*Research Article*

# An Empirical Likelihood Method for Semiparametric Linear Regression with Right Censored Data

## Kai-Tai Fang,[1] Gang Li,[2] Xuyang Lu,[2] and Hong Qin[3]

[1] *Beijing Normal University-Hong Kong Baptist University, United International College, Zhuhai 519085, China*
[2] *Department of Biostatistics, School of Public Health, University of California, Los Angeles, CA 90095-1772, USA*
[3] *Department of Statistics, Central China Normal University, Wuhan 430079, China*

Correspondence should be addressed to Gang Li; vli@ucla.edu

This paper develops a new empirical likelihood method for semiparametric linear regression with a completely unknown error distribution and right censored survival data. The method is based on the Buckley-James (1979) estimating equation. It inherits some appealing properties of the complete data empirical likelihood method. For example, it does not require variance estimation which is problematic for the Buckley-James estimator. We also extend our method to incorporate auxiliary information. We compare our method with the synthetic data empirical likelihood of Li and Wang (2003) using simulations. We also illustrate our method using Stanford heart transplantation data.

## 1. Introduction

Suppose that one observes right censored regression data consisting of $n$ i.i.d. triples $(X_i, Z_i, \delta_i) = (X_i, Y_i \wedge C_i, I[Y_i \leq C_i])$, $i = 1, \ldots, n$, where for subject $i$, $Y_i$ is a known monotone transformation of the survival time of interest, $C_i$ is the corresponding censoring time, and $X_i = (X_{i1}, \ldots, X_{ip})^\tau$ is a vector of $p$ covariates. We consider the problem of making inferences for the slope parameter of the semiparametric linear regression model as follows:

$$Y_i = X_i^\tau \beta + \epsilon_i, \quad i = 1, \ldots, n, \tag{1}$$

where $\beta = (\beta_1, \ldots, \beta_p)^\tau$ is a $p \times 1$ vector of unknown regression coefficients and $\epsilon_i$'s are independent and identically distributed random errors with an unknown distribution $F$. Assume that $\epsilon_i$ is independent of $(X_i, C_i)$ and that conditional on $X_i$, $Y_i$ and $C_i$ are independent, $i = 1, \ldots, n$. Because the error distribution $F$ is completely unknown, we assume no intercept term in model (1) in order for $F$ to be identifiable. Assume further that the covariance matrix of $X_i$ is positive definite. Model (1) provides a useful alternative to the popular Cox [1] model when the proportional hazards assumption does not hold. Furthermore, it makes no parametric assumption on the error distribution and is, thus, more flexible than

parametric accelerated failure time models (cf. Andersen et al. [2]).

Buckley and James [3] extended the least squares method to estimate the regression coefficient $\beta$ in model (1) with right censored data. The Buckley-James estimator can be calculated using an iterative algorithm. Its consistency and asymptotic normality have been established by Lai and Ying [4], Ritov [5], Tsiatis [6], Ying [7], and others. However, its asymptotic variance involves the unknown hazard function of $\epsilon$, and its derivatives whose nonparametric estimations are problematic (cf. Ritov [5]). To overcome this problem, several approaches have been studied in the literature. Jin et al. [3, 8], Wei et al. [9], and Lin and Geyer [10] considered rank-based inferences. Koul et al. [11] introduced a synthetic data approach which was further investigated by Zhou [12], Lai et al. [13], and others. Let $Y_{iG} = \delta_i Y_i / G(Y_i-)$, where $G$ is the survival function of the censoring time. It can be shown that $E(Y_{iG} \mid X_i) = E(Y_i \mid X_i)$ if $C$ is independent of $X$ and $Y$. Koul et al. [11] proposed to estimate the regression parameters by regressing $Y_{i\widehat{G}}$ on $X_i$, where $\widehat{G}$ is the Kaplan-Meier [14] estimate of $G$ and $Y_{i\widehat{G}}$ is referred to as a synthetic variable. Koul et al. [11] showed that the synthetic data estimator is asymptotically normal and its variance is simple to estimate.

However, it requires a strong assumption that the censoring time is independent of both the survival time and the covariates. It can also have poor small sample performance in the presence of heavy censoring (cf. Li and Wang [15]). Jing and Qin [16] and Li and Wang [15] independently developed empirical likelihood-based inference for $\beta$ using synthetic data. Their methods provide substantial improvement over the normal theory method of Koul et al. [11], but are still not very satisfactory for small samples (cf. Li and Wang [15]). The synthetic data approach makes an independence assumption between $C$ and $(Y, X)$. Li and Lu [17] showed that it can yield seriously biased parameter estimate if $C$ depends on $X$. Zhou and Li [18] developed a censored data EL method under a weaker censorship assumption that $Y$ and $C$ are conditionally independent given $X$. Their method also demonstrated better small sample performance than the synthetic data EL method when the errors $\epsilon_i$'s are independent and identically distributed (homogeneous). On the other hand, it is not easy to extend the method of Zhou and Li [18] to incorporate auxiliary information and to construct confidence regions for linear combinations of the regression coefficients, which are relatively easy using the synthetic data EL method (cf. Li and Wang [15]).

In this paper, we develop a new empirical likelihood method for linear regression with right censored data. We construct an estimated empirical likelihood based on the estimation equation for the Buckley-James [3] estimator. The approach inherits many appealing features of empirical likelihood. For example, the shape and orientation of a confidence region are entirely determined by data. Most importantly, it does not involve variance estimation. Our simulation shows that, under the assumptions of model (1), our method is far superior to the synthetic data empirical likelihood method of Li and Wang [15], with substantially shorter confidence intervals and higher coverage probabilities. Compared to Zhou and Li [18]'s method, our method may not be as efficient, but it is easier to compute. It can also be extended easily to incorporate auxiliary information and to construct confidence regions for linear combinations of the regression coefficients. More discussion of the pros and cons of our method in relation to existing methods are given later in Remark 2.

The use of empirical likelihood dates back at least to Thomas and Grunkemeier [19] who constructed confidence intervals for survival probabilities. The idea was later popularized by Owen [20, 21] who derived confidence regions for the mean of a random vector. Since the work of Owen [20, 21], there has been extensive developments of empirical likelihood methods for a wide range of applications including, among others, linear models [22–24], generalized linear models [25], quantile estimation [26], biased sample models [27], generalized estimating equations [28], dependent process model [29], partial linear models [30], mixture proportions [31], confidence bands for survival functions [32], confidence bands for quantile functions [33], censored cost regression [34], and confidence tubes for multiple quantile plots [35]. Some nice discussion of properties of empirical likelihood can be found in DiCiccio, Hall and Romano [36], Hall [37], and Hall and Scala [38], and others.

A comprehensive survey of empirical likelihood and further references can be found in Owen [39] and Li et al. [40].

In Section 2, we derive an estimated empirical likelihood for $\beta$ by combining the ideas of Owen [21, 22] and Buckley and James [3]. An adjustment factor is adopted so that the adjusted empirical likelihood has an asymptotic standard Chi-square distribution. We also discuss how to incorporate auxiliary information using empirical likelihood. Section 3 presents results from a simulation study to illustrate the performance of our method compared with the synthetic data empirical likelihood method. A real data example is also provided. Section 4 gives some concluding remarks. The proofs are collected in the appendix.

## 2. Empirical Likelihood Inference

*2.1. Empirical Likelihood Based on the Buckley-James Estimating Equation.* We motivate our procedure by first considering the case where $\mu_X = E(X)$ and $F$ are known.

We first give a review of the Buckley-James equation. It can be shown that, under model (1),

$$\beta = \left[E\left\{(X_i - \mu_X)X_i\right\}\right]^{-1} E\left\{(X_i - \mu_X)Y_i\right\}, \qquad (2)$$

or, equivalently,

$$E\left\{(X_i - \mu_X)(Y_i - X_i^\tau\beta)\right\} = 0. \qquad (3)$$

Because $Y_i$ is not always observable, we impute $Y_i$ by its conditional expectation given the observed data as follows:

$$
\begin{aligned}
Y_i^* &= E\left(Y_i \mid X_i, Z_i, \delta_i\right) \\
&= \delta_i Z_i + (1 - \delta_i) \\
&\quad \times \left\{X_i^\tau\beta + \frac{\int_{Z_i - X_i^\tau\beta}^{\infty} t\, dF(t)}{1 - F(Z_i - X_i^\tau\beta)}\right\}.
\end{aligned}
\qquad (4)
$$

Noting that $E(Y_i^* \mid X_i) = E(Y_i \mid X_i)$, we can replace $Y_i$ by $Y_i^*$ in (3). This leads to

$$E\left\{W_i(\beta, F, \mu_X)\right\} = 0, \qquad (5)$$

where

$$
\begin{aligned}
W_i(\beta, F, \mu_X) &= (X_i - \mu_X) \\
&\quad \times \left\{ \delta_i(Z_i - X_i^\tau\beta) \right. \\
&\qquad \left. + (1 - \delta_i)\frac{\int_{Z_i - X_i^\tau\beta}^{\infty} t\, dF(t)}{1 - F(Z_i - X_i^\tau\beta)}\right\}.
\end{aligned}
\qquad (6)
$$

Now, the problem of testing $H_0 : \beta = \beta_0$ is equivalent to testing

$$H_0 : E\left\{W_i(\beta_0, F, \mu_X)\right\} = 0, \qquad (7)$$

based on $n$ i.i.d. observations $W_i(\beta_0, F, \mu_X)$, $i = 1, \ldots, n$. This problem can be readily solved using the results of Owen [21]. Specifically, define

$$
\begin{aligned}
&l_n(\beta, F, \mu_X) \\
&= -2 \sup \left\{ \sum_{i=1}^{n} \log(np_i) \mid \sum_{i=1}^{n} p_i W_i(\beta, F, \mu_X) = 0, \right. \\
&\qquad\qquad \left. \sum_{i=1}^{n} p_i = 1, \ p_i \geq 0, \ 1 \leq i \leq n \right\}.
\end{aligned}
\tag{8}
$$

Owen [21] showed that

$$
l_n(\beta, F, \mu_X) = 2 \sum_{i=1}^{n} \log\left\{ 1 + \lambda^\tau W_i(\beta, F, \mu_X) \right\},
\tag{9}
$$

where $\lambda$ is the solution of the equation

$$
\frac{1}{n} \sum_{i=1}^{n} \frac{W_i(\beta, F, \mu_X)}{1 + \lambda^\tau W_i(\beta, F, \mu_X)} = 0.
\tag{10}
$$

Moreover, under $H_0$, $l_n(\beta_0, F, \mu_X)$ has an asymptotic central Chi-square distribution with $p$ degrees of freedom. Thus, one would reject $H_0$ if $l_n(\beta_0, F, \mu_X) > \chi^2_{p,\alpha}$, where $\chi^2_{p,\alpha}$ is the upper $\alpha$ quantile of the standard Chi-square distribution with $p$ degrees of freedom.

Now, we consider the case where both $F$ and $\mu_X$ are unknown. Define $e_i^\beta = Z_i - X_i^\tau \beta$, $i = 1, \ldots, n$. We estimate $F$ by

$$
\widehat{F}_n^\beta(t) = 1 - \prod_{i=1}^{n} \left[ \frac{n-i}{n-i+1} \right]^{I[e_{(i)}^\beta \leq t, \delta_{(i)} = 1]},
\tag{11}
$$

the Kaplan-Meier [14] estimator of $F$ based on $\{(e_i^\beta, \delta_i), i = 1, \ldots, n\}$, where $e_{(1)}^\beta \leq \cdots \leq e_{(n)}^\beta$ are the order statistics of the $e^\beta$-sample, and $\delta_{(i)}$ is the $\delta$ associated with $e_{(i)}^\beta$, $i = 1, \ldots, n$. In addition, we estimate $\mu_X$ by the sample mean $\overline{X} = n^{-1} \sum_{i=1}^{n} X_i$.

We propose to use $l_n(\beta_0, \widehat{F}_n^{\beta_0}, \overline{X})$ as a likelihood ratio statistic for testing $H_0$. However, we can no longer use Owen's [21] result for the null limiting distribution of $l_n(\beta_0, \widehat{F}_n^{\beta_0}, \overline{X})$ because $W_i(\beta, \widehat{F}_n^\beta, \overline{X})$'s are not i.i.d. The following theorem states that an adjustment factor is needed so that the limiting null distribution is standard Chi-squared.

**Theorem 1.** *Assume that the conditions listed in the appendix hold. Define that*

$$
c_n(\beta) = \frac{\operatorname{tr}\left( \widehat{\Sigma}_2^{-1}(\beta) S_n(\beta) \right)}{\operatorname{tr}\left( \widehat{\Sigma}_1^{-1}(\beta) S_n(\beta) \right)},
\tag{12}
$$

*where*

$$
S_n(\beta) = \left\{ \sum_{i=1}^{n} W_i\left(\beta, \widehat{F}_n^\beta, \overline{X}\right) \right\} \left\{ \sum_{i=1}^{n} W_i\left(\beta, \widehat{F}_n^\beta, \overline{X}\right) \right\}^\tau,
$$

$$
\widehat{\Sigma}_1(\beta) = \frac{1}{n} \sum_{i=1}^{n} W_i\left(\beta, \widehat{F}_n^\beta, \overline{X}\right) W_i\left(\beta, \widehat{F}_n^\beta, \overline{X}\right)^\tau,
$$

$$
\widehat{\Sigma}_2(\beta) = \frac{1}{n} \sum_{i=1}^{n} \int_0^\infty \left( u - \frac{\int_u^\infty v \, dF_n^\beta(v)}{1 - F_n^\beta(u)} \right)^2
$$
$$
\times \widehat{V}(u) \, dN_i(u),
\tag{13}
$$

$$
\widehat{V}(u) = \frac{\sum_{i=1}^{n} \left\{ X_i - \overline{X}(u) \right\} \left\{ X_i - \overline{X}(u) \right\}^\tau Y_i(u)}{\sum_{j=1}^{n} Y_j(u)},
$$

$$
\overline{X}(u) = \frac{\sum_{i=1}^{n} X_i Y_i(u)}{\sum_{j=1}^{n} Y_j(u)},
$$

$$
N_i(u) = I\left( e_i^\beta \leq u, \ \delta_i = 1 \right),
$$

$$
Y_i(u) = I\left( e_i^\beta \geq u \right).
$$

*Then, under $H_0 : \beta = \beta_0$,*

$$
c_n(\beta_0) l_n\left(\beta_0, \widehat{F}_n^{\beta_0}, \overline{X}\right) \xrightarrow{d} \chi^2_p,
\tag{14}
$$

*where $\chi^2_p$ is a standard Chi-square random variable with $p$ degrees of freedom.*

It follows immediately that an approximate $\alpha$-level test rejects $H_0$ if

$$
c_n(\beta_0) l_n\left(\beta_0, \widehat{F}_n^{\beta_0}, \overline{X}\right) > \chi^2_{p,\alpha}.
\tag{15}
$$

Moreover, an approximate $1 - \alpha$ confidence region for $\beta$ is given by

$$
\left\{ \beta : c_n(\beta) l_n\left(\beta, \widehat{F}_n^\beta, \overline{X}\right)(\beta) \leq \chi^2_{p,\alpha} \right\}.
\tag{16}
$$

*Remark 2.* Although both the above derived method and Zhou and Li [18] use the Buckley-James estimation equation, a sample version of (7), they are different in that we use the complete data likelihood, whereas Zhou and Li [18] uses the exact censored data likelihood to construct the EL. Similar to Li and Wang [15], the above method can be extended to incorporate auxiliary information and to obtain EL procedure for a subset, contrast, or linear combinations of the regression coefficients, which does not seem easy when using the method of Zhou and Li [18]. As an illustration, we show below how to extend our method to incorporate auxiliary information.

*2.2. Empirical Likelihood with Auxiliary Information.* Auxiliary population characteristics of the covariate $X$ are sometimes available in practice. Effective usage of the auxiliary

information can lead to more efficient inference (cf. Chen and Qin [41], Qin and Lawless [28] and Zhang [42, 43]). Here, we show how to use empirical likelihood to incorporate auxiliary information of $X$.

Assume that the available auxiliary information on $X$ is given in the form $Eg(X) = 0$, where $g(x) = (g_1(x), \ldots, g_r(x))^\tau, r \geq 1$, is a vector of $r$ known functions. To make use of the auxiliary information, we maximize

$$\prod_{i=1}^{n} p_i \tag{17}$$

subject to $\sum_{i=1}^{n} p_i = 1$, $\sum_{i=1}^{n} p_i g(X_i) = 0$, and $\sum_{i=1}^{n} p_i W_i(\beta, \widehat{F}_n^\beta, \overline{X}) = 0$.

Let $A_{ni}(\beta) = (g^\tau(X_i), W_i(\beta, \widehat{F}_n^\beta, \overline{X})^\tau)^\tau$. By the method of Lagrange multipliers, it can be shown that (17) is maximized at

$$p_{in} = \frac{1}{n} \frac{A_{ni}(\beta)}{1 + \zeta_n^\tau A_{ni}(\beta)}, \quad i = 1, \ldots, n, \tag{18}$$

where $\zeta_n$ satisfies the following equation

$$\frac{1}{n} \sum_{i=1}^{n} \frac{A_{ni}(\beta)}{1 + \zeta_n^\tau A_{ni}(\beta)} = 0. \tag{19}$$

Hence, the empirical log-likelihood ratio function for $\beta$ is given by

$$l_{n,AU}(\beta) = -2 \sum_{i=1}^{n} \log n p_{in} = 2 \sum_{i=1}^{n} \log \left(1 + \zeta_n^\tau A_{ni}(\beta)\right). \tag{20}$$

Similar to the previous section, an adjustment factor is needed for $l_{n,AU}(\beta)$ to have a standard Chi-square asymptotic distribution, as stated in the following theorem.

**Theorem 3.** *Assume that $V_1 = Eg(X_i)g^\tau(X_i)$ is positive definite and that $Eg(X)W_i(\beta, F, \mu_X)^\tau)$ exists. Define that $V_{n1}(\beta) = (1/n) \sum_{i=1}^{n} g(X_i)g^\tau(X_i)$,*

$$V_{n2}(\beta) = \frac{1}{n} \sum_{i=1}^{n} g(X_i) W_i\left(\beta, \widehat{F}_n^\beta, \overline{X}\right),$$

$$V_{n1,AU}(\beta) = \begin{pmatrix} V_{n1}(\beta), & V_{n2}(\beta) \\ V_{n2}^\tau(\beta), & \widehat{\Sigma}_1(\beta) \end{pmatrix}, \tag{21}$$

$$V_{n2,AU}(\beta) = \begin{pmatrix} V_{n1}(\beta), & 0 \\ 0 & \widehat{\Sigma}_2(\beta) \end{pmatrix},$$

*where $\widehat{\Sigma}_1(\beta)$ and $\widehat{\Sigma}_2(\beta)$ are defined in Theorem 1. Then, under the conditions of Theorem 1,*

$$c_{n,AU}(\beta_0) l_{n,AU}(\beta_0) \xrightarrow{d} \chi_{p+r}^2, \tag{22}$$

*as $n \to \infty$, where*

$$c_{n,AU}(\beta) = \frac{\text{tr}\left(V_{n2,AU}^{-1}(\beta) \Psi_n(\beta)\right)}{\text{tr}\left(V_{n1,AU}^{-1}(\beta) \Psi_n(\beta)\right)} \tag{23}$$

*and $\Psi_n(\beta) = (\sum_{i=1}^{n} A_{ni}(\beta))(\sum_{i=1}^{n} A_{ni}(\beta))^\tau$.*

## 3. Numerical Results

We carried out Monte Carlo simulations to examine the performance of our proposed empirical likelihood method based on the Buckley-James estimating equation (ELEE) in comparison to the empirical likelihood method based on synthetic data (ELSD) [15, 16]. We considered five models. In model A, the data were generated from $Y = 1 + X + \epsilon$, where $X$ and $\epsilon$ are independent normal random variables with mean 0 and variances 0.25, respectively, the censoring time $C$ is a normal random variable with mean $\mu$ and standard deviation 4. Model B is the same as model A except that $X \sim \text{Bernoulli}(0.5) - 0.5$. Model C assumes that $Y = X + \epsilon$, where $X \sim N(0, 0.5^2)$, $\epsilon \sim$ Weibull (shape = 1.843, scale = 1), and $C \sim N(\mu, 4^2)$. Model D is the same as model A except that $C \sim N(\mu + 2X, 15)$, allowing the censoring time to depend on $X$. Model E is the same as model A except that $\epsilon \sim N(0, X^2)$, allowing for heterogeneous errors. We adjust $\mu$ to produce different censoring rate (CR). We also vary the sample size $n$. The achieved confidence levels and average lengths of the ELEE and ELSD confidence intervals for the slope parameter are summarized in Table 1. Each entry in the table was computed using 3,000 Monte Carlo samples.

We see from Table 1 that under models A–D, the coverage probabilities of ELEE method are consistently close to or slightly above the nominal level, whereas the ELSD method can have severe under-coverage for small samples ($n = 50$) and large censoring rate (75%). Furthermore, the ELEE confidence intervals are much narrower than the ELSD confidence intervals. In particular, the ELSD method failed completely with unreasonably low coverage under model D when the censoring time is dependent on $X$. On the other hand, under model E with heterogeneous errors and independent censoring time, ELEE showed larger coverage probability errors than ELSD, as one would have expected. Thus, the ELEE method seems to dominate the ELSD method when the errors are homogeneous, but can be outperformed by ELSD in the presence of heterogeneous errors.

We now illustrate our method using the Stanford heart transplant data (Miller [44], Table 1). The data include the lengths of survival (in days) after transplantation, ages at time of transplant, and T5 mismatch scores for 69 patients who received heart transplants at Stanford and were followed to April 1, 1974. Twenty-four patients were still alive on April 1, 1974 and thus their survival times were censored. For illustration purpose, we considered two models, labeled as (I) and (II), respectively, where the dependent variable $Y$ is the logarithm to base 10 of the length of survival from transplantation. Specifically, model (I) regresses $Y$ on the mismatch score T5, and model (II) regresses $Y$ on age. As in Koul et al. [11], regression of survival on the mismatch score T5 was performed with nonrejection-related death being treated as censoring since the mismatch score is directed at the rejection phenomenon [44]. Table 2 reports the parameter estimates and 95% confidence intervals for the slope parameters using our empirical likelihood method based on the Buckley-James estimating equation (ELEE) and the empirical likelihood method based on synthetic data (ELSD) [15, 16].

Table 1: Comparison of the coverage probability (CP) and average width (Width) of two empirical likelihood confidence intervals for the slope parameter under four different models with various sample size ($n$) and censoring rate (CR). Here, ELEE is the proposed method, and ELSD is the method of Li and Wang [15]. Each entry is based on 3,000 Monte Carol samples.

| Model | $n$ | CR | Nominal level = 90% | | | | Nominal level = 95% | | | |
| | | | CP | | width | | CP | | Width | |
| | | | ELEE | ELSD | ELEE | ELSD | ELEE | ELSD | ELEE | ELSD |
| A | 50 | 0.75 | 0.94 | 0.77 | 1.59 | 2.18 | 0.98 | 0.84 | 1.95 | 2.70 |
| | 100 | 0.75 | 0.94 | 0.82 | 0.85 | 1.66 | 0.97 | 0.89 | 1.03 | 2.02 |
| | 500 | 0.75 | 0.91 | 0.88 | 0.31 | 0.81 | 0.96 | 0.94 | 0.37 | 0.97 |
| | 50 | 0.3 | 0.94 | 0.87 | 0.69 | 1.30 | 0.97 | 0.92 | 0.82 | 1.55 |
| | 100 | 0.3 | 0.93 | 0.89 | 0.45 | 0.94 | 0.97 | 0.94 | 0.53 | 1.12 |
| | 500 | 0.3 | 0.91 | 0.90 | 0.18 | 0.43 | 0.96 | 0.95 | 0.21 | 0.51 |
| | 50 | 0.1 | 0.95 | 0.88 | 0.59 | 1.10 | 0.98 | 0.93 | 0.70 | 1.30 |
| | 100 | 0.1 | 0.93 | 0.89 | 0.39 | 0.79 | 0.97 | 0.94 | 0.47 | 0.94 |
| | 500 | 0.1 | 0.90 | 0.90 | 0.16 | 0.36 | 0.95 | 0.95 | 0.19 | 0.43 |
| B | 50 | 0.75 | 0.93 | 0.83 | 1.20 | 2.05 | 0.95 | 0.88 | 1.40 | 2.50 |
| | 100 | 0.75 | 0.94 | 0.87 | 0.77 | 1.49 | 0.97 | 0.92 | 0.92 | 1.80 |
| | 500 | 0.75 | 0.93 | 0.89 | 0.30 | 0.67 | 0.96 | 0.94 | 0.36 | 0.80 |
| | 50 | 0.3 | 0.95 | 0.88 | 0.66 | 1.23 | 0.98 | 0.94 | 0.78 | 1.48 |
| | 100 | 0.3 | 0.94 | 0.90 | 0.44 | 0.88 | 0.97 | 0.95 | 0.52 | 1.05 |
| | 500 | 0.3 | 0.92 | 0.90 | 0.18 | 0.39 | 0.96 | 0.95 | 0.21 | 0.47 |
| | 50 | 0.1 | 0.94 | 0.89 | 0.58 | 1.08 | 0.97 | 0.95 | 0.69 | 1.29 |
| | 100 | 0.1 | 0.94 | 0.90 | 0.39 | 0.77 | 0.97 | 0.94 | 0.46 | 0.92 |
| | 500 | 0.1 | 0.91 | 0.91 | 0.16 | 0.35 | 0.96 | 0.95 | 0.19 | 0.41 |
| C | 50 | 0.75 | 0.93 | 0.77 | 1.49 | 2.01 | 0.96 | 0.83 | 2.01 | 2.46 |
| | 100 | 0.75 | 0.93 | 0.82 | 0.80 | 1.56 | 0.97 | 0.88 | 0.97 | 1.90 |
| | 500 | 0.75 | 0.92 | 0.87 | 0.29 | 0.76 | 0.96 | 0.93 | 0.35 | 0.92 |
| | 50 | 0.3 | 0.93 | 0.86 | 0.67 | 1.21 | 0.97 | 0.92 | 0.81 | 1.45 |
| | 100 | 0.3 | 0.93 | 0.88 | 0.44 | 0.88 | 0.97 | 0.93 | 0.53 | 1.05 |
| | 500 | 0.3 | 0.91 | 0.89 | 0.18 | 0.40 | 0.96 | 0.94 | 0.21 | 0.48 |
| | 50 | 0.1 | 0.94 | 0.87 | 0.60 | 1.01 | 0.97 | 0.93 | 0.71 | 1.21 |
| | 100 | 0.1 | 0.93 | 0.89 | 0.39 | 0.73 | 0.97 | 0.94 | 0.47 | 0.87 |
| | 500 | 0.1 | 0.92 | 0.90 | 0.16 | 0.33 | 0.96 | 0.95 | 0.19 | 0.39 |
| D | 50 | 0.75 | 0.95 | 0.68 | 1.74 | 2.45 | 0.97 | 0.77 | 1.87 | 2.98 |
| | 100 | 0.75 | 0.94 | 0.60 | 0.83 | 1.83 | 0.97 | 0.69 | 1.01 | 2.21 |
| | 500 | 0.75 | 0.92 | 0.12 | 0.30 | 0.89 | 0.96 | 0.18 | 0.36 | 1.06 |
| | 50 | 0.3 | 0.94 | 0.81 | 0.68 | 1.31 | 0.97 | 0.88 | 0.82 | 1.56 |
| | 100 | 0.3 | 0.93 | 0.76 | 0.45 | 0.94 | 0.97 | 0.84 | 0.53 | 1.12 |
| | 500 | 0.3 | 0.91 | 0.39 | 0.18 | 0.43 | 0.96 | 0.51 | 0.21 | 0.51 |
| | 50 | 0.1 | 0.94 | 0.86 | 0.59 | 1.09 | 0.97 | 0.92 | 0.71 | 1.30 |
| | 100 | 0.1 | 0.94 | 0.87 | 0.39 | 0.78 | 0.97 | 0.93 | 0.47 | 0.93 |
| | 500 | 0.1 | 0.91 | 0.78 | 0.16 | 0.36 | 0.95 | 0.86 | 0.19 | 0.42 |
| E | 50 | 0.75 | 0.78 | 0.76 | 1.52 | 2.25 | 0.85 | 0.83 | 1.79 | 2.75 |
| | 100 | 0.75 | 0.78 | 0.80 | 0.85 | 1.74 | 0.85 | 0.87 | 0.62 | 2.11 |
| | 500 | 0.75 | 0.73 | 0.85 | 0.34 | 0.88 | 0.81 | 0.92 | 0.40 | 1.06 |
| | 50 | 0.3 | 0.81 | 0.85 | 0.76 | 1.43 | 0.87 | 0.91 | 0.89 | 1.71 |
| | 100 | 0.3 | 0.79 | 0.87 | 0.53 | 1.05 | 0.86 | 0.92 | 0.62 | 1.26 |
| | 500 | 0.3 | 0.77 | 0.89 | 0.22 | 0.50 | 0.84 | 0.94 | 0.26 | 0.60 |
| | 50 | 0.1 | 0.81 | 0.86 | 0.69 | 1.24 | 0.87 | 0.92 | 0.81 | 1.48 |
| | 100 | 0.1 | 0.80 | 0.87 | 0.49 | 0.91 | 0.86 | 0.93 | 0.57 | 1.08 |
| | 500 | 0.1 | 0.76 | 0.89 | 0.20 | 0.42 | 0.83 | 0.94 | 0.23 | 0.50 |

Model A: $Y = 1 + X + \epsilon$, where $X \sim N(0, 0.5^2)$, $\epsilon \sim N(0, 0.5^2)$, and $C \sim N(\mu, 4^2)$; model B: $Y = 1 + X + \epsilon$, where $X \sim Bernoulli(0.5) - 0.5$, $\epsilon \sim N(0, 0.5^2)$, and $C \sim N(\mu, 4^2)$; model C: $Y = X + \epsilon$, where $X \sim N(0, 0.5^2)$, $\epsilon \sim$ Weibull (shape = 1.843, scale = 1), and $C \sim N(\mu, 4^2)$; model D (Dependent censoring): $Y = 1 + X + \epsilon$, where $X \sim N(0, 0.5^2)$, $\epsilon \sim N(0, 0.5^2)$, and $C \sim N(\mu + 2X, 15)$; model E: $Y = 1 + X + \epsilon$, where $X \sim N(0, 0.5^2)$, $\epsilon \sim N(0, X^2)$, and $C \sim N(\mu, 4^2)$.

Table 2: Empirical likelihood confidence interval estimates for heart transplant data. Nominal level = 95%.

| Model | Parameter | Parameter estimates | | Confidence intervals | |
|---|---|---|---|---|---|
| | | BJ | KSV | ELEE | ELSD |
| (I) | $\beta_{T5}$ | $-0.593$ | $0.258$ | $(-2.740, 0.645)$ | $(-0.596, 0.928)$ |
| (II) | $\beta_{\text{age}}$ | $-0.028$ | $0.055$ | $(-0.065, 0.032)$ | $(0.001, 0.128)$ |

Note: BJ refers to the Buckley-James [3] estimate, and KSV refers to the synthetic data estimate of Koul et al. [11]. ELEE is the proposed empirical likelihood method based on the Buckley-James estimating equation, and ELSD is the empirical likelihood method based on synthetic data [15, 16].

It is seen from Table 2 that both the point estimates and confidence intervals are quite different between the ELEE and ELSD methods for this data. For example, the KSV slope estimates of T5 and age are positive, which seems to contradict to the common belief that the survival time tends to be negatively correlated with T5 (the mismatch score) and age at diagnosis. For model (II), the ELSD confidence interval does not include zero and thus concludes a significant age effect at the 5% significant level, whereas the ELEE method does not produce a significant result. The ELSD results could be misleading in this example since the censoring time appears to depend on age and T5 as pointed out by Leurgans [45] and Li and Lu [17]. We have also examined the Cox-Snell residual plots for a number of parametric accelerated failure time (AFT) models. We found that the log-normal AFT model fits the data fairly well which indicates that the semiparametric AFT model should also fit the data well.

## 4. Discussion

This research adds a new tool to the toolbox of empirical likelihood (EL) methods for linear regression with right censored data. The three EL methods, namely, the method of Li and Wang [15], the method of Zhou and Li [18], and the method of this paper should be regarded as complementary, as opposed to competing, methods for linear regression with right censored survival data. Each of the three approaches has its own merits and shortcomings. None dominates the others in every aspect. So it is important to understand the pros and cons of these methods. First of all, if the errors in the regression model are i.i.d (homogeneous) and the censoring time is conditionally independent of the survival time given the covariates, then the method of Li and Zhou [18] and the method of this paper are expected to be superior to the Li and Wang [15] method, as demonstrated by simulations in Table 1 and Li and Zhou [18]. This is because the former two methods use the Buckley-James estimating equations which take advantage of the homogeneous error assumption to implicitly impute a censored observation from the model using all residuals, whereas the latter uses synthetic data that does not utilize the homogeneous error assumption. Furthermore, based on our limited experience, the Li and Zhou [18] seems to be more efficient than the method of this paper. This is not a surprise since, Li and Zhou method [18] uses the exact censored likelihood, whereas the ELEE method of this paper uses an approximate likelihood. On the other hand, Li and Zhou [18] only developed an EL procedure for the whole vector of regression coefficients and did not

discuss how to extend their method to incorporate auxiliary information which can be easily done using the approach of this paper as described in Section 2.2. Secondly, if the errors are heterogeneous, then the Li and Zhou [18] method and the method of this paper are expected to fail as indicated by Table 1 (model E), but the synthetic data method of Li and Wang [15] would still be asymptotically valid under the stronger censorship assumption that the censoring time is independent of both the survival time and the covariates. Finally, the method of Li and Wang [15] can fail completely when the censoring time is dependent on the covariate as shown in Table 1 (model D).

Some further extensions are warranted in future research. For example, the method of this paper can be extended to draw inference for linear combinations of the regression coefficients along the lines of Li and Wang ([15, Section 3]). It can also be extended to construct EL confidence regions based on estimating equations for a class of M-estimators described by Ritov [5]. Finally, this paper focuses only on EL methods for right censored data. It would be interesting to investigate how these EL methods compare to other methods such as the rank-based regression methods in future studies.

## Appendix

*Assumption A.1.* The covariate $X_i$ has compact support.

*Assumption A.2.* There exists a function $V(u)$ such that $\widehat{V}(u)$ converges uniformly to $V(u)$ in probability as $n \to \infty$.

The following lemmas are needed to prove Theorem 1.

**Lemma A.3.** *Let $\lambda(t)$ denote the hazard function of $F$. Then,*

$$n^{-1/2} \sum_{i=1}^{n} W_i \left( \beta, \widehat{F}_n^{\beta}, \overline{X} \right) \xrightarrow{\mathscr{L}} N \left( 0, \Sigma_2 \left( \beta \right) \right), \quad (A.1)$$

*where $\Sigma_2(\beta) = \int_0^\infty \{w(t)\}^2 V(t)\lambda(t)dt$ and*

$$w(t; F) = t - \frac{\int_t^\infty u dF(u)}{1 - F(t)}. \quad (A.2)$$

*Moreover, $\Sigma_2(\beta)$ can be consistently estimated by $\widehat{\Sigma}_2(\beta)$.*

*Proof.* By Proposition 4.1 of Ritov [5]

$$
\begin{aligned}
& n^{-1/2} \sum_{i=1}^{n} W_i \left( \beta, \widehat{F}_n^{\beta}, \overline{X} \right) \\
& = n^{-1/2} \sum_{i=1}^{n} \int w\left(t\right) \left\{ X_i - \overline{X}\left(t\right) \right\} dM_i\left(t\right) + o_p\left(1\right),
\end{aligned}
\tag{A.3}
$$

where $M_i(t) = N_i(t) - \int_0^t Y_i(u)\lambda(u)du$, $i = 1,\ldots,n$, are orthogonal locally square integrable martingales with respect to the filtration $\mathscr{F}_n(t) = \sigma[I(e_i^{\beta} \leq t), \delta_i I(e_i^{\beta} \leq t), X_i; i = 1,\ldots,n]$. This, together with Rebolledo's martingale central limit theorem (cf. Andersen et al. [2]), proves (A.1). □

**Lemma A.4.** *Under the conditions of Theorem 1,*

(a) $\max_{1 \leq i \leq n} \|W_i(\beta, \widehat{F}_n^{\beta}, \overline{X})\| = o_p(n^{1/2})$,

(b) $\lambda = O_p(n^{-1/2})$.

*Proof.* (a) By Assumption A.1, $\|X_i - \overline{X}\| \leq K$ for some constant $K$. Thus,

$$
\begin{aligned}
& \max_{1 \leq i \leq n} \left\| W_i \left( \beta, \widehat{F}_n^{\beta}, \overline{X} \right) \right\| \\
& \leq K \max_{1 \leq i \leq n} \left| \delta_i e_i^{\beta} \right| \\
& \quad + K \max_{1 \leq i \leq n} \left| (1 - \delta_i) \left\{ e_i^{\beta} - w\left(e_i^{\beta}; F\right) \right\} \right| \\
& \quad + K \max_{1 \leq i \leq n} \left| 1 - \delta_i \right| \left| w\left(e_i^{\beta}; F_n^{\beta}\right) - w\left(e_i^{\beta}; F\right) \right| \\
& \leq o\left(n^{1/2}\right) + o\left(n^{1/2}\right) + K \sup_t \left| w\left(t; F_n^{\beta}\right) - w\left(t; F\right) \right| \\
& = o\left(n^{1/2}\right) + O_p\left(n^{-1/2}\right) = o\left(n^{1/2}\right),
\end{aligned}
\tag{A.4}
$$

where we have used Lemma 3 of Owen [21] in the second step and Lemma 4.1 of Ritov [5] in the third step.

(b) Define that $\widetilde{\Sigma}_1(\beta) = (1/n) \sum_{i=1}^{n} W_i(\beta, F, \mu_X) W_i(\beta, F, \mu_X)^{\tau}$. Then, applying Lemma 4.1 of Ritov [5], it can be shown that

$$
\widehat{\Sigma}_1\left(\beta\right) = \widetilde{\Sigma}_1\left(\beta\right) + o_p\left(1\right).
\tag{A.5}
$$

Let $\lambda = \rho\theta$, where $\rho \geq 0$ and $\|\theta\| = 1$. Then, by Owen [21],

$$
\theta' \widetilde{\Sigma}_1\left(\beta\right) \theta \geq \sigma_p + o_p\left(1\right),
\tag{A.6}
$$

where $\sigma_p$ is the smallest eigenvalue of $E\{W_i(\beta, F, \mu) W_i(\beta, F, \mu)^{\tau}\}$. Let $e_j$ be the unit vector in the $j$th coordinate direction. By Lemma A.3,

$$
\left\| \frac{1}{n} \sum_{j=1}^{p} e_j' \sum_{i=1}^{n} W_{in} \right\| = O_p\left(n^{-1/2}\right).
\tag{A.7}
$$

The conclusion of (b) then follows from (10), (A.5)–(A.7), and the arguments used in the proof of (2.14) of Owen [21]. □

*Proof of Theorem 1.* For simplicity, we denote $W_i(\beta, \widehat{F}_n^{\beta}, \overline{X})$ by $W_i$ throughout the proof.

Applying Taylor's expansion to (9), we have

$$
\widetilde{l}_n\left(\beta\right) = 2 \sum_{i=1}^{n} \left( \lambda^{\tau} W_i - \frac{1}{2} (\lambda^{\tau} W_i)^2 \right) + r_n,
\tag{A.8}
$$

where

$$
\left| r_n \right| \leq C \sum_{i=1}^{n} (\lambda^{\tau} W_i)^3 \quad \text{in probability.}
\tag{A.9}
$$

Note that

$$
\left| r_n \right| \leq C \|\lambda\|^3 \max_{1 \leq i \leq n} \|W_i\| \sum_{i=1}^{n} \|W_i\|^2.
\tag{A.10}
$$

This, together with (a), (b), and the fact that

$$
\begin{aligned}
& \frac{1}{n} \sum_{i=1}^{n} \|W_i\|^2 \\
& \leq \frac{1}{n} \sum_{i=1}^{n} \left\| W_i\left(\beta, F, \mu_X\right) \right\|^2 \\
& \quad + \frac{1}{n} \sum_{i=1}^{n} \left\| W_i\left(\beta, \widehat{F}_n^{\beta}, \overline{X}\right) - W_i\left(\beta, F, \mu_X\right) \right\|^2 = O_p\left(1\right),
\end{aligned}
\tag{A.11}
$$

implies that

$$
\left| r_n \right| = o_p\left(1\right).
\tag{A.12}
$$

Note that

$$
\begin{aligned}
0 & = \frac{1}{n} \sum_{i=1}^{n} \frac{W_i}{1 + \lambda^{\tau} W_i} \\
& = \frac{1}{n} \sum_{i=1}^{n} W_i \left[ 1 - \lambda^{\tau} W_i + \frac{(\lambda^{\tau} W_i)^2}{1 + \lambda^{\tau} W_i} \right] \\
& = \frac{1}{n} \sum_{i=1}^{n} W_i - \left( \frac{1}{n} \sum_{i=1}^{n} W_i W_i^{\tau} \right) \lambda \\
& \quad + \frac{1}{n} \sum_{i=1}^{n} \frac{W_i (\lambda^{\tau} W_i)^2}{1 + \lambda^{\tau} W_i}.
\end{aligned}
\tag{A.13}
$$

By (10), (A.5), (A.13), and Lemma A.4, we get

$$
\lambda = \left( \sum_{i=1}^{n} W_i W_i^{\tau} \right)^{-1} \sum_{i=1}^{n} W_i + o_p\left(n^{-1/2}\right).
\tag{A.14}
$$

Again by (10), we have

$$
\begin{aligned}
0 & = \sum_{i=1}^{n} \frac{\lambda^{\tau} W_i}{1 + \lambda^{\tau} W_i} \\
& = \sum_{i=1}^{n} (\lambda^{\tau} W_i) - \sum_{i=1}^{n} (\lambda^{\tau} W_i)^2 \\
& \quad + \frac{1}{n} \sum_{i=1}^{n} \frac{(\lambda^{\tau} W_i)^3}{1 + \lambda^{\tau} W_i}.
\end{aligned}
\tag{A.15}
$$

Moreover, by Lemma A.4 and (A.11), we have

$$\frac{1}{n}\sum_{i=1}^{n}\frac{(\lambda^{\tau}W_i)^3}{1+\lambda^{\tau}W_i} = o_p(1).\tag{A.16}$$

It then follows from (A.15) and (A.16) that

$$\sum_{i=1}^{n}\lambda^{\tau}W_i = \sum_{i=1}^{n}(\lambda^{\tau}W_i)^2 + o_p(1).\tag{A.17}$$

Combining (A.8), (A.17) and (A.14) yields the following identity

$$c_n(\beta)l_n\left(\beta,\widehat{F}_n^{\beta},\overline{X}\right)$$
$$= \left(\frac{1}{\sqrt{n}}\sum_{i=1}^{n}W_i\right)^{\tau}\widehat{\Sigma}_2^{-1}(\beta)\left(\frac{1}{\sqrt{n}}\sum_{i=1}^{n}W_i\right) + o_p(1).\tag{A.18}$$

This, together with Lemma A.3, proves Theorem 1. □

The following lemma is needed to prove Theorem 3.

**Lemma A.5.** *Assume that the conditions of Theorems 1 and 3 hold. Then,*

$$\frac{1}{\sqrt{n}}\sum_{i=1}^{n}A_{ni}(\beta)\xrightarrow{\mathscr{L}}N\left(0,V_{2,AU}(\beta)\right),\tag{A.19}$$

*where*

$$V_{2,AU}(\beta) = \begin{pmatrix} V_1(\beta), & 0 \\ 0, & \Sigma_2(\beta) \end{pmatrix}.\tag{A.20}$$

*Proof.* This result is a direct consequence of Lemma A.3 and the following facts

$$\frac{1}{\sqrt{n}}\sum_{i=1}^{n}g(X_i)\xrightarrow{\mathscr{L}}N\left(0,V_1(\beta)\right),$$
$$\mathrm{Cov}\left(\frac{1}{\sqrt{n}}\sum g(X_i),\frac{1}{\sqrt{n}}\sum_{i=1}^{n}W_i\left(\beta,\widehat{F}_n^{\beta},\overline{X}\right)\right)\longrightarrow 0.\tag{A.21}$$

□

*Proof of Theorem 3.* The theorem can be proved using Lemma A.5 and along the lines of the proof of Theorem 1. We omit the details. □

## Acknowledgments

## References

[1] D. R. Cox, "Regression models and life-tables," *Journal of the Royal Statistical Society B*, vol. 34, no. 2, pp. 187–220, 1972.

[2] P. K. Andersen, O. Borgan, R. D. Gill, and N. Keiding, *Statistical Models Based on Counting Processes*, Springer, New York, NY, USA, 1993.

[3] J. Buckley and I. James, "Linear regression with censored data," *Biometrika*, vol. 66, no. 3, pp. 429–436, 1979.

[4] T. L. Lai and Z. Ying, "Large sample theory of a modified Buckley-James estimator for regression analysis with censored data," *Annals of Statistics*, vol. 19, pp. 1370–1402, 1991.

[5] Y. Ritov, "Estimation in a linear regression model with censored data," *Annals of Statistics*, vol. 18, pp. 303–328, 1990.

[6] A. A. Tsiatis, "Estimating regression parameters using linear rank tests for censored data," *Annals of Statistics*, vol. 18, pp. 354–372, 1990.

[7] Z. Ying, "A large sample study of rank estimation for censored regression data," *Annals of Statistics*, vol. 21, p. 7699, 1993.

[8] Z. Jin, D. Y. Lin, L. J. Wei, and Z. Ying, "Rank-based inference for the accelerated failure time model," *Biometrika*, vol. 90, no. 2, pp. 341–353, 2003.

[9] L. J. Wei, Z. Ying, and D. Y. Lin, "Linear regression analysis of censored survival data based on rank tests," *Biometrika*, vol. 77, no. 4, pp. 845–851, 1990.

[10] D. Y. Lin and C. J. Geyer, "Computational methods for semiparametric linear regression with censored data," *Journal of Computational and Graphical Statistics*, vol. 1, pp. 77–90, 1992.

[11] H. Koul, V. Susarla, and J. van Ryzin, "Regression analysis with randomly right-censored data," *Annals of Statistics*, vol. 9, pp. 1276–1288, 1981.

[12] M. Zhou, "Asymptotic normality of the synthetic estimator for censored survival data," *Annals of Statistics*, vol. 20, pp. 1002–1021, 1992.

[13] T. L. Lai, Z. L. Ying, and Z. K. Zheng, "Asymptotic normality of a class of adaptive statistics with applications to synthetic data methods for censored regression," *Journal of Multivariate Analysis*, vol. 52, no. 2, pp. 259–279, 1995.

[14] E. Kaplan and P. Meier, "Nonparametric estimation from incomplete observations," *Journal of the American Statistical Association*, vol. 53, no. 282, pp. 457–481, 1958.

[15] G. Li and Q. H. Wang, "Empirical likelihood regression analysis for right censored data," *Statistica Sinica*, vol. 13, no. 1, pp. 51–68, 2003.

[16] G. Qin and B. Y. Jing, "Empirical likelihood for censored linear regression," *Scandinavian Journal of Statistics*, vol. 28, no. 4, pp. 661–673, 2001.

[17] G. Li and X. Lu, "Comments on: a review on empirical likelihood methods for regression," *Test*, vol. 18, no. 3, pp. 463–467, 2009.

[18] M. Zhou and G. Li, "Empirical likelihood analysis of the Buckley-James estimator," *Journal of Multivariate Analysis*, vol. 99, no. 4, pp. 649–664, 2008.

[19] D. R. Thomas and G. L. Grunkemeier, "Confidence interval estimation of survival probabilities for censored data," *Journal of the American Statistical Association*, vol. 70, pp. 865–871, 1975.

[20] A. B. Owen, "Empirical likelihood ratio confidence intervals for a single functional," *Biometrika*, vol. 75, no. 2, pp. 237–249, 1988.

[21] A. Owen, "Empirical likelihood ratio confidence regions," *Annals of Statistics*, vol. 18, no. 1, pp. 90–120, 1990.

[22] A. Owen, "Empirical likelihood for linear models," *Annals of Statistics*, vol. 19, pp. 1725–1747, 1991.

[23] S. X. Chen, "On the accuracy of empirical likelihood confidence regions for linear regression model," *Annals of the Institute of Statistical Mathematics*, vol. 45, no. 4, pp. 621–637, 1993.

[24] S. X. Chen, "Empirical likelihood confidence intervals for linear regression coefficients," *Journal of Multivariate Analysis*, vol. 49, no. 1, pp. 24–40, 1994.

[25] E. Kolaczyk, "Empirical likelihood for generalized linear models," *Statistica Sinica*, vol. 4, pp. 199–218, 1994.

[26] S. X. Chen and P. Hall, "Smoothed empirical likelihood confidence intervals for quantiles," *Annals of Statistics*, vol. 21, no. 3, pp. 1166–1181, 1993.

[27] J. Qin, "Empirical likelihood in biased sample problems," *Annals of Statistics*, vol. 21, no. 3, pp. 1182–1196, 1993.

[28] J. Qin and J. F. Lawless, "Empirical likelihood and general estimating equations," *Annals of Statistics*, vol. 22, pp. 300–325, 1994.

[29] Y. Kitamura, "Empirical likelihood methods with weakly dependent processes," *Annals of Statistics*, vol. 25, no. 5, pp. 2084–2102, 1997.

[30] Q. H. Wang and B. Y. Jing, "Empirical likelihood for partial linear models with fixed design," *Statistics & Probability Letters*, vol. 41, pp. 425–433, 1999.

[31] J. Qin, "Empirical likelihood ratio based confidence intervals for mixture proportions," *Annals of Statistics*, vol. 27, no. 4, pp. 1368–1384, 1999.

[32] M. Hollander, I. W. McKeague, and J. Yang, "Likelihood ratio-based confidence bands for survival functions," *Journal of the American Statistical Association*, vol. 92, no. 437, pp. 215–226, 1997.

[33] G. Li, M. Hollander, I. W. McKeague, and J. Yang, "Nonparametric likelihood ratio confidence bands for quantile functions from incomplete survival data," *Annals of Statistics*, vol. 24, no. 2, pp. 628–640, 1996.

[34] X. H. Zhou, G. S. Qin, H. Z. Lin, and G. Li, "Inferences in censored cost regression models with empirical likelihood," *Statistica Sinica*, vol. 16, no. 4, pp. 1213–1232, 2006.

[35] J. H. J. Einmahl and I. W. McKeague, "Confidence tubes for multiple quantile plots via empirical likelihood," *Annals of Statistics*, vol. 27, no. 4, pp. 1348–1367, 1999.

[36] T. J. DiCiccio, P. Hall, and J. P. Romano, "Empirical likelihood is Bartlett-correctable," *Annals of Statistics*, vol. 19, no. 2, pp. 1053–1061, 1991.

[37] P. Hall, *The Bootstrap and Edgeworth Expansion*, Springer, New York, NY, USA, 1992.

[38] P. Hall and B. La Scala, "Methodology and algorithms of empirical likelihood," *International Statistical Review*, vol. 58, no. 2, pp. 109–127, 1990.

[39] A. Owen, *Empirical Likelihood*, Chapman & Hall, London, UK, 2001.

[40] G. Li, R. Li, and M. Zhou, "Empirical likelihood in survival analysis," in *Contemporary Multivariate Analysis and Experimental Design*, J. Fan and G. Li, Eds., pp. 336–350, World Scientific, 2005.

[41] J. Chen and J. Qin, "Empirical likelihood estimation for finite populations and the effective usage of auxiliary information," *Biometrika*, vol. 80, no. 1, pp. 107–116, 1993.

[42] B. Zhang, "$M$-estimation and quantile estimation in the presence of auxiliary information," *Journal of Statistical Planning and Inference*, vol. 44, no. 1, pp. 77–94, 1995.

[43] B. Zhang, "Confidence intervals for a distribution function in the presence of auxiliary information," *Computational Statistics and Data Analysis*, vol. 21, no. 3, pp. 327–342, 1996.

[44] R. G. Miller, "Least squares regression with censored data," *Biometrika*, vol. 63, no. 3, pp. 449–464, 1976.

[45] S. Leurgans, "Linear models, random censoring and synthetic data," *Biometrika*, vol. 74, no. 2, pp. 301–309, 1987.