



Regression under Cox's model for recall-based time-to-event data in observational studies



Sedigheh Mirzaei Salehabadi*, Debasis Sengupta

Applied Statistical Unit, Indian Statistical Institute, Kolkata, 700108, India

ARTICLE INFO

Article history:

Received 3 March 2015

Received in revised form 17 June 2015

Accepted 5 July 2015

Available online 13 July 2015

Keywords:

Informative censoring

Interval censoring

Proportional hazards

Related risk regression model

Retrospective study

Turnbull estimator

ABSTRACT

In some retrospective observational studies, the subject is asked to recall the age at a particular landmark event. The resulting data may be partially incomplete because of the inability of the subject to recall. This type of incompleteness may be regarded as interval censoring, where the censoring is likely to be informative. The problem of fitting Cox's relative risk regression model to such data is considered. While a partial likelihood is not available, a method of semi-parametric inference of the regression parameters as well as the baseline distribution is proposed. Monte Carlo simulations show reasonable performance of the regression parameters, compared to Cox estimators of the same parameters computed from the complete version of the data. The proposed method is illustrated through the analysis of data on age at menarche from an anthropometric study of adolescent and young adult females in Kolkata, India.

© 2015 Elsevier B.V. All rights reserved.

1. Introduction

Retrospective studies on a landmark event can produce dichotomous data on the current status of an individual (whether or not the event has occurred till the day of observation). From the perspective of the time to event, these data can be regarded as left or right censored. In some retrospective studies, the subject is asked to recall the time of the landmark event, in case it has already taken place. Such retrospective data can be incomplete because of the possibility that the time is forgotten. Sometimes the subject may be able to specify only a range for the time-to-event. For some other subjects, the event may be found not to have happened till the time of visit. Thus, data arising from this kind of retrospective studies are interval-censored. However, the chance of recall may depend on the time span between the occurrence of the event and the time of interview. For instance, between two young adult females interviewed at the same age, the one having experienced menarche more recently may have a higher chance of recalling the date. Thus, the censoring mechanism in this set-up is likely to be informative. Mirzaei et al. (2015) and Mirzaei and Sengupta (2015) have proposed parametric and nonparametric methods of likelihood based inference, when the data are subjected to informative interval censoring of this type. They have shown theoretically as well as through simulation that estimators of survival function that ignore the informative nature of censoring can have large bias even when the sample size is large (Mirzaei et al., 2015). On the other hand, the problem of regression with the above type of censored data has not been addressed yet.

The relative risk regression model, also known as the proportional hazards model, is widely used in the analysis of event time data with covariates. The method of analysis proposed by Cox (1972) can accommodate right-censored data which are usual in survival problems, and left-truncated data which arise when there are delayed entries in a cohort

* Corresponding author.

E-mail addresses: sedigheh_r@isical.ac.in (S. Mirzaei Salehabadi), sdebasis@isical.ac.in (D. Sengupta).

(Breslow et al., 1983). Other models which are used for more complex observation schemes include the accelerated failure time (AFT) model (Wei, 1992), the additive hazard regression model (Klein and Moeschberger, 2003), proportional odds ratio model (Dabrowska and Doksum, 1988) and so on (Vonta, 1996). There has also been some work on more general regression models for survival data, such as single index regression models (Chaudhuri, 2007) and models with random effect/frailty (Wienke, 2010).

Discrete-time regression models for right-truncated data have been developed and applied in the analysis of AIDS incidence and induction time distributions (Kalbfleisch and Lawless, 1991; Gross and Huber-Carol, 1992). Finkelstein (1986) and Finkelstein et al. (1993) discussed methods for fitting a discrete proportional hazards model for the case where the data are either interval-censored or right-truncated. In both cases, a score test was developed for testing the hypothesis of a zero regression coefficient. Tu et al. (1993) discussed a discrete proportional hazards model and an associated EM algorithm for data that are censored as well as truncated. Alioum and Commenges (1996) discussed a method for fitting a relative risk regression model for arbitrarily interval censored data. Their method assumes the censoring to be non-informative.

DeMasi et al. (1997) and Tanaka and Rao (2005) considered the regression problem for informatively censored data. Their model treats informative censoring as a type of risk in a competing risks setup, where the subject experiences two types of mutually exclusive events. This set-up is not meant to model the informative censoring found in recall data.

In this paper, we consider regression under Cox's model for the special type of informatively censored data arising from uncertainly recalled event time in a retrospective study. In Section 2, we develop a semiparametric maximum likelihood estimator of the regression coefficients under the model. In Section 3, we report results of a simulation study of the performance of the proposed maximum likelihood estimator. Section 4 illustrates this method with data on menarcheal age of adolescent and young adult females, collected during the course of a project undertaken by the Indian Statistical Institute, Kolkata. Proofs of all the results are given in Appendix.

2. Model and inference

2.1. Model

Consider a subject having time of occurrence of the landmark event T_i , which is a single sample from a distribution F_i with density f_i and support $[t_{\min}, t_{\max}]$, for $i = 1, \dots, n$. Let these subjects be interviewed at times $S_1, \dots, S_n \in [t_{\min}, t_{\max}]$, respectively. Suppose U_i is the unobservable time that the i th subject would take to forget the epoch of his/her landmark event. For the sake of simplicity, we assume that when the subject forgets the epoch, there is no recollection of an approximate range of time also. There are observable indicators δ_i and ε_i of the events $T_i \leq S_i$ and $U_i > S_i - T_i \geq 0$, respectively. We assume that U_1, \dots, U_n are samples from a distribution with distribution function π , and that these are independent of both T_i and S_i . It follows that, given S_i and T_i , the non-recall probability depends on the time elapsed since the landmark event as

$$P(\varepsilon_i = 0 | T_i = t, S_i = s) = \pi(s - t), \quad s > t. \quad (1)$$

According to this model, the likelihood, conditional on the ages at interview, is

$$\prod_{i=1}^n [\bar{F}_i(S_i)]^{1-\delta_i} \left[\{f_i(T_i)(1 - \pi(S_i - T_i))\}^{\varepsilon_i} \left(\int_0^{S_i} f_i(u)\pi(S_i - u)du \right)^{1-\varepsilon_i} \right]^{\delta_i}. \quad (2)$$

Here the informativeness of the censoring mechanism is captured through the function π . If π is a constant, then the likelihood (2) becomes a multiple of the likelihood for non-informatively left- or right censored data with possibility of no censoring. As a further special case, if $\pi = 1$, then the likelihood (2) simplifies to the likelihood for dichotomous data. If $\pi = 0$, i.e., there is perfect recall with probability 1, then $\varepsilon_i = 1$ for all i such that $\delta_i = 1$, and the likelihood reduces to that for right-censored data.

Let Z_i be the r -dimensional vector of covariates, assumed to be independent of both S_i and U_i . Note that the distribution of T_i would depend on Z_i . Under Cox's relative risk regression model, the probability of the individual i , with covariate vector Z_i , having the event after time t is

$$\bar{F}_i(t) = [\bar{F}_0(t)]^{\exp(\beta^T Z_i)}, \quad (3)$$

where \bar{F}_0 is the baseline survival function, assumed to have a density.

2.2. Identifiability

Before embarking on developing a method of estimation, we need to check the identifiability of β , F_0 and π . By substituting (3) in the likelihood (2), after dropping the subscript i for simplicity and following Theorem 1 of Mirzaei et al. (2015), one can show that a typical factor in the product likelihood is equal to the conditional density of the observable vector (V, δ) , given S and Z , where $V = (S - T)\varepsilon$. The conditional density is written alternatively as

$$h(v, \delta | s, z; \beta) = \begin{cases} \bar{F}_0(s)^{\exp(\beta^T z)} & \text{if } v = 0 \text{ and } \delta = 0, \\ \int_0^s -\frac{d}{du} \left(\bar{F}_0(u)^{\exp(\beta^T z)} \right) \pi(s - u) du & \text{if } v = 0 \text{ and } \delta = 1, \\ -\frac{d}{dv} \left(\bar{F}_0(s - v)^{\exp(\beta^T z)} \right) (1 - \pi(v)) & \text{if } v > 0 \text{ and } \delta = 1, \\ 0 & \text{otherwise.} \end{cases} \tag{4}$$

Among the unknown parameters β , F_0 and π , the interest lies mainly in β , and possibly in F_0 . We address the question as to whether β , F_0 and π are identifiable from h , in the next theorem.

Theorem 1. *Suppose, for any number τ in the support of F_0 , the model (4) holds for some $s > \tau$, and for $z = 0$ as well as for any r linearly independent values of the vector z , where r is the dimension of z . Then the parameters β , F_0 and π are identifiable from h under this model.*

We now proceed with the estimation problem, after assuming that the distribution of the covariate vector ensures the identifiability of all the unknown parameters.

2.3. Piecewise constant non-recall probability

The integral contained in the likelihood (2) makes it difficult to maximize. For the sake of mathematical tractability, we now assume that π is a piecewise constant function of the form

$$\pi(x) = b_1 I(x_1 < x \leq x_2) + b_2 I(x_2 < x \leq x_3) + \dots + b_k I(x_k < x < \infty), \tag{5}$$

where the jump points satisfy the relations $0 = x_1 < x_2 < \dots < x_k$, and the coefficients satisfy the relations $0 < b_1 < b_2 < \dots < b_k \leq 1$, so that π is a non-decreasing function.

In view of (5), the likelihood (2) reduces to

$$L = \prod_{i=1}^n [\bar{F}_i(S_i)]^{1-\delta_i} \left[\left\{ f_i(T_i) \left(1 - \sum_{l=1}^k b_l I(W_{l+1}(S_i) < T_i \leq W_l(S_i)) \right) \right\}^{\varepsilon_i} \cdot \left[\sum_{l=1}^k b_l (F_i(W_l(S_i)) - F_i(W_{l+1}(S_i))) \right]^{1-\varepsilon_i} \right]^{\delta_i}, \tag{6}$$

where $W_l(S_i) = (S_i - x_l) \vee t_{\min}$ for $l = 1, \dots, k$ and $W_{k+1}(S_i) = t_{\min}$, $i = 1, 2, \dots, n$. Note that

$$t_{\min} = W_{k+1}(S_i) \leq W_k(S_i) \leq W_{k-1}(S_i) \leq \dots \leq W_1(S_i). \tag{7}$$

Depending on the value of S_i , some of the above inequalities may in fact be equalities. Specifically, if l is an index such that $S_i - x_{l+1} \leq t_{\min} < S_i - x_l$ then $t_{\min} = W_{k+1}(S_i) = \dots = W_{l+1}(S_i)$. The remaining inequalities would be strict.

The likelihood (6) can be rewritten as

$$L = \prod_{i=1}^n [\bar{F}_0(S_i)^{\exp(\beta^T Z_i)}]^{1-\delta_i} \left[\left\{ \left(\bar{F}_0(T_i)^{\exp(\beta^T Z_i)} - \bar{F}_0(T_i)^{\exp(\beta^T Z_i)} \right) \cdot \left(1 - \sum_{l=1}^k b_l I(W_{l+1}(S_i) < T_i \leq W_l(S_i)) \right) \right\}^{\varepsilon_i} \cdot \left[\sum_{l=1}^k b_l \left(\bar{F}_0(W_{l+1}(S_i))^{\exp(\beta^T Z_i)} - \bar{F}_0(W_l(S_i))^{\exp(\beta^T Z_i)} \right) \right]^{1-\varepsilon_i} \right]^{\delta_i}, \tag{8}$$

which does not involve any integration, and is mathematically more tractable than (2).

2.4. Maximum likelihood estimation

We assume that the parameters k and x_1, x_2, \dots, x_k of the function π are known. The likelihood (8) involves the unspecified baseline survival function \bar{F}_0 , apart from other unknown parameters β and $\eta = (b_1, b_2, \dots, b_k)^T$. Since the factors in the product likelihood are probabilities of various intervals, it is clear that estimation of \bar{F}_0 would amount to assignment of probabilities to various intervals that would act as basic units. If an appropriate set of intervals (including single points that can be regarded as degenerate intervals) are identified, then the values of \bar{F}_0 at the requisite points can be expressed as sums of probabilities of these intervals. Therefore, we need a systematic identification of these intervals.

The likelihood (8) involves probabilities assigned to intervals of the type $[t, t_{\max}]$ and $(t, t_{\max}]$, as per the baseline probability distribution \bar{F}_0 . Since these intervals have overlap, we express them as unions of smaller, disjoint intervals. Let $\mathcal{I}_1, \mathcal{I}_2$ and \mathcal{I}_3 be sets of indices i (between 1 and n) that satisfy the conditions $\delta_i = 0, \delta_i \varepsilon_i = 1$ and $\delta_i(1 - \varepsilon_i) = 1$, respectively. The set \mathcal{I}_1 contains indices of subjects for whom the event is yet to happen till the time of observation, \mathcal{I}_2 contains indices

of subjects who have experienced the event and remember the date of occurrence, while \mathcal{I}_3 is the set of indices of subjects who have experienced the event but forgotten the date. Consider the intervals

$$\begin{aligned} A_i &= (S_i, t_{\max}] \quad \text{for } i \in \mathcal{I}_1; \\ A_i &= [T_i, t_{\max}] \quad \text{for } i \in \mathcal{I}_2; \\ A'_i &= (T_i, t_{\max}] \quad \text{for } i \in \mathcal{I}_2; \\ A_{il} &= \begin{cases} (W_l(S_i), t_{\max}], & l = 1, \dots, k, \\ [W_l(S_i), t_{\max}], & l = k + 1, \end{cases} \quad \text{for } i \in \mathcal{I}_2 \cup \mathcal{I}_3 \end{aligned} \tag{9}$$

and the sets

$$\begin{aligned} \mathcal{A}_1 &= \{A_i : i \in \mathcal{I}_1\}; \\ \mathcal{A}_2 &= \{A_i \setminus A'_i : i \in \mathcal{I}_2\}; \\ \mathcal{A}_3 &= \{A'_i : i \in \mathcal{I}_2\}; \\ \mathcal{A}_4 &= \{A_{i(l+1)} \setminus A_{il} : 1 \leq l \leq k \text{ and } i \in \mathcal{I}_3\}. \end{aligned} \tag{10}$$

As the baseline distribution is absolutely continuous, the elements of \mathcal{A}_2 and \mathcal{A}_3 are all distinct with probability 1. Let n_2 be the number of elements of \mathcal{I}_2 . The elements of \mathcal{A}_2 are singletons; we arrange them in increasing order, and denote them as B_1, B_2, \dots, B_{n_2} . We also arrange the elements of \mathcal{A}_3 in the corresponding order and denote them as $B_{n_2+1}, B_{n_2+2}, \dots, B_{2n_2}$. We then collect the unique elements of $\mathcal{A}_1 \cup \mathcal{A}_4$ that are distinct from $B_1, B_2, \dots, B_{2n_2}$, and denote them as $B_{2n_2+1}, B_{2n_2+2}, \dots, B_M$. Observe that the collection B_1, B_2, \dots, B_M consist of the distinct elements of $\mathcal{A}_1 \cup \mathcal{A}_2 \cup \mathcal{A}_3 \cup \mathcal{A}_4$, arranged in a particular order. Denote the non-empty subsets of the index set $\{1, 2, \dots, M\}$ by $s_1, s_2, \dots, s_{2^M-1}$. Define

$$I_r = \left\{ \bigcap_{i \in s_r} B_i \right\} \cap \left\{ \bigcap_{i \notin s_r} B_i^c \right\} \quad \text{for } r = 1, 2, \dots, 2^M - 1. \tag{11}$$

Some of the I_r 's may be empty sets, denoted here by ϕ . Let

$$\mathcal{C} = \{s_r : I_r \neq \phi, 1 \leq r \leq 2^M - 1\}. \tag{12}$$

It can be verified that all the non-empty I_r 's are distinct and disjoint. Let \mathcal{A} be the class of all sets I_r such that $s_r \in \mathcal{C}$.

Note that each of the intervals B_1, \dots, B_M is a union of disjoint sets that are members of \mathcal{A} . For any Borel set A , suppose $P_0(A)$ is the probability assigned to A as per the baseline probability distribution corresponding to the survival function \bar{F}_0 . Let $p_r = P_0(I_r)$, for $I_r \in \mathcal{A}$. Then the likelihood (8) reduces to

$$\begin{aligned} L &= \prod_{i \in \mathcal{I}_1} \left(\sum_{\substack{r: I_r \subseteq A_i \\ s_r \in \mathcal{C}}} p_r \right)^{\exp(\beta^T Z_i)} \times \prod_{i \in \mathcal{I}_2} \left(1 - \sum_{l=1}^k b_l I(T_i \in A_{i(l+1)} \setminus A_{il}) \right) \cdot \left[\left(\sum_{\substack{r: I_r \subseteq A_i \\ s_r \in \mathcal{C}}} p_r \right)^{\exp(\beta^T Z_i)} - \left(\sum_{\substack{r: I_r \subseteq A'_i \\ s_r \in \mathcal{C}}} p_r \right)^{\exp(\beta^T Z_i)} \right] \\ &\times \prod_{i \in \mathcal{I}_3} \left[\sum_{l=1}^k b_l \left\{ \left(\sum_{\substack{r: I_r \subseteq A_{i(l+1)}} \\ s_r \in \mathcal{C}}} p_r \right)^{\exp(\beta^T Z_i)} - \left(\sum_{\substack{r: I_r \subseteq A_{il}} \\ s_r \in \mathcal{C}}} p_r \right)^{\exp(\beta^T Z_i)} \right\} \right]. \end{aligned} \tag{13}$$

Thus, maximizing the likelihood (8) is equivalent to maximizing the likelihood (13) with respect to β, η and the set of p_r 's with $s_r \in \mathcal{C}$. The p_r 's are nuisance parameters when the main objective is to estimate β . The number of these parameters can be very high. This problem is simplified if it can be shown algebraically that some of the estimated p_r 's are zero. With this goal, we consider the following subsets of \mathcal{C} .

$$\begin{aligned} \mathcal{C}_1 &= \{s : s \in \mathcal{C}; \text{ there is another element } s' \in \mathcal{C}, \text{ such that } s \subset s'\}, \\ \mathcal{C}_2 &= \{s : s \in \mathcal{C}; \text{ there is another element } s' \in \mathcal{C}, \text{ such that} \\ &\quad s' \setminus (s \cap s') \text{ consists of a singleton } j \text{ and } s \setminus (s \cap s') = \{j + n_2\}\}, \\ \mathcal{C}_0 &= \mathcal{C} \setminus (\mathcal{C}_1 \cup \mathcal{C}_2). \end{aligned} \tag{14}$$

The next result shows that the maximization of the likelihood can be restricted to \mathcal{C}_0 .

Theorem 2. For fixed values of β and η , maximizing the likelihood (13) with respect to p_r for $s_r \in \mathcal{C}$ is almost surely equivalent to maximizing it with respect to p_r for $s_r \in \mathcal{C}_0$, i.e.,

$$\max_{p_r: p_r \in [0, 1], \sum_{s_r \in \mathcal{C}} p_r = 1} L = \max_{p_r: p_r \in [0, 1], \sum_{s_r \in \mathcal{C}_0} p_r = 1} L.$$

Let us relabel the intervals I_j , $s_j \in \mathcal{C}_0$, by J_1, J_2, \dots, J_v . Further, let $q_j = P(J_j)$ for $j = 1, 2, \dots, v$. **Theorem 2** implies that maximizing the likelihood (13) is almost surely equivalent to maximizing

$$\begin{aligned}
 L(q_1, \dots, q_v, \eta, \beta) &= \prod_{i \in \mathcal{I}_1} \left(\sum_{j: J_j \subseteq A_i} q_j \right)^{\exp(\beta^T Z_i)} \times \prod_{i \in \mathcal{I}_2} \left(1 - \sum_{l=1}^k b_l I(T_i \in A_{i(l+1)} \setminus A_{il}) \right) \\
 &\cdot \left[\left(\sum_{j: J_j \subseteq A_i} q_j \right)^{\exp(\beta^T Z_i)} - \left(\sum_{j: J_j \subseteq A'_i} q_j \right)^{\exp(\beta^T Z_i)} \right] \\
 &\times \prod_{i \in \mathcal{I}_3} \left[\sum_{l=1}^k b_l \left\{ \left(\sum_{j: J_j \subseteq A_{i(l+1)}} q_j \right)^{\exp(\beta^T Z_i)} - \left(\sum_{j: J_j \subseteq A_{il}} q_j \right)^{\exp(\beta^T Z_i)} \right\} \right]. \tag{15}
 \end{aligned}$$

with respect to $q_1, q_2, \dots, q_v, \eta$ and β , subject to the restriction $\sum_{j=1}^v q_j = 1$.

In order to maximize the likelihood (15), we need to identify the sets J_j , $j = 1, \dots, v$, that is, the intervals I_j , $s_j \in \mathcal{C}_0$, defined through (11) and (14). This identification involves elaborate combinatorial calculations. In fact, simulations reported in Mirzaei and Sengupta (2015) show (in the case of nonparametric estimation in the absence of covariates) that these calculations consume much more computational time than the actual maximization. They have shown that the set of times of exact recall can serve as a readily available and approximate support of the estimated baseline distribution, so that the computational speed can be enhanced several times without sacrificing the quality of the solution substantially. In the next section, we prove a similar result for the regression problem.

2.5. Approximate MLE

Let $\mathcal{A}_0 = \{J_1, J_2, \dots, J_v\}$, and $\mathcal{A}_2 = \{\{T_i\}, i \in \mathcal{I}_2\}$ as already defined in (10). Further, let n_i be the cardinality of \mathcal{I}_i , $i = 1, 2, 3$. The task of maximizing the likelihood (15) can be simpler for large n_2 , as the following result shows.

Theorem 3. *The set \mathcal{A}_2 is contained in the set \mathcal{A}_0 almost surely. Further, if the inspection times take values from a finite set and the range of values of $\beta^T Z_i$ in (15) is bounded, then the probability of \mathcal{A}_0 being equal to \mathcal{A}_2 goes to one as $n_2 \rightarrow \infty$.*

One can form a computationally simpler estimator on the basis of Theorem 3. According to this theorem, the maximum likelihood estimator has mass only at points of exact recall of the event, when n_2 is large. In such a case, the likelihood (15) involves J_j 's that are singletons only. Therefore, irrespective of the value of n_2 , one can maximize (15) with respect to point masses corresponding to the times of exact recall.

Formally, let t_1, \dots, t_{n_2} be the ordered set of distinct ages at event that have been exactly recalled, and $q_1^*, \dots, q_{n_2}^*$ be the probability masses allocated to them. Maximizing the likelihood (15), subject to the constraint that $q_j = 0$ whenever $J_j \notin \mathcal{A}_2$, is equivalent to maximizing the following approximate likelihood subject to $\sum_{j=1}^{n_2} q_j^* = 1$ and $q_j^* \geq 0$:

$$\begin{aligned}
 L_a(q_1^*, \dots, q_{n_2}^*, \eta, \beta) &= \prod_{i \in \mathcal{I}_1} \left(\sum_{j: t_j \geq t_{m_i}} q_j^* \right)^{\exp(\beta^T Z_i)} \prod_{i \in \mathcal{I}_2} \left(1 - \sum_{l=1}^k b_l I(T_i \in A_{i(l+1)} \setminus A_{il}) \right) \cdot \left[\left(\sum_{j: t_j \geq t_{m_i}} q_j^* \right)^{\exp(\beta^T Z_i)} - \left(\sum_{j: t_j > t_{m_i}} q_j^* \right)^{\exp(\beta^T Z_i)} \right] \\
 &\times \prod_{i \in \mathcal{I}_3} \left[\sum_{l=1}^k b_l \left\{ \left(\sum_{j: t_j \geq t_{m_{i(l+1)}}} q_j^* \right)^{\exp(\beta^T Z_i)} - \left(\sum_{j: t_j \geq t_{m_{il}}} q_j^* \right)^{\exp(\beta^T Z_i)} \right\} \right] \tag{16}
 \end{aligned}$$

where $m_i = \inf\{j : t_j \in A_i\}$ for $i \in \mathcal{I}_1 \cup \mathcal{I}_2$, and, $m_{il} = \inf\{j : t_j \in A_{il}\}$, $l = 1, 2, \dots, k$ for $i \in \mathcal{I}_3$.

In order to remove the range restriction on the parameters $q_1^*, \dots, q_{n_2}^*$, we use the reparametrization $\gamma_d = \log(-\log(\sum_{j: t_j \geq t_d} q_j^*))$, $d = 1, 2, \dots, n_2$. Thus, the approximate likelihood (16) can be expressed as

$$\begin{aligned}
 L_a(\gamma, \eta, \beta) &= \prod_{i \in \mathcal{I}_1} \left(e^{-e^{Z_i \beta + \gamma_{m_i}}} \right) \times \prod_{i \in \mathcal{I}_2} \left(1 - \sum_{l=1}^k b_l I(T_i \in A_{i(l+1)} \setminus A_{il}) \right) \cdot \left[\left(e^{-e^{Z_i \beta + \gamma_{m_i}}} \right) - \left(e^{-e^{Z_i \beta + \gamma_{m_i+1}}} \right) \right] \\
 &\times \prod_{i \in \mathcal{I}_3} \left[\sum_{l=1}^k b_l \left\{ \left(e^{-e^{Z_i \beta + \gamma_{m_{i(l+1)}}}} \right) - \left(e^{-e^{Z_i \beta + \gamma_{m_{il}}} } \right) \right\} \right] \tag{17}
 \end{aligned}$$

where $\gamma = (\gamma_1, \dots, \gamma_{n_2+1})$ and $\gamma_{n_2+1} = \infty$. The log of the above expression simplifies to

$$\ell_a(\gamma, \eta, \beta) = \sum_{i=1}^n \log \left[\sum_{j=1}^{n_2} \alpha_{ij} \left[e^{(-e^{z_i\beta+\gamma_j})} - e^{(-e^{z_i\beta+\gamma_{j+1}})} \right] \right] \tag{18}$$

where, for $j = 1, 2, \dots, n_2$, we have

$$\alpha_{ij} = \begin{cases} I(J_j \subseteq A_i) & \text{if } i \in \mathcal{I}_1, \\ 1 - \sum_{l=1}^k b_l \cdot I(T_i \in A_{i(l+1)} \setminus A_{il}) \cdot I(J_j \subseteq A_i \setminus A'_i) & \text{if } i \in \mathcal{I}_2, \\ \sum_{l=1}^k b_l \cdot I(J_j \subseteq A_{i(l+1)} \setminus A_{il}) & \text{if } i \in \mathcal{I}_3. \end{cases} \tag{19}$$

We obtain the approximate maximum likelihood estimate (AMLE) of the parameters γ , η and β by maximizing the above approximate log-likelihood. The first and second derivatives of $\ell_a(\gamma, \eta, \beta)$ with respect to γ and β are as given below.

$$\begin{aligned} \frac{\partial \ell_a(\gamma, \eta, \beta)}{\partial b_l} &= \sum_{i=1}^n \frac{\sum_{j=1}^{n_2} g_{ij}}{\sum_t \alpha_{it} g_{it}} \xi_{ijl} \\ \frac{\partial \ell_a(\gamma, \eta, \beta)}{\partial \gamma_j} &= \sum_{i=1}^n \mu_{ij} h_{ij} \quad \text{for } j = 1, \dots, n_2, \\ \frac{\partial \ell_a(\gamma, \eta, \beta)}{\partial \beta_k} &= \frac{\sum_{i=1}^n \sum_{j=1}^{n_2} \alpha_{ij} [h_{ij} - h_{ij+1}] x_{ik}}{\sum_t \alpha_{it} g_{it}} \quad \text{for } k = 1, 2, \dots, r, \end{aligned} \tag{20}$$

where $g_{ij} = [e^{(-e^{z_i\beta+\gamma_j})} - e^{(-e^{z_i\beta+\gamma_{j+1}})}]$, $\mu_{ij} = (\alpha_{ij-1} - \alpha_{ij}) / \sum_t \alpha_{it} g_{it}$,

$$\xi_{ijl} = -I(T_i \in A_{i(l+1)} \setminus A_{il}) \cdot I(J_j \subseteq A_i \setminus A'_i) I(i \in \mathcal{I}_2) + I(J_j \subseteq A_{i(l+1)} \setminus A_{il}) \cdot I(i \in \mathcal{I}_3),$$

and

$$h_{ij} = (e^{-e^{z_i\beta+\gamma_j}}) (-e^{z_i\beta+\gamma_j}) \quad \text{for } j = 1, 2, \dots, n_2.$$

$$\begin{aligned} -\partial^2 \ell_a / \partial \gamma_j \partial \gamma_k &= \begin{cases} \sum_{i=1}^n \mu_{ij} \mu_{ik} h_{ij} h_{ik} & \text{if } j \neq k, \\ \sum_{i=1}^n -\mu_{ij} d_{ij} + (\mu_{ij} h_{ij})^2 & \text{if } j = k, \end{cases} \\ -\partial^2 \ell_a / \partial \beta_k \partial \beta_l &= -\sum_{i=1}^n x_{ik} x_{il} \left\{ \frac{\sum_j \alpha_{ij} (d_{ij} - d_{ij+1})}{\sum_t \alpha_{it} g_{it}} - \left[\frac{\sum_j \alpha_{ij} (h_{ij} - h_{ij+1})}{\sum_t \alpha_{it} g_{it}} \right]^2 \right\}, \\ -\partial^2 \ell_a / \partial \beta_k \partial \gamma_j &= \sum_{i=1}^n -x_{ik} \left[\mu_{ij} d_{ij} - \mu_{ij} h_{ij} \frac{\sum_r \alpha_{ir} (h_{ir} - h_{ir+1})}{\sum_t \alpha_{it} g_{it}} \right], \\ -\partial^2 \ell_a / \partial b_r \partial b_l &= \sum_{i=1}^n \left(\frac{\sum_{j=1}^{n_2} g_{ij}}{\sum_t \alpha_{it} g_{it}} \right)^2 \xi_{ijr}, \\ -\partial^2 \ell_a / \partial b_l \partial \gamma_j &= \sum_{i=1}^n \frac{[\xi_{i(j-1)l} - \xi_{ijl}] \sum_t \alpha_{it} g_{it} - [\alpha_{ij-1} - \alpha_{ij}] \sum_t g_{it} \xi_{ijl}}{\left(\sum_t \alpha_{it} g_{it} \right)^2}, \\ -\partial^2 \ell_a / \partial b_l \partial \beta_k &= \sum_{i=1}^n \frac{\sum_{j=1}^{n_2} [h_{ij} - h_{ij+1}] x_{ik} \sum_t \alpha_{it} g_{it} - \sum_{j=1}^{n_2} \alpha_{ij} [h_{ij} - h_{ij+1}] x_{ik} \sum_t g_{it}}{\left(\sum_t \alpha_{it} g_{it} \right)^2} \cdot \xi_{ijl}, \end{aligned}$$

where $d_{ij} = [e^{-e^{z_i\beta+\gamma_j}}](-e^{z_i\beta+\gamma_j}) + [e^{-e^{-z_i\beta+\gamma_j}}](-e^{-z_i\beta+\gamma_j})^2$ for $j = 1, 2, \dots, n_2$. A Newton–Raphson iteration can be used to compute the AMLEs $\hat{\gamma}, \hat{\eta}, \hat{\beta}$. The corresponding AMLE of the baseline distribution function \hat{F}_0 is

$$\hat{F}_0(t) = \sum_{j:t_j \geq t} \hat{q}_j^* \tag{21}$$

3. Simulation study of small sample performance

For the purpose of simulation, we generate samples of time-to-event from a relative risk regression model with survival function $\bar{F}_i(t) = [\bar{F}_0(t)]^{\exp(\beta^T Z_i)}$, where the baseline distribution function $F_0(t)$ is Weibull with shape and scale parameters $\alpha = 11$ and $\beta = 13$, respectively, and discard the samples lying outside the interval [8, 16]. This truncated distribution has median 11.57. The vector of covariates, $Z = (Z_1, Z_2)$, consists of a binary variable, taking values 1 and 0 with probabilities 0.25 and 0.75, and a continuous variable having the uniform distribution over the interval [0, 5]. We choose the vector of regression coefficients as $\beta = (\beta_1, \beta_2) = (1.5, 1.5)$. The ‘time of interview’ is generated from the discrete uniform distribution over the set of integers {7, 8, . . . , 21}. These choices are in line with the data analytic example of the next section, where the time to landmark event is the age at menarche in years. As for the forgetting probability π , we use (5) with $k = 7, x_1 = 0, x_2 = 1.7, x_3 = 3.4, x_4 = 5.1, x_5 = 6.8, x_6 = 8.5$ and $x_7 = 10.2$ and the vector parameter $\eta = (b_1, b_2, \dots, b_7) = (0.01, 0.15, 0.15, 0.15, 0.15, 0.15, 0.15)$.

The approximate log-likelihood (18) is maximized with respect to γ, β and η , subject to the constraint that the elements of η are nonnegative and in non-decreasing order.

As a benchmark of performance, one can consider the hypothetical situation when all the event times are perfectly recalled, that is, the data are right censored. In this case, one can use the MLE obtained by maximizing Cox’s partial likelihood. We refer to this estimator based on ‘complete recall’ data as the ‘complete recall MLE’. On the other hand, if one uses only the ‘current status’ information, namely whether the event of interest has happened till the time of interview, then the corresponding likelihood is

$$\prod_{i=1}^n [\bar{F}_0^{\exp(\beta^T Z_i)}(S_i)]^{1-\delta_i} \cdot [1 - \bar{F}_0^{\exp(\beta^T Z_i)}(S_i)]^{\delta_i},$$

which can be maximized with respect to β and the values of \bar{F}_0 at the possible times of inspection (namely, the integers 7–21). We refer to this estimator as the ‘current status MLE’. Another option is to use the recalled event time whenever available, but to disregard the informativeness of the censoring. A penalized version of the corresponding likelihood is maximized in the function `shr` of the `SmoothHazard` package of R, which fits the Cox model by using an approximation of the hazard function by a linear combination of M-splines. We refer to this estimator as the ‘SmoothHazard MLE’.

We now compare the performance of the proposed AMLE of the regression coefficients with the three estimators described above. Table 1 shows the bias, the standard deviation (Stdev) and the mean squared error (MSE) of the estimated regression coefficients. The results reported here are based on 500 simulation runs for sample sizes $n = 50, 200$ and 1000 . It is clear that the standard deviation of the proposed AMLE, as well as its mean square error, is larger than those of the (hypothetical) ‘complete recall MLE’, but smaller than the ‘current status MLE’. The gap between the performances of the first two estimators becomes small as the sample size increases, though the gap between the AMLE and the ‘current status MLE’ does not reduce as much. The ‘SmoothHazard MLE’ has a persistent bias even when n is large. This outcome is expected, as the estimator is based on the assumption that the censoring is non-informative. Thus, neither the ‘current status MLE’ nor the ‘SmoothHazard MLE’ is able to successfully utilize the information contained in the recalled time-to-event data, while the proposed AMLE is able to do so.

Figs. 1 and 2 show the plots of the empirical bias and the empirical standard deviation of the estimated baseline survival functions, for $n = 50, 200$ and 1000 . It is clear that the empirical bias as well as the empirical standard deviation of the estimated baseline survival function become smaller as the sample size increases.

We now turn to the problem of testing for the significance of the estimators of the regression coefficients. The standard theory of parametric estimation generally does not hold for an infinite dimensional nuisance parameter. However, in the case of the Cox regression model for randomly right censored data, it has been shown that an asymptotic theory based on partial likelihood works in an analogous manner to that based on the asymptotic theory of parametric likelihood (Andersen and Gill, 1982), and that the partial likelihood may be viewed as the full likelihood maximized with respect to the baseline hazard subject to a piecewise linear constraint (Johansen, 1983). We now run some simulations to check whether the likelihood (18) with the nuisance parameters \bar{F}_0 replaced by the estimator (21) can be used similarly to obtain an approximate test of significance of the regression coefficients, even though there is no asymptotic theory as yet to justify such an approximation.

The ‘score vector’ (borrowing terminology of parametric likelihood theory) based on $\frac{\partial \ell_a(\gamma, \eta, \beta)}{\partial \beta}$, can be written as

$$U = \frac{\sum_{i=1}^n \sum_{j=1}^{n_2} \alpha_{ij} \left(\hat{F}(t_j) \log(\hat{F}(t_j)) - \hat{F}(t_{j+1}) \log(\hat{F}(t_{j+1})) \right) Z_i}{\sum_l \alpha_{il} g_{il}} \tag{22}$$

Table 1
Bias, Stdev and MSE of estimated regression coefficients.

Estimator	Property	n = 50		n = 200		n = 1000	
		β_1	β_2	β_1	β_2	β_1	β_2
Complete Recall MLE	Bias	0.2981	-0.0734	0.0089	0.0037	-0.0026	0.0008
	Stdev	0.8321	0.8499	0.1293	0.5048	0.0848	0.2271
	MSE	0.7812	0.7277	0.0168	0.2548	0.0072	0.0515
Proposed AMLE	Bias	0.2593	-0.0547	0.0105	-0.0047	0.0083	-0.0011
	Stdev	1.3145	1.2913	0.1739	0.5057	0.0885	0.2272
	MSE	1.7904	1.6658	0.0303	0.2553	0.0079	0.0516
Current Status MLE	Bias	-0.392	-0.3316	-0.026	-0.0367	-0.0032	0.0010
	Stdev	1.4048	1.3405	0.3225	0.9847	0.2353	0.6113
	MSE	2.1271	1.9069	0.1047	0.9709	0.0553	0.3737
Smooth-Hazard MLE	Bias	-0.170	3.551	-0.310	2.841	0.191	1.782
	Stdev	0.740	1.739	0.322	0.850	0.123	0.219
	MSE	0.569	15.648	0.198	8.780	0.0505	3.250

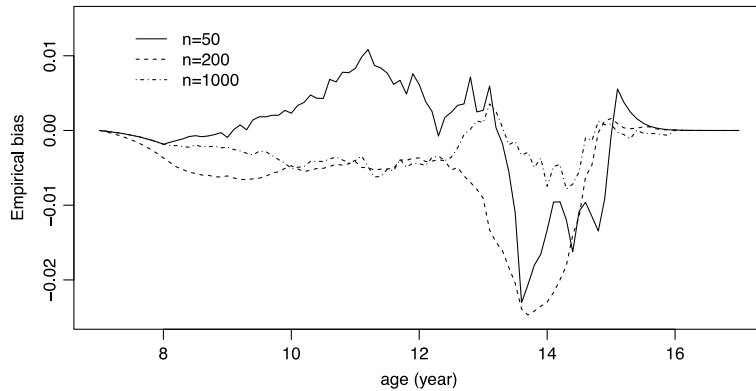


Fig. 1. Empirical bias of the estimated baseline survival function with $n = 50, 200$ and 1000 .

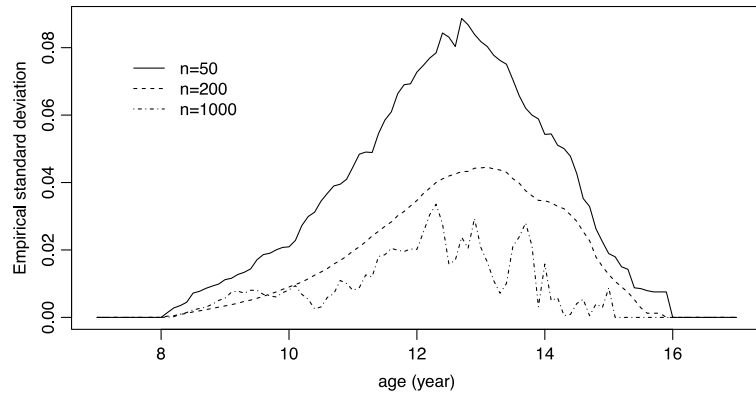


Fig. 2. Empirical standard deviation of estimated baseline survival function with $n = 50, 200$ and 1000 .

The relevant part of the ‘information matrix’ is $V = A_{22} - A_{21}A_{11}^{-1}A_{12}$, where

$$A = \begin{bmatrix} A_{11} & A_{12} \\ A_{21} & A_{22} \end{bmatrix},$$

$$A_{11} = - \begin{bmatrix} \partial^2 \ell_a / \partial \gamma \partial \gamma^T & \partial^2 \ell_a / \partial \gamma \partial \eta^T \\ \partial^2 \ell_a / \partial \eta \partial \gamma^T & \partial^2 \ell_a / \partial \eta \partial \eta^T \end{bmatrix},$$

$$A_{12} = - (\partial^2 \ell_a / \partial \gamma \partial \beta^T \quad \partial^2 \ell_a / \partial \eta \partial \beta^T) = A_{21}^T,$$

and

$$A_{22} = -\partial^2 \ell_a / \partial \beta \partial \beta^T,$$

Table 2The empirical type I error probability of test $H_0 : \beta = 0$.

Asymptotic	$n = 50$	$n = 200$	$n = 1000$
Type I error	0.041	0.038	0.022

Table 3Estimated regression coefficients and their p -values.

Covariates	Estimated value	p -value
Father passed high school	0.091	0.0036
Mother passed high school	0.249	0.0061
Monthly family expenditure (Rupees)	0.0002	0.0047

the quantities being estimated at $\beta = 0$, $\gamma = \hat{\gamma}_0$ and $\eta = \hat{\eta}_0$, the restricted AMLE's at $\beta = 0$. The hypothesis $\beta = 0$ may be tested by taking $U^T V^{-1} U$ as an approximate χ^2 statistic with 2 degrees of freedom. In order to check the behavior of this statistic, we generate data of sizes $n = 50, 200$ and 1000 for 1000 runs, under the null hypothesis when the baseline distribution function $F_0(t)$ is the Weibull distribution with shape and scale parameters $\alpha = 11$ and $\beta = 13$, respectively, truncated to the interval $[8, 16]$. The vector of covariates, $Z = (Z_1, Z_2)$, consists of a binary variable, taking values 1 and 0 with probabilities 0.25 and 0.75, and a continuous variable having the uniform distribution over the interval $[0, 5]$. Table 2 shows the value of empirical type I error probability of this test for different sizes of data. It can be seen that the error probability is less than 0.05. This indicates that the 'score test' is somewhat conservative, and more so for larger sample size.

4. Data analysis

In this section, we illustrate the use of the proposed method with data collected from an anthropometric study conducted by the Biological Anthropology Unit of the Indian Statistical Institute in and around the city of Kolkata, India from 2005 to 2011 (ISI, 2012, p. 108). In this retrospective data set, individuals aged between 7 and 21 years were surveyed through stratified sampling among students of educational institutions sampled at the first stage. The subjects, stratified by age year, were interviewed on or around their birthdays. The data set contains age, measurement of body dimensions, menarcheal status, age at menarche (if recalled), and some socioeconomic information. We used a part of this data set, and regarded the onset of menarche as the landmark event. Whenever the subject reported having had menarche but could not recall the date exactly, we regarded it as a case of non-recall.

There are many studies concerning the effects of socioeconomic factors on the measures of body shape (anthropometric indices or ratios) and physical maturation (e.g., biological parameters of the adolescent growth spurt) of children. Some of the important factors which affect age at menarche (maturation in girls) are diet and physical activities which can be directly related to parents' education and monthly family expenditure (Khan et al., 1996; Padez, 2003; Aryeetey et al., 2011). We considered three socioeconomic variables: two binary variables indicating whether the father or the mother of the subject had passed high school, and a real variable representing monthly family expenditure in Indian Rupees (indexed with respect to 2008 as base year).

We considered a subset of the original data, consisting of 673 respondents who came from a nuclear family and were the only child of their respective parents. Among 673 samples, 241 individuals did not have menarche, 147 individuals had menarche and recalled the date of its onset, while 285 individuals had menarche but could not recall the date. There were 492 individuals with father having passed high school and 420 individuals with mother having passed high school. The median of monthly family expenditure was Rupees 7808. As for the forgetting probability π , we modeled it over the interval 0–13 years (maximum possible separation between menarcheal age and age at observation in the sample). We used a piecewise constant model, with $k = 8$ and equal length of the intervals over which the probability is constant. The computational method for AMLE was as described in Section 3.

Table 3 shows the estimated regression coefficients and the corresponding p -values. It is found that all the coefficients are significant at the 1% level. The vector of the three regression coefficients has p -value 0.00093.

Fig. 3 shows a plot of the estimated survival functions of four hypothetical subjects with covariate profiles described below.

CASE (a) Neither parent passed high school, monthly family income is equal to the median income of the group (Rs. 7808). We represent this case as $Z = (0, 0, 7808)$.

CASE (b) Only the father passed high school, monthly family income is equal to the median income of the group. We represent this case as $Z = (1, 0, 7808)$.

CASE (c) Both the parents passed high school, monthly family income is equal to the median income of the group. We represent this case as $Z = (1, 1, 7808)$.

CASE (d) Both the parents passed high school, monthly family income is equal to Rupees 10,000. We represent this case as $Z = (1, 1, 10000)$.

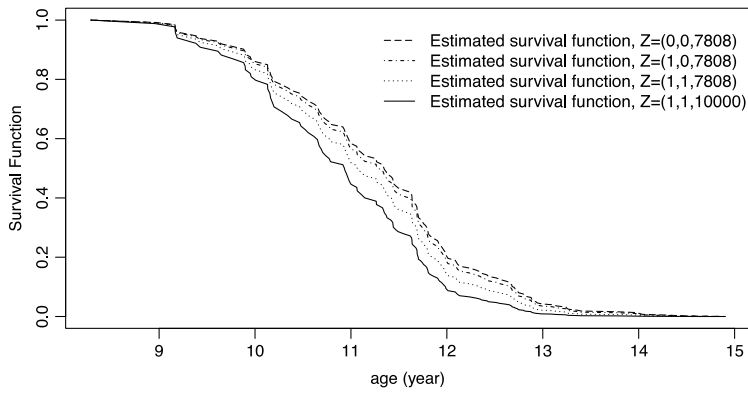


Fig. 3. Estimated survival function in different cases.

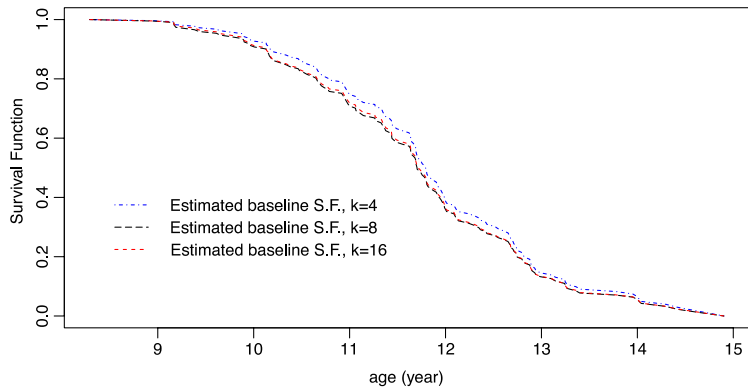


Fig. 4. Estimated baseline survival function with different k .

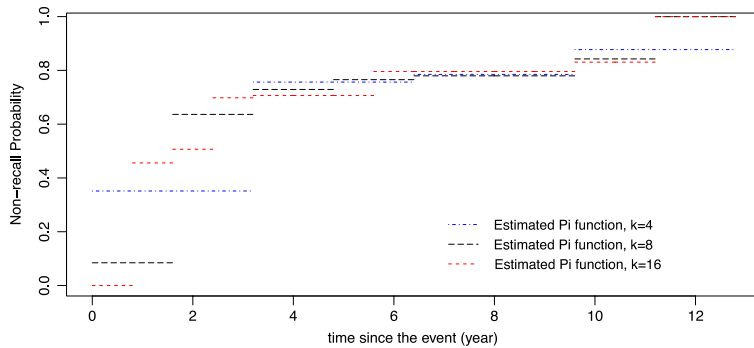


Fig. 5. Estimated π function with different k .

It is clear that the fact of any parent having passed high school is associated with earlier maturation. In particular, the mother’s educational status is found to account for a greater reduction of the survival function. Also, even a small increase in the monthly family expenditure is found to have a considerable impact on the survival function of the age at menarche.

The chosen value of k was obtained after considering a coarser and a finer partition for the piecewise constant model of π . Specifically, the range 0–13 years was split experimentally into k equal intervals, with $k = 4, 8$ and 16 , and the resulting estimated baseline survival functions were compared. Fig. 4 shows plots of the estimated baseline survival function for different values of k . It is seen that by increasing k from 4 to 8, one observes a substantial change in the estimated baseline survival function, though the change is much less when k is increased from 8 to 16. The integrated mean square difference between baseline survival functions (scaled by the integral of the square of the function for the lower value of k) is 0.92 when one compares $k = 4$ with $k = 8$. The same criterion produces the value 0.021 when the comparison is between the curves for $k = 8$ and $k = 16$. We have chosen $k = 8$, as the alternative choice $k = 16$ does not produce a substantially different estimate of the baseline survival function. Fig. 5 shows the estimated function π for different values of k . Once again, the estimates of π for $k = 8$ and $k = 16$ differ much less than those for $k = 4$ and $k = 8$.

5. Concluding remarks

In this paper, we have presented a method for fitting the Cox regression model to recall-based time-to-event data with covariates, where there is informative censoring. Simulation results indicate that the estimators of the regression coefficients are reasonable, though there is no proof of consistency of these estimators as of now. It may be recalled that there is no known proof of consistency of MLEs of the Cox regression parameters even in the case of non-informatively interval-censored data. Some results are available in the special case of status data with fixed and multiple inspection times (and in particular, for the further special case of current status data) (Huang, 1996; Yu et al., 2006; Liu and Shen, 2009). The problem of establishing consistency may be considered in future.

Fitting of a semi-parametric regression model is generally a more complex inferential problem than that of estimating only a distribution. The complexity in the present case is even greater because the informative interval censoring model leads to a large number of nuisance parameters, including the probability masses allocated, as per the baseline distribution of the Cox model, to intersections of different intervals. The tasks of formation of these intervals and tracking of their probability masses are greatly simplified by the approximation inspired by Theorem 3. The Cox regression model appears to be suited to the formulation of the approximate likelihood through masses at the times of exactly recalled events. It is this matching of the models that makes the AMLE computationally tractable. A different approach may be needed for other regression models.

The proposed approach can be adapted to handle left truncated data. Assuming that there is a time of left truncation associated with each observation, each term in the likelihood would have to be divided by the upper tail probability at the point of truncation. It can be shown that the simplification given through Theorem 2 will continue to hold, since the shift of mass envisaged in the proof of that theorem does not alter the factors in the denominator. The objective function (18) would then be replaced by

$$\ell_a(\gamma, \eta, \beta) = \sum_{i=1}^n \log \left[\sum_{j=1}^{n_2} \alpha_{ij} \left[e^{(-e^{z_i \beta + \gamma_j})} - e^{(-e^{z_i \beta + \gamma_{j+1}})} \right] \right] - \log \left[\sum_{j=1}^{n_2} \psi_{ij} e^{(-e^{z_i \beta + \gamma_j})} \right], \quad (23)$$

where ψ_{ij} 's are known constants like α_{ij} 's. The optimization problem is therefore similar.

The data set analyzed in Section 4 also contains 'partial' recall data relating to the week/month/year of menarche. In this paper, we have regarded a date as not recalled at all even when a range of possible dates is available. Apart from simplification of the problem, this strategy also minimizes errors in recall, which has been recognized as a problematic issue with recall data (Rabe-Hesketh et al., 2001; Wen and Chen, 2014). If reliability of partial recall data is not an issue, one might look for more sophisticated modeling to handle it. The work presented in this paper can be used as a point of departure for analysis under such models. Another direction of future research could be extension of this model to include frailty. Methodology for other forms of regression models, such as the accelerated failure time model, may be developed also.

Acknowledgments

This research is partially sponsored by the project "Physical growth, body composition and nutritional status of the Bengal school aged children, adolescent, and young adults of Calcutta, India: Effects of socioeconomic factors on secular trends", funded by the Neys Van Hoogstraten Foundation of the Netherlands, and the project "Optimization and Reliability Modeling" funded by the Indian Statistical Institute, Kolkata. The authors thank Professor Parasmani Dasgupta, leader of the first project, for making the data available for this research and Professor Biswabrata Pradhan, leader of the second project, for facilitating useful discussions. The investigation into the issue of identifiability was enriched by a discussion with Professor Probal Chaudhuri.

Appendix

A.1. Proof of Theorem 1

Without loss of generality, we assume that a possible value of z is the 0 vector (this can be always achieved through a shift of origin, the effect of which can be absorbed through the baseline survival function). For the sake of contradiction, let us assume there are two values of the triplet (β, \bar{F}_0, π) , say $(\beta_1, \bar{F}_{01}, \pi_1)$ and $(\beta_2, \bar{F}_{02}, \pi_2)$, such that their substitutions on the right hand side of (4) produce the same function. Then we have, for all z and s and all positive $v < s$,

$$\begin{aligned} -\frac{d}{dv} \left(\bar{F}_{01}(s-v)^{\exp(\beta_1^T z)} \right) (1 - \pi_1(v)) &= h(v, 1|z; \beta) \\ &= -\frac{d}{dv} \left(\bar{F}_{02}(s-v)^{\exp(\beta_2^T z)} \right) (1 - \pi_2(v)). \end{aligned}$$

Hence,

$$\frac{\frac{d}{dv} \left(\bar{F}_{01}(s-v)^{\exp(\beta_1^T z)} \right)}{\frac{d}{dv} \left(\bar{F}_{02}(s-v)^{\exp(\beta_2^T z)} \right)} = \frac{1 - \pi_2(v)}{1 - \pi_1(v)} \quad \forall z, s, v < s, \quad (\text{A.1})$$

i.e.,

$$\frac{\exp(\beta_1^T z) \bar{F}_{01}(s-v)^{\exp(\beta_1^T z)-1} f_{01}(s-v)}{\exp(\beta_2^T z) \bar{F}_{02}(s-v)^{\exp(\beta_2^T z)-1} f_{02}(s-v)} = \frac{1-\pi_2(v)}{1-\pi_1(v)} \quad \forall z, s, v < s. \tag{A.2}$$

In particular, the above identity holds for $z = 0$, i.e.,

$$\frac{f_{01}(s-v)}{f_{02}(s-v)} = \frac{1-\pi_2(v)}{1-\pi_1(v)} \quad \forall s, v < s. \tag{A.3}$$

After combining the above equation with (A.2), we obtain

$$\frac{\bar{F}_{01}(s-v)^{\exp(\beta_1^T z)-1}}{\bar{F}_{02}(s-v)^{\exp(\beta_2^T z)-1}} = \exp((\beta_2 - \beta_1)^T z) \quad \forall z, s, v < s. \tag{A.4}$$

By taking the limit of the left hand side as v goes to s , we obtain $\exp((\beta_2 - \beta_1)^T z) = 1$, i.e.,

$$\beta_2^T z = \beta_1^T z \quad \forall z. \tag{A.5}$$

Since the above equation holds for r linearly independent values of the vector z (as assumed in the statement of the theorem), we have

$$\beta_1 = \beta_2. \tag{A.6}$$

It follows from Eqs. (A.4) and (A.5) that

$$\left[\frac{\bar{F}_{01}(s-v)}{\bar{F}_{02}(s-v)} \right]^{\exp(\beta_1^T z)-1} = 1 \quad \forall z, s, v < s. \tag{A.7}$$

Therefore,

$$\bar{F}_{01} = \bar{F}_{02}. \tag{A.8}$$

From Eqs. (A.1), (A.6) and (A.8), we have

$$\pi_1(v) = \pi_2(v) \quad \forall v, \tag{A.9}$$

i.e.,

$$(\beta_1, \bar{F}_{01}, \pi_1) = (\beta_2, \bar{F}_{02}, \pi_2),$$

which is a contradiction.

A.2. Proof of Theorem 2

Since \mathcal{C} is the union of disjoint sets \mathcal{C}_0 and $\mathcal{C}_1 \cup \mathcal{C}_2$, we can rewrite the likelihood (13) as

$$\begin{aligned} L &= \prod_{i \in \mathcal{I}_1} \left(\sum_{\substack{r: I_r \subseteq A_i \\ s_r \in \mathcal{C}_0}} p_r + \sum_{\substack{r: I_r \subseteq A_i \\ s_r \in \mathcal{C}_1 \cup \mathcal{C}_2}} p_r \right)^{\exp(\beta^T Z_i)} \times \prod_{i \in \mathcal{I}_2} \left(1 - \sum_{l=1}^k b_l I(T_i \in A_{il}) \right) \\ &\cdot \left[\left(\sum_{\substack{r: I_r \subseteq A_i \\ s_r \in \mathcal{C}_0}} p_r + \sum_{\substack{r: I_r \subseteq A_i \\ s_r \in \mathcal{C}_1 \cup \mathcal{C}_2}} p_r \right)^{\exp(\beta^T Z_i)} - \left(\sum_{\substack{r: I_r \subseteq A'_i \\ s_r \in \mathcal{C}_0}} p_r + \sum_{\substack{r: I_r \subseteq A'_i \\ s_r \in \mathcal{C}_1 \cup \mathcal{C}_2}} p_r \right)^{\exp(\beta^T Z_i)} \right] \\ &\times \prod_{i \in \mathcal{I}_3} \left[\sum_{l=1}^k b_l \left\{ \left(\sum_{\substack{r: I_r \subseteq A_{i(l+1)} \\ s_r \in \mathcal{C}_0}} p_r + \sum_{\substack{r: I_r \subseteq A_{i(l+1)} \\ s_r \in \mathcal{C}_1 \cup \mathcal{C}_2}} p_r \right)^{\exp(\beta^T Z_i)} - \left(\sum_{\substack{r: I_r \subseteq A_{il} \\ s_r \in \mathcal{C}_0}} p_r + \sum_{\substack{r: I_r \subseteq A_{il} \\ s_r \in \mathcal{C}_1 \cup \mathcal{C}_2}} p_r \right)^{\exp(\beta^T Z_i)} \right\} \right]. \tag{A.10} \end{aligned}$$

For every $s_r \in \mathcal{C}_2$, there exists a unique $s_{r^*} \in \mathcal{C}_0$ such that $s_{r^*} \setminus s_r \cap s_r = \{j_r\}$ and $s_r \setminus s_{r^*} \cap s_r = \{n_2 + j_r\}$ for some integer j_r in between 1 and n_2 , where n_2 is as defined after (10). If any probability mass is shifted from I_r to I_{r^*} , the likelihood (A.10) can possibly be affected only through terms that involve the sets B_{j_r} and $B_{n_2+j_r}$, defined in (11). Given the fact that the baseline distribution is absolutely continuous, there is almost surely a unique $i_r \in \mathcal{I}_2$ such that $B_{j_r} = A_{i_r} \setminus A'_{i_r} = \{T_{i_r}\}$ and

$B_{n_2+j_r} = A'_{i_r}$. The individual indexed by i_r is the only one whose contribution to the likelihood is affected by the change. For this individual, $I_{r^*} \subset B_{j_r} \subseteq A_{i_r}$, but $I_{r^*} \not\subseteq A'_{i_r}$. On the other hand, $I_r \subseteq B_{n_2+j_r} = A'_{i_r} \subseteq A_{i_r}$. Therefore, the first exponentiated term in the second line of (A.10) remains the same after the shift of mass, while there is a reduction in the subtracted term in that line. The likelihood increases as a result.

We now turn to shifting of probability mass out of I_r , where $s_r \in \mathcal{C}_1$. For any such s_r , define the non-empty set $\mathcal{C}_{s_r} = \{s' : s' \in \mathcal{C}_0, s_r \subset s'\}$. If \mathcal{C}_{s_r} is a singleton, we denote the only member by s_{r^*} . If \mathcal{C}_{s_r} is not a singleton, we denote by s_{r^*} that member which satisfies the condition: 'for all $\beta \in \cup_{j: s_j \in \mathcal{C}_{s_r}; s_j \neq s_{r^*}} I_j$, there is a real number $\alpha \in I_{r^*}$ such that $\alpha < \beta$ '. Thus, for every $s_r \in \mathcal{C}_1$, we have a uniquely defined $s_{r^*} \in \mathcal{C}_0$.

If p_r is increased at the expense of p_{r^*} , the likelihood (A.10) can possibly change only through terms that involve sets B_j such that $j \in s_{r^*} \setminus s_r$. We shall show that for an individual i , whose contribution to the likelihood involves such sets, that contribution generally increases due to the said shift of probability mass. In a particular case (Case (iii) below), where this shift cannot be proved to increase the likelihood, there is another way of shifting mass out of p_r that would definitely increase the likelihood.

CASE (i) Let $j \in s_{r^*} \setminus s_r$ and $B_j = A_{i_j}$ for some $i_j \in \mathcal{I}_1$. Any shift of probability mass from I_r to I_{r^*} would increase the contribution of the i_j th individual to the likelihood, since $I_{r^*} \subseteq A_{i_j}$ but $I_r \not\subseteq A_{i_j}$.

CASE (ii) Let $j \in s_{r^*} \setminus s_r$ and $B_j = A_{i_j} \setminus A'_{i_j}$ for some $i_j \in \mathcal{I}_2$. In this case, $I_{r^*} \subseteq A_{i_j}$ but $I_{r^*} \not\subseteq A'_{i_j}$. By construction, $B_{n_2+j} = A'_{i_j}$, which is disjoint with B_j . In order that I_{r^*} is not a null set, we must have $n_2 + j \notin s_r$. It follows that I_r is not contained in B_j or B_{n_2+j} . Thus, $I_r \not\subseteq A_{i_j}$ and $I_r \not\subseteq A'_{i_j}$. Clearly, a transfer of probability mass from I_r to I_{r^*} would increase the contribution of the i_j th individual to the likelihood.

CASE (iii) Let $j \in s_{r^*} \setminus s_r$ and $B_j = A'_{i_j}$ for some $i_j \in \mathcal{I}_2$. Since $j \notin s_r$, we have $I_r \subseteq B_j^c = [t_{\min}, T_{i_j})$. Therefore, for each of the intervals B_l with $l \in s_r$, $B_l \cap [t_{\min}, T_{i_j}) \neq \emptyset$. On the other hand, since $I_{r^*} \neq \emptyset$, we have $B_l \cap (T_{i_j}, t_{\max}] \neq \emptyset$ for $l \in s_r$. It follows that each of the intervals B_l , $l \in s_r$, contains a left- and a right-neighborhood of the point T_{i_j} . Consequently, T_{i_j} is contained in these intervals. Hence, the set $s_{r^\dagger} = \{l : T_{i_j} \in B_l\}$ is a superset of s_r contained in \mathcal{C}_0 , with $I_{r^\dagger} = \{T_{i_j}\} \neq \emptyset$. As argued in Case (ii), a transfer of probability mass from I_r to I_{r^\dagger} would increase the contribution of the i_j th individual to the likelihood.

CASE (iv) Let $j \in s_{r^*} \setminus s_r$ and $B_j = A_{i_j(l+1)} \setminus A_{i_j l}$ for some $l \in \{1, \dots, k\}$ and some $i_j \in \mathcal{I}_3$. A transfer of probability mass from I_r to I_{r^*} would increase the contribution of the i_j th individual to the likelihood. This is because of the fact that $I_{r^*} \subseteq A_{i_j(l+1)}$ and $I_{r^*} \not\subseteq A_{i_j l}$, whereas I_r is not contained in either of these sets.

It transpires that maximization of L can be achieved even in the presence of the constraint $p_r = 0$ for $s_r \in \mathcal{C}_1 \cup \mathcal{C}_2$. Thus, L can be fully maximized over the restricted set $\{p_r : s_r \in \mathcal{C}_0\}$.

A.3. Proof of Theorem 3

Let $i \in \mathcal{I}_2$ and the index j_i be such that $s_{j_i} = \{j : T_i \in B_j\}$. Since each time-to-event has an absolutely continuous distribution, the recalled times T_i , $i \in \mathcal{I}_2$ are distinct with probability 1. Therefore, $\{T_i\} \in \{B_1, B_2, \dots, B_{n_2}\}$ almost surely. It follows that $T_i \in I_{j_i} \subseteq \{T_i\}$, i.e., $I_{j_i} = \{T_i\}$ with probability 1. It is also easy to see that s_{j_i} does not belong to \mathcal{C}_1 or \mathcal{C}_2 , with probability 1. Therefore, $s_{j_i} \in \mathcal{C}_0$ and hence $\mathcal{A}_2 \subseteq \mathcal{A}_0$ almost surely.

Let $J_j \in \mathcal{A}_0 \setminus \mathcal{A}_2$. There is an index r such that $I_r = J_j \neq \emptyset$ and $s_r \in \mathcal{C}_0$, even though $I_r \neq \{T_i\}$ for any $i \in \mathcal{I}_2$. We shall show that the existence of I_r implies an event with probability going to zero.

It is easy to see that $i \notin s_r$ for $i = 1, 2, \dots, n_2$. Thus, I_r can be written as

$$I_r = \left\{ \bigcap_{i \in s_r} B_i \right\} \cap \left\{ \bigcap_{i \notin s_r} B_i^c \right\} = I'_r \setminus \{T_i, i \in \mathcal{I}_2\},$$

where

$$I'_r = L_r \cap R_r, \quad L_r = \left\{ \bigcap_{i \in s_r} B_i \right\}, \quad R_r = \left\{ \bigcap_{i \notin s_r, i > n_2} B_i^c \right\}.$$

If there is an $i \in \mathcal{I}_2$ such that $T_i \in I'_r$, then the index set $s_{r^*} = s_r \cup \{i\}$ corresponds to the non-null set $I_{r^*} = \{T_i\}$. It follows that $s_r \in \mathcal{C}_1$, which leads to the contradictory conclusion $s_r \notin \mathcal{C}_0$. Therefore $T_i \notin I'_r$ for any $i \in \mathcal{I}_2$.

We now show that an upper bound of the probability of the above event goes to zero as $n_2 \rightarrow \infty$. Since the set L_r is obtained as an intersection of sets of the form $(S_i, t_{\max}]$, $(T_i, t_{\max}]$, $(W_l(S_i), t_{\max}]$ or $[t_{\min}, t_{\max}]$, the intersection itself must be an interval of the form $(l_r, t_{\max}]$. On the other hand, since the set R_r is obtained as an intersection of sets that are complements of sets of the above type, the intersection itself must be an interval of the form $[t_{\min}, m_r]$. Thus, the set I'_r is the interval $(l_r, m_r]$. By the argument given in the preceding paragraph, neither l_r nor m_r is equal to T_i for any $i \in \mathcal{I}_2$ (otherwise s_r would not be in \mathcal{C}_0). Therefore, both l_r and m_r are of the form S_i for some $i \in \mathcal{I}_1$ or of the form $S_i - x_l$ for some $i \in \mathcal{I}_3$ and some $l \in \{1, \dots, k\}$.

Let $w_1 < w_2 < \dots < w_K$ be the feasible values of S_i and $S_i - x_l$ (where $1 \leq l \leq k$) that are strictly between t_{\min} and t_{\max} . Since the baseline distribution is absolutely continuous, we have $1 > \bar{F}(w_1) > \bar{F}(w_2) > \dots > \bar{F}(w_K) > 0$. The values of l_r and m_r are taken from the set w_1, w_2, \dots, w_K .

The probability of the event “ $T_i \notin I'_r$ for any $i \in \mathcal{I}_2$ ” is

$$\prod_{i \in \mathcal{I}_2} \left[1 - \left\{ \bar{F}^{\exp(\beta^T Z_i)}(l_r) - \bar{F}^{\exp(\beta^T Z_i)}(m_r) \right\} \right] = \prod_{i \in \mathcal{I}_2} \left[1 - \left\{ \bar{F}^B(l_r) \right\}^{\exp(\beta^T Z_i)/B} + \left\{ \bar{F}^B(m_r) \right\}^{\exp(\beta^T Z_i)/B} \right],$$

where B is an upper bound on $\exp(\beta^T Z_i)$. Since $u^{\exp(\beta^T Z_i)/B}$ is a strictly concave function of u , we have

$$1 - (1 - u_2 + u_1)^{\exp(\beta^T Z_i)/B} < u_2^{\exp(\beta^T Z_i)/B} - u_1^{\exp(\beta^T Z_i)/B}$$

for $0 < u_1 < u_2 < 1$. Using this inequality for $u_1 = \bar{F}^B(m_r)$ and $u_2 = \bar{F}^B(l_r)$, we have

$$\begin{aligned} \prod_{i \in \mathcal{I}_2} \left[1 - \left\{ \bar{F}^{\exp(\beta^T Z_i)}(l_r) - \bar{F}^{\exp(\beta^T Z_i)}(m_r) \right\} \right] &< \prod_{i \in \mathcal{I}_2} \left[1 - \bar{F}^B(l_r) + \bar{F}^B(m_r) \right]^{\exp(\beta^T Z_i)/B} \\ &< \left[1 - \bar{F}^B(l_r) + \bar{F}^B(m_r) \right]^{n_2 L/B}, \end{aligned}$$

where L is a lower bound on $\exp(\beta^T Z_i)$. Since $[1 - \bar{F}^B(w_{j_1}) + \bar{F}^B(w_{j_2})] \in (0, 1)$ for any j_1 and j_2 with $1 \leq j_1 < j_2 \leq K$, we have $[1 - \bar{F}^B(l_r) + \bar{F}^B(m_r)] \in (0, 1)$. Therefore, the last expression goes to zero as $n_2 \rightarrow \infty$. This completes the proof.

References

- Alioum, A., Commenges, D., 1996. A proportional hazards model for arbitrarily censored and truncated data. *Biometrics* 52, 512–524.
- Andersen, P.K., Gill, R.D., 1982. Cox's regression model for counting processes: A large sample study. *Ann. Statist.* 10, 1100–1120.
- Aryeetey, R., Ashinyo, A., Adjuik, M., 2011. Age at menarche among basic level school girls in Medina, Accra. *African J. Reprod. Health* 103, 103–110.
- Breslow, N.E., Lubin, J.H., Marek, P., Langholz, B., 1983. Multiplicative models and cohort analysis. *J. Amer. Statist. Assoc.* 78, 1–12.
- Chaudhuri, P., 2007. On single index regression models for multivariate survival time data. In: Nair, V. (Ed.), *Advances in Statistical Modeling and Inference: Essays in Honor of Kjell A. Doksum*. pp. 223–232.
- Cox, D.R., 1972. Regression models and life-tables. *J. R. Stat. Soc. Ser. B* 34, 187–220. With discussion by F. Downton, Richard Peto, D.J. Bartholomew, D.V. Lindley, P.W. Glassborow, D.E. Barton, Susannah Howard, B. Benjamin, John J. Gart, L.D. Meshalkin, A.R. Kagan, M. Zelen, R.E. Barlow, Jack Kalbfleisch, R.L. Prentice and Norman Breslow, and a reply by D.R. Cox.
- Dabrowska, D.M., Doksum, K.A., 1988. Estimation and testing in a two-sample generalized odds-rate model. *J. Amer. Statist. Assoc.* 83, 744–749.
- DeMasi, R.A., Qaqish, B., Sen, P.K., 1997. Statistical models and asymptotic results for multivariate failure time data with generalized competing risks. *Sankhyā Ser. A* 59, 408–434.
- Finkelstein, D.M., 1986. A proportional hazards model for interval-censored failure time data. *Biometrics* 42, 845–854.
- Finkelstein, D.M., Moore, D.F., Schoenfeld, D.A., 1993. A proportional hazards model for truncated AIDS data. *Biometrics* 49, 731–740.
- Gross, S.T., Huber-Carol, C., 1992. Regression models for truncated survival data. *Scand. J. Stat.* 19, 193–213.
- Huang, J., 1996. Efficient estimation for the proportional hazards model with interval censoring. *Ann. Statist.* 24, 540–568.
- ISI, 2012. Annual Report of the Indian Statistical Institute 2011–12. Indian Statistical Institute. Available at: URL <http://library.isical.ac.in/jspui/handle/10263/5345?mode=full>.
- Johansen, S., 1983. An extension of Cox's regression model. *Internat. Statist. Rev.* 51, 165–174.
- Kalbfleisch, J.D., Lawless, J.F., 1991. Regression models for right-truncated data with applications to aids incubation times and reporting lags. *Statist. Sinica* 81, 19–32.
- Khan, A.D., Schroeder, D.G., Martorell, R., Haas, J.D., Rivera, J., 1996. Early childhood determinants of age at menarche in rural Guatemala. *Am. J. Hum. Biol.* 8, 717–723.
- Klein, J.P., Moeschberger, M.L., 2003. *Survival Analysis: Techniques for Censored and Truncated Data*. Springer-Verlag, New York.
- Liu, H., Shen, Y., 2009. A semiparametric regression cure model for interval-censored data. *J. Amer. Statist. Assoc.* 104, 1168–1178.
- Mirzaei, S.S., Das, R., Sengupta, D., 2015. Parametric estimation of menarcheal age distribution based on recall data. *Scand. J. Stat.* <http://dx.doi.org/10.1111/sjos.12107>.
- Mirzaei Salehabadi, S., Sengupta, D., 2015. Nonparametric estimation of time-to-event distribution based on recall data in observational studies. *Lifetime Data Anal.* <http://dx.doi.org/10.1007/s10985-015-9338-8>.
- Padez, C., 2003. Age at menarche of schoolgirls in Maputo, Mozambique. *Ann. Hum. Biol.* 30, 487–495.
- Rabe-Hesketh, S., Yang, S., Pickles, A., 2001. Multilevel models for censored and latent responses. *Stat. Methods Med. Res.* 10., 409–427.
- Tanaka, Y., Rao, P.V., 2005. A proportional hazards model for informatively censored survival times. *J. Statist. Plann. Inference* 129, 253–262.
- Tu, X.M., Meng, X.L., Pagano, M., 1993. The AIDS epidemic: Estimating survival after AIDS diagnosis from surveillance data. *J. Amer. Statist. Assoc.* 88, 26–36.
- Vonta, F., 1996. Efficient estimation in a non-proportional hazards model in survival analysis. *Scand. J. Stat.* 23, 49–61.
- Wei, L.J., 1992. The accelerated failure time model: A useful alternative to the Cox regression model in survival analysis (with discussion). *Stat. Med.* 11, 1871–1879.
- Wen, C.C., Chen, Y.H., 2014. Functional inference for interval-censored data in proportional odds model with covariate measurement error. *Statist. Sinica* 24.
- Wienke, A., 2010. *Frailty Models in Survival Analysis*. Chapman and Hall /CRC.
- Yu, Q., Wong, G.Y.C., Kong, F., 2006. Consistency of the semi-parametric MLE in linear regression models with interval-censored data. *Scand. J. Stat.* 33, 367–378.