boldsymbol{\beta}_0 + \boldsymbol{\gamma}_0 = 0$  and to model (5) by reparameterizing  $\boldsymbol{\beta}_0 + \boldsymbol{\gamma}_0$  as a new parameter. Hence, model (6) has great flexibility and, more importantly, is able to simultaneously investigate the time-accelerated effect of covariate  $\mathbf{Z}$  and the proportional hazards effect of covariate  $\mathbf{X}$ . Noting that the parametric and nonparametric components are bundled together in  $\Lambda_0(te^{\boldsymbol{\beta}_0^T \mathbf{Z}})$ , the theoretical development and numerical implementation of model (6) are very challenging. In contrast to the conventional martingale-based estimating equations proposed by Chen and Wang [4] and Chen and Jewell [3], we propose an intuitive spline-based sieve maximum likelihood estimation procedure for model (6) to improve the estimation efficiency. The numerical implementation of the proposed method can be achieved through the conventional gradient-based search algorithm, such as the Newton–Raphson algorithm. The variance estimates can be obtained from the inverse of Fisher’s information matrix, and thus achieves semiparametric efficiency.

The main contributions of this paper are threefold. First, we proposed a sieve MLE method for the general accelerated hazards model in which the nonparametric function and the regression parameter are entangled with each other. The asymptotic properties of the resulting estimators are established and the estimator for the regression parameter achieves the semiparametric efficiency bound. Second, compared with the weighted estimating equation approach where the optimal weight depends on the form of the baseline function and thus it is challenging to find such an optimal weight and difficult to implement in practice, the proposed sieve MLE method is easier to be carried out. Third, the standard error estimates are obtained directly by either inverting the observed information matrix of all the parameters or inverting the efficient information matrix of the regression parameters, and both methods are more computationally tractable compared with the resampling techniques.

The rest of this article is organized as follows. We propose the sieve maximum likelihood estimation procedure in Section 2 and establish the asymptotic properties of the resultant estimators in Section 3, while proofs are presented in Section 7. We conduct simulation studies to assess the proposed method with finite samples in Section 4. As an illustration, a real data set is analyzed in Section 5. Some concluding remarks are provided in Section 6.

## 2. Sieve maximum likelihood estimation

Let  $T$  be the failure time, let  $C$  be the censoring time, denote  $a \wedge b$  as the minimum of  $a$  and  $b$ , and let  $I(\cdot)$  be the indicator function. We observe the data  $\{Y_i \equiv T_i \wedge C_i, \Delta_i \equiv I(T_i \leq C_i), \mathbf{Z}_i, \mathbf{X}_i\}$ ,  $i = 1, \dots, n$ , which are independent and identically distributed (i.i.d.) copies of  $\{Y \equiv T \wedge C, \Delta \equiv I(T \leq C), \mathbf{Z}, \mathbf{X}\}$ . Covariates  $\mathbf{Z}$  and  $\mathbf{X}$  may share the same components. Assume that  $T$  and  $C$  are conditionally independent given covariates  $\mathbf{Z}$  and  $\mathbf{X}$ . Under model (6), the conditional survival and density functions of  $T$  given both  $\mathbf{Z}$  and  $\mathbf{X}$  are  $S(t|\mathbf{Z}, \mathbf{X}) = \exp\{-\Lambda_0(t e^{\beta_0^T \mathbf{Z}}) e^{\gamma_0^T \mathbf{X}}\}$  and  $f(t|\mathbf{Z}, \mathbf{X}) = S(t|\mathbf{Z}, \mathbf{X}) \lambda_0(t e^{\beta_0^T \mathbf{Z}}) e^{\beta_0^T \mathbf{Z} + \gamma_0^T \mathbf{X}}$ , respectively, where  $\lambda_0(t) = d\Lambda_0(t)/dt$  is the baseline hazard function. The likelihood function of parameters  $(\beta, \gamma, \lambda)$  based on the observed data can be derived as

$$\prod_{i=1}^n [\{\lambda(Y_i e^{\beta^T \mathbf{Z}_i}) e^{(\beta^T \mathbf{Z}_i + \gamma^T \mathbf{X}_i)}\}^{\Delta_i} \exp\{-\Lambda(Y_i e^{\beta^T \mathbf{Z}_i}) e^{\gamma^T \mathbf{X}_i}\}],$$

where  $\lambda(t) = d\Lambda(t)/dt$ . The log-likelihood function is given by

$$l_n(\beta, \gamma, \lambda) = n^{-1} \sum_{i=1}^n \left[ \Delta_i \{\beta^T \mathbf{Z}_i + \gamma^T \mathbf{X}_i + \log \lambda(Y_i e^{\beta^T \mathbf{Z}_i})\} - \int_0^{Y_i e^{\beta^T \mathbf{Z}_i}} \lambda(s) ds e^{\gamma^T \mathbf{X}_i} \right]. \tag{7}$$

To overcome the nonnegative constraint on  $\lambda(\cdot)$ , let  $g(t) = \log \lambda(t)$  and then (7) is recast as

$$l_n(\beta, \gamma, g) = n^{-1} \sum_{i=1}^n \left[ \Delta_i \{\beta^T \mathbf{Z}_i + \gamma^T \mathbf{X}_i + g(Y_i e^{\beta^T \mathbf{Z}_i})\} - \int_0^{Y_i e^{\beta^T \mathbf{Z}_i}} \exp\{g(s)\} ds e^{\gamma^T \mathbf{X}_i} \right]. \tag{8}$$

In what follows, we propose a spline-based method to estimate the function  $g$ . Denote  $b = \sup_{y, \mathbf{z}, \beta} y \exp(\beta^T \mathbf{z})$ , then  $0 < b < \infty$  under conditions C1 and C2 listed in Section 3. Let  $0 \equiv t_0 < t_1 < \dots < t_{K_n} < t_{K_n+1} \equiv b$  be a partition of  $[0, b]$  with  $K_n = O(n^v)$  and  $\max_{0 \leq j \leq K_n} |t_{j+1} - t_j| = O(n^{-v})$  for  $v \in (0, 0.5)$ . Denote the set of partition points by  $T_{K_n} = \{t_1, \dots, t_{K_n}\}$ , and let  $\mathcal{S}_n(T_{K_n}, K_n, p)$  be the space of polynomial splines of order  $p$  defined in Schumaker ([17], page 108, Definition 4.1). According to Schumaker ([17], page 117, Corollary 4.10), there exists a local basis  $\{B_j : 1 \leq j \leq q_n\}$  with  $q_n = K_n + p$  such that for any  $s \in \mathcal{S}_n(T_{K_n}, K_n, p)$ , we can write

$$s(t) = \mathbf{a}^T \mathbf{B}(t) = \sum_{j=1}^{q_n} a_j B_j(t),$$

where  $\mathbf{a} = (a_1, \dots, a_{q_n})^T$  and  $\mathbf{B} = (B_1, \dots, B_{q_n})^T$ . Under some suitable smoothness assumptions,  $g_0$ , the true function of  $g$ , can be well approximated by some function in  $\mathcal{S}_n(T_{K_n}, K_n, p)$ .

Let  $\mathcal{B} \subseteq \mathbb{R}^{d_1}$  and  $\mathcal{T} \subseteq \mathbb{R}^{d_2}$  denote the parameter spaces of  $\boldsymbol{\beta}$  and  $\boldsymbol{\gamma}$ , respectively, where  $d_1$  and  $d_2$  are their corresponding dimensions. As a result, we seek a member of  $\mathcal{S}_n(T_{K_n}, K_n, p)$  along with values of  $(\boldsymbol{\beta}, \boldsymbol{\gamma}) \in \mathcal{B} \times \mathcal{T}$  that maximizes the log-likelihood function. Specifically, we define  $(\widehat{\boldsymbol{\beta}}_n, \widehat{\boldsymbol{\gamma}}_n, \widehat{\mathbf{a}}_n)$  to be the parameter values that maximize

$$l_n(\boldsymbol{\beta}, \boldsymbol{\gamma}, \mathbf{a}) = n^{-1} \sum_{i=1}^n \left[ \Delta_i \{ \boldsymbol{\beta}^T \mathbf{Z}_i + \boldsymbol{\gamma}^T \mathbf{X}_i + \mathbf{a}^T \mathbf{B}(Y_i e^{\boldsymbol{\beta}^T \mathbf{Z}_i}) \} - \int_0^{Y_i e^{\boldsymbol{\beta}^T \mathbf{Z}_i}} \exp\{ \mathbf{a}^T \mathbf{B}(s) \} ds e^{\boldsymbol{\gamma}^T \mathbf{X}_i} \right].$$

### 3. Asymptotic properties

Denote the true parameter  $\boldsymbol{\theta}_0 = (\boldsymbol{\alpha}_0, g_0)$  with  $\boldsymbol{\alpha}_0 = (\boldsymbol{\beta}_0^T, \boldsymbol{\gamma}_0^T)^T$ . To establish the asymptotic properties of the proposed estimator  $(\widehat{\boldsymbol{\alpha}}_n, \widehat{g}_n)$  with  $\widehat{\boldsymbol{\alpha}}_n = (\widehat{\boldsymbol{\beta}}_n^T, \widehat{\boldsymbol{\gamma}}_n^T)^T$  and  $\widehat{g}_n(t) = \widehat{\mathbf{a}}_n^T \mathbf{B}(t)$ , we need the following regularity conditions.

- C1. The parameter spaces  $\mathcal{B}$  and  $\mathcal{T}$  are both compact and contain the true parameters  $\boldsymbol{\beta}_0$  and  $\boldsymbol{\gamma}_0$  as their interior points, respectively.
- C2. The domain of the covariate  $\mathbf{V} \equiv (\mathbf{Z}^T, \mathbf{X}^T)^T$ , denoted by  $\mathcal{V}$ , is a bounded subset of  $\mathbb{R}^d$ , where  $d = d_1 + d_2$ , and both  $E(\mathbf{Z}\mathbf{Z}^T)$  and  $E(\mathbf{X}\mathbf{X}^T)$  are nonsingular.
- C3. For  $i = 1, 2$ , assume that  $\boldsymbol{\beta}_i \in \mathcal{B}$ ,  $\boldsymbol{\gamma}_i \in \mathcal{T}$ , and  $\log \lambda_i(\cdot) \in \mathcal{G}^p$ , and denote  $\Lambda_i(t) = \int_0^t \lambda_i(s) ds$ . If  $\Lambda_1(te^{\boldsymbol{\beta}_1^T \mathbf{z}})e^{\boldsymbol{\gamma}_1^T \mathbf{x}} = \Lambda_2(te^{\boldsymbol{\beta}_2^T \mathbf{z}})e^{\boldsymbol{\gamma}_2^T \mathbf{x}}$  for any  $t \in [0, b]$  and  $\mathbf{v} = (\mathbf{z}^T, \mathbf{x}^T)^T \in \mathcal{V}$ , then  $\boldsymbol{\beta}_1 = \boldsymbol{\beta}_2$ ,  $\boldsymbol{\gamma}_1 = \boldsymbol{\gamma}_2$ , and  $\lambda_1 = \lambda_2$ .
- C4. Let  $\epsilon_0 = Y e^{\boldsymbol{\beta}_0^T \mathbf{Z}}$ . There exists a truncation time  $\tau < \infty$  such that, for some positive constant  $\delta_0$ ,  $P(\epsilon_0 > \tau | \mathbf{V}) \geq \delta_0$  almost surely with respect to the probability measure of  $\mathbf{V}$ .
- C5. The conditional density of  $C$  given  $\mathbf{V}$  and its derivative are uniformly bounded for all possible values of  $\mathbf{V}$ .
- C6. Let  $\mathcal{G}^p$  denote the collection of bounded functions  $g$  on  $[0, b]$  with bounded derivatives  $g^{(j)}$ ,  $j = 1, \dots, k$ , such that the  $k$ th derivative  $g^{(k)}$  satisfies the following Lipschitz continuity condition,

$$|g^{(k)}(s) - g^{(k)}(t)| \leq L|s - t|^m \quad \text{for } s, t \in [0, b],$$

where  $k$  is a positive integer and  $m \in (0, 1]$  such that  $p = m + k \geq 3$ , and  $L < \infty$  is a constant. The true log baseline hazard function  $g_0(\cdot) = \log \lambda_0(\cdot)$  belongs to  $\mathcal{G}^p$ . For notational simplicity, we may also use  $g'$  and  $g''$  to denote the first and second derivatives of  $g$ , respectively.

- C7. For some  $\eta \in (0, 1)$ ,  $u^T \text{Var}(\mathbf{V} | \epsilon_0, \Delta = 1) u \geq \eta u^T E(\mathbf{V}\mathbf{V}^T | \epsilon_0, \Delta = 1) u$  almost surely for all  $u \in \mathbb{R}^d$ .  $E(\Delta \mathbf{W}\mathbf{W}^T)$  is nonsingular, where  $\mathbf{W} = (\{1 + \epsilon_0 g'_0(\epsilon_0)\} \mathbf{Z}^T, \mathbf{X}^T)^T$ .

C8. Let  $M(t) = \Delta I(Ye^{\beta_0^T \mathbf{Z}} \leq t) - \int_0^t I(Ye^{\beta_0^T \mathbf{Z}} \geq s)e^{g_0(s)} e^{\gamma_0^T \mathbf{X}} ds$ ,

$$l_{\beta_0}^*(\mathcal{O}) = \int_0^b \left[ \mathbf{Z} - \frac{E\{\mathbf{Z}I(Ye^{\beta_0^T \mathbf{Z}} \geq t)e^{\gamma_0^T \mathbf{X}}\}}{E\{I(Ye^{\beta_0^T \mathbf{Z}} \geq t)e^{\gamma_0^T \mathbf{X}}\}} \right] \{1 + tg'_0(t)\} dM(t),$$

$$l_{\gamma_0}^*(\mathcal{O}) = \int_0^b \left[ \mathbf{X} - \frac{E\{\mathbf{X}I(Ye^{\beta_0^T \mathbf{Z}} \geq t)e^{\gamma_0^T \mathbf{X}}\}}{E\{I(Ye^{\beta_0^T \mathbf{Z}} \geq t)e^{\gamma_0^T \mathbf{X}}\}} \right] dM(t),$$

$$l_{\alpha_0}^*(\mathcal{O}) = (l_{\beta_0}^*(\mathcal{O})^T, l_{\gamma_0}^*(\mathcal{O})^T)^T, \quad \mathcal{I}(\alpha_0) = E\{l_{\alpha_0}^*(\mathcal{O})^{\otimes 2}\},$$

where  $\mathcal{O} = (Y, \Delta, \mathbf{Z}, \mathbf{X})$  and  $\mathbf{a}^{\otimes 2} = \mathbf{a}\mathbf{a}^T$  for a column vector  $\mathbf{a}$ . Assume that  $\mathcal{I}(\alpha_0)$  is nonsingular.

Conditions C1–C2 and C4–C5 are common assumptions in the context of survival analysis. Condition C3 is required to guarantee the identifiability of model (6). Obviously, the model is unidentifiable if and only if  $\Lambda_0(t) = c_1 t^{c_2}$  for some positive constants  $c_1$  and  $c_2$  (Chen and Jewell [3]). Condition C6 requires  $p \geq 3$  to guarantee the desirable control of the spline approximation error rates of the first and second derivatives of  $g_0$ . Condition C7 is a technical assumption and can be justified in many applications. This assumption is also imposed by Wellner and Zhang [23] for the panel count data model and Ding and Nan [7] for the censored linear regression model. Condition C8 is a natural assumption that essentially requires the semiparametric efficiency information matrix to be invertible.

Following Ding and Nan [7], we define

$$\mathcal{H}^p = \{\xi(\cdot, \beta) : \xi(t, \mathbf{z}, \beta) = g(\psi(t, \mathbf{z}, \beta)), g \in \mathcal{G}^p, t \in [0, b], \mathbf{z} \in \mathcal{Z}, \beta \in \mathcal{B}\},$$

where

$$\psi(t, \mathbf{z}, \beta) = te^{(\beta - \beta_0)^T \mathbf{z}}.$$

Here  $\xi$  is a composite function of  $g$  composed with  $\psi$ , and  $\xi(t, \mathbf{z}, \beta_0) = g(t)$ . We equip the functional space  $\mathcal{H}^p$  with the norm  $\|\cdot\|_2$  defined as

$$\|\xi(\cdot, \beta)\|_2 = \left[ \int_{\mathcal{Z}} \int_0^b \{g(te^{(\beta - \beta_0)^T \mathbf{z}})\}^2 d\Lambda_0(t) dF_{\mathbf{Z}}(\mathbf{z}) \right]^{1/2}$$

for any  $\xi(\cdot, \beta) \in \mathcal{H}^p$ , where  $F_{\mathbf{Z}}(\mathbf{z})$  is the cumulative distribution function of  $\mathbf{Z}$ . For any  $\theta_i = (\beta_i, \gamma_i, \xi_i(\cdot, \beta_i))$ ,  $i = 1, 2$ , in the space  $\Theta^p \equiv \mathcal{B} \times \mathcal{T} \times \mathcal{H}^p$ , define the distance,

$$d(\theta_1, \theta_2) = (|\beta_1 - \beta_2|^2 + |\gamma_1 - \gamma_2|^2 + \|\xi_1(\cdot, \beta_1) - \xi_2(\cdot, \beta_2)\|_2^2)^{1/2},$$

where  $|\cdot|$  is the Euclidean norm. Let  $\mathcal{G}_n^p = \mathcal{S}_n(T_{K_n}, K_n, p)$ ,

$$\mathcal{H}_n^p = \{\xi(\cdot, \beta) : \xi(t, \mathbf{z}, \beta) = g(\psi(t, \mathbf{z}, \beta)), g \in \mathcal{G}_n^p, t \in [0, b], \mathbf{z} \in \mathcal{Z}, \beta \in \mathcal{B}\},$$

and  $\Theta_n^p = \mathcal{B} \times \mathcal{T} \times \mathcal{H}_n^p$ . It is easy to see that  $\Theta_n^p \subseteq \Theta_{n+1}^p \cdots \subseteq \Theta^p$  for  $n \geq 1$ . Note that the sieve estimator  $\hat{\theta}_n = (\hat{\beta}_n, \hat{\gamma}_n, \hat{\xi}_n(\cdot, \hat{\beta}_n))$  is the maximizer of the empirical log-likelihood over

the sieve space  $\Theta_n^p$ , where  $\widehat{\xi}_n(t, \mathbf{z}, \widehat{\boldsymbol{\beta}}_n) = \widehat{g}_n(te^{\widehat{\boldsymbol{\beta}}_n - \boldsymbol{\beta}_0} \mathbf{z})$ . The following theorem provides the convergence rate of the proposed estimator  $\widehat{\boldsymbol{\theta}}_n$  to the true parameter  $\boldsymbol{\theta}_0 = (\boldsymbol{\beta}_0, \boldsymbol{\gamma}_0, \xi_0(\cdot, \boldsymbol{\beta}_0)) = (\boldsymbol{\beta}_0, \boldsymbol{\gamma}_0, g_0)$ .

**Theorem 1.** *Suppose that conditions C1–C7 hold and  $(2p + 2)^{-1} < v < (2p)^{-1}$ , then*

$$d(\widehat{\boldsymbol{\theta}}_n, \boldsymbol{\theta}_0) = O_p(n^{-\min\{pv, (1-v)/2\}}).$$

The proof of Theorem 1 is provided in Section 7 by verifying the conditions of Theorem 1 in Shen and Wong [19]. Theorem 1 implies that, if  $v = (2p + 1)^{-1}$ ,  $d(\widehat{\boldsymbol{\theta}}_n, \boldsymbol{\theta}_0) = O_p(n^{-p/(1+2p)})$ , which is the optimal convergence rate in the nonparametric setting. Although the overall convergence rate is slower than  $n^{-1/2}$ , the proposed estimator for the regression parameter  $\boldsymbol{\alpha}_0$  is still asymptotically normal at the rate of  $n^{-1/2}$  and attains the semiparametric efficiency bound. We summarize these asymptotic results in the following theorem.

**Theorem 2.** *Suppose that conditions C1–C8 hold and  $(2p + 2)^{-1} < v < (2p)^{-1}$ , then  $n^{1/2}(\widehat{\boldsymbol{\alpha}}_n - \boldsymbol{\alpha}_0)$  converges in distribution to a mean zero normal random vector with covariance matrix  $\mathcal{I}^{-1}(\boldsymbol{\alpha}_0)$  equal to the semiparametric efficiency bound of  $\boldsymbol{\alpha}_0$ .*

The proof of Theorem 2 is also presented in Section 7 by checking the conditions in Theorem 6.1 of Wellner and Zhang [23], which relies heavily on the empirical process theory. A consistent estimator for the limiting covariance matrix is summarized by the following theorem.

**Theorem 3.** *Let  $\widehat{\boldsymbol{\tau}}_{\boldsymbol{\alpha}_0}^*(\mathcal{O}) = (\widehat{\boldsymbol{\tau}}_{\boldsymbol{\beta}_0}^*(\mathcal{O})^T, \widehat{\boldsymbol{\tau}}_{\boldsymbol{\gamma}_0}^*(\mathcal{O})^T)^T$ , where*

$$\begin{aligned} \widehat{\boldsymbol{\tau}}_{\boldsymbol{\beta}_0}^*(\mathcal{O}) &= \int_0^b \{\mathbf{Z} - \bar{\mathbf{Z}}(t, \widehat{\boldsymbol{\beta}}_n, \widehat{\boldsymbol{\gamma}}_n)\} \{1 + t\widehat{g}'_n(t)\} d\widehat{M}(t), \\ \widehat{\boldsymbol{\tau}}_{\boldsymbol{\gamma}_0}^*(\mathcal{O}) &= \int_0^b \{\mathbf{X} - \bar{\mathbf{X}}(t, \widehat{\boldsymbol{\beta}}_n, \widehat{\boldsymbol{\gamma}}_n)\} d\widehat{M}(t), \\ \bar{\mathbf{Z}}(t, \boldsymbol{\beta}, \boldsymbol{\gamma}) &= \frac{P_n\{\mathbf{Z}I(Ye^{\boldsymbol{\beta}^T \mathbf{Z}} \geq t)e^{\boldsymbol{\gamma}^T \mathbf{X}}\}}{P_n\{I(Ye^{\boldsymbol{\beta}^T \mathbf{Z}} \geq t)e^{\boldsymbol{\gamma}^T \mathbf{X}}\}}, \\ \bar{\mathbf{X}}(t, \boldsymbol{\beta}, \boldsymbol{\gamma}) &= \frac{P_n\{\mathbf{X}I(Ye^{\boldsymbol{\beta}^T \mathbf{Z}} \geq t)e^{\boldsymbol{\gamma}^T \mathbf{X}}\}}{P_n\{I(Ye^{\boldsymbol{\beta}^T \mathbf{Z}} \geq t)e^{\boldsymbol{\gamma}^T \mathbf{X}}\}}, \\ \widehat{M}(t) &= \Delta I(Ye^{\widehat{\boldsymbol{\beta}}_n^T \mathbf{Z}} \leq t) - \int_0^t I(Ye^{\widehat{\boldsymbol{\beta}}_n^T \mathbf{Z}} \geq s) \exp\{\widehat{g}_n(s)\} e^{\widehat{\boldsymbol{\gamma}}_n^T \mathbf{X}} ds, \end{aligned}$$

and  $P_n$  is the empirical measure with respect to  $\mathcal{O}$ . Suppose that conditions in Theorem 2 hold, then  $P_n\{\widehat{\boldsymbol{\tau}}_{\boldsymbol{\alpha}_0}^*(\mathcal{O})^{\otimes 2}\}$  converges to  $\mathcal{I}(\boldsymbol{\alpha}_0)$  in probability.

## 4. Simulation studies

We conducted simulation studies to assess the proposed sieve MLE for finite samples. We simulated covariates  $\mathbf{Z}$  and  $\mathbf{X}$  independently from the Bernoulli distribution with success probability 0.5 and then generated the survival times  $T$  from model (6). We set the true parameter values  $\boldsymbol{\beta}_0 = 1.5$  and  $\boldsymbol{\gamma}_0 = 0.5$  and considered four different baseline hazard functions for  $\lambda_0(t)$ : (i)  $\lambda_0(t) = 1/(1+t)$ ; (ii)  $\lambda_0(t) = (t-0.5)^2$ ; (iii)  $\lambda_0(t) = \log(1+t)$ ; and (iv)  $\lambda_0(t) = 1 + \cos(5t + 10)$ . For each case, we generated censoring times  $C$  from  $\text{Unif}(c_1, c_2)$  with truncation at  $\tau = c_2 - 1$  to achieve censoring rates of 20% and 40%, respectively. We considered the sample size  $n = 200$  and 400.

In the implementation of the sieve MLE, we chose the cubic B-spline and took the data-adaptive interior knots as the median of  $\{Y_i e^{\boldsymbol{\beta}^T \mathbf{Z}_i} : i = 1, \dots, n\}$  with a given  $\boldsymbol{\beta}$  in cases (i)–(iii) and the 20th, 40th, 60th, and 80th quantiles in case (iv). In particular, we adopted the following procedure to obtain the sieve MLE.

- (1) Choose initial values  $(\tilde{\boldsymbol{\beta}}^{(0)}, \tilde{\boldsymbol{\gamma}}^{(0)}, \tilde{\mathbf{a}}^{(0)})$  and set  $k = 0$ .
- (2) At step  $k + 1$ , obtain  $\tilde{\mathbf{a}}^{(k+1)}$  by solving  $\partial l_n(\tilde{\boldsymbol{\beta}}^{(k)}, \tilde{\boldsymbol{\gamma}}^{(k)}, \mathbf{a})/\partial \mathbf{a} = \mathbf{0}$  using the Newton–Raphson algorithm with the initial value  $\tilde{\mathbf{a}}^{(k)}$  until the maximum componentwise difference between the two consecutive values is less than  $10^{-3}$ .
- (3) Obtain  $(\tilde{\boldsymbol{\beta}}^{(k+1)}, \tilde{\boldsymbol{\gamma}}^{(k+1)})$  by solving  $\partial l_n(\boldsymbol{\beta}, \boldsymbol{\gamma}, \tilde{\mathbf{a}}^{(k+1)})/\partial(\boldsymbol{\beta}, \boldsymbol{\gamma}) = \mathbf{0}$  using the Newton–Raphson algorithm with initial value  $(\tilde{\boldsymbol{\beta}}^{(k)}, \tilde{\boldsymbol{\gamma}}^{(k)})$  until the maximum componentwise difference between the two consecutive values is less than  $10^{-3}$ .
- (4) Repeat steps (2) and (3) until the maximum componentwise differences of two consecutive values are less than  $10^{-3}$ . The resultant estimators, denoted by  $(\hat{\boldsymbol{\beta}}_n, \hat{\boldsymbol{\gamma}}_n, \hat{\mathbf{a}}_n)$ , are taken as the sieve MLE.

Table 1 summarizes the estimates from 1000 replications for the censoring rate of 20% with  $n = 200$  and 400, respectively. The column labeled “EST” is the average value of the estimates, “SE” is the sample standard error of the estimates, “ESE<sub>1</sub>” is the average of standard error estimates by inverting the observed information matrix of all parameters including the basis spline coefficients, and “CP<sub>1</sub>” is the corresponding coverage proportion of 95% confidence intervals. We also present the column “ESE<sub>2</sub>”, which is the average of standard error estimates by inverting the estimated information matrix of the regression parameter  $\boldsymbol{\alpha}_0$  based on Theorem 3 and list the column “CP<sub>2</sub>”, which is the corresponding coverage proportion of 95% confidence intervals. The column “MSE” refers to the average value of the mean squared errors.

Clearly, the proposed sieve MLE method performs well under all of the four different baseline hazard functions. The parameter estimates are virtually unbiased for both  $\boldsymbol{\beta}$  and  $\boldsymbol{\gamma}$ , and the bias decreases as the sample size increases. The estimated standard errors by inverting the observed information matrix of all parameters or those by inverting the information matrix based on the efficient score function agree well with the sample standard errors. The coverage probabilities are around the nominal level 95% for all cases. The estimated baseline hazard function using the B-spline approximation under  $n = 200$  and  $n = 400$  are presented in Figure 1. It can be seen that the estimated baseline hazard functions are reasonably close to the corresponding true curves. We also explored the situation with a censoring rate of 40%. The corresponding results of



**Table 1.** Simulation results under the proposed accelerated hazards model with a censoring rate of 20%

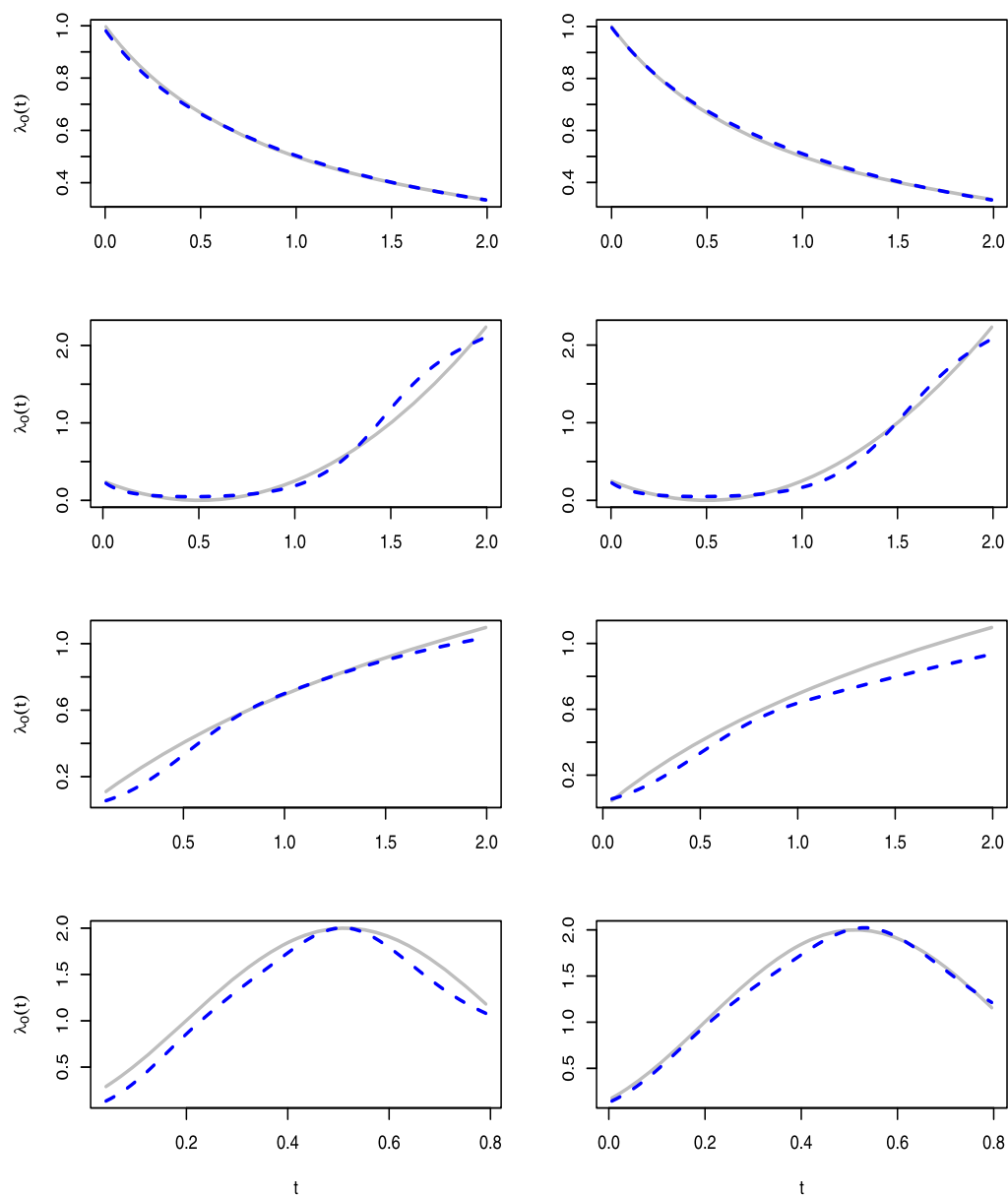
$\lambda_0(\cdot)$	$n$	True value	Sieve MLE						
			EST	SE	ESE <sub>1</sub>	CP <sub>1</sub>	ESE <sub>2</sub>	CP <sub>2</sub>	MSE <sub><math>\times 10^2</math></sub>
(i)	200	$\beta(1.5)$	1.496	0.238	0.232	0.939	0.225	0.936	5.657
		$\gamma(0.5)$	0.511	0.162	0.160	0.947	0.160	0.949	2.628
	400	$\beta(1.5)$	1.503	0.162	0.163	0.948	0.161	0.944	2.606
		$\gamma(0.5)$	0.504	0.114	0.112	0.950	0.112	0.950	1.292
(ii)	200	$\beta(1.5)$	1.505	0.046	0.043	0.935	0.040	0.922	0.210
		$\gamma(0.5)$	0.514	0.168	0.165	0.950	0.167	0.952	2.817
	400	$\beta(1.5)$	1.506	0.031	0.030	0.932	0.028	0.915	0.099
		$\gamma(0.5)$	0.503	0.114	0.116	0.954	0.117	0.954	1.293
(iii)	200	$\beta(1.5)$	1.501	0.099	0.098	0.955	0.096	0.947	0.984
		$\gamma(0.5)$	0.514	0.166	0.163	0.945	0.166	0.950	2.774
	400	$\beta(1.5)$	1.503	0.068	0.069	0.952	0.069	0.950	0.466
		$\gamma(0.5)$	0.498	0.111	0.114	0.961	0.115	0.962	1.236
(iv)	200	$\beta(1.5)$	1.493	0.115	0.115	0.911	0.099	0.902	1.311
		$\gamma(0.5)$	0.515	0.169	0.162	0.944	0.163	0.942	2.865
	400	$\beta(1.5)$	1.498	0.079	0.077	0.924	0.072	0.913	0.623
		$\gamma(0.5)$	0.504	0.117	0.114	0.954	0.114	0.951	1.358

\*EST, the average value of the parameter estimates; SE, the sample standard error of the estimates; ESE<sub>1</sub>, the estimate of the standard error by inverting the information matrix of all parameters; CP<sub>1</sub>, the corresponding coverage probability of 95% confidence intervals; ESE<sub>2</sub>, the estimate of the standard error by inverting the information matrix based on the efficient score function; CP<sub>2</sub>, the corresponding coverage probability of 95% confidence intervals; MSE, the mean squared errors of the parameter estimates.

the estimates for the regression parameters based on 1000 replications are presented in Table 2, from which similar conclusions can be drawn as before. Moreover, the estimated baseline hazard functions are plotted in Figure 2, which deteriorate slightly compared with those in Figure 1.

### 5. Application

As an illustration, we applied the proposed general class of accelerated hazards models to a study of bone marrow transplantation with 137 patients of acute leukemia (Copelan *et al.* [5] and Klein and Moeschberger [11]). The disease-free survival time, including the time to relapse, death, or the end of study, is of primary interest. Patients were followed for approximate 7.2 years, of whom around 39.4% were censored. Several potential risk factors were measured at the time of transplantation. Patients were classified into three risk categories based on their disease status: 38 patients with acute lymphoblastic leukemia (ALL), 54 patients with acute myelocytic leukemia (AML) low risk, and 45 patients with AML high risk. Both patients and donors' ages and the waiting times from diagnosis to transplantation were recorded. The AML patients with their French–American–British (FAB) classification of grade 4 or 5 based on standard morphological criteria were also considered as a covariate in our regression model. Patients were either given



**Figure 1.** True baseline hazard function (solid line) and its estimate (dashed line) using the B-spline approximation under  $n = 200$  (left panel) and  $n = 400$  (right panel) with a censoring rate of 20%. From top to bottom, the plots correspond to cases (i) to (iv) for the baseline hazard functions.

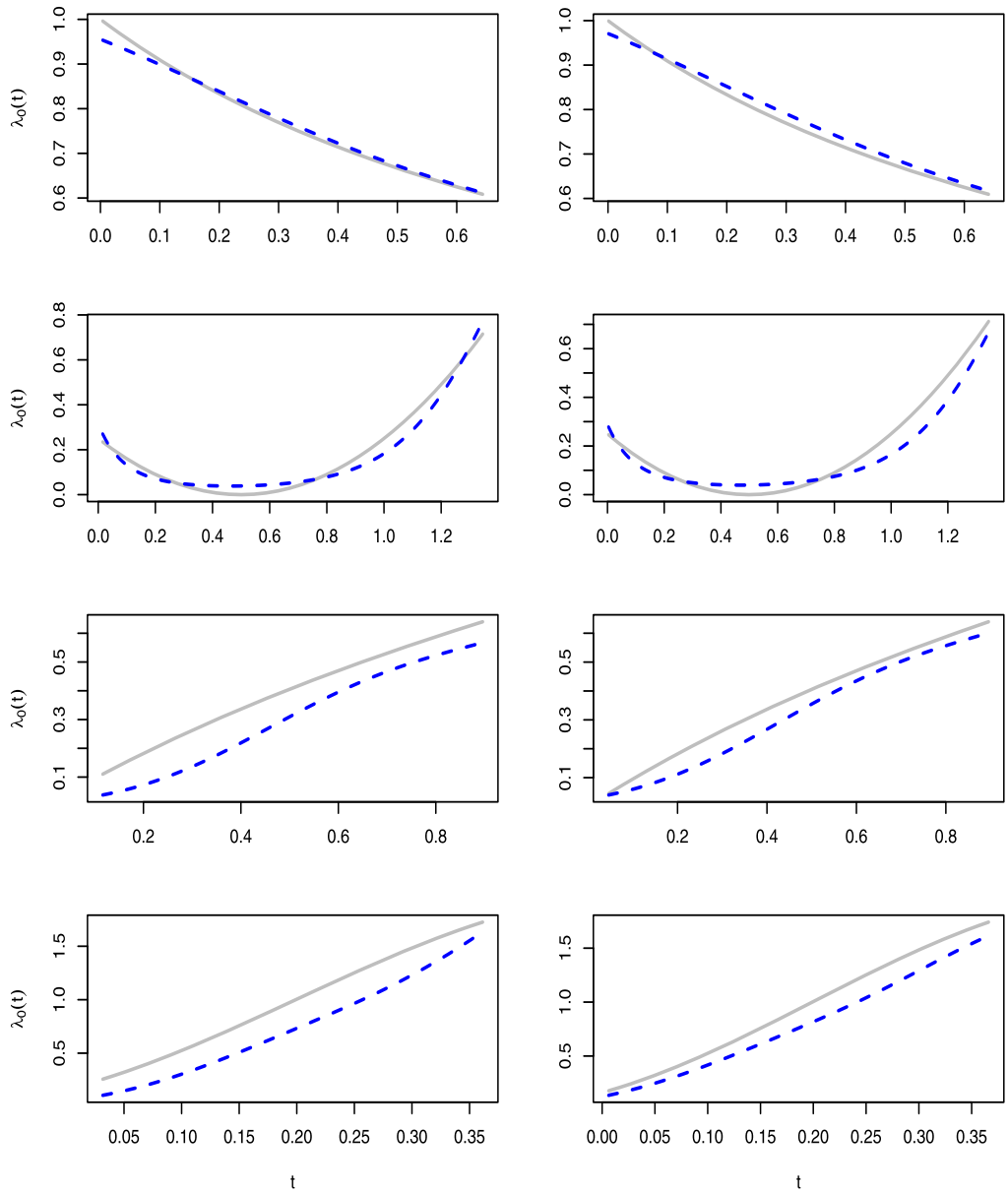
**Table 2.** Simulation results under the proposed accelerated hazards model with a censoring rate of 40%

$\lambda_0(\cdot)$	$n$	True value	Sieve MLE						
			EST	SE	ESE <sub>1</sub>	CP <sub>1</sub>	ESE <sub>2</sub>	CP <sub>2</sub>	MSE <sub><math>\times 10^2</math></sub>
(i)	200	$\beta(1.5)$	1.501	0.258	0.262	0.952	0.250	0.949	6.625
		$\gamma(0.5)$	0.518	0.188	0.185	0.950	0.185	0.949	3.539
	400	$\beta(1.5)$	1.504	0.180	0.180	0.942	0.178	0.936	3.241
		$\gamma(0.5)$	0.504	0.130	0.130	0.941	0.130	0.941	1.679
(ii)	200	$\beta(1.5)$	1.501	0.069	0.063	0.921	0.061	0.916	0.472
		$\gamma(0.5)$	0.523	0.196	0.189	0.947	0.194	0.949	3.864
	400	$\beta(1.5)$	1.498	0.048	0.044	0.925	0.042	0.913	0.235
		$\gamma(0.5)$	0.509	0.130	0.132	0.957	0.135	0.962	1.696
(iii)	200	$\beta(1.5)$	1.511	0.142	0.138	0.943	0.139	0.943	2.025
		$\gamma(0.5)$	0.520	0.196	0.189	0.952	0.191	0.953	3.843
	400	$\beta(1.5)$	1.506	0.099	0.097	0.946	0.099	0.948	0.977
		$\gamma(0.5)$	0.501	0.129	0.132	0.956	0.133	0.958	1.666
(iv)	200	$\beta(1.5)$	1.504	0.127	0.146	0.936	0.117	0.935	1.605
		$\gamma(0.5)$	0.516	0.196	0.187	0.947	0.187	0.943	3.843
	400	$\beta(1.5)$	1.508	0.090	0.097	0.948	0.086	0.933	0.811
		$\gamma(0.5)$	0.504	0.126	0.131	0.963	0.131	0.961	1.577

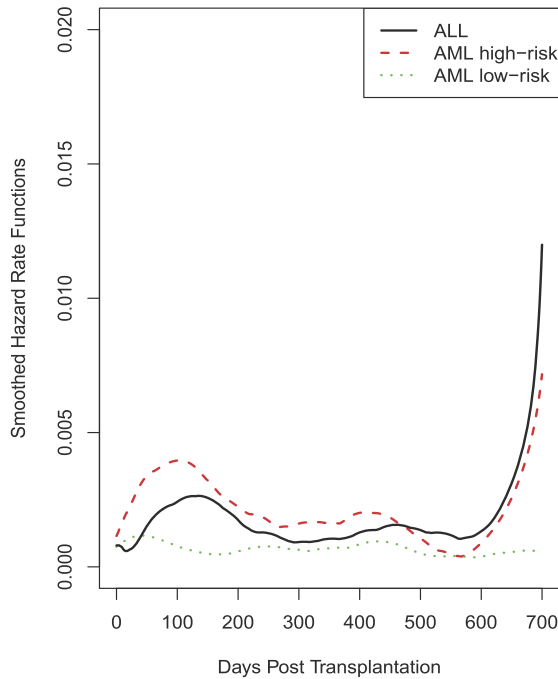
\*EST, the average value of the parameter estimates; SE, the sample standard error of the estimates; ESE<sub>1</sub>, the estimate of the standard error by inverting the information matrix of all parameters; CP<sub>1</sub>, the corresponding coverage probability of 95% confidence intervals; ESE<sub>2</sub>, the estimate of the standard error by inverting the information matrix based on the efficient score function; CP<sub>2</sub>, the corresponding coverage probability of 95% confidence intervals; MSE, the mean squared errors of the parameter estimates.

a graft-versus-host prophylactic combining methotrexate (MTX) with cyclosporine and possibly methylprednisolone or given only a combination of cyclosporine and methylprednisolone. In our analysis, we used  $X_1 = 1$  to indicate the patient with AML low risk and  $X_1 = 0$  otherwise,  $X_2 = 1$  to indicate the patient with AML high risk and  $X_2 = 0$  otherwise,  $X_3$  to denote the patient’s age centered by 28 years,  $X_4$  to denote the donor’s age centered by 28 years,  $X_5 = 1$  to indicate the AML patient with FAB grade 4 or 5 and  $X_5 = 0$  otherwise, and  $X_6$  to denote the patient’s waiting time from diagnosis to transplantation centered by 9 months, and  $X_7 = 1$  to indicate the patient treated with MTX and  $X_7 = 0$  otherwise.

To make a preliminary investigation of whether the hazards of the patients within three different risk categories were identical at the beginning of study, we plotted the kernel-smoothed hazard rate functions with bandwidth 100 days in Figure 3. It can be observed that the smoothed hazards of patients with ALL and AML low risk are almost the same at the initiation of the study. However, the smoothed hazard of patients with AML high risk appears to be slightly higher than those of the other two at time  $t = 0$ , and increases to a higher level and then lies between those of patients with ALL and patients with AML low risk during the later follow-up of the study. Figure 3 reveals that the hazards may not be proportional from the very beginning of the study, which results in the usual proportional hazards assumption questionable. Intuitively, it is more appealing



**Figure 2.** True baseline hazard function (solid line) and its estimate (dashed line) using the B-spline approximation under  $n = 200$  (left panel) and  $n = 400$  (right panel) with a censoring rate of 40%. From top to bottom, the plots correspond to cases (i) to (iv) for the baseline hazard functions.



**Figure 3.** Smoothed hazard rate functions for patients with ALL, AML high-risk, and AML low-risk, respectively.

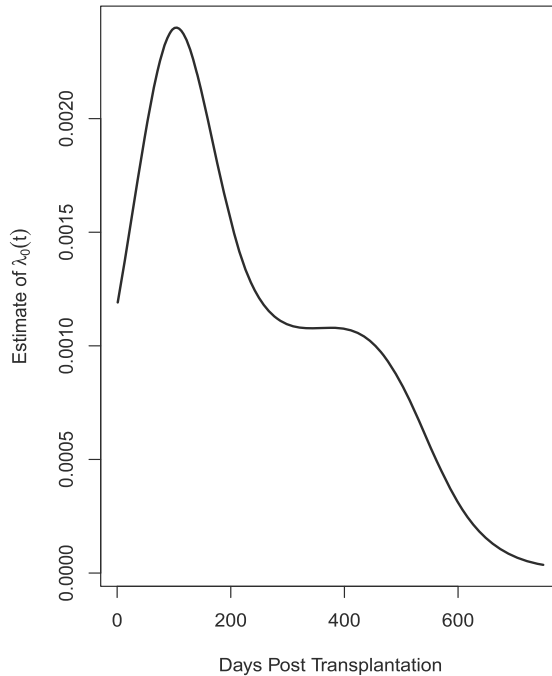
to consider the time scaled effects of risk status categories as well as the proportional effects of all the risk factors by employing our proposed general accelerated hazards model through setting  $\mathbf{Z} = (X_1, X_2)^T$  and  $\mathbf{X} = (X_1, \dots, X_7)^T$ . We applied the proposed sieve MLE with smoothing splines to fit the data. For comparison, we fitted the Cox proportional hazards model to the data without taking into consideration of the time scaled effects of risk status categories. The estimation results are summarized in Table 3. For the regression parameters in the proportional hazards component, all of the three considered methods agree in general: Patients without FAB grade 4 or 5 and those with AML low risk were associated with lower hazard risks and thus led to longer progression-free survival times, while the effects of other covariates were not significant. For the time scaled effects of risk status categories, the sieved MLE method shows that patients with AML low risk had significant decelerated hazard risks while the scaled time effect of patients with AML high risk was not significant.

Figure 4 exhibits the estimate of the baseline hazard function  $\lambda_0(t)$  using the proposed sieve MLE method. Clearly, patients at the beginning of post transplantation would suffer from the drastically increasing risk due to the incompatibility between the donor and patient, and then the hazard gradually decreased with time.

**Table 3.** Analysis results of the bone marrow transplantation data

Estimation	Z <sub>1</sub> (AML L)	Z <sub>2</sub> (AML H)	X <sub>1</sub> (AML L)	X <sub>2</sub> (AML H)	X <sub>3</sub> (P Age)	X <sub>4</sub> (D Age)	X <sub>5</sub> (FAB)	X <sub>6</sub> (Wait T)	X <sub>7</sub> (MTX)
Sieve MLE									
EST	-0.651	-0.128	-0.716	-0.033	0.009	0.000	0.804	-0.011	0.348
ESE <sub>1</sub>	0.045	0.214	0.370	0.380	0.020	0.018	0.276	0.012	0.252
p-value <sub>1</sub>	<0.001	0.548	0.053	0.930	0.650	0.994	0.004	0.329	0.166
ESE <sub>2</sub>	0.119	0.105	0.365	0.371	0.020	0.016	0.269	0.010	0.240
p-value <sub>2</sub>	<0.001	0.220	0.050	0.929	0.652	0.993	0.003	0.262	0.147
Cox model									
EST			-1.051	-0.188	0.012	-0.001	0.812	-0.011	0.294
ESE			0.368	0.359	0.020	0.018	0.275	0.011	0.250
p-value			0.004	0.600	0.530	0.940	0.003	0.310	0.240

\*EST, the parameter estimates; ESE<sub>1</sub>, the estimate of the standard error by inverting the information matrix of all parameters; ESE<sub>2</sub>, the estimate of the standard error by inverting the information matrix based on the efficient score function; ESE, the estimate of the standard error by inverting the information matrix based on the Cox proportional hazards regression model.



**Figure 4.** Estimated baseline hazard  $\lambda_0(t)$  using the proposed sieve MLE method for the bone marrow transplantation data.

### 6. Remark

The general accelerated hazards model enables us to evaluate the time scaled effects and the proportional hazards effects of covariates simultaneously. However, it is difficult in practice to classify the risk factors rigorously into either the time scaled or the proportional hazards components of the model. It often depends on the objectives of the study, the interest of the investigator, and the underlying biological process. If there is no such biological information as guidance, some data-driven methods could be used for the classification of covariates. For example, when the number of covariates is small, all the possible models from different combinations of covariates in the time scaled and the proportional hazards parts can be considered. To facilitate the selection of the models, some criteria for evaluating the goodness of model fitting should be considered. When the number of covariates is moderately large, this exhaustive method could be time-consuming, while similar automatic structure discovery procedures as presented in Zhang, Cheng and Liu [26] may warrant further research.

### 7. Proofs of theorems

Before proving the theorems presented in Section 3, we introduce some useful lemmas. Define

$$\mathcal{H} = \left\{ h : h(\cdot, \boldsymbol{\beta}) = \frac{\partial \xi_\eta(\cdot, \boldsymbol{\beta})}{\partial \eta} \Big|_{\eta=0} = w(\psi(\cdot, \boldsymbol{\beta})), \xi_\eta \in \mathcal{H}^p \right\}.$$

**Lemma 1.** *Denote*

$$l(\boldsymbol{\beta}, \boldsymbol{\gamma}, \xi(\cdot, \boldsymbol{\beta}); \mathcal{O}) = \Delta \{ \boldsymbol{\beta}^T \mathbf{Z} + \boldsymbol{\gamma}^T \mathbf{X} + g(Ye^{\boldsymbol{\beta}^T \mathbf{Z}}) \} - \int_0^{Ye^{\boldsymbol{\beta}^T \mathbf{Z}}} \exp\{g(s) + \boldsymbol{\gamma}^T \mathbf{X}\} ds.$$

*Under conditions C1, C2, C4 and C6,  $l$  has bounded and continuous first and second derivatives with respect to  $\boldsymbol{\beta} \in \mathcal{B}$ ,  $\boldsymbol{\gamma} \in \mathcal{T}$ , and  $\xi(\cdot, \boldsymbol{\beta}) \in \mathcal{H}^p$ .*

**Proof.** After some algebraic calculations, we have

$$\begin{aligned} l'_\beta(\boldsymbol{\beta}, \boldsymbol{\gamma}, \xi(\cdot, \boldsymbol{\beta}); \mathcal{O}) &= \Delta \mathbf{Z} \{ 1 + Ye^{\boldsymbol{\beta}^T \mathbf{Z}} g'(Ye^{\boldsymbol{\beta}^T \mathbf{Z}}) \} \\ &\quad - \mathbf{Z} Y \exp\{g(Ye^{\boldsymbol{\beta}^T \mathbf{Z}}) + \boldsymbol{\beta}^T \mathbf{Z} + \boldsymbol{\gamma}^T \mathbf{X}\}, \\ l'_\gamma(\boldsymbol{\beta}, \boldsymbol{\gamma}, \xi(\cdot, \boldsymbol{\beta}); \mathcal{O}) &= \Delta \mathbf{X} - \mathbf{X} \int_0^{Ye^{\boldsymbol{\beta}^T \mathbf{Z}}} \exp\{g(s) + \boldsymbol{\gamma}^T \mathbf{X}\} ds, \\ l'_\xi(\boldsymbol{\beta}, \boldsymbol{\gamma}, \xi(\cdot, \boldsymbol{\beta}); \mathcal{O})[h(\cdot, \boldsymbol{\beta})] &= \Delta w(Ye^{\boldsymbol{\beta}^T \mathbf{Z}}) - \int_0^{Ye^{\boldsymbol{\beta}^T \mathbf{Z}}} \exp\{g(s) + \boldsymbol{\gamma}^T \mathbf{X}\} w(s) ds, \\ l''_{\beta\beta}(\boldsymbol{\beta}, \boldsymbol{\gamma}, \xi(\cdot, \boldsymbol{\beta}); \mathcal{O}) &= \Delta \mathbf{Z} \mathbf{Z}^T Y^2 e^{2\boldsymbol{\beta}^T \mathbf{Z}} g''(Ye^{\boldsymbol{\beta}^T \mathbf{Z}}) + \Delta \mathbf{Z} \mathbf{Z}^T Y e^{\boldsymbol{\beta}^T \mathbf{Z}} g'(Ye^{\boldsymbol{\beta}^T \mathbf{Z}}) \end{aligned}$$

$$\begin{aligned}
 & -\mathbf{Z}\mathbf{Z}^T Y^2 g'(Ye^{\beta^T \mathbf{Z}}) \exp\{g(Ye^{\beta^T \mathbf{Z}}) + 2\beta^T \mathbf{Z} + \boldsymbol{\gamma}^T \mathbf{X}\} \\
 & -\mathbf{Z}\mathbf{Z}^T Y \exp\{g(Ye^{\beta^T \mathbf{Z}}) + \beta^T \mathbf{Z} + \boldsymbol{\gamma}^T \mathbf{X}\}, \\
 l''_{\boldsymbol{\gamma}\boldsymbol{\gamma}}(\boldsymbol{\beta}, \boldsymbol{\gamma}, \xi(\cdot, \boldsymbol{\beta}); \mathcal{O}) &= -\mathbf{X}\mathbf{X}^T \int_0^{Ye^{\beta^T \mathbf{Z}}} \exp\{g(s) + \boldsymbol{\gamma}^T \mathbf{X}\} ds, \\
 l''_{\boldsymbol{\beta}\boldsymbol{\gamma}}(\boldsymbol{\beta}, \boldsymbol{\gamma}, \xi(\cdot, \boldsymbol{\beta}); \mathcal{O}) &= -Y\mathbf{Z}\mathbf{X}^T \exp\{g(Ye^{\beta^T \mathbf{Z}}) + \beta^T \mathbf{Z} + \boldsymbol{\gamma}^T \mathbf{X}\}, \\
 l''_{\beta\xi}(\boldsymbol{\beta}, \boldsymbol{\gamma}, \xi(\cdot, \boldsymbol{\beta}); \mathcal{O})[h(\cdot, \boldsymbol{\beta})] &= \Delta Y\mathbf{Z}e^{\beta^T \mathbf{Z}} w'(Ye^{\beta^T \mathbf{Z}}) \\
 & \quad - Y\mathbf{Z}w(Ye^{\beta^T \mathbf{Z}}) \exp\{g(Ye^{\beta^T \mathbf{Z}}) + \beta^T \mathbf{Z} + \boldsymbol{\gamma}^T \mathbf{X}\}, \\
 l''_{\boldsymbol{\gamma}\xi}(\boldsymbol{\beta}, \boldsymbol{\gamma}, \xi(\cdot, \boldsymbol{\beta}); \mathcal{O})[h] &= -\mathbf{X} \int_0^{Ye^{\beta^T \mathbf{Z}}} \exp\{g(s) + \boldsymbol{\gamma}^T \mathbf{X}\} w(s) ds, \\
 l''_{\xi\xi}(\boldsymbol{\beta}, \boldsymbol{\gamma}, \xi(\cdot, \boldsymbol{\beta}); \mathcal{O})[h_1, h_2] &= - \int_0^{Ye^{\beta^T \mathbf{Z}}} \exp\{g(s) + \boldsymbol{\gamma}^T \mathbf{X}\} w_1(s)w_2(s) ds,
 \end{aligned}$$

where  $h(\cdot, \boldsymbol{\beta}) = w(\psi(\cdot, \boldsymbol{\beta}))$ ,  $h_1 = w_1(\psi(\cdot, \boldsymbol{\beta}))$ ,  $h_2(\cdot, \boldsymbol{\beta}) = w_2(\psi(\cdot, \boldsymbol{\beta})) \in \mathcal{H}$ . Under conditions C1, C2, C4 and C6, all the above derivatives are continuous and bounded.

Employing Corollary 6.21 in Schumaker [17], we directly have the following lemma with its proof omitted. □

**Lemma 2.** For  $g_0 \in \mathcal{G}^p$ , there exists a function  $g_{0n} \in \mathcal{G}_n^p$  such that

$$\|g_{0n} - g_0\|_\infty = O(n^{-pv}),$$

where  $\|\cdot\|_\infty$  is the sup-norm.

**Lemma 3.** Let  $\boldsymbol{\theta}_{0n} = (\boldsymbol{\beta}_0, \boldsymbol{\gamma}_0, \xi_{0n}(\cdot, \boldsymbol{\beta}_0))$  with  $\xi_{0n}(\cdot, \boldsymbol{\beta}_0) = g_{0n}$  defined in Lemma 2, and  $\mathcal{F}_n = \{l(\boldsymbol{\theta}; \mathcal{O}) - l(\boldsymbol{\theta}_{0n}; \mathcal{O}) : \boldsymbol{\theta} \in \boldsymbol{\Theta}_n^p\}$ . If conditions C1–C4 and C6 hold, then the  $\varepsilon$ -bracketing number associated with  $\|\cdot\|_\infty$  for  $\mathcal{F}_n$ , denoted by  $N_{[\cdot]}(\varepsilon, \mathcal{F}_n, \|\cdot\|_\infty)$ , is bounded by  $(1/\varepsilon)^{cq_n+d}$ , i.e.,

$$N_{[\cdot]}(\varepsilon, \mathcal{F}_n, \|\cdot\|_\infty) \lesssim (1/\varepsilon)^{cq_n+d}$$

for a constant  $c$ . Hereafter, we use the symbol  $\lesssim$  to denote that the left-hand side is bounded above by a constant times the right-hand side.

**Proof.** Denote the ceiling of  $x$  by  $\lceil x \rceil$ . By the calculation in Shen and Wong ([19], page 597), for any  $\varepsilon > 0$ , there exists a set of brackets

$$\{[g_i^L, g_i^U] : i = 1, \dots, \lceil (1/\varepsilon)^{cq_n} \rceil\}$$

such that for any  $g \in \mathcal{G}_n^p$ ,  $g_i^L(t) \leq g(t) \leq g_i^U(t)$  over  $t \in [0, b]$  for some  $1 \leq i \leq \lceil (1/\varepsilon)^{cq_n} \rceil$ , where  $\|g_i^U - g_i^L\|_\infty \leq \varepsilon$  and  $c$  is a constant. Since  $\mathcal{B}$  and  $\mathcal{T}$  are both compact under condition



C1,  $\mathcal{B}$  and  $\mathcal{T}$  can be covered by  $\lceil c_2(1/\varepsilon)^{d_1} \rceil$  and  $\lceil c_3(1/\varepsilon)^{d_2} \rceil$  balls with radius  $\varepsilon$ , respectively. Thus, for any  $\boldsymbol{\beta} \in \mathcal{B}$  and  $\boldsymbol{\gamma} \in \mathcal{T}$ , there exist  $\boldsymbol{\beta}_\ell$  for some  $1 \leq \ell \leq \lceil c_2(1/\varepsilon)^{d_1} \rceil$  and  $\boldsymbol{\gamma}_k$  for some  $1 \leq k \leq \lceil c_3(1/\varepsilon)^{d_2} \rceil$  such that  $\|\boldsymbol{\beta} - \boldsymbol{\beta}_\ell\| \leq \varepsilon$  and  $\|\boldsymbol{\gamma} - \boldsymbol{\gamma}_k\| \leq \varepsilon$ . Hence,  $\|\boldsymbol{\beta}^\top \mathbf{Z} - \boldsymbol{\beta}_\ell^\top \mathbf{Z}\| \leq c_4 \varepsilon$  and  $\|\boldsymbol{\gamma}^\top \mathbf{X} - \boldsymbol{\gamma}_k^\top \mathbf{X}\| \leq c_4 \varepsilon$  for some constant  $c_4$  under condition C2. Define

$$m_{i,\ell,k}^L(\mathcal{O}) = \Delta \left\{ \boldsymbol{\beta}_\ell^\top \mathbf{Z} + \boldsymbol{\gamma}_k^\top \mathbf{X} - 2c_4 \varepsilon + g_i^L(Y e^{\boldsymbol{\beta}_\ell^\top \mathbf{Z} + \xi_{i\ell}}) \right\} \\ - \int_0^{Y e^{\boldsymbol{\beta}_\ell^\top \mathbf{Z} + c_4 \varepsilon}} \exp\{g_i^U(s) + \boldsymbol{\gamma}_k^\top \mathbf{X} + c_4 \varepsilon\} ds - l(\boldsymbol{\theta}_{0n}; U)$$

and

$$m_{i,\ell,k}^U(\mathcal{O}) = \Delta \left\{ \boldsymbol{\beta}_\ell^\top \mathbf{Z} + \boldsymbol{\gamma}_k^\top \mathbf{X} + 2c_4 \varepsilon + g_i^U(Y e^{\boldsymbol{\beta}_\ell^\top \mathbf{Z} + \xi_{iu}}) \right\} \\ - \int_0^{Y e^{\boldsymbol{\beta}_\ell^\top \mathbf{Z} - c_4 \varepsilon}} \exp\{g_i^L(s) + \boldsymbol{\gamma}_k^\top \mathbf{X} - c_4 \varepsilon\} ds - l(\boldsymbol{\theta}_{0n}; \mathcal{O}),$$

where

$$g_i^L(e^{\boldsymbol{\beta}_\ell^\top \mathbf{Z} + \xi_{i\ell}}) = \min_{|s| \leq c_4 \varepsilon} g_i^L(e^{\boldsymbol{\beta}_\ell^\top \mathbf{Z} + s}) \quad \text{and} \quad g_i^U(e^{\boldsymbol{\beta}_\ell^\top \mathbf{Z} + \xi_{iu}}) = \max_{|s| \leq c_4 \varepsilon} g_i^U(e^{\boldsymbol{\beta}_\ell^\top \mathbf{Z} + s}).$$

After some calculations, we have  $|m_{i,\ell,k}^U(\mathcal{O}) - m_{i,\ell,k}^L(\mathcal{O})| \lesssim \varepsilon$  and for any  $m(\boldsymbol{\theta}; \mathcal{O}) \in \mathcal{F}_n$ , there exist some  $i, \ell$ , and  $k$  such that  $m(\boldsymbol{\theta}; \mathcal{O}) \in [m_{i,\ell,k}^L(\mathcal{O}), m_{i,\ell,k}^U(\mathcal{O})]$ . Therefore, we have

$$N_{[\cdot]}(\varepsilon, \mathcal{F}_n, \|\cdot\|_\infty) \lesssim (1/\varepsilon)^{cqn} (1/\varepsilon)^{d_1} (1/\varepsilon)^{d_2} = (1/\varepsilon)^{cqn+d}$$

for a constant  $c$ , which completes the proof.  $\square$

**Lemma 4.** *Let*

$$(w_1^*(t), \dots, w_{d_1}^*(t))^\top = \left\{ 1 + t g_0'(t) \right\} \frac{E\{\mathbf{Z} I(Y e^{\boldsymbol{\beta}_0^\top \mathbf{Z}} \geq t) e^{\boldsymbol{\gamma}_0^\top \mathbf{X}}\}}{E\{I(Y e^{\boldsymbol{\beta}_0^\top \mathbf{Z}} \geq t) e^{\boldsymbol{\gamma}_0^\top \mathbf{X}}\}}$$

and

$$(w_{d_1+1}^*(t), \dots, w_d^*(t))^\top = \frac{E\{\mathbf{X} I(Y e^{\boldsymbol{\beta}_0^\top \mathbf{Z}} \geq t) e^{\boldsymbol{\gamma}_0^\top \mathbf{X}}\}}{E\{I(Y e^{\boldsymbol{\beta}_0^\top \mathbf{Z}} \geq t) e^{\boldsymbol{\gamma}_0^\top \mathbf{X}}\}}.$$

If conditions C1–C6 hold, then there exist  $w_{j_n}^* \in \mathcal{G}_n^2$  ( $j = 1, \dots, d$ ) such that  $\|w_{j_n}^* - w_j^*\|_\infty = O(n^{-2v})$ ,  $j = 1, \dots, d$ .

**Proof.** Note that  $E\{I(Y e^{\boldsymbol{\beta}_0^\top \mathbf{Z}} \geq t) | \mathbf{Z}, \mathbf{X}\} = S_{C|\mathbf{Z}, \mathbf{X}}(te^{-\boldsymbol{\beta}_0^\top \mathbf{Z}}) \exp\{-\Lambda_0(t) \exp(\boldsymbol{\gamma}_0^\top \mathbf{X})\}$ , where  $S_{C|\mathbf{Z}, \mathbf{X}}(\cdot)$  is the conditional survival function of  $C$  given  $\mathbf{Z}$  and  $\mathbf{X}$ . It can be shown that the first and second derivatives of  $w_j^*$  are bounded under conditions C1–C6. Thus, according to Corollary 6.21 of Schumaker [17], the conclusion of this lemma follows.  $\square$

**Lemma 5.** Let  $h_j^*(\cdot, \boldsymbol{\beta}) = w_j^*(\psi(\cdot, \boldsymbol{\beta}))$  with  $h_j^*(t, \mathbf{z}, \boldsymbol{\beta}_0) = w_j^*(t)$  where  $w_j^*$  is defined in Lemma 4,  $j = 1, \dots, d$ . For  $\eta > 0$ , denote

$$\mathcal{F}_{jn}(\eta) = \{l'_\xi(\boldsymbol{\theta}; \mathcal{O})[h_j^* - h_j] : \boldsymbol{\theta} \in \Theta_n^p, h_j \in \mathcal{H}_n^p, d(\boldsymbol{\theta}, \boldsymbol{\theta}_0) \leq \eta, \|h_j^* - h_j\|_\infty \leq \eta\}.$$

If conditions C1–C6 hold, then  $N_{[\cdot]}(\varepsilon, \mathcal{F}_{jn}(\eta), \|\cdot\|_\infty) \lesssim (\eta/\varepsilon)^{c_1 q_n + d}$  for a constant  $c_1$ .

**Lemma 6.** Define

$$\mathcal{F}_{jn}^\beta(\eta) = \{l'_{\beta_j}(\boldsymbol{\theta}; \mathcal{O}) - l'_{\beta_j}(\boldsymbol{\theta}_0; \mathcal{O}) : \boldsymbol{\theta} \in \Theta_n^p, d(\boldsymbol{\theta}, \boldsymbol{\theta}_0) \leq \eta, \|g'(\psi(\cdot, \boldsymbol{\beta})) - g'(\psi(\cdot, \boldsymbol{\beta}_0))\|_2 \leq \eta\},$$

$$j = 1, \dots, d_1,$$

$$\mathcal{F}_{jn}^\gamma(\eta) = \{l'_{\gamma_j}(\boldsymbol{\theta}; \mathcal{O}) - l'_{\gamma_j}(\boldsymbol{\theta}_0; \mathcal{O}) : \boldsymbol{\theta} \in \Theta_n^p, d(\boldsymbol{\theta}, \boldsymbol{\theta}_0) \leq \eta\}, \quad j = 1, \dots, d_2,$$

and

$$\mathcal{F}_{jn}^\xi(\eta) = \{l'_\xi(\boldsymbol{\theta}; \mathcal{O})[h_j^*(\cdot, \boldsymbol{\beta})] - l'_\xi(\boldsymbol{\theta}_0; \mathcal{O})[h_j^*(\cdot, \boldsymbol{\beta}_0)] : \boldsymbol{\theta} \in \Theta_n^p, d(\boldsymbol{\theta}, \boldsymbol{\theta}_0) \leq \eta\},$$

$$j = 1, \dots, d,$$

where  $l'_{\beta_j}(\boldsymbol{\theta}; \mathcal{O})$  and  $l'_{\gamma_j}(\boldsymbol{\theta}; \mathcal{O})$  are the  $j$ th element of  $l'_\beta(\boldsymbol{\theta}; \mathcal{O})$  and  $l'_\gamma(\boldsymbol{\theta}; \mathcal{O})$ , respectively, and  $h_j^*$  is defined in Lemma 5. Suppose that conditions C1–C6 hold, then

$$N_{[\cdot]}(\varepsilon, \mathcal{F}_{jn}^\beta(\eta), \|\cdot\|_\infty) \lesssim (\eta/\varepsilon)^{c_1 q_n + d},$$

$$N_{[\cdot]}(\varepsilon, \mathcal{F}_{jn}^\gamma(\eta), \|\cdot\|_\infty) \lesssim (\eta/\varepsilon)^{c_2 q_n + d},$$

$$N_{[\cdot]}(\varepsilon, \mathcal{F}_{jn}^\xi(\eta), \|\cdot\|_\infty) \lesssim (\eta/\varepsilon)^{c_3 q_n + d},$$

for some constants  $c_1, c_2$ , and  $c_3$ .

The proofs of Lemmas 5 and 6 are similar to that of Lemma 3 and thus omitted here for the sake of space. The detailed proofs are available as supplementary materials from the authors.

**Proof of Theorem 1.** To obtain the convergence rate of the proposed estimator, we need to verify conditions C1–C3 of Theorem 1 in Shen and Wong [19]. Some algebraic calculations yield that

$$E\{l(\boldsymbol{\beta}, \boldsymbol{\gamma}, \xi(\cdot, \boldsymbol{\beta}); \mathcal{O})\} = E[\Delta\{\boldsymbol{\beta}^T \mathbf{Z} + \boldsymbol{\gamma}^T \mathbf{X} + g(Ye^{\boldsymbol{\beta}^T \mathbf{Z}})\}] - E[\Delta \exp\{(g(Ye^{\boldsymbol{\beta}^T \mathbf{Z}}) + \boldsymbol{\beta}^T \mathbf{Z} + \boldsymbol{\gamma}^T \mathbf{X}) - \{g_0(Ye^{\boldsymbol{\beta}_0^T \mathbf{Z}}) + \boldsymbol{\beta}_0^T \mathbf{Z} + \boldsymbol{\gamma}_0^T \mathbf{X}\})\}]$$

and

$$E\{l(\boldsymbol{\beta}_0, \boldsymbol{\gamma}_0, \xi_0(\cdot, \boldsymbol{\beta}_0); \mathcal{O})\} - E\{l(\boldsymbol{\beta}, \boldsymbol{\gamma}, \xi(\cdot, \boldsymbol{\beta}); \mathcal{O})\}$$

$$= E[\Delta\{\exp\{(g(Ye^{\boldsymbol{\beta}^T \mathbf{Z}}) + \boldsymbol{\beta}^T \mathbf{Z} + \boldsymbol{\gamma}^T \mathbf{X}) - \{g_0(Ye^{\boldsymbol{\beta}_0^T \mathbf{Z}}) + \boldsymbol{\beta}_0^T \mathbf{Z} + \boldsymbol{\gamma}_0^T \mathbf{X}\})\}]$$

$$\begin{aligned}
& -1 - (\{g(Ye^{\beta^T \mathbf{Z}}) + \beta^T \mathbf{Z} + \boldsymbol{\gamma}^T \mathbf{X}\} - \{g_0(Ye^{\beta_0^T \mathbf{Z}} + \beta_0^T \mathbf{Z} + \boldsymbol{\gamma}_0^T \mathbf{X})\}) \\
& \geq \frac{1}{2} E\{\Delta(\{g(Ye^{\beta^T \mathbf{Z}}) + \beta^T \mathbf{Z} + \boldsymbol{\gamma}^T \mathbf{X}\} - \{g_0(Ye^{\beta_0^T \mathbf{Z}}) + \beta_0^T \mathbf{Z} + \boldsymbol{\gamma}_0^T \mathbf{X}\})^2\} \\
& \geq \frac{1}{2} E\{\Delta(\{g(Ye^{\beta^T \mathbf{Z}}) - g_0(Ye^{\beta_0^T \mathbf{Z}})\} + (\beta - \beta_0)^T \mathbf{Z} + (\boldsymbol{\gamma} - \boldsymbol{\gamma}_0)^T \mathbf{X})^2\}.
\end{aligned} \tag{9}$$

Using the Taylor expansion, we have

$$\begin{aligned}
& E\{\Delta(\{g(Ye^{\beta^T \mathbf{Z}}) - g_0(Ye^{\beta_0^T \mathbf{Z}})\} + (\beta - \beta_0)^T \mathbf{Z} + (\boldsymbol{\gamma} - \boldsymbol{\gamma}_0)^T \mathbf{X})^2\} \\
& = E\{\Delta(\{g(Ye^{\beta^T \mathbf{Z}}) - g(Ye^{\beta_0^T \mathbf{Z}})\} + \{g(Ye^{\beta_0^T \mathbf{Z}}) - g_0(Ye^{\beta_0^T \mathbf{Z}})\} \\
& \quad + (\beta - \beta_0)^T \mathbf{Z} + (\boldsymbol{\gamma} - \boldsymbol{\gamma}_0)^T \mathbf{X})^2\} \\
& = E(\Delta[g'(\varepsilon_0)\varepsilon_0(\beta - \beta_0)^T \mathbf{Z} + O(|\beta - \beta_0|^2) + \{g(\varepsilon_0) - g_0(\varepsilon_0)\} \\
& \quad + (\beta - \beta_0)^T \mathbf{Z} + (\boldsymbol{\gamma} - \boldsymbol{\gamma}_0)^T \mathbf{X}]^2) \\
& = E(\Delta[\{g(\varepsilon_0) - g_0(\varepsilon_0)\} + (\boldsymbol{\alpha} - \boldsymbol{\alpha}_0)^T \mathbf{W}]^2) + o(d^2(\boldsymbol{\theta}, \boldsymbol{\theta}_0)),
\end{aligned} \tag{10}$$

where  $\mathbf{W} = (\{g'_0(\varepsilon_0)\varepsilon_0 + 1\}\mathbf{Z}^T, \mathbf{X}^T)^T$ . Obviously,

$$\begin{aligned}
& E(\Delta[\{g(\varepsilon_0) - g_0(\varepsilon_0)\} + (\boldsymbol{\alpha} - \boldsymbol{\alpha}_0)^T \mathbf{W}]^2) \\
& \geq E[\Delta\{g(\varepsilon_0) - g_0(\varepsilon_0)\}^2] + E[\Delta\{(\boldsymbol{\alpha} - \boldsymbol{\alpha}_0)^T \mathbf{W}\}^2] \\
& \quad - 2|E[\Delta\{g(\varepsilon_0) - g_0(\varepsilon_0)\}(\boldsymbol{\alpha} - \boldsymbol{\alpha}_0)^T \mathbf{W}]|.
\end{aligned} \tag{11}$$

On the other hand, it follows from the Cauchy–Schwarz inequality and condition C7 that

$$\begin{aligned}
& |E[\Delta\{g(\varepsilon_0) - g_0(\varepsilon_0)\}(\boldsymbol{\alpha} - \boldsymbol{\alpha}_0)^T \mathbf{W}]|^2 \\
& = |E[\Delta\{g(\varepsilon_0) - g_0(\varepsilon_0)\}E\{(\boldsymbol{\alpha} - \boldsymbol{\alpha}_0)^T \mathbf{W} | \varepsilon_0, \Delta = 1\}]|^2 \\
& \leq E[\Delta\{g(\varepsilon_0) - g_0(\varepsilon_0)\}^2]E[E\{(\boldsymbol{\alpha} - \boldsymbol{\alpha}_0)^T \mathbf{W} | \varepsilon_0, \Delta = 1\}]^2 \\
& \leq (1 - \eta)E[\Delta\{g(\varepsilon_0) - g_0(\varepsilon_0)\}^2]E[\Delta\{(\boldsymbol{\alpha} - \boldsymbol{\alpha}_0)^T \mathbf{W}\}^2].
\end{aligned} \tag{12}$$

Note that

$$d^2(\boldsymbol{\theta}, \boldsymbol{\theta}_0) \lesssim E[\Delta\{g(\varepsilon_0 e^{(\beta - \beta_0)^T \mathbf{Z}}) - g_0(\varepsilon_0)\}^2] + |\beta - \beta_0|^2 + |\boldsymbol{\gamma} - \boldsymbol{\gamma}_0|^2 \lesssim d^2(\boldsymbol{\theta}, \boldsymbol{\theta}_0) \tag{13}$$

and

$$\begin{aligned}
& E[\Delta\{g(\varepsilon_0 e^{(\beta - \beta_0)^T \mathbf{Z}}) - g_0(\varepsilon_0)\}^2] + |\beta - \beta_0|^2 + |\boldsymbol{\gamma} - \boldsymbol{\gamma}_0|^2 \\
& \lesssim E[\Delta\{g(\varepsilon_0) - g_0(\varepsilon_0)\}^2] + |\beta - \beta_0|^2 + |\boldsymbol{\gamma} - \boldsymbol{\gamma}_0|^2 \\
& \lesssim E[\Delta\{g(\varepsilon_0 e^{(\beta - \beta_0)^T \mathbf{Z}}) - g_0(\varepsilon_0)\}^2] + |\beta - \beta_0|^2 + |\boldsymbol{\gamma} - \boldsymbol{\gamma}_0|^2
\end{aligned} \tag{14}$$

under conditions C1–C4. Therefore, it follows from (10)–(14) that

$$\begin{aligned}
 & E\{\Delta(\{g(Ye^{\beta^T\mathbf{Z}}) - g_0(Ye^{\beta_0^T\mathbf{Z}})\} + (\boldsymbol{\beta} - \boldsymbol{\beta}_0)^T\mathbf{Z} + (\boldsymbol{\gamma} - \boldsymbol{\gamma}_0)^T\mathbf{X})^2\} \\
 & \geq \{1 - (1 - \eta)^{1/2}\}(E[\Delta\{g(\varepsilon_0) - g_0(\varepsilon_0)\}^2] + E[\Delta\{(\boldsymbol{\alpha} - \boldsymbol{\alpha}_0)^T\mathbf{W}\}^2]) \\
 & \gtrsim E[\Delta\{g(\varepsilon_0) - g_0(\varepsilon_0)\}^2] + |\boldsymbol{\beta} - \boldsymbol{\beta}_0|^2 + |\boldsymbol{\gamma} - \boldsymbol{\gamma}_0|^2 \\
 & \gtrsim d^2(\boldsymbol{\theta}, \boldsymbol{\theta}_0).
 \end{aligned}
 \tag{15}$$

Hence, condition C1 in Theorem 1 of Shen and Wong [19] holds from (9) and (15).

Next, we verify their condition C2. Note that under our conditions C2 and C6,

$$\begin{aligned}
 & \{l(\boldsymbol{\beta}, \boldsymbol{\gamma}, \xi(\cdot, \boldsymbol{\beta}); \mathcal{O}) - l(\boldsymbol{\beta}_0, \boldsymbol{\gamma}_0, \xi_0(\cdot, \boldsymbol{\beta}_0); \mathcal{O})\}^2 \\
 & \lesssim \{(\boldsymbol{\beta} - \boldsymbol{\beta}_0)^T\mathbf{Z}\}^2 + \{(\boldsymbol{\gamma} - \boldsymbol{\gamma}_0)^T\mathbf{X}\}^2 \\
 & \quad + \Delta\{g(Ye^{\beta^T\mathbf{Z}}) - g(Ye^{\beta_0^T\mathbf{Z}})\}^2 + \Delta\{g(Ye^{\beta_0^T\mathbf{Z}}) - g_0(Ye^{\beta_0^T\mathbf{Z}})\}^2 \\
 & \quad + \int_0^Y [\exp\{g(se^{\beta^T\mathbf{Z}}) + \boldsymbol{\beta}^T\mathbf{Z} + \boldsymbol{\gamma}^T\mathbf{X}\} - \exp\{g_0(se^{\beta_0^T\mathbf{Z}}) + \boldsymbol{\beta}_0^T\mathbf{Z} + \boldsymbol{\gamma}_0^T\mathbf{X}\}]^2 ds.
 \end{aligned}
 \tag{16}$$

Using the Taylor expansion, we have for any  $\boldsymbol{\theta} \in \Theta^p$ ,

$$E[\Delta\{g(Ye^{\beta^T\mathbf{Z}}) - g(Ye^{\beta_0^T\mathbf{Z}})\}^2] \lesssim |\boldsymbol{\beta} - \boldsymbol{\beta}_0|^2,
 \tag{17}$$

and

$$\begin{aligned}
 & E \int_0^Y [\exp\{g(se^{\beta^T\mathbf{Z}}) + \boldsymbol{\beta}^T\mathbf{Z} + \boldsymbol{\gamma}^T\mathbf{X}\} - \exp\{g_0(se^{\beta_0^T\mathbf{Z}}) + \boldsymbol{\beta}_0^T\mathbf{Z} + \boldsymbol{\gamma}_0^T\mathbf{X}\}]^2 ds \\
 & \lesssim E \int_0^Y \exp\{g_0(se^{\beta_0^T\mathbf{Z}}) + \boldsymbol{\beta}_0^T\mathbf{Z} + \boldsymbol{\gamma}_0^T\mathbf{X}\} \{g(se^{\beta_0^T\mathbf{Z}}) - g_0(se^{\beta_0^T\mathbf{Z}})\}^2 ds \\
 & \quad + |\boldsymbol{\beta} - \boldsymbol{\beta}_0|^2 + |\boldsymbol{\gamma} - \boldsymbol{\gamma}_0|^2 \\
 & = E[\Delta\{g(Ye^{\beta_0^T\mathbf{Z}}) - g_0(Ye^{\beta_0^T\mathbf{Z}})\}^2] + |\boldsymbol{\beta} - \boldsymbol{\beta}_0|^2 + |\boldsymbol{\gamma} - \boldsymbol{\gamma}_0|^2 \lesssim d^2(\boldsymbol{\theta}, \boldsymbol{\theta}_0).
 \end{aligned}
 \tag{18}$$

Thus, combining (16)–(18), we obtain that

$$E\{l(\boldsymbol{\beta}, \boldsymbol{\gamma}, \xi(\cdot, \boldsymbol{\beta}); \mathcal{O}) - l(\boldsymbol{\beta}_0, \boldsymbol{\gamma}_0, \xi_0(\cdot, \boldsymbol{\beta}_0); \mathcal{O})\}^2 \lesssim d^2(\boldsymbol{\theta}, \boldsymbol{\theta}_0),$$

which implies condition C2 in Theorem 1 of Shen and Wong [19].

As  $\widehat{\boldsymbol{\theta}}_n$  maximizes the log-likelihood  $P_n l(\boldsymbol{\theta}; \mathcal{O})$  over the sieve space  $\Theta_n^p$ ,  $\widehat{\boldsymbol{\theta}}_n$  satisfies inequality (1.1) in Shen and Wong [19] with  $\eta_n = 0$ . It follows from Lemma 2 that there exists a  $\xi_{0n}(\cdot, \boldsymbol{\beta}_0) \in \mathcal{H}_n^p$  such that  $\|\xi_{0n} - \xi_0\|_\infty = O(n^{-pv})$ . The Kullback–Leibler distance between

$\theta_0 = (\beta_0, \gamma_0, \xi_0(\cdot, \beta_0))$  and  $\theta_{0n} = (\beta_0, \gamma_0, \xi_{0n}(\cdot, \beta_0))$  is given by

$$\begin{aligned} K(\theta_0, \theta_{0n}) &= P\{l(\theta_0; \mathcal{O}) - l(\theta_{0n}; \mathcal{O})\} \\ &\lesssim \|\xi_0(\cdot, \beta_0) - \xi_{0n}(\cdot, \beta_0)\|_2^2 \\ &\lesssim \|\xi_{0n}(\cdot, \beta_0) - \xi_0(\cdot, \beta_0)\|_\infty^2 \\ &= O(n^{-2pv}). \end{aligned}$$

Thus, it follows from Theorem 1 of Shen and Wong [19] that

$$d(\widehat{\theta}_n, \theta_0) = O_p(n^{-\min(pv, (1-v)/2)}),$$

which completes the proof of Theorem 1. □

**Proof of Theorem 2.** Employing Theorem 2.1 of Ding and Nan [7], it suffices to verify the following conditions to prove Theorem 2.

- A1.  $d(\widehat{\theta}_n, \theta_0) = O_p(n^{-\rho})$  for some  $\rho > 0$ .
- A2.  $Pl_\alpha(\alpha_0, \xi_0(\cdot, \beta_0); \mathcal{O}) = \mathbf{0}$  and  $Pl_\xi(\alpha_0, \xi_0(\cdot, \beta_0); \mathcal{O})[h] = \mathbf{0}$  for all  $h \in \mathcal{H}$ .
- A3. There exists an  $\mathbf{h}^* = (h_1^*, \dots, h_d^*)^\top$ , where  $h_j^* \in \mathcal{H}$  for  $j = 1, \dots, d$ , such that

$$Pl''_{\alpha\xi}(\alpha_0, \xi_0(\cdot, \beta_0); \mathcal{O})[h] - Pl''_{\xi\xi}(\alpha_0, \xi_0(\cdot, \beta_0); \mathcal{O})[\mathbf{h}^*, h] = \mathbf{0},$$

for all  $h \in \mathcal{H}$ . Furthermore, the matrix  $P\{l''_{\alpha\alpha}(\alpha_0, \xi_0(\cdot, \beta_0); \mathcal{O}) - l''_{\xi\alpha}(\alpha_0, \xi_0(\cdot, \beta_0); \mathcal{O})[\mathbf{h}^*]\}$  is non-singular.

- A4.  $P_n l'_\alpha(\widehat{\alpha}_n, \widehat{\xi}_n(\cdot, \widehat{\beta}_n); \mathcal{O}) = o_p(n^{-1/2})$  and  $P_n l'_\xi(\widehat{\alpha}_n, \widehat{\xi}_n(\cdot, \widehat{\beta}_n); \mathcal{O})[\mathbf{h}^*] = o_p(n^{-1/2})$ .
- A5. Let  $G_n = n^{1/2}(P_n - P)$ . For any  $c > 0$ ,

$$\sup_{d(\theta, \theta_0) \leq cn^{-\rho}, \theta \in \Theta_n^p} |G_n l'_\alpha(\alpha, \xi(\cdot, \beta); \mathcal{O}) - G_n l'_\alpha(\alpha_0, \xi_0(\cdot, \beta_0); \mathcal{O})| = o_p(1)$$

and

$$\begin{aligned} &\sup_{d(\theta, \theta_0) \leq cn^{-\rho}, \theta \in \Theta_n^p} |G_n l'_\xi(\alpha, \xi(\cdot, \beta); \mathcal{O})[\mathbf{h}^*(\cdot, \beta)] - G_n l'_\xi(\alpha_0, \xi_0(\cdot, \beta_0); \mathcal{O})[\mathbf{h}^*(\cdot, \beta)]| \\ &= o_p(1). \end{aligned}$$

- A6. For some  $\zeta > 1$  satisfying that  $\rho\zeta > 1/2$  and for any  $\theta$  in a neighborhood of  $\theta_0$ ,  $\{\theta : d(\theta, \theta_0) \leq cn^{-\rho}, \theta \in \Theta_n^p\}$  say,

$$\begin{aligned} &|Pl'_\alpha(\alpha, \xi(\cdot, \beta); \mathcal{O}) - Pl'_\alpha(\alpha_0, \xi_0(\cdot, \beta_0); \mathcal{O}) - Pl''_{\alpha\alpha}(\alpha_0, \xi_0(\cdot, \beta_0); \mathcal{O})(\alpha - \alpha_0) \\ &\quad - Pl''_{\alpha\xi}(\alpha_0, \xi_0(\cdot, \beta_0); \mathcal{O})[\xi(\cdot, \beta) - \xi_0(\cdot, \beta_0)]| \\ &= O(d^\zeta(\theta, \theta_0)) \end{aligned}$$

and

$$\begin{aligned} & |Pl'_{\xi}(\boldsymbol{\alpha}, \xi(\cdot, \boldsymbol{\beta}); \mathcal{O})[\mathbf{h}^*(\cdot, \boldsymbol{\beta})] - Pl'_{\xi}(\boldsymbol{\alpha}_0, \xi_0(\cdot, \boldsymbol{\beta}_0); \mathcal{O})[\mathbf{h}^*(\cdot, \boldsymbol{\beta}_0)] \\ & \quad - Pl''_{\xi\boldsymbol{\alpha}}(\boldsymbol{\alpha}_0, \xi_0(\cdot, \boldsymbol{\beta}_0); \mathcal{O})[\mathbf{h}^*(\cdot, \boldsymbol{\beta}_0)](\boldsymbol{\alpha} - \boldsymbol{\alpha}_0) \\ & \quad - Pl'_{\xi\xi}(\boldsymbol{\alpha}_0, \xi_0(\cdot, \boldsymbol{\beta}_0); \mathcal{O})[\mathbf{h}^*(\cdot, \boldsymbol{\beta}), \xi(\cdot, \boldsymbol{\beta}) - \xi_0(\cdot, \boldsymbol{\beta}_0)]| \\ & = O(d^{\zeta}(\boldsymbol{\theta}, \boldsymbol{\theta}_0)). \end{aligned}$$

We briefly discuss these conditions prior to verification of each of them. The convergence rate in A1, established in Theorem 1, is a prerequisite while proving the asymptotic normality. Condition A2 evaluates the score function (at the population level) at the true value, while A3 states the least favorable direction and the nonsingular information matrix along such a direction; both of them are standard in the maximum likelihood theory. Condition A4 assesses the score function (at the sample level) at the estimator and the stochastic equicontinuity in A5 can typically be verified by either the Donsker property or the maximal inequality (van der Vaart and Wellner [21]). Finally, the Taylor expansion results in A6.

First, A1 holds by choosing  $\rho = \min(pv, (1 - v)/2)$  from Theorem 1. Using the fact of zero-mean score functions, it is easy to show A2 holds.

Next, we find  $\mathbf{h}^* = (h_1^*, \dots, h_d^*)^T$  with  $\mathbf{h}^*(t, z, \boldsymbol{\beta}_0) = \mathbf{w}^*(t)$  such that A3 holds. Denote  $\mathbf{h}^*_{\boldsymbol{\beta}} = (h_1^*, \dots, h_{d_1}^*)^T$ ,  $\mathbf{h}^*_{\boldsymbol{\gamma}} = (h_{d_1+1}^*, \dots, h_d^*)^T$ ,  $\mathbf{w}^*_{\boldsymbol{\beta}} = (w_1^*, \dots, w_{d_1}^*)^T$ , and  $\mathbf{w}^*_{\boldsymbol{\gamma}} = (w_{d_1+1}^*, \dots, w_d^*)^T$ . For any  $h \in \mathcal{H}$ ,

$$\begin{aligned} & Pl''_{\boldsymbol{\alpha}\xi}(\boldsymbol{\alpha}_0, \xi_0(\cdot, \boldsymbol{\beta}_0); \mathcal{O})[h] - Pl''_{\xi\xi}(\boldsymbol{\alpha}_0, \xi_0(\cdot, \boldsymbol{\beta}_0); \mathcal{O})[\mathbf{h}^*, h] \\ & = P \left( \begin{aligned} & l''_{\boldsymbol{\beta}\xi}(\boldsymbol{\beta}_0, \boldsymbol{\gamma}_0, \xi_0(\cdot, \boldsymbol{\beta}_0); \mathcal{O})[h] - l''_{\xi\xi}(\boldsymbol{\beta}_0, \boldsymbol{\gamma}_0, \xi_0(\cdot, \boldsymbol{\beta}_0); \mathcal{O})[\mathbf{h}^*_{\boldsymbol{\beta}}, h] \\ & l''_{\boldsymbol{\gamma}\xi}(\boldsymbol{\beta}_0, \boldsymbol{\gamma}_0, \xi_0(\cdot, \boldsymbol{\beta}_0); \mathcal{O})[h] - l''_{\xi\xi}(\boldsymbol{\beta}_0, \boldsymbol{\gamma}_0, \xi_0(\cdot, \boldsymbol{\beta}_0); \mathcal{O})[\mathbf{h}^*_{\boldsymbol{\gamma}}, h] \end{aligned} \right). \end{aligned}$$

Some calculation entails that

$$\begin{aligned} & Pl''_{\boldsymbol{\beta}\xi}(\boldsymbol{\beta}_0, \boldsymbol{\gamma}_0, \xi_0(\cdot, \boldsymbol{\beta}_0); \mathcal{O})[h] - Pl''_{\xi\xi}(\boldsymbol{\beta}_0, \boldsymbol{\gamma}_0, \xi_0(\cdot, \boldsymbol{\beta}_0); \mathcal{O})[\mathbf{h}^*_{\boldsymbol{\beta}}, h] \\ & = E\{\Delta Y \mathbf{Z} e^{\boldsymbol{\beta}_0^T \mathbf{Z}} \mathbf{Z}' w'(Y e^{\boldsymbol{\beta}_0^T \mathbf{Z}})\} - E[Y \mathbf{Z} w(Y e^{\boldsymbol{\beta}_0^T \mathbf{Z}}) \exp\{g_0(Y e^{\boldsymbol{\beta}_0^T \mathbf{Z}}) + \boldsymbol{\beta}_0^T \mathbf{Z} + \boldsymbol{\gamma}_0^T \mathbf{X}\}] \\ & \quad + E\left[\int_0^{Y e^{\boldsymbol{\beta}_0^T \mathbf{Z}}} \exp\{g_0(s) + \boldsymbol{\gamma}_0^T \mathbf{X}\} \mathbf{w}^*_{\boldsymbol{\beta}}(s) w(s) ds\right]. \end{aligned}$$

In what follows, we calculate the above expectations using the ordinary properties of conditional expectation. We denote the conditional survival function of  $T$  given  $\mathbf{Z}$  and  $\mathbf{X}$  by  $S_{T|\mathbf{Z}, \mathbf{X}}(\cdot|\mathbf{Z}, \mathbf{X})$  and the corresponding conditional density function by  $f_{T|\mathbf{Z}, \mathbf{X}}(\cdot|\mathbf{Z}, \mathbf{X})$ . After some tedious but straightforward calculations, we have

$$\begin{aligned} & E\{\Delta Y \mathbf{Z} e^{\boldsymbol{\beta}_0^T \mathbf{Z}} \mathbf{Z}' w'(Y e^{\boldsymbol{\beta}_0^T \mathbf{Z}})|C, \mathbf{Z}, \mathbf{X}\} \\ & = \mathbf{Z} C w(C e^{\boldsymbol{\beta}_0^T \mathbf{Z}}) f_{T|\mathbf{Z}, \mathbf{X}}(C|\mathbf{Z}, \mathbf{X}) - \mathbf{Z} \int_0^C w(t e^{\boldsymbol{\beta}_0^T \mathbf{Z}}) f_{T|\mathbf{Z}, \mathbf{X}}(t|\mathbf{Z}, \mathbf{X}) dt \end{aligned}$$

$$\begin{aligned}
 & -\mathbf{Z} \int_0^C t w(te^{\beta_0^T \mathbf{Z}}) [g_0'(te^{\beta_0^T \mathbf{Z}}) e^{\beta_0^T \mathbf{Z}} - \exp\{g_0(te^{\beta_0^T \mathbf{Z}}) + \beta_0^T \mathbf{Z} + \boldsymbol{\gamma}_0^T \mathbf{X}\}] f_{T|\mathbf{Z}, \mathbf{X}}(t|\mathbf{Z}, \mathbf{X}) dt, \\
 & E[\boldsymbol{\gamma} \mathbf{Z} w(Ye^{\beta_0^T \mathbf{Z}}) \exp\{g_0(Ye^{\beta_0^T \mathbf{Z}}) + \beta_0^T \mathbf{Z} + \boldsymbol{\gamma}_0^T \mathbf{X}\} | C, \mathbf{Z}, \mathbf{X}] \\
 & = \mathbf{Z} \int_0^C t w(te^{\beta_0^T \mathbf{Z}}) \exp\{g_0(te^{\beta_0^T \mathbf{Z}}) + \beta_0^T \mathbf{Z} + \boldsymbol{\gamma}_0^T \mathbf{X}\} f_{T|\mathbf{Z}, \mathbf{X}}(t|\mathbf{Z}, \mathbf{X}) dt \\
 & \quad + \mathbf{Z} C w(Ce^{\beta_0^T \mathbf{Z}}) f_{T|\mathbf{Z}, \mathbf{X}}(C|\mathbf{Z}, \mathbf{X}),
 \end{aligned}$$

and

$$\begin{aligned}
 & E \left[ \int_0^{Ye^{\beta_0^T \mathbf{Z}}} \exp\{g_0(s) + \boldsymbol{\gamma}_0^T \mathbf{X}\} \mathbf{w}_\beta^*(s) h(s) ds \mid C, \mathbf{Z}, \mathbf{X} \right] \\
 & = \int_0^C \mathbf{w}_\beta^*(se^{\beta_0^T \mathbf{Z}}) w(se^{\beta_0^T \mathbf{Z}}) f_{T|\mathbf{Z}, \mathbf{X}}(s|\mathbf{Z}, \mathbf{X}) ds.
 \end{aligned}$$

Thus, we obtain that

$$\begin{aligned}
 & Pl''_{\beta\xi}(\boldsymbol{\beta}_0, \boldsymbol{\gamma}_0, \xi_0(\cdot, \boldsymbol{\beta}_0); \mathcal{O})[h] - Pl''_{\xi\xi}(\boldsymbol{\beta}_0, \boldsymbol{\gamma}_0, \xi_0(\cdot, \boldsymbol{\beta}_0); \mathcal{O})[\mathbf{h}_\beta^*, h] \\
 & = E \left[ \int_0^C \{ \mathbf{w}_\beta^*(se^{\beta_0^T \mathbf{Z}}) - \mathbf{Z}(1 + se^{\beta_0^T \mathbf{Z}} g_0'(se^{\beta_0^T \mathbf{Z}})) \} w(se^{\beta_0^T \mathbf{Z}}) f_{T|\mathbf{Z}, \mathbf{X}}(s|\mathbf{Z}, \mathbf{X}) ds \right] \\
 & = E \left[ \int_0^{Ce^{\beta_0^T \mathbf{Z}}} \{ \mathbf{w}_\beta^*(t) - \mathbf{Z}(1 + t g_0'(t)) \} w(t) f_{T|\mathbf{Z}, \mathbf{X}}(te^{-\beta_0^T \mathbf{Z}}|\mathbf{Z}, \mathbf{X}) e^{-\beta_0^T \mathbf{Z}} dt \right] \\
 & = E \left[ \int_0^{+\infty} I(Ce^{\beta_0^T \mathbf{Z}} \geq t) \{ \mathbf{w}_\beta^*(t) - \mathbf{Z}(1 + t g_0'(t)) \} w(t) f_{T|\mathbf{Z}, \mathbf{X}}(te^{-\beta_0^T \mathbf{Z}}|\mathbf{Z}, \mathbf{X}) e^{-\beta_0^T \mathbf{Z}} dt \right] \\
 & = \int_0^{+\infty} [E\{I(Ce^{\beta_0^T \mathbf{Z}} \geq t) \exp\{-\Lambda_0(t)e^{\boldsymbol{\gamma}_0^T \mathbf{X}}\} e^{\boldsymbol{\gamma}_0^T \mathbf{X}}\} \mathbf{w}_\beta^*(t) \\
 & \quad - (1 + t g_0'(t)) E\{\mathbf{Z} I(Ce^{\beta_0^T \mathbf{Z}} \geq t) \exp(-\Lambda_0(t)e^{\boldsymbol{\gamma}_0^T \mathbf{X}}) e^{\boldsymbol{\gamma}_0^T \mathbf{X}}\}] w(t) \lambda_0(t) dt.
 \end{aligned}$$

Therefore, we take  $\mathbf{h}_\beta^*$  with

$$\mathbf{h}_\beta^*(t, z, \boldsymbol{\beta}_0) = \mathbf{w}_\beta^*(t) = \frac{\{1 + t g_0'(t)\} E[\mathbf{Z} I(Ce^{\beta_0^T \mathbf{Z}} \geq t) \exp\{-\Lambda_0(t)e^{\boldsymbol{\gamma}_0^T \mathbf{X}}\} e^{\boldsymbol{\gamma}_0^T \mathbf{X}}]}{E[I(Ce^{\beta_0^T \mathbf{Z}} \geq t) \exp\{-\Lambda_0(t)e^{\boldsymbol{\gamma}_0^T \mathbf{X}}\} e^{\boldsymbol{\gamma}_0^T \mathbf{X}}]},$$

which makes  $Pl''_{\beta g}(\boldsymbol{\beta}_0, \boldsymbol{\gamma}_0, g_0; \mathcal{O})[h] - Pl''_{gg}(\boldsymbol{\beta}_0, \boldsymbol{\gamma}_0, g_0; \mathcal{O})[\mathbf{h}_\beta^*, h] = \mathbf{0}$  for any  $h \in \mathcal{H}$ .

Based on similar but simpler calculations, we also have that

$$\begin{aligned}
 & Pl''_{\boldsymbol{\gamma} g}(\boldsymbol{\beta}_0, \boldsymbol{\gamma}_0, g_0; \mathcal{O})[h] - Pl''_{gg}(\boldsymbol{\beta}_0, \boldsymbol{\gamma}_0, g_0; \mathcal{O})[\mathbf{h}_\boldsymbol{\gamma}^*, h] \\
 & = \int_0^{+\infty} [E\{I(Ce^{\beta_0^T \mathbf{Z}} \geq t) \exp\{-\Lambda_0(t)e^{\boldsymbol{\gamma}_0^T \mathbf{X}}\} e^{\boldsymbol{\gamma}_0^T \mathbf{X}}\} \mathbf{w}_\boldsymbol{\gamma}^*(t)
 \end{aligned}$$

$$- E \{ \mathbf{X} I ( C e^{\beta_0^T \mathbf{Z}} \geq t ) \exp \{ - \Lambda_0 ( t ) e^{\gamma_0^T \mathbf{X}} \} e^{\gamma_0^T \mathbf{X}} \} w ( t ) \lambda_0 ( t ) dt .$$

Thus, we take  $\mathbf{h}_\gamma^*$  with

$$\mathbf{h}_\gamma^* ( t, z, \beta_0 ) = \mathbf{w}_\gamma^* ( t ) = \frac{ E [ \mathbf{X} I ( C e^{\beta_0^T \mathbf{Z}} \geq t ) \exp \{ - \Lambda_0 ( t ) e^{\gamma_0^T \mathbf{X}} \} e^{\gamma_0^T \mathbf{X}} ] }{ E [ I ( C e^{\beta_0^T \mathbf{Z}} \geq t ) \exp \{ - \Lambda_0 ( t ) e^{\gamma_0^T \mathbf{X}} \} e^{\gamma_0^T \mathbf{X}} ] }$$

such that  $Pl''_{\gamma g}(\beta_0, \gamma_0, g_0; \mathcal{O})[h] - Pl''_{gg}(\beta_0, \gamma_0, g_0; \mathcal{O})[\mathbf{h}_\gamma^*, h] = \mathbf{0}$  for any  $h \in \mathcal{H}$ .

Note that

$$\begin{aligned} & P ( Y e^{\beta_0^T \mathbf{Z}} \geq t | C, \mathbf{Z}, \mathbf{X} ) \\ &= P ( T e^{\beta_0^T \mathbf{Z}} \geq t, T \leq C | C, \mathbf{Z}, \mathbf{X} ) + P ( C e^{\beta_0^T \mathbf{Z}} \geq t, T > C | C, \mathbf{Z}, \mathbf{X} ) \\ &= P ( t e^{-\beta_0^T \mathbf{Z}} \leq T \leq C | C, \mathbf{Z}, \mathbf{X} ) I ( C \geq t e^{-\beta_0^T \mathbf{Z}} ) + P ( T > C | C, \mathbf{Z}, \mathbf{X} ) I ( C \geq t e^{-\beta_0^T \mathbf{Z}} ) \\ &= S_{T|Z, X} ( t e^{-\beta_0^T \mathbf{Z}} | \mathbf{Z}, \mathbf{X} ) I ( C \geq t e^{-\beta_0^T \mathbf{Z}} ) \\ &= I ( C \geq t e^{-\beta_0^T \mathbf{Z}} ) \exp \{ - \Lambda_0 ( t ) e^{\gamma_0^T \mathbf{X}} \} . \end{aligned}$$

Then,  $\mathbf{w}_\beta^*$  and  $\mathbf{w}_\gamma^*$  can be simplified as

$$\mathbf{w}_\beta^* ( t ) = \frac{ \{ 1 + t g_0' ( t ) \} E \{ \mathbf{Z} I ( \varepsilon_0 \geq t ) e^{\gamma_0^T \mathbf{X}} \} }{ E \{ I ( \varepsilon_0 \geq t ) e^{\gamma_0^T \mathbf{X}} \} }$$

and

$$\mathbf{w}_\gamma^* ( t ) = \frac{ E \{ \mathbf{X} I ( \varepsilon_0 \geq t ) e^{\gamma_0^T \mathbf{X}} \} }{ E \{ I ( \varepsilon_0 \geq t ) e^{\gamma_0^T \mathbf{X}} \} } .$$

Hence, we have found  $\mathbf{h}^* = ( h_1^*, \dots, h_d^* )^T$  such that for any  $h \in \mathcal{H}$ ,

$$Pl''_{\alpha \xi}(\alpha_0, \xi_0(\cdot, \beta_0); \mathcal{O})[h] - Pl''_{\xi \xi}(\alpha_0, \xi_0(\cdot, \beta_0); \mathcal{O})[\mathbf{h}^*, h] = \mathbf{0} .$$

Furthermore, we obtain

$$l'_\beta(\beta_0, \gamma_0, \xi_0(\cdot, \beta_0); \mathcal{O}) - l'_\xi(\beta_0, \gamma_0, \xi_0(\cdot, \beta_0); \mathcal{O})[\mathbf{h}_\beta^*] = l_{\beta_0}^* (\mathcal{O})$$

and

$$l'_\gamma(\beta_0, \gamma_0, \xi_0(\cdot, \beta_0); \mathcal{O}) - l'_\xi(\beta_0, \gamma_0, \xi_0(\cdot, \beta_0); \mathcal{O})[\mathbf{h}_\gamma^*] = l_{\gamma_0}^* (\mathcal{O}) ,$$

which are the efficient score functions for  $\beta_0$  and  $\gamma_0$ , respectively. We can also show that

$$P \{ l''_{\alpha \alpha}(\alpha_0, \xi_0(\cdot, \beta_0); \mathcal{O}) - l''_{\xi \alpha}(\alpha_0, \xi_0(\cdot, \beta_0); \mathcal{O})[\mathbf{h}^*] \} = - P l_{\alpha_0}^* (\mathcal{O})^{\otimes 2} ,$$

which is the negative information matrix for  $\alpha_0$ . Thus, it is invertible under condition C8. Hence, A3 holds.



Using Lemmas 4–6, the Taylor expansion, the maximal inequality in Lemma 3.4.2 of van der Vaart and Wellner [21] or Theorem 11.3 of Kosorok [12], and the Markov inequality, we can show that assumptions A4–A6 hold for  $\rho = \min\{pv, (1 - v)/2\}$  and

$$\zeta = \min\{2(p - 1)v, 1/2 + (p - 5/2)v, 1 - v\} / \min\{pv, (1 - v)/2\} > 1.$$

Therefore, by Theorem 6.1 of Wellner and Zhang [23], we have

$$n^{1/2}(\widehat{\alpha}_n - \alpha_0) = \{\mathcal{I}(\alpha_0)\}^{-1} n^{1/2} P_n l_{\alpha_0}^*(\theta_0; \mathcal{O}) + o_p(1) \rightarrow_d N(0, \{\mathcal{I}(\alpha_0)\}^{-1}),$$

where  $l_{\alpha_0}^*(\theta_0; \mathcal{O}) = l'_{\alpha}(\theta_0; \mathcal{O}) - l'_{\xi}(\theta_0; \mathcal{O})[\mathbf{h}^*]$  is the efficient score function for  $\alpha_0$ . This completes the proof of Theorem 2. □

**Proof of Theorem 3.** Define

$$\mathbf{w}_n^*(t) = (\{1 + t\widehat{g}'_n(t)\}\bar{\mathbf{Z}}^T(t; \widehat{\beta}_n, \widehat{\gamma}_n), \bar{\mathbf{X}}^T(t; \widehat{\beta}_n, \widehat{\gamma}_n))^T.$$

Then we have

$$\widehat{l}_{\alpha_0}^*(\widehat{\theta}_n; \mathcal{O}) = l'_{\alpha}(\widehat{\theta}_n; \mathcal{O}) - l'_{\xi}(\widehat{\theta}_n; \mathcal{O})[\mathbf{h}_n^*].$$

Let

$$\begin{aligned} \mathcal{I}^{jk}(\theta_0) &= P[\{l'_{\alpha_j}(\theta_0; \mathcal{O}) - l'_{\xi}(\theta_0; \mathcal{O})[h_j^*]\}\{l'_{\alpha_k}(\theta_0; \mathcal{O}) - l'_{\xi}(\theta_0; \mathcal{O})[h_k^*]\}] \\ &\equiv PA^{jk}(\theta_0; \mathcal{O}) \end{aligned}$$

and

$$\begin{aligned} \widehat{\mathcal{I}}_n^{jk}(\theta_0) &= P_n[\{l'_{\alpha_j}(\widehat{\theta}_n; \mathcal{O}) - l'_{\xi}(\widehat{\theta}_n; \mathcal{O})[h_{jn}^*]\}\{l'_{\alpha_k}(\widehat{\theta}_n; \mathcal{O}) - l'_{\xi}(\widehat{\theta}_n; \mathcal{O})[h_{kn}^*]\}] \\ &\equiv P_n A_n^{jk}(\widehat{\theta}_n; \mathcal{O}) \end{aligned}$$

for  $j, k = 1, \dots, d$ . It suffices to show that  $P_n A_n^{jk}(\widehat{\theta}_n; \mathcal{O})$  converges to  $PA^{jk}(\theta_0; \mathcal{O})$  in probability. Note that

$$\begin{aligned} &P_n A_n^{jk}(\widehat{\theta}_n; \mathcal{O}) - PA^{jk}(\theta_0; \mathcal{O}) \\ &= P_n\{A_n^{jk}(\widehat{\theta}_n; \mathcal{O}) - A^{jk}(\theta_0; \mathcal{O})\} + (P_n - P)A^{jk}(\theta_0; \mathcal{O}). \end{aligned}$$

Clearly,  $(P_n - P)A^{jk}(\theta_0; \mathcal{O}) = o_p(1)$ . On the other hand, under conditions C2 and C6, we have

$$\begin{aligned} &P\{A_n^{jk}(\widehat{\theta}_n; \mathcal{O}) - A^{jk}(\theta_0; \mathcal{O})\}^2 \\ &\lesssim P\{l'_{\alpha_j}(\widehat{\theta}_n; \mathcal{O}) - l'_{\alpha_j}(\theta_0; \mathcal{O})\}^2 + P\{l'_{\alpha_k}(\widehat{\theta}_n; \mathcal{O}) - l'_{\alpha_k}(\theta_0; \mathcal{O})\}^2 \\ &\quad + P\{l'_{\xi}(\widehat{\theta}_n; \mathcal{O})[h_j^* - h_{jn}^*]\}^2 + P\{l'_{\xi}(\widehat{\theta}_n; \mathcal{O})[h_k^* - h_{kn}^*]\}^2 \\ &\quad + P\{l'_{\xi}(\widehat{\theta}_n; \mathcal{O})[h_j^*] - l'_{\xi}(\theta_0; \mathcal{O})[h_j^*]\}^2 + P\{l'_{\xi}(\widehat{\theta}_n; \mathcal{O})[h_k^*] - l'_{\xi}(\theta_0; \mathcal{O})[h_k^*]\}^2. \end{aligned}$$

It is easy to show that

$$\begin{aligned} P\{l'_{\alpha_j}(\widehat{\theta}_n; \mathcal{O}) - l'_{\alpha_j}(\theta_0; \mathcal{O})\}^2 &\lesssim \|\widehat{g}'_n - g'_0\|_2^2 + d^2(\widehat{\theta}_n, \theta_0), \quad j = 1, \dots, d_1, \\ P\{l'_{\alpha_j}(\widehat{\theta}_n; \mathcal{O}) - l'_{\alpha_j}(\theta_0; \mathcal{O})\}^2 &\lesssim d^2(\widehat{\theta}_n, \theta_0), \quad j = d_1 + 1, \dots, d, \\ P\{l'_{\xi}(\widehat{\theta}_n; \mathcal{O})[h_j^* - h_{jn}^*]\}^2 &\lesssim d^2(\widehat{\theta}_n, \theta_0) + |\widehat{\beta} - \beta_0|, \quad j = 1, \dots, d, \\ P\{l'_{\xi}(\widehat{\theta}_n; \mathcal{O})[h_j^*] - l'_{\xi}(\theta_0; \mathcal{O})[h_j^*]\}^2 &\lesssim d^2(\widehat{\theta}_n, \theta_0), \quad j = 1, \dots, d. \end{aligned}$$

Thus, it follows from Theorem 1 that  $P_n\{A_n^{jk}(\widehat{\theta}_n; \mathcal{O}) - A^{jk}(\theta_0; \mathcal{O})\} = o_p(1)$ , which completes the proof of the theorem.  $\square$

## Acknowledgments

The authors would like to thank the Editor, Professor Eric Moulines, the Associate Editor and the two referees for their insightful comments and constructive suggestions that greatly improved the paper. Zhao's research was partly supported by the Research Grant Council of Hong Kong (504011, 503513), The Hong Kong Polytechnic University, and the National Natural Science Foundation of China (11371299); Wu's research was partly supported by the National Natural Science Foundation of China (11201350); and Yin's research was partly supported by the Research Grant Council of Hong Kong (17125814).

## References

- [1] Buckley, J. and James, I. (1979). Linear regression with censored data. *Biometrika* **66** 429–436.
- [2] Chen, X. (2007). Large sample sieve estimation of semi-nonparametric models. In *Handbook of Econometrics, Part B* (J.J. Heckman and E.E. Leamer, eds.) **6** 5549–5632. North-Holland: Elsevier.
- [3] Chen, Y.Q. and Jewell, N.P. (2001). On a general class of semiparametric hazards regression models. *Biometrika* **88** 687–702. [MR1859402](#)
- [4] Chen, Y.Q. and Wang, M.-C. (2000). Analysis of accelerated hazards models. *J. Amer. Statist. Assoc.* **95** 608–618. [MR1803176](#)
- [5] Copelan, E.A., Biggs, J.C., Thompson, J.M., Crilley, P., Szer, J., Klein, J.P., Kapoor, N., Avalos, B.R., Cunningham, I., Atkinson, K., Downs, K., Harmon, G.S., Daly, M.B., Brodsky, I., Bulova, S.I. and Tutschka, P.J. (1991). Treatment for acute myelocytic leukemia with allogeneic bone marrow transplantation following preparation with Bu/Cy. *Blood* **78** 838–843.
- [6] Cox, D.R. (1972). Regression models and life-tables. *J. R. Stat. Soc. Ser. B. Stat. Methodol.* **34** 187–220. With discussion by F. Downton, Richard Peto, D. J. Bartholomew, D. V. Lindley, P. W. Glassborow, D. E. Barton, Susannah Howard, B. Benjamin, John J. Gart, L. D. Meshalkin, A. R. Kagan, M. Zelen, R. E. Barlow, Jack Kalbfleisch, R. L. Prentice and Norman Breslow, and a reply by D. R. Cox. [MR0341758](#)
- [7] Ding, Y. and Nan, B. (2011). A sieve  $M$ -theorem for bundled parameters in semiparametric models, with application to the efficient estimation in a linear model for censored data. *Ann. Statist.* **39** 3032–3061. [MR3012400](#)

- [8] Huang, J. (1999). Efficient estimation of the partly linear additive Cox model. *Ann. Statist.* **27** 1536–1563. [MR1742499](#)
- [9] Huang, J. and Rossini, A.J. (1997). Sieve estimation for the proportional-odds failure-time regression model with interval censoring. *J. Amer. Statist. Assoc.* **92** 960–967. [MR1482126](#)
- [10] Jin, Z., Lin, D.Y., Wei, L.J. and Ying, Z. (2003). Rank-based inference for the accelerated failure time model. *Biometrika* **90** 341–353. [MR1986651](#)
- [11] Klein, J.P. and Moeschberger, M.L. (2003). *Survival Analysis: Techniques for Censored and Truncated Data*, 2nd ed. New York: Springer.
- [12] Kosorok, M.R. (2008). *Introduction to Empirical Processes and Semiparametric Inference*. Springer Series in Statistics. New York: Springer. [MR2724368](#)
- [13] Lai, T.L. and Ying, Z. (1991). Rank regression methods for left-truncated and right-censored data. *Ann. Statist.* **19** 531–556. [MR1105835](#)
- [14] Prentice, R.L. (1978). Linear rank tests with right censored data. *Biometrika* **65** 167–179. [MR0497517](#)
- [15] Ritov, Y. (1990). Estimation in a linear regression model with censored data. *Ann. Statist.* **18** 303–328. [MR1041395](#)
- [16] Rossini, A.J. and Tsiatis, A.A. (1996). A semiparametric proportional odds regression model for the analysis of current status data. *J. Amer. Statist. Assoc.* **91** 713–721. [MR1395738](#)
- [17] Schumaker, L.L. (1981). *Spline Functions: Basic Theory*. Pure and Applied Mathematics. New York: Wiley. A Wiley-Interscience Publication. [MR0606200](#)
- [18] Shen, X. (1998). Proportional odds regression and sieve maximum likelihood estimation. *Biometrika* **85** 165–177. [MR1627289](#)
- [19] Shen, X. and Wong, W.H. (1994). Convergence rate of sieve estimates. *Ann. Statist.* **22** 580–615. [MR1292531](#)
- [20] Tsiatis, A.A. (1990). Estimating regression parameters using linear rank tests for censored data. *Ann. Statist.* **18** 354–372. [MR1041397](#)
- [21] van der Vaart, A.W. and Wellner, J.A. (1996). *Weak Convergence and Empirical Processes: With Applications to Statistics*. Springer Series in Statistics. New York: Springer. [MR1385671](#)
- [22] Wei, L.J., Ying, Z. and Lin, D.Y. (1990). Linear regression analysis of censored survival data based on rank tests. *Biometrika* **77** 845–851. [MR1086694](#)
- [23] Wellner, J.A. and Zhang, Y. (2007). Two likelihood-based semiparametric estimation methods for panel count data with covariates. *Ann. Statist.* **35** 2106–2142. [MR2363965](#)
- [24] Zeng, D. and Lin, D.Y. (2007). Efficient estimation for the accelerated failure time model. *J. Amer. Statist. Assoc.* **102** 1387–1396. [MR2412556](#)
- [25] Zeng, D., Yin, G. and Ibrahim, J.G. (2005). Inference for a class of transformed hazards models. *J. Amer. Statist. Assoc.* **100** 1000–1008. [MR2201026](#)
- [26] Zhang, H.H., Cheng, G. and Liu, Y. (2011). Linear or nonlinear? Automatic structure discovery for partially linear models. *J. Amer. Statist. Assoc.* **106** 1099–1112. [MR2894767](#)

Received April 2014 and revised March 2016