# AN ADDITIVE RISKS REGRESSION MODEL FOR MIDDLE-CENSORED LIFETIME DATA

## P.G. Sankaran[1], Prasad S.[2]

## Abstract

Middle-censoring refers to data arising in situations where the exact lifetime of study subjects becomes unobservable if it happens to fall in a random censoring interval. In the present paper we propose a semiparametric additive risks regression model for analysing middle-censored lifetime data arising from an unknown population. We estimate the regression parameters and the unknown baseline survival function by two different methods. The first method uses the martingale-based theory and the second method is an iterative method. We report simulation studies to assess the finite sample behaviour of the estimators. Then, we illustrate the utility of the model with a real life data set. The paper ends with a conclusion.

**Key words:** additive risks model, counting process, martingales, middle-censoring.

## 1. Introduction

Middle-censoring introduced by Jammalamadaka & Mangalam (2003) occurs in situations where a data point becomes unobservable if it falls inside a random censoring interval. In such situations, the exact values are available for some individuals and for others, random censoring intervals are observed. To be more precise, let $T$ be the random variable representing the lifetime of interest and let $(U,V)$ be a bivariate random variable, representing the censoring interval, such that $P(U < V) = 1$. Under the middle-censored set-up, the exact lifetime $T$ becomes unobservable if $T \in (U,V)$, and in such instances we only observe the censoring interval $(U,V)$. Otherwise we observe $T$. We may find several such situations in survival studies and reliability applications. For example, in a prognostic study, the patients under observation may be withdrawn from the study for a short period of time for some unforeseen reasons and may return to the study with a changed status of event of interest. In reliability applications, it may happen that a failure of equipment occurs during a period of time when we accidentally fail to observe the study subjects. In

[1]Department of Statistics, Cochin University of Science and Technology, Kerala, India. E-mail: sankaran.p.g@gmail.com
[2]Department of Statistics, Cochin University of Science and Technology, Kerala, India. E-mail: hariprasadtvpm@gmail.com(Corresponding Author).

such contexts we only observe a censorship indicator and the interval of censorship.

As was pointed out by Jammalamadaka & Mangalam (2003), one can observe that the left censored data and right censored data are in fact special cases of this more general censoring scheme, by suitable choices of the interval, and also that such a censoring scheme is not complementary to the usual double censoring discussed in Klein & Moeschberger (2005) and Sun (2006).

Jammalamadaka & Mangalam (2003) pointed out various applications of middle-censoring scheme and developed a nonparametric maximum likelihood estimator (NPMLE) of the distribution function of the random variable. They proved that the NPMLE is always a self-consistent estimator (SCE) (Tarpey & Flury, 1996). Some rigorous treatments of this censoring scheme are found in Jammalamadaka & Iyer (2004), Iyer et al. (2008), Mangalam et al. (2008), Jammalamadaka & Mangalam (2009), Shen (2010, 2011), Davarzani & Parsian (2011) and Davarzani et al. (2015).

In survival studies, covariates or explanatory variables are usually used to represent heterogeneity in a population. The main objective in such situations is to understand and exploit the relationship between the lifetime and covariates. To this end we generally employ regression models. In the presence of covariates, Sankaran & Prasad (2014) discussed a parametric proportional hazards regression model for the analysis of middle-censored lifetime data. Jammalamadaka & Leong (2015) analysed discrete middle-censored data in the presence of covariates with an accelerated failure time regression model. Recently, Jammalamadaka et al. (2016) developed an iterative algorithm for analysing a semiparametric proportional hazards regression model under middle-censoring scheme, while Bennett et al. (2017) considered a parametric accelerated failure time regression model under this censoring scheme.

One extensively used semiparametric regression model is the well-known proportional hazards (PH) model by Cox (1972). It is a multiplicative hazards model in the sense that if $T$ has a baseline hazard function $h_0(t)$ and if $\mathbf{z}$ is a $p \times 1$ vector of the recorded covariates then the hazard function of $T$ conditional on $\mathbf{z}$ is modelled as

$$h(t|\mathbf{z}) = h_0(t)\exp(\mathbf{z}^\top \boldsymbol{\theta}),$$

where $\boldsymbol{\theta} = (\theta_1, \theta_2, ..., \theta_p)^\top$ is the vector of regression coefficients and $h_0(t)$ is left arbitrary. Here, $a^\top$ represents the transpose of vector $a$ . In this model the effect of the covariates is acting multiplicatively on the baseline hazard function. But it is well known that in many occasions the PH model does not fit a given lifetime data well. One important alternative to the PH model is the additive risks (AR) model introduced by Aalen (1989) and later studied by Lin & Ying (1994). The model

associates the conditional hazard function with the covariates by

$$h(t|\mathbf{z}) = h_0(t) + \mathbf{z}^\top \boldsymbol{\theta}. \tag{1}$$

In contrast to the PH model, the AR model given in (1) specifies that the hazard rate associated with a given set of covariates is the sum of the baseline hazard function and the regression function of covariates. This kind of model assumption is particularly useful in tumorigenicity experiments that investigate the dose effect on tumor risk, since the excess risk is often the quantity of interest (Breslow & Day, 1987). For a comprehensive review on properties and inference procedures of model (1), one may refer to Aranda-Ordaz (1983), Cox & Oakes (1984), Thomas (1986), Breslow & Day (1980), and Lin & Ying (1994). For a nonparametric treatment of model (1) one may refer to Aalen (1980, 1989). Model (1) is further explored in the context of left truncated current status data by Wang et al. (2015).

   In the present work, we aim at estimating the unknown baseline survival function $S_0(t)$ of a continuous type lifetime variate $T$, which is subject to middle-censoring, and estimation of the unknown regression coefficients under model (1). We propose two different inference methods in Section 2. Simulation studies to assess the performance of the estimators under both methods for practical sample sizes are carried out and the results are compared in Section 3. The utility of the methods are illustrated with the help of a real life example in Section 4. Finally, some important conclusions are provided in Section 5.

## 2. Inference Procedure

Let the lifetime variate $T$ admit an absolutely continuous cumulative distribution function (cdf) $F_0(t)$. Assume that $T$ is middle-censored by the random censoring interval $(U,V)$ having bivariate cdf given by $G(u,v) = P(U \le u, V \le v)$. Let us further assume that under model (1), $T$ is independent of $(U,V)$, given the covariate $\mathbf{z}$. Thus, one can observe the vector $(X, \delta, \mathbf{z})$, where

$$X = \begin{cases} T & \text{if } \delta = 1 \\ (U,V) & \text{if } \delta = 0, \end{cases}$$

and $\delta = I(X = T)$ is the uncensoring indicator. Now, we state an important assumption regarding the identifiability of the cdf $F_0(t)$. Let $[a,b]$, $a \le b$ be any arbitrary interval in the support of $T$. Define, for $r \in [a,b]$,

$$A_0(r) = G(r-, \infty) - G(r-, r) = P(U < r < V). \tag{2}$$

Now, consider a situation where $A_0(r) = 1$ for $r \in [a, b]$ for which $F_0(b) > F_0(a-)$. *i.e.* censoring occurs with probability 1 on this interval where $F_0$ has a positive mass. Consequently, there will not be any exact observation in this interval, making it impossible to distinguish two distributions which are identical outside $[a, b]$ but differ only on $[a, b]$. To overcome this issue we make the following assumption.
A1: The probability defined in (2) is strictly less than one.

In the following we describe two different estimation methods: one makes use of the classic martingale theory and the other by using of an iterative method.

## 2.1. Martingale Method

Here we provide an inference procedure based on the martingale feature associated with the observed data. First a partial likelihood function is developed under model (1), similar to the one for the Cox PH model (Kalbfleisch & Prentice, 2011). Then, the stochastic integral representation of the score function derived from the partial likelihood function is used to infer about the unknown regression coefficient.

The observed data consists of $n$ independent and identically distributed replicates $(X_i, \mathbf{z}_i, \delta_i)$ of $(X, \mathbf{z}, \delta), 1 \leq i \leq n$. When the lifetime is subject to middle-censoring, we shall define the counting process corresponding to the $i$'th individual as $N_i(t) = I(X_i \leq t, \delta_i = 1), t \geq 0$, which indicates whether the event occurred at time $t$, for $i = 1, 2, ..., n$. The at-risk process may be similarly defined as $R_i(t) = I(X_i \geq t, \delta_i = 1) + I(U_i \geq t, \delta_i = 0)$ which is a 0-1 predictable process, where the value 1 indicates whether the $i$'th individual is at risk at time $t$, for $i = 1, 2, ..., n$, *i.e.*, whether it is uncensored and waiting for a possible event at the epoch $t$. Denote the filtration $\sigma\{N_i(u), R_i(u+), \mathbf{z}_i : i = 1, 2, ..., n; 0 \leq u \leq t\}$ by $\mathscr{F}_t$. Under model (1) the conditional cumulative hazard rate for the $i$'th individual is given by $H(t|\mathbf{z}_i) = H_0(t) + \mathbf{z}_i^\top \theta t$, where $H_0(t) = \int_0^t h_0(a) da$ is the baseline cumulative hazard function. Model (1) assumes that

$$E[N_i(t)|\mathscr{F}_{t-}] = (h_0(t) + \theta^\top \mathbf{z}_i) R_i(t) dt,$$

and the intensity function corresponding to the counting process $N_i(t)$ can thus be written as $R_i(t) dH(t|\mathbf{z}_i) = R_i(t)\{dH_0(t) + \mathbf{z}_i^\top \theta dt\}$. With this, the counting process can be uniquely decomposed so that for every $i$ and $t$,

$$N_i(t) = M_i(t) + \int_0^t R_i(a) \, dH(a|\mathbf{z}_i), \tag{3}$$

where $M_i(\cdot)$ is a local square integrable martingale (Andersen & Gill, 1982). From (3), we have

$$dN_i(t) = dM_i(t) + R_i(t)dH(t|\mathbf{z}_i), \tag{4}$$

so that

$$\sum_{i=1}^{n} dM_i(t) = \sum_{i=1}^{n} [dN_i(t) - R_i(t)(dH_0(t) + \boldsymbol{\theta}^{\top} \mathbf{z}_i dt)] = 0. \tag{5}$$

To estimate $\boldsymbol{\theta}$, let us now consider the partial likelihood function suggested by Cox (1972) and further discussed in Cox (1975). It is defined as

$$L(\boldsymbol{\theta}) = \prod_{i=1}^{k} \frac{h(t_{(i)}|\mathbf{z}_{(i)})}{\sum_{l=1}^{n} R_l(t_{(i)}) h(t_{(i)}|\mathbf{z}_l)}, \tag{6}$$

where $t_{(1)}, t_{(2)}, ..., t_{(k)}$ are the $k$ observed exact lifetimes which are arranged in increasing order of magnitude. The motivation for (6) is that when we have the information that an event occurs at time point $t$ and that the at-risk set is $R(t)$, the right-hand side of (6) is precisely the probability that it is individual $i \in R(t)$, who registered the event. Since $T$ is assumed to be of continuous type, the possibility of ties is ruled out. However, (6) is not a usual likelihood, as it is not obtained from the probability of some observable events. A detailed discussion on this is available in Lawless (2011). Under the model assumption (1), we can rewrite (6) as

$$L(\boldsymbol{\theta}) = \prod_{i=1}^{k} \frac{h_0(t_{(i)}) + \mathbf{z}_{(i)}^{\top} \boldsymbol{\theta}}{\sum_{l=1}^{n} R_l(t_{(i)}) \left( (h_0(t_{(i)}) + \mathbf{z}_l^{\top} \boldsymbol{\theta}) \right)}. \tag{7}$$

The value of $\boldsymbol{\theta}$ that maximizes (7) can be obtained by maximizing

$$l(\boldsymbol{\theta}) = \log(L(\boldsymbol{\theta})) = \sum_{i=1}^{k} \Bigg[ \log\Big( h_0(t_{(i)}) + \mathbf{z}_{(i)}^{\top} \boldsymbol{\theta} \Big) -$$

$$\log\Bigg( \sum_{l=1}^{n} R_l(t_{(i)}) \Big( (h_0(t_{(i)}) + \mathbf{z}_l^{\top} \boldsymbol{\theta}) \Big) \Bigg) \Bigg]. \tag{8}$$

In terms of the counting process defined earlier, we can rewrite (8) as

$$C(\boldsymbol{\theta}) = \sum_{i=1}^{n} \int_0^{\infty} \log\Big( h_0(s) + \mathbf{z}_i^{\top} \boldsymbol{\theta} \Big) dN_i(s) -$$

$$\int_0^{\infty} \log\Bigg( \sum_{l=1}^{n} R_l(s) \big( h_0(s) + \mathbf{z}_l^{\top} \boldsymbol{\theta} \big) \Bigg) d\bar{N}(s), \tag{9}$$

where $\bar{N}(s) = \sum_{i=1}^n N_i(s)$. The score function is simply the derivative of (9) with respect to $\theta$, and is given by

$$U(\theta) = \sum_{i=1}^n \int_0^\infty \left( h_0(s) + \mathbf{z}_i^\top \theta \right)^{-1} \mathbf{z}_i \mathrm{d}N_i(s) -$$
$$\int_0^\infty \left( \sum_{l=1}^n R_l(s) \left( (h_0(s) + \mathbf{z}_l^\top \theta) \right) \right)^{-1} \left( \sum_{l=1}^n R_l(s)\mathbf{z}_l \right) \mathrm{d}\bar{N}(s). \qquad (10)$$

Using the idea of Lin & Ying (1994), we propose to estimate the true regression coefficient $\theta_0$ from the following estimating equation, which is obtained by an algebraic simplification of (10).

$$U(\theta) = \sum_{i=1}^n \int_0^\infty \mathbf{z}_i \{ \mathrm{d}N_i(t) - R_i(t)d\hat{H}_0(\theta,t) - R_i(t)\mathbf{z}_i^\top \theta \mathrm{d}t \},$$

which is equivalent to

$$U(\theta) = \sum_{i=1}^n \int_0^\infty \{ \mathbf{z}_i - \bar{\mathbf{z}} \} \{ \mathrm{d}N_i(t) - R_i(t)\mathbf{z}_i^\top \theta \mathrm{d}t \}, \qquad (11)$$

where $\bar{\mathbf{z}} = \sum_{i=1}^n \mathbf{z}_i R_i(t) / \sum_{i=1}^n R_i(t)$, with the convention that $0/0 = 0$. The identity (11) is based on a simple fact that when $\theta_0$ is the true parameter value, $U(\theta_0)$ is a martingale integral and therefore has mean zero. Note that (11) is linear in $\theta$ and the resulting estimator takes an explicit form given by

$$\hat{\theta} = \left[ \sum_{i=1}^n \int_0^\infty [\mathbf{z}_i - \bar{\mathbf{z}}]^{\otimes 2} R_i(t) \mathrm{d}t \right]^{-1} \sum_{i=1}^n \int_0^\infty [\mathbf{z}_i - \bar{\mathbf{z}}] \mathrm{d}N_i(t), \qquad (12)$$

where $a^{\otimes 2} = aa^\top$. Since $M_i(t)$ is a martingale, we have $\sum_{i=1}^n \mathrm{d}M_i(t) = 0$. Thus, from the representation given in (3), a Breslow type estimator (Breslow, 1972) for the cumulative hazard function $H_0(t)$ can be obtained as

$$\hat{H}_o(\hat{\theta},t) = \int_0^t \frac{\sum_{i=1}^n \{ \mathrm{d}N_i(a) - R_i(a)\mathbf{z}_i^\top \hat{\theta} \mathrm{d}a \}}{\sum_{i=1}^n R_i(a)}. \qquad (13)$$

This naturally leads to the following estimator of conditional survival function $S(t|\mathbf{z})$.

$$\hat{S}(t|\mathbf{z}) = \exp\{ -\hat{H}_0(\hat{\theta},t) - \mathbf{z}^\top \hat{\theta} t \}. \qquad (14)$$

An algebraic manipulation of (4) yields

$$U(\theta) = \sum_{i=1}^{n} \int_{0}^{\infty} (\mathbf{z}_i - \bar{\mathbf{z}}) dM_i(t), \tag{15}$$

which is a martingale integral. It follows from standard counting process theory that $n^{-1/2}U(\theta_0)$ converges weakly to a $p$-variate normal with mean zero and a covariance matrix that can be estimated consistently by

$$A = \frac{1}{n} \sum_{i=1}^{n} \int_{0}^{\infty} (\mathbf{z}_i - \bar{\mathbf{z}})^{\otimes 2} dN_i(t). \tag{16}$$

Also, the random vector $n^{1/2}(\hat{\theta} - \theta_0)$ converges weakly to a $p$-variate normal distribution with mean zero and a covariance matrix that can be consistently estimated by $B^{-1}AB^{-1}$, where

$$B = \frac{1}{n} \sum_{i=1}^{n} \int_{0}^{\infty} R_i(t)(\mathbf{z}_i - \bar{\mathbf{z}})^{\otimes 2} dt. \tag{17}$$

Specifically, $(B^{-1}AB^{-1})^{-\frac{1}{2}}(\hat{\theta} - \theta_0)$ converges in distribution to $N(0, I)$. It can be observed that neither $A$ nor $B$ involves the regression parameters. The estimator (13) provides the basis for estimating survival probabilities. Using standard counting process techniques, it follows that the process $\sqrt{n}(\hat{H}_0(\hat{\theta}, t) - H_0(t))$ converges weakly to a zero mean Gaussian process, whose covariance function at $(t, s), t \geq s$ can be estimated consistently by

$$\int_{0}^{s} \frac{n \sum_{i=1}^{n} dN_i(a)}{(\sum_{1}^{n} R_i(a))^2} + C'(t)B^{-1}AB^{-1}C(s) - C'(t)B^{-1}D(s) - C'(s)B^{-1}D(t),$$

where $C(t) = \bar{\mathbf{z}}t$ and $D(t) = \int_{0}^{t} \frac{\sum_{1}^{n}(\mathbf{z}_i - \bar{\mathbf{z}})dN_i(a)}{\sum_{1}^{n} R_i(a)}$ with $k'(a) = dk(a)/da$.
Using functional delta method (Andersen et al., 2012), it follows that the process $\sqrt{n}(\hat{S}(t|\mathbf{z}) - S(t|\mathbf{z}))$ converges weakly to a zero-mean Gaussian process, whose covariance function at $(t, s), t \geq s$ can be estimated consistently by

$$\hat{S}(t|\mathbf{z})\hat{S}(s|\mathbf{z})\Big( \int_{0}^{s} \frac{n \sum_{i=1}^{n} dN_i(a)}{(\sum_{1}^{n} R_i(a))^2}$$
$$+ W'(t,\mathbf{z})B^{-1}AB^{-1}W(s,\mathbf{z}) + W'(t,\mathbf{z})B^{-1}D(s) + W'(s,\mathbf{z})B^{-1}D(t)\Big),$$

where $W(t, \mathbf{z}) = (\mathbf{z} - \bar{\mathbf{z}})t$.

## 2.2. The Iterative Method

In this section, an iterative method is proposed for estimating the unknown baseline survival function $S_0(t)$ of the lifetime variate $T$ and the regression coefficient vector $\theta$ under model (1). Assume that $T$ is middle-censored by the random censoring interval $(U, V)$ such that, given the covariate $\mathbf{z}$, $T$ and $(U, V)$ are independently distributed. Let the observed data be as before. For convenience let us assume that the first $n_1$ observations are exact lifetimes, and the remaining $n_2$ are censored intervals, with $n_1 + n_2 = n$. Now, the likelihood corresponding to the observed data, excluding the normalizing constant, can be written as

$$L(\theta) = \prod_{i=1}^{n_1} f(t_i|\mathbf{z}_i) \cdot \prod_{i=n_1+1}^{n_1+n_2} \left( S(u_i|\mathbf{z}_i) - S(v_i|\mathbf{z}_i) \right). \tag{18}$$

Under the model assumption given in (1), the conditional survival function is obtained as

$$S(t|z) = S_0(t)\exp(-\theta^\top \mathbf{z}t), \tag{19}$$

where $S_0(t) = \exp(-H_0(t))$. Thus, the density function of $T$ given $\mathbf{z}$ is given by

$$f(t|\mathbf{z}) = \exp(-\theta^\top \mathbf{z}t)\left(\theta^\top \mathbf{z}S_0(t) - S_0'(t)\right). \tag{20}$$

Therefore, (18) becomes

$$L(\theta) = \prod_{i=1}^{n_1}\exp(-\theta^\top \mathbf{z}_i t_i)\left(\theta^\top \mathbf{z}_i S_0(t_i) - S_0'(t_i)\right) \times$$
$$\prod_{i=n_1+1}^{n_1+n_2} \left( S_0(u_i)\exp(-\theta^\top \mathbf{z}_i u_i) - S_0(v_i)\exp(-\theta^\top \mathbf{z}_i v_i)\right). \tag{21}$$

The log-likelihood is given by

$$l(\theta) = \log L(\theta) = \sum_{i=1}^{n_1}\left( -\theta^\top \mathbf{z}_i t_i + \log(\theta^\top \mathbf{z}_i S_0(t_i) - S_0'(t_i))\right) +$$
$$\sum_{i=n_1+1}^{n_1+n_2} \log\left( S_0(u_i)\exp(-\theta^\top \mathbf{z}_i u_i) - S_0(v_i)\exp(-\theta^\top \mathbf{z}_i v_i)\right), \tag{22}$$

and its partial derivative with respect to $\theta_r$, for $r = 1, 2, ..., p$, is given by

$$\frac{\partial l(\boldsymbol{\theta})}{\partial \theta_r} = \sum_{i=1}^{n_1} z_{ir}(t_i + (\boldsymbol{\theta}^\top \mathbf{z}_i S_0(t_i) - S_0'(t_i))^{-1} S_0(t_i)) + \sum_{i=n_1+1}^{n_1+n_2} z_{ir}\Big(S_0(u_i)\exp(-\boldsymbol{\theta}^\top \mathbf{z}_i u_i) -$$

$$S_0(v_i)\exp(-\boldsymbol{\theta}^\top \mathbf{z}_i v_i)\Big)^{-1} \Big(v_i S_0(v_i)\exp(-\boldsymbol{\theta}^\top \mathbf{z}_i v_i) - u_i S_0(u_i)\exp(-\boldsymbol{\theta}^\top \mathbf{z}_i u_i)\Big), \tag{23}$$

where $z_{ir}$ is the $r$'th component in the covariate vector corresponding to $i$'th individual. Note that (23) involves both unknown quantities $\boldsymbol{\theta}$ and $S_0(t)$ and explicit solution for $\boldsymbol{\theta}$ cannot be obtained directly from it. We provide an iterative algorithm to estimate the maximum likelihood estimates of these two quantities, where at each iteration a better update is obtained. To begin with the algorithm we consider the SCE of the baseline survival function as an initial approximation.

In the case of middle-censored data, Jammalamadaka & Mangalam (2003) showed that the NPMLE of $S_0(t)$ is always an SCE, which takes the form

$$\hat{S}_0(t) = 1 - \frac{1}{n} \sum_{i=1}^{n} \left\{ \delta_i I(T_i \leq t) + (1-\delta_i)I(V_i \leq t) + (1-\delta_i)I(t \in (U_i, V_i)) \right.$$

$$\left. \frac{\hat{F}_0(t) - \hat{F}_0(U_i)}{\hat{F}_0(V_i-) - \hat{F}_0(U_i)} \right\}. \tag{24}$$

Now, we give the algorithm in the following few steps.

**Step 1.** Set the vector $\boldsymbol{\theta} = 0$.

**Step 2.** At the first iteration, find the SCE $S_0^{(1)}(t)$ of $S_0(t)$ using (24) and substitute this in (23) and solve $\partial l(\boldsymbol{\theta})/\partial \theta_r = 0, r = 1, 2, ..., p$ to get the estimator $\boldsymbol{\theta}^{(1)}$ of $\boldsymbol{\theta}$.

**Step 3.** Find $\tilde{t}_i^{(1)} = S_0^{(1)-1}\left(S_0^{(1)}(t_i)\exp(-\boldsymbol{\theta}^{(1)\top}\mathbf{z}_i t_i)\right)$ and similarly find $\tilde{u}_i^{(1)}$ and $\tilde{v}_i^{(1)}$ as our updated observations at the first iteration.

**Step 4.** At the $j$'th iteration ($j > 1$), use $\tilde{t}_i^{(j-1)}, i = 1, 2, ..., n_1$ and $(\tilde{u}_i^{(j-1)}, \tilde{v}_i^{(j-1)}), i = n_1 + 1, ..., n$ as our data points in (24) and obtain $S_0^{(j)}(t)$. Substitute $S_0^{(j)}(t)$ in (23) and solve $\partial l(\boldsymbol{\theta})/\partial \theta_r = 0, r = 1, 2, ..., p$ to obtain the $j$'th iterated update $\boldsymbol{\theta}^{(j)}$ of $\boldsymbol{\theta}$.

**Step 5.** Repeat Step 4 until convergence is met, say when $\|\boldsymbol{\theta}^{(k)} - \boldsymbol{\theta}^{(k+1)}\| < 0.0001$ and $\sup_t \left\{ \left| S_0^{(k)}(t) - S_0^{(k+1)}(t) \right| \right\} < 0.001$, for some finite positive integer $k$.

Note that Step 3 in the algorithm is justified, because if $a_i = S_0^{(1)}(t_i)\exp(-\boldsymbol{\theta}^{(1)\top}\mathbf{z}_i t_i)$, then the $a_i$ 's have a uniform distribution over $[0, 1]$. Therefore, to scale these back to baseline distribution we need to find $\tilde{t}_i = \inf\{t : S_0^{(1)}(t) \leq a_i\}$. Thus, the correct choice is $\tilde{t}_i = S_0^{(1)-1}(a_i) = S_0^{(1)-1}\left(S_0^{(1)}(t_i)\exp(-\boldsymbol{\theta}^{(1)\top}\mathbf{z}_i t_i)\right)$.

We now define our parameter space to be $(\Theta, \Phi)$, where $\Theta \subseteq \mathbb{R}_p$ contains $\boldsymbol{\theta}$ and

$\Phi = \{\phi(t) : [0, \infty] \to [0, 1]$ and $\phi(\cdot)$ is absolutely continuous and nonincreasing$\}$ contains $S_0(t)$. Let us name the estimator obtained for $\theta$ as $\hat{\theta}_{(n)}$ and that for $S_0(t)$ as $\hat{S}_{0(n)}(t)$. Besides the identifiability condition A1, the following conditions are also assumed to hold for establishing the consistency property.

A2: Conditional on $\mathbf{z}$, $T$ is independent of $(U, V)$.

A3: The joint distribution of $(U, V, \mathbf{z})$ does not depend on the true parameter $(\theta^0, S_0^0(t))$.

A4: The covariate space is bounded. That is, there exist some finite $M > 0$ such that $P\{\|\mathbf{z}\| \leq M\} = 1$, where $\|\cdot\|$ is the usual metric on $\mathbb{R}_p$.

A5: Distribution of $\mathbf{z}$ is not concentrated on any proper affine subspace of $\mathbb{R}_p$.

***Theorem:*** Suppose that $\Theta \in \mathbb{R}_p$ is bounded and assumptions (A1) to (A5) hold. Then, the estimator $(\hat{\theta}_{(n)}, \hat{S}_{0(n)}(t))$ is consistent for the true parameter $(\theta^0, S_0^0(t))$ in the sense that if we define a metric $d : \Theta \times \Phi \to \mathbb{R}$ by

$$d\big((\theta_1, S_{01}(t)), (\theta_2, S_{02}(t))\big) = \|\theta_1 - \theta_2\| + \int |S_{01}(t) - S_{02}(t)| dF_0(t) +$$

$$\left[ \int \big((S_{01}(u) - S_{02}(u))^2 + (S_{01}(v) - S_{02}(v))^2\big) dG(u, v) \right]^{\frac{1}{2}}, \quad (25)$$

where $\theta_1, \theta_2 \in \Theta$ and $S_{01}(t), S_{02}(t) \in \Phi$, then $d\big((\hat{\theta}_{(n)}, \hat{S}_{0(n)}(t)), (\theta^0, S_0^0(t))\big) \to 0$ almost surely (a.s.).

***Proof:***
In the following discussion we denote $Y_i = (X_i, \delta_i)$. Let the probability function of $Y = (X, \delta)$ be given by

$$p(y; \theta, S_0(t)) = \prod_{i=1}^{n} f(t_i|\mathbf{z}_i)^{\delta_i} [S_0(u_i)\exp(-\theta^\top \mathbf{z}_i u_i) - S_0(v_i)\exp(-\theta^\top \mathbf{z}_i v_i)]^{1-\delta_i} \times$$

$$g(u_i, v_i|\mathbf{z}_i) q(\mathbf{z}_i), \quad (26)$$

where $g$ is the joint density of $(U, V)$, conditional on $\mathbf{z}$ and $q$ is the density of $\mathbf{z}$. Using $(A2)$ and $(A3)$, the log-likelihood function scaled by $1/n$ for the sample $(y_i, \mathbf{z}_i), i = 1, 2, ..., n$, up to terms not depending on $(\theta^0, S_0^0(t))$ is

$$l(\theta, S_0(t)) = \frac{1}{n} \sum_{i=1}^{n} \big\{ \delta_i \log f(t_i|\mathbf{z}_i) + (1 - \delta_i) \log [S_0(u_i)\exp(-\theta^\top \mathbf{z}_i u_i) -$$

$$S_0(v_i)\exp(-\theta^\top \mathbf{z}_i v_i)] \big\}. \quad (27)$$

We write $p_n(y) = p(y; \hat{\theta}_{(n)}, \hat{S}_{0(n)}(t))$ and $p_0(y) = p(y; \theta^0, S_0^0(t))$ where $(\hat{\theta}_{(n)}, \hat{S}_{0(n)}(t))$ is the MLE that maximizes the likelihood function over $\Theta \times \Phi$ and $(\theta^0, S_0^0(t)) \in$

$\Theta \times \Phi$. Therefore,

$$\sum_{i=1}^{n} \log p_n(Y_i) \geq \sum_{i=1}^{n} \log p_0(Y_i)$$

and hence

$$\sum_{i=1}^{n} \log \frac{p_n(Y_i)}{p_0(Y_i)} \geq 0.$$

By the concavity of the function $x \mapsto \log x$, for any $0 < \alpha < 1$,

$$\frac{1}{n} \sum_{i=1}^{n} \log \left( (1-\alpha) + \alpha \frac{p_n(Y_i)}{p_0(Y_i)} \right) \geq 0. \tag{28}$$

The left hand side can be written as

$$\int \log \left( (1-\alpha) + \alpha \frac{p_n(Y_i)}{p_0(Y_i)} \right) d(\mathbb{P}_n - \mathbb{P})(Y) + \int \log \left( (1-\alpha) + \alpha \frac{p_n(Y_i)}{p_0(Y_i)} \right) d\mathbb{P}(Y), \tag{29}$$

where $\mathbb{P}_n$ is the empirical measure of $Y$ and $\mathbb{P}$ is the joint probability measure of $Y$. Let us assume that the sample space $\Omega$ consists of all infinite sequences $Y_1, Y_2, ...,$ along with the usual sigma field generated by the product topology on $\prod_1^{\infty}(\mathbb{R}^3 \times \{0,1\})$ and the product measure $\mathbf{P}$. For $p$ defined in (26) let us define a class of functions $\mathscr{P} = \left\{ p(y, \theta, S_0(t)) : (\theta, S_0(t)) \in (\Theta \times \Phi) \right\}$ and a class of functions $\mathscr{H} = \left\{ \log(1 - \alpha + \alpha p/p_0) : p \in \mathscr{P} \right\}$, where $p_0 = p(y, \theta^0, S_0^0(t))$. Then, it follows from Huang & Wellner (1995) that $\mathscr{H}$ is a Donsker class. With this and Glivenko-Cantelli theorem, there exists a set $\Omega_0 \in \Omega$ with $\mathbf{P}(\Omega_0) = 1$ such that for every $\omega \in \Omega_0$, the first term of (29) converges to zero. Now, fix a point $\omega \in \Omega_0$ and write $\hat{\theta}_{(n)} = \hat{\theta}_{(n)}(\omega)$ and $\hat{S}_{0(n)}(\cdot) = \hat{S}_{0(n)}(\cdot, \omega)$. By our assumption $\Theta$ is bounded, and hence for any subsequence of $\hat{\theta}_{(n)}$, we can find a subsequence converging to $\theta_* \in \Theta^C$, the closure of $\Theta$. Also, by Helly's selection theorem, for any subsequence of $\hat{S}_{0(n)}(t)$, we can find a further subsequence converging to some nonincreasing function $S_{0*}(t)$. Choose the convergent subsequence of $\hat{\theta}_{(n)}$ and the convergent subsequence of $\hat{S}_{0(n)}(t)$ so that they have the same indices, and without loss of generality, assume that $\hat{\theta}_{(n)}$ converges to $\theta_*$ and that $\hat{S}_{0(n)}(t)$ converges to $S_{0*}(t)$. Let $p_*(y) = p(y, \theta_*, S_{0*}(t))$. By the bounded convergence theorem, the second term of (29) converges to

$$\int \log \left( (1-\alpha) + \alpha \frac{p_*(y)}{p_0(y)} \right) d\mathbb{P}(y)$$

and by (28) this is nonnegative. But by Jensen's inequality, it must be non-positive.

Therefore, it must be zero and it follows that

$$p_*(y) = p_0(y) \quad \mathbb{P}-\text{almost surely.}$$

This implies

$$S_{0*}(t) = S_0^0(t) \quad F_0 - \text{almost surely.}$$

Therefore, by bounded convergence theorem,

$$\int |\hat{S}_{0(n)}(t) - S_0^0(t)| \mathrm{d}F_0(t) \to 0. \tag{30}$$

Also,

$$S_{0*}(u)\exp(-\boldsymbol{\theta}_*^\top \mathbf{z}u) = S_0^0(u)\exp(-\boldsymbol{\theta}^{0^\top} \mathbf{z}u) \quad \mathbb{P}-\text{almost surely}$$

and

$$S_{0*}(v)\exp(-\boldsymbol{\theta}_*^\top \mathbf{z}v) = S_0^0(v)\exp(-\boldsymbol{\theta}^{0^\top} \mathbf{z}v) \quad \mathbb{P}-\text{almost surely.}$$

This together with (A5) imply that there exist $\mathbf{z}_1 \neq \mathbf{z}_2$ such that for some $c > 0$,

$$S_{0*}(c)\exp(-\boldsymbol{\theta}_*^\top \mathbf{z}_1 c) = S_0^0(c)\exp(-\boldsymbol{\theta}^{0^\top} \mathbf{z}_1 c)$$

and

$$S_{0*}(c)\exp(-\boldsymbol{\theta}_*^\top \mathbf{z}_2 c) = S_0^0(c)\exp(-\boldsymbol{\theta}^{0^\top} \mathbf{z}_2 c).$$

Since $S_{0*}(c) > 0$ and $S_0^0(c) > 0$, this implies $(\boldsymbol{\theta}_* - \boldsymbol{\theta}^{0^\top})(\mathbf{z}_1 - \mathbf{z}_2) = 0$. Again, by (A5), the collection of such $\mathbf{z}_1$ and $\mathbf{z}_2$ has positive probability and there exist at least $p$ such pairs that constitute a full rank $p \times p$ matrix. It follows that $\boldsymbol{\theta}_* = \boldsymbol{\theta}^0$. This in turn implies that

$$S_{0*}(u) = S_0^0(u) \quad \text{and} \quad S_{0*}(v) = S_0^0(v) \quad G - \text{almost surely.}$$

Therefore, by bounded convergence theorem,

$$\int \left( (\hat{S}_{0(n)}(u) - S_0^0(u))^2 + (\hat{S}_{0(n)}(v) - S_0^0(v))^2 \right) \mathrm{d}G(u, v) \to 0. \tag{31}$$

Equations (30) and (31) together with $\boldsymbol{\theta}_* = \boldsymbol{\theta}^0$ hold for all $\omega \in \Omega_0$ with $\mathbf{P}(\Omega_0) = 1$. This completes the proof.

Table 1: Absolute bias, MSE and bootstrap coverage probability (BCP) of the estimator of $\theta$ under Method-1 and Method-2 with mild (10%) censoring

| $\lambda$ | $\theta$ | Method | $n=30$ | | | $n=50$ | | | $n=75$ | | |
|---|---|---|---|---|---|---|---|---|---|---|---|
| | | | Bias | MSE | BCP | Bias | MSE | BCP | Bias | MSE | BCP |
| 0.1 | 0.25 | 1 | 0.0033 | 0.0008 | 0.903 | 0.0061 | 0.0054 | 0.900 | 0.0091 | 0.0067 | 0.898 |
| | | 2 | 0.0347 | 0.0011 | 0.940 | 0.0380 | 0.0051 | 0.937 | 0.0396 | 0.0073 | 0.934 |
| 1.0 | 0.5 | 1 | 0.0104 | 0.0009 | 0.895 | 0.0134 | 0.0021 | 0.893 | 0.0163 | 0.0069 | 0.889 |
| | | 2 | 0.0373 | 0.0018 | 0.928 | 0.0405 | 0.0063 | 0.924 | 0.0454 | 0.0078 | 0.920 |
| 2.5 | -0.50 | 1 | 0.0077 | 0.0019 | 0.921 | 0.0089 | 0.0037 | 0.916 | 0.0108 | 0.0073 | 0.915 |
| | | 2 | 0.0247 | 0.0012 | 0.926 | 0.0259 | 0.0049 | 0.925 | 0.0307 | 0.0067 | 0.921 |
| 4.0 | -0.01 | 1 | 0.0336 | 0.0017 | 0.924 | 0.0366 | 0.0029 | 0.922 | 0.0410 | 0.0055 | 0.918 |
| | | 2 | 0.0448 | 0.0013 | 0.934 | 0.0484 | 0.0062 | 0.931 | 0.0507 | 0.0106 | 0.929 |

Table 2: Absolute bias, MSE and bootstrap coverage probability (BCP) of the estimator of $\theta$ under Method-1 and Method-2 with moderate (20%) censoring

| $\lambda$ | $\theta$ | Method | $n=30$ | | | $n=50$ | | | $n=75$ | | |
|---|---|---|---|---|---|---|---|---|---|---|---|
| | | | Bias | MSE | BCP | Bias | MSE | BCP | Bias | MSE | BCP |
| 0.1 | 0.25 | 1 | 0.0047 | 0.0024 | 0.901 | 0.0085 | 0.0095 | 0.897 | 0.0109 | 0.0114 | 0.893 |
| | | 2 | 0.0366 | 0.0021 | 0.938 | 0.0401 | 0.0063 | 0.936 | 0.0424 | 0.0111 | 0.930 |
| 1.0 | 0.5 | 1 | 0.0121 | 0.0027 | 0.894 | 0.0152 | 0.0075 | 0.891 | 0.0197 | 0.0088 | 0.886 |
| | | 2 | 0.0385 | 0.0028 | 0.927 | 0.0418 | 0.0102 | 0.922 | 0.0493 | 0.0126 | 0.918 |
| 2.5 | -0.5 | 1 | 0.0091 | 0.0031 | 0.919 | 0.0101 | 0.0052 | 0.914 | 0.0139 | 0.0114 | 0.910 |
| | | 2 | 0.0265 | 0.0025 | 0.925 | 0.0278 | 0.0078 | 0.923 | 0.0344 | 0.0115 | 0.917 |
| 4.0 | -0.01 | 1 | 0.0346 | 0.0036 | 0.923 | 0.0402 | 0.0061 | 0.918 | 0.0435 | 0.0077 | 0.916 |
| | | 2 | 0.0465 | 0.0032 | 0.932 | 0.0512 | 0.0083 | 0.927 | 0.053 | 0.0118 | 0.924 |

Table 3: Absolute bias, MSE and bootstrap coverage probability (BCP) of the estimator of $\theta$ under Method-1 and Method-2 with heavy (30%) censoring

| $\lambda$ | $\theta$ | Method | $n=30$ | | | $n=50$ | | | $n=75$ | | |
|---|---|---|---|---|---|---|---|---|---|---|---|
| | | | Bias | MSE | BCP | Bias | MSE | BCP | Bias | MSE | BCP |
| 0.1 | 0.25 | 1 | 0.0057 | 0.0042 | 0.900 | 0.0118 | 0.0107 | 0.894 | 0.0151 | 0.0129 | 0.889 |
| | | 2 | 0.0384 | 0.0031 | 0.937 | 0.0415 | 0.0079 | 0.934 | 0.0441 | 0.0147 | 0.925 |
| 1.0 | 0.5 | 1 | 0.0141 | 0.0041 | 0.892 | 0.0163 | 0.0121 | 0.889 | 0.0222 | 0.0143 | 0.882 |
| | | 2 | 0.0405 | 0.0044 | 0.925 | 0.0429 | 0.0139 | 0.919 | 0.0539 | 0.0174 | 0.916 |
| 2.5 | -0.5 | 1 | 0.0104 | 0.0050 | 0.918 | 0.0143 | 0.0097 | 0.909 | 0.0170 | 0.0151 | 0.904 |
| | | 2 | 0.0283 | 0.0038 | 0.923 | 0.0296 | 0.0090 | 0.918 | 0.0372 | 0.0156 | 0.915 |
| 4.0 | -0.01 | 1 | 0.0361 | 0.0056 | 0.921 | 0.0432 | 0.0088 | 0.916 | 0.0462 | 0.0101 | 0.913 |
| | | 2 | 0.0484 | 0.0044 | 0.931 | 0.0541 | 0.0096 | 0.925 | 0.0561 | 0.0131 | 0.921 |

## Remark 2.1

The asymptotic distributions of the estimators $\hat{\theta}_{(n)}$ and $\hat{S}_{0(n)}(t)$ do not seem to be easy to establish under the iterative method. We consider this as a problem for future research.

## Remark 2.2

A likelihood ratio test can be carried out to test the significance of regression coefficients. The null hypothesis $H_0 : \theta = 0$ can be tested against $H_1 : \theta \neq 0$, where 0 is the null vector of the same order, with the test statistic $-2\log\frac{L(0)}{L(\hat{\theta})}$, which follows $\chi^2_{(p)}$ distribution. The test results in rejecting the null hypothesis for small P-values.

## 3. Simulation Studies

A simulation study is carried out to assess the finite sample properties of the estimators. We consider the exponential distribution with mean $\lambda^{-1}$ as the distribution of lifetime variable $T$. Also, we choose independent exponential distributions with fixed means $\lambda_1^{-1}$ and $\lambda_2^{-1}$ as the distributions for the censoring random variate $U$ and the interval of censorship $V - U$ respectively, and these two distributions are assumed to be independent of $T$. We consider a single covariate $z$ in the present study, which is generated from uniform distribution over $[0, 10]$ and let $\theta$ be the corresponding regression coefficient. Under the AR model in (1), the survival function of $T$ given $z$ may be written as

$$S(t|z) = S_0(t) \exp(-\theta z t), \tag{32}$$

where $S_0(t) = \exp(-\lambda t)$. It can be observed that (32) is the survival function corresponding to an exponential variate with mean $(\lambda + \theta z)^{-1}$. A large number of observations are generated from (32) for fixed values of $\lambda$ and $\theta$. Now corresponding to each observation on $T$, a random censoring interval is generated from $(U, V)$, where the distribution parameters are fixed as $\lambda_1^{-1} = 20$ and $\lambda_2^{-1} = 10$. If we find $T \notin (U, V)$ then $T$ is selected in the sample, otherwise we choose the interval as the observation. As we generate large number of observations we can now choose a sample of required size $n$. We consider three different censoring rates: 10% (mild), 20% (moderate) and 30% (heavy) for our inference. The martingale-based inference procedure, denoted as Method-1, and iterative inference procedure, denoted as Method-2, which are described in Section 2, are employed to obtain the estimates of $S_0(t)$ and $\theta$ and using 1000 iterations for various choices of $\lambda$ and $\theta$. The absolute bias and mean squared error (MSE) are computed and are given in Table 1 to Table 3. Also in each case, a 95% bootstrap confidence interval for regression parameter is computed. The proportion of times the true parameter value lies in such intervals is called bootstrap coverage probabilities (BCP). They are also reported in Table 1 to Table 3. It is evident that both bias and MSE are small in each case and they decrease as the sample size increases. The bootstrap coverage probabilities are found fairly large, close to one. Further, as the censoring rate increases the bias and MSE increase, while the BCP decreases. Also, for each combination of parameter values, and with sample size 75, we shall find out a cubic polynomial estimate of the form $S_0(t) = c_0 + c_1 t + c_2 t^2 + c_3 t^3$ with each of its coefficients being the average of corresponding coefficients obtained for all the iterations, for the baseline survival function. These estimated survival curves corresponding to both methods are plotted in Figure 1 to Figure 3, where continuous curve represents the true baseline
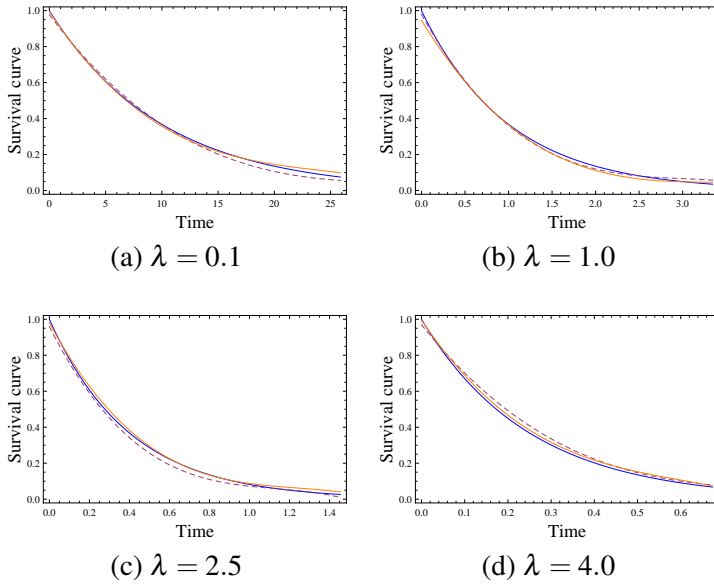
Figure 1: Plots of baseline survival curve and its estimates under Method-1 (dashed curve) and Method-2 (dotted curve) with mild (10%) censoring
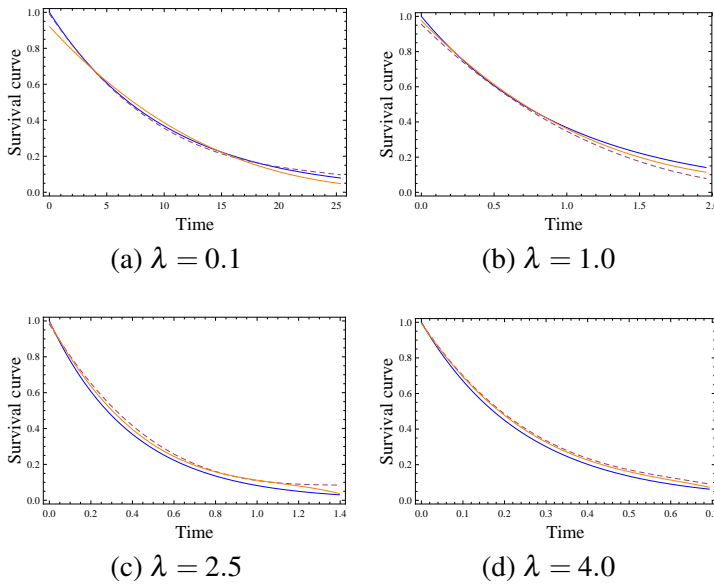


Figure 2: Plots of baseline survival curve and its estimates under Method-1 (dashed curve) and Method-2 (dotted curve) with moderate (20%) censoring
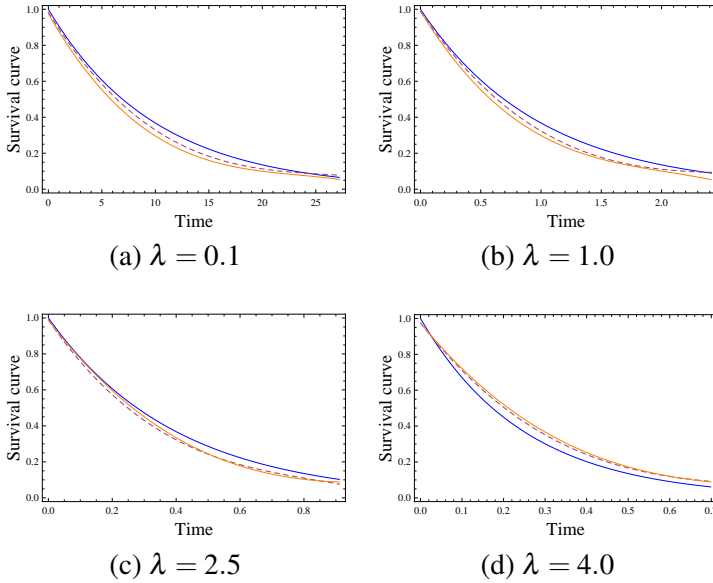
(a) $\lambda = 0.1$

(b) $\lambda = 1.0$

(c) $\lambda = 2.5$

(d) $\lambda = 4.0$

Figure 3: Plots of baseline survival curve and its estimates under Method-1 (dashed curve) and Method-2 (dotted curve) with heavy (30%) censoring

survival function and dashed curve represents the corresponding the estimate under Method-1 and dotted curve represents the corresponding estimate under Method-2. We see that both the estimated curves are close to the true curve.

## 4. Illustrative Data Analysis

The proposed methods are applied to a real life data studied by Ichida et al. (1993). The data deals with an evaluation of a protocol change in disinfectant practices in a medical center where patients are suffering from burn wounds. The control of infection is the major concern in burn management and the study aims at comparing two different controlling methods: routine bathing care method and body cleansing method. The time (in days) until a patient develops staphylococcus infection is considered as the lifetime variable. Although the original study involves several covariates, for the illustration purpose we consider two of them, namely treatment ($z_1$), which is coded as 1-for routine bathing and 2-for body cleansing, and percentage of total surface area burned ($z_2$). Let $\theta_1$ and $\theta_2$ respectively be the unknown regression coefficients. A random censoring interval $(U, V)$, where $U$ and $V - U$ are independent exponential variates with means $\lambda_1^{-1} = 20$ and $\lambda_2^{-1} = 10$ is generated first. Then, an individual from among all exact 48 lifetimes is selected at random and if lifetime of the patient happens to fall in the generated censoring interval, that

Table 4: Estimates of coefficients of survival curve and regression coefficients under Method-1 and Method-2.

| Method | $S_0(t)$ | | | | $\theta$ | |
|---|---|---|---|---|---|---|
| | $c_0$ | $c_1$ | $c_2$ | $c_3$ | $\theta_1$ | $\theta_2$ |
| 1 | 0.93665 | -0.04872 | 0.000635 | -9.223e-6 | 0.0112 | 0.1005 |
| 2 | 0.96574 | -0.05991 | 0.00121 | -9.256e-6 | 0.00895 | 0.1760 |

lifetime is assumed to have middle censored and that interval is considered as the corresponding observation. Otherwise the lifetime is maintained. This process is repeated until around 25% of the observations are censored. The data resulted consists of twelve censored observations. We apply the two methods of estimation given in Section 2 and obtained the estimates of the baseline survival function of the form $S_0 = c_0 + c_1 t + c_2 t^2 + c_3 t^3$ and the regression coefficient $\theta$. The estimated values, under both methods, of the coefficients of survival curves as well as regression coefficients are listed in Table 4. To test the significance of the covariate effect under the iterative method, we consider the null hypothesis $H_0 : \theta = 0$, where $\theta = (\theta_1, \theta_2)$ and 0 is null vector of the same order, and we use the likelihood ratio test described in Remark 2.2. The P-value of 0.008 indicates that the covariate effects are significant.

Now, we check the overall fit of the model by using Cox-Snell residuals (Cox & Snell, 1968). Suppose that the AR model given in (1) is fitted to the data. If the model assumption is correct then the probability integral transform of the true death time $T$ assumes a uniform distribution over $[0, 1]$ or equivalently the random variable $H(T_j | \mathbf{z}_j)$, which is the true cumulative hazard function corresponding to (1), has an exponential distribution with hazard rate 1. Then, the Cox-Snell residuals are defined to be the fitted cumulative hazard function values $\hat{r}_j = \hat{H}_0(t_j) + \mathbf{z}_j^\top \hat{\theta} t_j$ with the estimated parameters. If the model is reasonable and the estimates of the parameters are close to the true values, then these quantities should look like a censored sample from unit exponential distribution. To check whether the $r_j$'s behave as a sample from the unit exponential distribution we compute the Nelson-Aalen estimator of the cumulative hazard rate of $r_j$'s. If the unit exponential distribution fits the data, then this estimator should be approximately equal to the cumulative hazard rate of the unit exponential distribution. Thus, a plot of $r_j$'s versus their estimated cumulative hazard rates should be a straight line through origin and with a slope of 1. Figure 2 shows the plots so obtained under both the models. The curves are close to the straight line indicating AR assumption is reasonable.
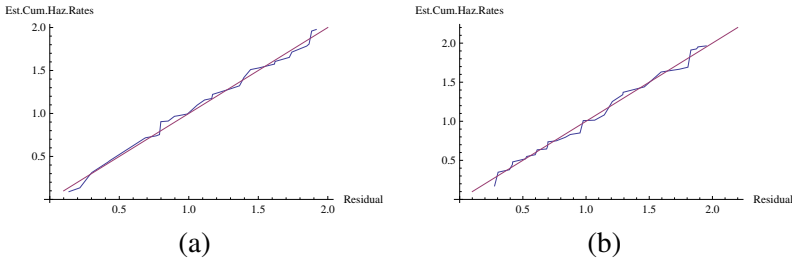
Figure 4: Plot of $r_j$'s against estimated cumulative hazard rates under (a) Method-1 and (b) Method-2.

## 5. Conclusion

The present study discussed the semiparametric regression problem for the analysis of middle-censored lifetime data. We considered two different methods of estimation, one making use of martingale-based theory and the other based on an iterative method for which a maximization procedure for finding the NPMLE is developed. Large sample properties including consistency and weak convergence of the estimators were established under the martingale-based method. Consistency of estimators was proved under the iterative method, whereas their weak convergence do not appear to be easy to establish, although one can perhaps extend the ideas used in (Huang & Wellner, 1995). Simulation studies showed that the inference procedures were efficient. The model was applied to a real data set. Although we considered time-fixed covariates in this work, the procedure can easily be extended to the case of time-varying covariates, as in the work of Lin & Ying (1994). The middle-censored data has a connection with mixed interval-censored (MIC) data (Yu et al., 2001). Although both sampling schemes differ in character, the observed data from MIC will reduce to data from middle-censoring, when there are no left censored or right censored observations. For a detailed discussion on this interrelationship one may refer to Shen (2011).

## Acknowledgements

# References

Aalen, O. O. (1980). A model for nonparametric regression analysis of counting processes. *Mathematical Statistics and Probability Theory*, 1–25.

Aalen, O. O. (1989). A linear regression model for the analysis of life times. *Statistics in Medicine*, 8(8), 907–925.

Andersen, P. K., Borgan, O., Gill, R. D., & Keiding, N. (2012). *Statistical Models based on Counting Processes*. Springer Science & Business Media.

Andersen, P. K. & Gill, R. D. (1982). Cox's regression model for counting processes: a large sample study. *The annals of statistics*, 1100–1120.

Aranda-Ordaz, F. J. (1983). An extension of the proportional-hazards model for grouped data. *Biometrics*, 109–117.

Bennett, N., Iyer, S. K., & Jammalamadaka, S. R. (2017). Analysis of gamma and Weibull lifetime data under a general censoring scheme and in the presence of covariates. *Communications in Statistics - Theory and Methods*, 46(5), 2277–2289.

Breslow, N. & Day, N. (1980). Conditional logistic regression for matched sets. *Statistical Methods in Cancer Research*, 1, 248–279.

Breslow, N. E. (1972). Discussion of paper of D. R. Cox. *Journal of Royal Statistical Society Series B*, 34, 216–7.

Breslow, N. E. & Day, N. E. (1987). *Statistical Methods in Cancer Research*, volume 2. International Agency for Research on Cancer, Lyon.

Cox, D. R. (1972). Regression models and life tables (with discussion). *Journal of the Royal Statistical Society*, 34, 187–220.

Cox, D. R. (1975). Partial likelihood. *Biometrika*, 62(2), 269–276.

Cox, D. R. & Oakes, D. (1984). *Analysis of Survival Data*, volume 21. CRC Press.

Cox, D. R. & Snell, E. J. (1968). A general definition of residuals. *Journal of the Royal Statistical Society. Series B (Methodological)*, 248–275.

Davarzani, N. & Parsian, A. (2011). Statistical inference for discrete middle-censored data. *Journal of Statistical Planning and Inference*, 141(4), 1455–1462.

Davarzani, N., Parsian, A., & Peeters, R. (2015). Statistical inference on middle-censored data in a dependent setup. *Journal of Statistical Theory and Practice*, *9*(3), 646–657.

Huang, J. & Wellner, J. A. (1995). Efficient estimation for the proportional hazards model with "case 2" interval censoring. *Technical Report No. 290, Department of Statistics, University of Washington, Seattle, USA*.

Ichida, J., Wassell, J., Keller, M., & Ayers, L. (1993). Evaluation of protocol change in burn-care management using the Cox proportional hazards model with time-dependent covariates. *Statistics in Medicine*, *12*(3-4), 301–310.

Iyer, S. K., Jammalamadaka, S. R., & Kundu, D. (2008). Analysis of middle-censored data with exponential lifetime distributions. *Journal of Statistical Planning and Inference*, *138*(11), 3550–3560.

Jammalamadaka, S. R. & Iyer, S. K. (2004). Approximate self consistency for middle-censored data. *Journal of Statistical Planning and Inference*, *124*(1), 75–86.

Jammalamadaka, S. R. & Leong, E. (2015). Analysis of discrete lifetime data under middle-censoring and in the presence of covariates. *Journal of Applied Statistics*, *42*(4), 905–913.

Jammalamadaka, S. R. & Mangalam, V. (2003). Nonparametric estimation for middle-censored data. *Journal of Nonparametric Statistics*, *15*(2), 253–265.

Jammalamadaka, S. R. & Mangalam, V. (2009). A general censoring scheme for circular data. *Statistical Methodology*, *6*(3), 280–289.

Jammalamadaka, S. R., Prasad, S. N., & Sankaran, P. G. (2016). A semi-parametric regression model for analysis of middle censored lifetime data. *Statistica*, *76*(1), 27.

Kalbfleisch, J. D. & Prentice, R. L. (2011). *The Statistical Analysis of Failure Time Data*, volume 360. John Wiley & Sons.

Klein, J. P. & Moeschberger, M. L. (2005). *Survival Analysis: Techniques for Censored and Truncated Data*. Springer Science & Business Media.

Lawless, J. F. (2011). *Statistical Models and Methods for Lifetime Data*, volume 362. John Wiley & Sons.

Lin, D. & Ying, Z. (1994). Semiparametric analysis of the additive risks model. *Biometrika*, *81*(1), 61–71.

Mangalam, V., Nair, G. M., & Zhao, Y. (2008). On computation of NPMLE for middle-censored data. *Statistics & Probability Letters*, *78*(12), 1452–1458.

Sankaran, P. G. & Prasad, S. (2014). Weibull regression model for analysis of middle-censored lifetime data. *Journal of Statistics and Management Systems*, *17*(5-6), 433–443.

Shen, P. (2010). An inverse-probability-weighted approach to the estimation of distribution function with middle-censored data. *Journal of Statistical Planning and Inference*, *140*(7), 1844–1851.

Shen, P. (2011). The nonparametric maximum likelihood estimator for middle-censored data. *Journal of Statistical Planning and Inference*, *141*(7), 2494–2499.

Sun, J. (2006). *The Statistical Analysis of Interval-censored Failure Time Data*. Springer Science & Business Media.

Tarpey, T. & Flury, B. (1996). Self-consistency: a fundamental concept in statistics. *Statistical Science*, 229–243.

Thomas, D. C. (1986). Use of auxiliary information in fitting nonproportional hazards models. *Modern Statistical Methods in Chronic Disease Epidemiology*, 197–210.

Wang, P., Tong, X., Zhao, S., & Sun, J. (2015). Regression analysis of left-truncated and case I interval-censored data with the additive hazards model. *Communications in Statistics - Theory and Methods*, *44*(8), 1537–1551.

Yu, Q., Wong, G. Y., & Li, L. (2001). Asymptotic properties of self-consistent estimators with mixed interval-censored data. *Annals of the Institute of Statistical Mathematics*, *53*(3), 469–486.