# Quantile Regression for Survival Data

**Limin Peng**[1]

[1]Department of Biostatistics and Bioinformatics, Emory University, Atlanta, USA, 30322;

## Abstract

Quantile regression offers a useful alternative strategy for analyzing survival data. Compared to traditional survival analysis methods, quantile regression allows for comprehensive and flexible evaluations of covariate effects on a survival outcome of interest, while providing simple physical interpretations on the time scale. Moreover, many quantile regression methods enjoy easy and stable computation. These appealing features make quantile regression a valuable practical tool for delivering in-depth analyses of survival data. In this paper, I review a comprehensive set of statistical methods for performing quantile regression with different types of survival data. This review covers various survival scenarios, including randomly censored data, data subject to left truncation or censoring, competing risks and semi-competing risks data, and recurrent events data. Two real examples are presented to illustrate the utility of quantile regression for practical survival data analyses.

### Keywords

Quantile regression; estimating equation; randomly censored data; competing risks data; semi-competing risks data; recurrent events data

## 1. Background and motivation

The problem of analyzing survival (or time-to-event) data arises in a number of scientific fields. For example, an event time of interest can be survival time of a cancer patient recorded in a medical study, time to high school dropout studied by sociologists, "survival" time of new business addressed in economic research, or lifetime of a part under stress evaluated in an engineering reliability study. A common feature of survival data is that they often contain incomplete time-to-event information due to censoring or truncation. Censoring occurs when an event time is known to have occurred only within certain intervals. Truncation is defined as a condition which excludes certain subjects from the study population. Statistical methods for analyzing survival data need to appropriately account for various forms of censoring and truncation.

To evaluate the association between a survival outcome and a set of explanatory variables (or covariates), the accelerated failure time (AFT) model has been extensively studied in literature as a counterpart of linear regression in survival analysis (Miller 1976, Buckley &

lpeng@emory.edu.

James 1979, Prentice 1978, Wei & Gail 1983, Tsiatis 1990, Ritov 1990, Wei et al. 1990, among others). Consider an event time $T$ and a $p \times 1$ covariate vector $\widetilde{Z}$. The AFT model regresses a survival response, $Y \doteq \log(T)$, or another monotone transformation of $T$, over $\widetilde{Z}$; that is,

$$Y = \widetilde{Z}^\top \boldsymbol{b} + \epsilon,$$

where $\boldsymbol{b}$ is a vector of unknown regression coefficients and $\epsilon$ is an error term with an unknown distribution independent of $\widetilde{Z}$. The assumption that $\epsilon$ is independent of $\widetilde{Z}$ confines covariates to affect only the location of the distribution of $Y$. However, this is often too restrictive in real applications. For example, an analysis of a dialysis dataset presented in Section 5.1 suggests that the symptom severity of restless leg syndrome (RLS) may only impact the lower range but not the upper range of the survival function of dialysis patients. The AFT model, which assumes pure location shift effects, would fail to accommodate such an inhomogeneous effect of RLS.

An alternative regression strategy for survival data is to use the Cox proportional hazards (PH) model (Cox 1972, Andersen & Gill 1982) Specifically, the Cox PH model relates the conditional hazard function of $T$, $\lambda(t \mid \widetilde{Z}) \doteq \lim_{\Delta \to 0} \Delta^{-1} \Pr\{t \in (t, t + \Delta] \mid T \geqslant t, \widetilde{Z}\}$, to covariates in $\widetilde{Z}$ multiplicatively without specifying a parametric baseline hazard; that is,

$$\lambda(t \mid \widetilde{Z}) = \lambda_0(t) \exp\{\widetilde{Z}^\top \boldsymbol{b}\},$$

where $\lambda_0(t)$ is an unspecified baseline hazard function, and $\boldsymbol{b}$ is an unknown regression coefficient vector. The Cox PH model is widely used in the practice of survival analysis. Nevertheless, the Cox regression model requires covariate-specific hazards be proportional. This key assumption essentially exerts a location shift model for a monotonically transformed survival function of $T$. This limits the applications of the Cox PH model in scenarios with inhomogeneous covariate effects as exemplified above. Under the Cox PH model, covariate effects are formulated on the conditional hazard function of $T$, which lacks a physical interpretation on event times (Reid 1994).

Quantile regression (Koenker & Bassett 1978) offers a natural remedy for accommodating heterogeneous covariate effects, while retaining straightforward physical interpretations. A comprehensive review of quantile regression methodology was provided by Koenker (2017). Quantile regression has received increased attention in survival analysis because event times themselves are often of scientific interest, and quantiles are more flexible and robust quantitative tools for characterizing event times than mean-based devices. For example, in the presence of censoring with bounded support, mean survival time may not be identifiable, while quantiles may be identifiable.

For a survival time $T$, a standard quantile regression model assumes that the $\tau$-th conditional quantile of $Y \doteq \log(T)$ given $Z = \left(1, \widetilde{Z}^\top\right)^\top$, defined as $Q_Y(\tau|Z) \doteq \inf\{t : \Pr(Y \leqslant t|Z) \geqslant \tau)$, is linearly related to covariates in $Z$. That is,

$$Q_Y(\tau \mid \boldsymbol{Z}) = \boldsymbol{Z}^\top \boldsymbol{\beta}_0(\tau), \ \tau \in [\tau_L, \tau_U], \tag{1}$$

where $0 \leqslant \tau_L \leqslant \tau_U < 1$ and $\boldsymbol{\beta}_0(\tau)$ is a $(p+1) \times 1$ vector of unknown regression coefficients. A non-intercept coefficient in $\boldsymbol{\beta}_0(\tau)$ represents a covariate effect on the $\tau$-th conditional quantile of $\log(T)$. By allowing $\boldsymbol{\beta}_0(\tau)$ to change with $\tau$, quantile regression permits varying covariate effects on different segments of the response distribution. This feature renders the flexibility to accommodate inhomogeneous covariate effects

When $\tau_L = \tau_U$, model (1) is referred to as a locally linear quantile regression model because it only asserts "local" linearity between the conditional quantile of $\log(T)$ and $\boldsymbol{Z}$ at a single quantile level. When $\tau_L < \tau_U$, model (1) imposes a "global" linearity for the conditional quantile of $\log(T)$ throughout the $\tau$-interval $[\tau_L, \tau_U]$, and thus is referred to as a globally linear quantile regression model. A globally linear quantile regression model can provide a platform for investigating the dynamic pattern of covariate effects, while paying the price of imposing a stronger model assumption compared to a locally linear quantile regression model.

It is easy to see that the AFT model is a special case of model (1) with $\boldsymbol{\beta}_0(\tau) = \{ Q_\epsilon(\tau), \boldsymbol{b}^\top \}^\top$, where $Q_\epsilon(\tau)$ denotes the $\tau$th-percentile of $\epsilon$. Under the standard Cox PH model, $Q_{\log \Lambda_0(T)}(\tau \mid \boldsymbol{Z}) = \log\left\{ -\log(1-\tau) \right\} + \widetilde{\boldsymbol{Z}}^\top \boldsymbol{b}$, where $\Lambda_0(t) = \int_0^t \lambda_0(u)du$. Given $\Lambda_0(\cdot)$ is an unknown function, the quantile interpretation of the Cox PH model is vague. Moreover, many interesting forms of heterogeneous covariate effects are excluded by the restrictive forms of $Q_T(\tau|Z)$ designated by the AFT model and the Cox PH model. In contrast, quantile regression modeling offers straightforward physical interpretations as well as greater flexibility to accommodate heterogeneous associations between covariates and the survival response. This serves as the key motivation for considering quantile regression as an alternative approach to analyzing survival data.

In this paper, we present a comprehensive methodological framework that has been developed to perform quantile regression with survival data. Due to space limit, the review is rather selective but yet cover a wide range of survival scenarios. More specifically, section 2 is focused on the standard survival setting with randomly censored data. Section 3 includes discussions of quantile regression methods applicable to more complex survival settings that involve left truncation or left censoring, competing risks and semi-competing risks. Recent method developments for recurrent events data are presented in Section 4. Two real data examples are presented in section 5 to illustrate the practical utility of quantile regression methods for survival analyses. Section 6 concludes this paper with a brief summary and a few remarks.

## 2. Quantile regression for randomly censored data

Let $T$ denote time to event subject to right censoring by $C$, and let $Y = \log(T)$. Define $\widetilde{T} = T \wedge C$ and $= I(T \leqslant C)$, where $\wedge$ is the minimum operator. The observed data consist of

$n$ i.i.d. replicates of $(\widetilde{T}, \ , \mathbf{Z})$, denoted by $(\widetilde{T}_i, \ _i, \mathbf{Z}_i)$, $i = 1, \ldots, n$. Define $\widetilde{Y} = \log(\widetilde{T})$, $\widetilde{Y}_i = \log(\widetilde{T}_i)$, $U = \log(C)$, and $U_i = \log(C_i)$.

### 2.1.  Random right censoring with *C* always known

In the absence of censoring, the regression quantile $\boldsymbol{\beta}_0(\tau)$ in model (1) is defined as the minimizer of the standard quantile loss function, $\sum_{i=1}^{n} \rho_\tau \{Y_i - (\mathbf{Z}_i^\top \boldsymbol{b}) \wedge U_i\}$, with respect to $\boldsymbol{b}$, where $\rho_\tau(x) = x\{\tau - I(x < 0)\}$ (Koenker & Bassett 1978).

When censoring is present and the censoring time $C$ is fixed at prespecified values, by utilizing the fact that $Q_{\widetilde{Y}}(\tau \mid \mathbf{Z}) = \{\mathbf{Z}^\top \boldsymbol{\beta}_0(\tau)\} \wedge U$, Powell (1984, 1986) proposed an adaptation of the standard quantile loss function, which led to an estimator of $\boldsymbol{\beta}_0(\tau)$ given by $\arg\min_{\boldsymbol{b}} r(\boldsymbol{b}, \tau)$, where

$$r(\boldsymbol{b}, \tau) = \sum_{i=1}^{n} \rho_\tau \left\{ \widetilde{Y}_i - \left( \mathbf{Z}_i^\top \boldsymbol{b} \right) \wedge U_i \right\}.$$

This estimation method is directly applicable to a more general case where $C$ is always known but not necessarily fixed, and is independent of $T$ given $\mathbf{Z}$. Unlike the standard quantile loss function, $r(\boldsymbol{b}, \tau)$ is not convex in $\boldsymbol{b}$ and thus it may have multiple local minima. Several authors, for example, Fitzenberger (1997), Buchinsky & Hahn (1998), Chernozhukov & Hong (2001), contributed strategies to improve the numerical performance of Powell's estimator. An implementation of Powell's method is available in the *crq* function in the R package *quantreg* (Koenker et al. 2019).

### 2.2.  Unconditionally random right censoring

Censoring time $C$ is not always observed in most survival settings. Under the assumption that $T$ and $C$ are independent and $C$ is independent of $\mathbf{Z}$ (i.e. unconditionally random right censoring), Ying et al. (1995) proposed to estimate $\boldsymbol{\beta}_0(\tau)$ by solving the following estimating equation,

$$n^{-1/2} \sum_{i=1}^{n} \mathbf{Z}_i \left[ \frac{I\left\{ \widetilde{Y}_i - \mathbf{Z}_i^\top \beta(\tau) > 0 \right\}}{\widehat{G}\left\{ \mathbf{Z}_i^\top \beta(\tau) \right\}} - (1 - \tau) \right] = 0, \tag{2}$$

where $\widehat{G}(\cdot)$ is the Kaplan-Meier estimate for $G(\cdot)$, which denotes the survival function of $U$. Since equation (2) is not continuous and may not have an exact zero-crossing, Ying et al. (1995) suggested obtaining the estimator of $\boldsymbol{\beta}_0(\tau)$ by minimizing the $L_2$ norm of the estimating function in (2). The resulting objective function however may have multiple minima.

Employing the inverse probability of censoring weighting (IPCW) technique (Robins & Rotnitzky 1992), Zhou (2006) studied an alternative estimating equation for $\boldsymbol{\beta}_0(\tau)$, which is given by

$$n^{-1/2} \sum_{i=1}^{n} Z_i \left[ \frac{I\left\{\tilde{Y}_i \leqslant Z_i^\top \beta(\tau), \Delta_i = 1\right\}}{\hat{G}(\tilde{Y}_i)} - \tau \right] = 0. \tag{3}$$

The estimating function in (3) is monotone (Fygenson & Ritov 1994). Consequently, the solution to equation (3) can be reformulated as the minimizer of a convex $L_1$-type convex function of $b$,

$$\sum_{i=1}^{n} \left\{ I(\Delta_i = 1) \left| \frac{\tilde{Y}_i}{\hat{G}(\tilde{Y}_i)} - b^\top \frac{Z_i}{\hat{G}(\tilde{Y}_i)} \right| \right\} + \left| M^* - b^\top \sum_{l=1}^{n} \left( -\frac{Z_l I(\Delta_l = 1)}{\hat{G}(\tilde{Y}_l)} + 2 Z_l \tau \right) \right|,$$

where $M^*$ is an extremely large positive number selected to bound $b^\top \sum_{l=1}^{n} \left( -\frac{Z_l I(\Delta_l = 1)}{\hat{G}(\tilde{Y}_l)} + 2 Z_l \tau \right)$ for all $b$'s in a compact parameter space. This minimization problem can be readily solved by the *l1fit()* function in S-PLUS or the *rq()* function in the R package *quantreg* (Koenker et al. 2019).

## 2.3. Conditionally random right censoring

Conditionally random right censoring, which assumes $C$ is independent of $T$ given $Z$, is the most commonly adopted random censoring mechanism. This censoring mechanism is less restrictive than those considered in Sections 2.1 and 2.2.

In this section, we first consider a globally linear quantile regression model,

$$Q_Y(\tau \mid Z) = Z^\top \beta_0(\tau), \ \tau \in [0, \tau_U], \tag{4}$$

and then discuss a locally linear quantile regression model with $\tau_L = \tau_U$.

**2.3.1.    Self-consistent approaches.—**By adapting the idea of redistributing censoring probability in the self-consistent Kaplan-Meier estimator (Efron 1967), Portnoy (2003) made the first attempt to estimate the globally linear quantile regression model (4) under the conditionally random right censoring assumption. The initial iterative self-consistent algorithm (Portnoy 2003) was simplified into a grid-based sequential estimation procedure (Neocleous et al. 2006), and the corresponding asymptotic studies was conducted by Portnoy & Lin (2010).

The grid-based estimation procedure of Neocleous et al. (2006) is outlined as follows. First, define a grid of $\tau$-values, $\mathscr{G}_n$, as $0 = \tau_0 < \tau_1 < \dots < \tau_{M_n} = \tau_U$. Let $\|\mathscr{G}_n\| = \max\{\tau_k - \tau_{k-1}: k = 1, \dots, M_n\}$. Without further mentioning, $\mathscr{G}_n$ will be adopted throughout Section 2.3. Assuming no censoring occurs below the $\tau_1$-th conditional quantile of $T$, one can obtain an estimate for $\beta_0(\tau_1)$ from applying uncensored quantile regression. Next, one can estimate $\beta_0(\tau_{k+1})$ sequentially for $k = 1, 2, \dots, M_n$ by minimizing

$$\sum_{i \in K^c} \rho_{\tau_{k+1}}(\widetilde{Y}_i - Z_i^\top b) + \sum_{i \in K}\left\{w_{k+1,i} \cdot \rho_{\tau_{k+1}}(\widetilde{Y}_i - Z_i^\top b)\right.$$
$$\left. + (1 - w_{k+1,i})\rho_{\tau_{k+1}}(Y^* - Z_i^\top b)\right\}, \tag{5}$$

where $Y^*$ is an extremely large value and $K$ denotes the set of indices of censored observations that have been previously crossed (i.e. $C_i \leqslant Z_i^\top \widehat{\beta}(\tau)$). The weight $w_{k+1,i}$ takes the form $(\tau_{k+1} - \tau_l)/(1 - \tau_l)$ to approximate the conditional probability, $\Pr(C_i < T_i < \exp\{Z_i\beta_0(\tau_{k+1})\}|C_i < T_i, Z_i)$ based on the estimates for $\beta_0(\tau_1), \ldots, \beta_0(\tau_k)$.

Peng (2012) proposed alternative formulations of the self-consistent approach based on stochastic integral equations. First, using stochastic integral formulation and applying stochastic integral by parts, Efron (1967)'s self-consistent estimating equation for $F_Y(t)$ in the one-sample case can be re-expressed as

$$F_Y(t) = n^{-1} \sum_{i=1}^{n}\left[N_i(t) + R_i(t)\{1 - F_Y(t)\}\int_0^t \frac{R_i(u)}{\{1 - F_Y(u)\}^2}dF_Y(u)\right],$$

where $N_i(t) = I(\widetilde{Y}_i \leqslant t, \Delta_i = 1)$, $R_i(t) = I(\widetilde{Y}_i \leqslant t, \Delta_i = 0)$, and $F_Y(t) = \Pr(Y \leqslant t)$.

With $t$ replaced by $Z_i^\top \beta(\tau)$, this equation evolves into an estimating equation for $\beta_0(\tau)$,

$$n^{1/2}S_n^{(SC)}(\beta, \tau) \doteq n^{-1/2}\sum_{i=1}^{n}Z_i\left[N_i\{Z_i^\top\beta(\tau)\} + R_i\{Z_i^\top\beta(\tau)\}(1 - \tau\right.$$
$$\left.)\int_0^\tau \frac{R_i\{Z_i^\top\beta(u)\}}{(1-u)^2}du - \tau\right] = 0. \tag{6}$$

Peng (2012) further justified several asymptotically equivalent variants of the estimating equation (6), one of which takes the form of

$$n^{1/2}S_n^{(MSC)}(\beta, \tau) \doteq n^{-1/2}\sum_{i=1}^{n}Z_i\left[N_i\{Z_i^\top\beta(\tau)\} + R_i\{Z_i^\top\beta(\tau)\}\frac{\tau - \psi_i(\beta,\tau)}{1 - \psi_i(\beta,\tau)} - \tau\right]$$
$$= 0. \tag{7}$$

Here $\psi_i(\beta, \tau) = \sup\{A_i(\beta, \tau)\} \cdot I(A_i(\beta, \tau) \text{ is not empty}) + \tau I(A_i(\beta, \tau) \text{ is empty})$ with $A_i(\beta, \tau) = \left\{u : 0 \leqslant u < \tau, \ Z_i^\top\beta(u-) \leqslant \widetilde{Y}_i \leqslant Z_i^\top\beta(u)\right\}$. The estimation procedure derived from equation (7) is identical to Neocleous et al. (2006)'s procedure except for the estimation of $\beta_0(\tau_1)$. That is, Neocleous et al. (2006)'s procedure estimates $\beta_0(\tau_1)$ as the minimizer of (5) with $w_{1,i} = 1$, while estimating $\beta_0(\tau_1)$ based on equation (7) is equivalent to minimizing (5) with $w_{1,i} = 0$.

Large sample studies for the estimator $\widehat{\beta}_{SC}(\tau)$ obtained from solving equation (6) are facilitated by the stochastic integral equation representation of $S_n^{(SC)}(\beta, \tau)$. Specifically,

under certain regularity conditions and given $\lim_{n \to \infty} \|\mathscr{G}_n\| = 0$,

$\sup_{\tau \in [\nu, \tau_U]} \|\widehat{\beta}_{SC}(\tau) - \beta_0(\tau)\| \to_p 0$, where $0 < \nu < \tau_U$. If $n^{1/2} \lim_{n \to \infty} \|\mathscr{G}_n\| = 0$ is further

satisfied, then $n^{1/2}\left\{\widehat{\beta}_{SC}(\tau) - \beta_0(\tau)\right\}$ converges weakly to a Gaussian process for $\tau \in [\nu, \tau_U]$.

Estimator defined based on (7) is asymptotically equivalent to $\widehat{\beta}_{SC}(\cdot)$.

**2.3.2. Martingale-based approach.**—Peng & Huang (2008) proposed to utilize the
martingale structure underlying randomly censored data to construct an estimating equation
for model (4). Define $\Lambda_Y(t|\mathbf{Z}) = -\log\{1 - \Pr(Y \leq t|\mathbf{Z})\}$, $N(t) = I(\widetilde{Y} \leq t, \Delta = 1)$, and $M(t) = N(t) - \Lambda_Y(t \wedge Y|\mathbf{Z})$. Let $N_i(t)$ and $M_i(t)$ be sample analogs of $N(t)$ and $M(t)$ respectively, $i = 1, \ldots, n$. Note that $M_i(t)$ is the martingale process associated with the counting process $N_i(t)$.
Thus, $E\{M_i(t)|\mathbf{Z}_i\} = 0$ for all $t > 0$, and

$E\left\{\sum_{i=1}^{n} \mathbf{Z}_i\left[N_i\left\{\mathbf{Z}_i^\top\beta_0(\tau)\right\} - \Lambda_Y\left\{\mathbf{Z}_i^\top\beta_0(\tau) \wedge \widetilde{Y}_i \mid \mathbf{Z}_i\right\}\right]\right\} = 0$ for $\tau[0, \tau_U]$. Since $\mathbf{Z}_i^\top\beta_0(\tau)$ is

monotone in $\tau \in [0, \tau_U]$, we have $\Lambda_Y\left\{\mathbf{Z}_i^\top\beta_0(\tau) \wedge \widetilde{Y}_i \mid \mathbf{Z}_i\right\} = \int_0^\tau I\left\{\widetilde{Y}_i \geq \mathbf{Z}_i^\top\beta_0(u)\right\}dH(u)$, where

$H(x) = -\log(1 - x)$. These findings suggest a stochastic integral based estimating equation,

$$n^{1/2}\mathbf{S}_n^{(PH)}(\beta, \tau) \doteq n^{-1/2}\sum_{i=1}^{n} \mathbf{Z}_i\left[N_i\left\{\mathbf{Z}_i^\top\beta(\tau)\right\} - \int_0^\tau I\left\{\widetilde{Y}_i \geq \mathbf{Z}_i^\top\beta(u)\right\}dH(u)\right] = 0. \quad (8)$$

An estimator of $\beta_0(\tau)$, denoted by $\widehat{\beta}_{PH}(\tau)$, can be obtained through approximating the
stochastic solution to equation (8). Specifically, let $\widehat{\beta}_{PH}(\tau)$ be a cadlag (i.e. right continuous
with finite left limits) step function of $\tau$ that jumps only at the grid points of $\mathscr{G}_n$. The
procedure to obtain $\widehat{\beta}_{PH}(\tau)$ follows.

1. Set $\exp\left\{\mathbf{Z}_i^\top\widehat{\beta}_{PH}(\tau_0)\right\} = 0$ for all $i$. Set $k = 0$.

2. Given $\exp\left\{\mathbf{Z}_i^\top\widehat{\beta}_{PH}(\tau_l)\right\}$ for $l \leq k$, obtain $\widehat{\beta}_{PH}(\tau_{k+1})$ as the minimizer of the
   following $L_1$-type convex objective function:

   $$l_{k+1}(\mathbf{h}) = \sum_{i=1}^{n}\left|\Delta_i\widetilde{Y}_i - \Delta_i\mathbf{Z}_i^\top\mathbf{h}\right| + \left|Y^* - \sum_{l=1}^{n}\left(-\Delta_l\mathbf{Z}_l^\top\mathbf{h}\right)\right|$$
   $$+ \left|Y^* - \sum_{r=1}^{n}\left[\left(2\mathbf{Z}_r^\top\mathbf{h}\right)\sum_{l=0}^{k} I\left\{\widetilde{Y}_r \geq \mathbf{Z}_r^\top\widehat{\beta}_{PH}(\tau_l)\right\}\{H(\tau_{l+1}) - H(\tau_l)\}\right]\right|,$$

   where $Y^*$ is an extremely large value.

3. Replace $k$ by $k + 1$ and repeat step 2 until $k = M_n$ or no feasible solution can be
   found for minimizing $l_k(\mathbf{h})$.

Peng & Huang (2008) established the uniform consistency and weak convergence of
$\widehat{\beta}_{PH}(\cdot)$. Moreover, Peng (2012) showed that $\widehat{\beta}_{PH}(\cdot)$ is asymptotically equivalent to the
self-consistent estimator $\widehat{\beta}_{SC}(\cdot)$ in that $\sup_{\tau \in [\nu, \tau_U]}\left\|n^{1/2}\left\{\widehat{\beta}_{PH}(\tau) - \widehat{\beta}_{SC}(\tau)\right\}\right\| = o_p(1)$ with $0 < \nu < \tau_U$. This theoretical result is consistent with the numerical results reported in Koenker (2008) and Peng (2012).

The *crq()* function in the R package *quantreg* (Koenker et al. 2019) provides an implementation of $\widehat{\beta}_{PH}(\tau)$ based on an algorithm slightly different from the one presented above. An asymptotically equivalent grid-free estimation procedure for model (4) was developed by Huang (2010).

**2.3.3.    Data augmentation approach.**—Yang et al. (2018) employed a variation of the data augmentation algorithm to tackle the estimation of model (4) with $\tau_U = 1$. The basic idea is to apply the general principle of data augmentation (Tanner & Wong 1987), and employ an alternating process between imputation of censored values from the quantile functions and refitting of the quantile model using the imputed values. More specifically, the algorithm starts with a set of initial values, $\widehat{\beta}^{(0)}(\tau_k)(k = 1, ..., M_n)$, obtained by parallel quantile regression estimators or existing quantile regression estimators. For $h = 1, ..., H$, draw $Y_i^{(h)}$ from the quantile process approximated by $Z_i^T\widehat{\beta}^{(h-1)}(\tau_k)$ conditional on the set of possible values for $Y_i$. Then, based on a pairwise bootstrapping sample of size $n$ from $\left\{Z_i, Y_i^{(h)}\right\}_{i=1}^n$, obtain updated estimates $\widehat{\beta}^{(h)}(\tau_k)$ via standard uncensored quantile regression. Lastly, take the final estimates as $\widehat{\beta}(\tau) = H^{-1}\sum_{h=1}^H \widehat{\beta}^{(h)}(\tau)$.

An appealing feature of Yang et al. (2018)'s approach is that it can handle different forms of censoring, including random censoring, double censoring, and interval censoring. As demonstrated by Monte Carlo simulations, the resulting estimator can achieve significant efficiency gains over the existing methods. The algorithm of Yang et al. (2018) is implemented by the R function *DArq()*.

**2.3.4.    Adjusted loss function methods.**—Assume a locally linear quantile regression model, which is model (1) with $\tau_L = \tau_U$ equal to a prespecified $\tau$, i.e.

$$Q_Y(\tau \mid Z) = Z^\top\beta_0(\tau).$$                                                  (9)

To account for random censoring in the estimation of model (9), Wang & Wang (2009) proposed to modify the standard quantile loss function by twisting the idea of the self-consistent Kaplan-Meier estimator (Efron 1967). That is, one may redistribute the probability mass associated with each censored case, $\Pr(T_i > C_i|C_i, Z_i)$, to the right through a local weighting scheme by $w_i(F_0)$, where

$$w_i(F_0) = \begin{cases} 1 & \Delta_i = 1 \text{ or } F_0(C_i \mid Z_i) > \tau \\ \dfrac{\tau - F_0(C_i \mid Z_i)}{1 - F_0(C_i \mid Z_i)} & \Delta_i = 0 \text{ and } F_0(C_i \mid Z_i) < \tau \end{cases}$$

with $F_0(t|z) = \Pr(T > t|Z = z)$. Suppose $F_0(t|z)$ is known. An estimator of $\beta_0(\tau)$ in model (9) can be obtained by minimizing the following objective function of $\beta$:

$$n^{-1}\sum_{i=1}^n \left[ w_i(F_0)\rho_\tau\left(\widetilde{Y}_i - Z_i^\top\beta\right) + \{1 - w_i(F_0)\}\rho_\tau\left(Y^* - Z_i^\top\beta\right)\right].$$             (10)

In practice, $F_0(t|z)$ is usually unknown. In this case, Wang & Wang (2009) proposed to minimize the objective function (10) with $F_0(\cdot)$ replaced by $\widehat{F}(\cdot)$, the local Kaplan-Meier estimator, namely,

$$\widehat{F}(t \mid z) = 1 - \prod_{j=1}^{n} \left\{ 1 - \frac{B_{nj}(z)}{\sum_{k=1}^{n} I(\widetilde{Y}_k \geqslant \widetilde{Y}_j) B_{nk}(z)} \right\}^{N_j(t)} .$$

Here $B_{nk}(z)$ is a sequence of nonnegative weights adding up to 1, for example, Nadaraya-Watson's type weight, $B_{nk}(z) = K\left(\frac{z - z_k}{h_n}\right) / \sum_{i=1}^{n} K\left(\frac{z - z_i}{h_n}\right)$, where $K(\cdot)$ is a density kernel function and $h_n$ is a positive bandwidth converging to 0 as $n \to \infty$. The resulting estimator is shown to be consistent and asymptotically normal with root $n$ rate under regularity conditions.

De Backer et al. (2019) and De Backer et al. (2020) investigated different strategies for adjusting the standard quantile loss function in order to accommodate randomly censored data. More specifically, letting $G_U(u|\mathbf{Z}) = \Pr(U > u|\mathbf{Z})$, De Backer et al. (2019) noted that the derivative of $\phi_\tau(a; Y, G_U(\cdot \mid \mathbf{Z})) \doteq (Y - a)\{\tau - I(Y \leqslant a)\} - (1 - \tau)\int_0^a \{1 - G_U(s \mid \mathbf{Z})\}ds$ with respect to $a$ equals $-\{I(Y > a) - G_U(a|\mathbf{Z})(1 - \tau)\}$, which, conditional on $\mathbf{Z}$, has expectation zero with $a = \mathbf{Z}^\top \beta_0(\tau)$ under model (9). This key fact leads to an adjusted loss function for censored quantile regression,

$$\sum_{i=1}^{n} \phi_\tau\left(\mathbf{Z}_i^\top \boldsymbol{\beta}; Y_i, \widehat{G}_U(\cdot \mid \mathbf{Z}_i)\right), \tag{11}$$

where $\widehat{G}_U(\cdot \mid z)$ is a consistent estimator of $G_U(\cdot|z)$. When $C$ is independent $\mathbf{Z}$, $\widehat{G}_U(\cdot \mid z)$ can be obtained through the Kaplan-Meier estimator of the survival distribution of $C$. Without assuming the independence between $C$ and $\mathbf{Z}$, $\widehat{G}_U(\cdot \mid z)$ can be obtained through semiparametric modeling of $C$ given $\mathbf{Z}$, or by directly using Beran's conditional Kaplan-Meier estimator (Beran 1981). De Backer et al. (2019) developed a numerically robust MM algorithm to solve the minimization of the non-convex adjusted loss function (11). Following a different view, De Backer et al. (2020) proposed to estimate model (9) based on a minimum distance loss function, given by $\sum_{i=1}^{n} \left\{ 1 - \widehat{F}\left(\mathbf{Z}_i^\top \beta(\tau) \mid \mathbf{Z}_i\right) - \tau \right\}^2$. De Backer et al. (2020) further suggested using a smooth double kernel version of $\widehat{F}(\cdot \mid \mathbf{Z}_i)$. They also discussed how to handle high-dimensional covariates by employing the effective dimension reduction technique (Li et al. 1999, Xia et al. 2010). Desirable asymptotic properties, consistency and asymptotic normality, were established for these estimators of $\boldsymbol{\beta}_0(\tau)$ in model (9).

## 2.4. Inference procedures

### 2.4.1. Variance estimation.—Bootstrapping procedures have been justified and commonly used to make inferences under quantile regression with either uncensored response or censored survival response. For example, to estimate the asymptotic variance of

the estimators discussed in sections 2.1–2.3, one may use resampling methods that follow the idea of Parzen & Ying (1994) or apply the standard bootstrapping procedures that use resampling with replacement (Koenker 2005, Peng & Huang 2008).

Alternative methods without involving resampling have been developed for variance estimation under quantile regression. A main challenge is how to estimate the unknown densities involved in the formulas for asymptotic variances. Under random right censoring with known censoring time or unconditionally random censoring, Huang (2002)'s technique can be directly applied to avoid smoothing-based density estimation, which may be unstable with small or moderate sample sizes. Specifically, let $\widehat{\beta}(\tau)$ denote an estimator of $\beta_0(\tau)$, and $S_n\{\widehat{\beta}(\tau), \tau\}$ denote the estimating function associated with $\widehat{\beta}(\tau)$, for example, the left-hand side of (2) and (3). Generally it can be shown that $S_n\{\beta_0(\tau), \tau\}$ converges to a mean-zero multivariate normal distribution with covariance matrix $\Sigma(\tau)$, which may be consistently estimated by $\widehat{\Sigma}(\tau)$. The following are the main steps to obtain a sample-based variance estimator:

**A.1**  Find a symmetric and nonsingular $(p+1) \times (p+1)$ matrix $E_n(\tau) \doteq \{e_{n,1}(\tau), \ldots, e_{n,p+1}(\tau)\}$ such that $\widehat{\Sigma}(\tau) = \{E_n(\tau)\}^2$.

**A.2**  Calculate $D_n(\tau) = \left(S_n^{-1}\{e_{n,1}(\tau), \tau\} - \widehat{\beta}(\tau), \ldots, S_n^{-1}\{e_{n,p+1}(\tau), \tau\} - \widehat{\beta}(\tau)\right)$, where $S_n^{-1}(e, \tau)$ is defined as the solution to $S_n(b, \tau) - e = 0$.

**A.3**  Estimate the asymptotic variance matrix of $n^{1/2}\left\{\widehat{\beta}(\tau) - \beta_0(\tau)\right\}$ by $n\{D_n(\tau)\}^{\otimes 2}$.

Under conditionally random censoring, the self-consistent estimators and the martingale-based estimator for model (4) take much more complex forms than those developed under the stronger censoring mechanism with either known censoring time or unconditionally independent censoring. To estimate the asymptotic variances of these estimators, it requires much more sophisticated twists of Huang (2002)'s technique to address the challenge associated with unknown densities. A sample-based variance estimation procedure for Peng & Huang (2008)'s estimator is available through adapting Sun et al. (2016)'s sample-based inference procedure for recurrent events data to the setting with randomly censored data.

**2.4.2.  Second-stage inference.**—Globally linear quantile regression model (4) provides a platform to explore the varying pattern of covariate effects across different quantile levels. Second-stage inference can be performed to address such interests. For example, one may estimate a functional of $\beta_0(\cdot)$, say $\Psi(\beta_0)$, to provide a meaningful summary of covariate effects over a range of $\tau$. It is often of interest to determine whether some covariates have constant effects so that a simpler model may be considered. In this case, the problem can be formulated as testing the null hypothesis $H_{0,j} : \beta_0^{(j)}(\tau) = \rho_0, \tau \in [\tau_L, \tau_U]$, where the superscript $(j)$ indicates the $j$th component of a vector, and $\rho_0$ is an unspecified constant, $j = 2, \ldots,$ or $p+1$. Of note, accepting $H_{0,j}$ for all $j \in \{2, \ldots, p+1\}$ may indicate the adequacy of an AFT model. Peng & Huang (2008) presented second-stage inference procedures for estimating $\Psi(\beta_0)$ and testing $H_0$ under model (4), which can be readily adapted to many other quantile regression settings.

# 3.   Quantile regression in complex survival settings

In practice, survival data often involve complications beyond random censoring, such as truncation, competing risks or semi-competing risks. Various methods have developed for quantile regression in more complex survival scenarios. In this section, we present a set of quantile regression methods developed for analyzing for doubly censored data with left truncation, competing risks data, and semi-competing risks data.

## 3.1.   Quantile regression with doubly censored data with left truncation

Ji et al. (2012) proposed an extension of Peng & Huang (2008)'s method to handle doubly censored data subject to left truncation. Such survival scenarios often arise in observational studies, where the event of interest can occur before study entry. Let $T$ denote the event time of interest and $C$ denote time to random right censoring. In addition, let $L$ denote left censoring time, always observed, and $A$ denote left truncation time. Define $X = L \vee (T \wedge C)$ and     as the censoring indicator which equals 1 if $L < T \leqslant C$, 2 if $T \leqslant L$, and 3 if $T > C$, where $\vee$ is the maximum operator. When $X$ is subject to left truncation by $A$, the observed data include $n$ i.i.d. replicates of $(X', L', A', \delta', Z)$, denoted by $\{(X'_i, L'_i, A'_i, \delta'_i, Z_i)\}^n_{i=1}$, where $\{X', L', A', \delta', Z'\}$ follows the conditional distribution of $\{X, L, A, \delta, Z\}$ given $X \geqslant A$. It is assumed that $(L, C, A)$ is independent of $T$ given $Z$. We refer to such data as doubly censored data with left truncation. With $L = 0$, the data reduce to the usual randomly left truncated right censored data.

To estimate model (4) with doubly censored data subject to left truncation, an estimating equation can be constructed based on the martingale structure underlying the observed survival data, namely, $M'(t) = N'(t) - \int_0^t R'(s) d\Lambda_Y(s \mid Z)$, where $N'(t) = I(\log(X') \leqslant t, \delta' = 1)$, $R'(t) = I\{\log(L' \vee A') < t \leqslant \log(X')\}$ denoting an at-risk process, and $\Lambda_Y(\cdot | Z)$ denotes the cumulative hazard function of $Y \doteq \log(T)$ given $Z$. It can be shown that $M'(t)$ is a martingale process. This fact suggests an estimating equation for $\beta_0(\cdot)$,

$$n^{1/2} S'_n(\beta, \tau) \doteq n^{-1/2} \sum_{i=1}^n Z_i \left( N'_i \left\{ Z_i^\top \beta(\tau) \right\} - \int_0^\tau R' \left\{ Z_i^\top \beta(u) \right\} dH(u) \right) = 0. \qquad (12)$$

To obtain an estimator of $\beta_0(\tau)$ based on equation (12), denoted by $\hat{\beta}_{PH, *}(\tau)$, one may follow the algorithm for $\hat{\beta}_{PH}(\tau)$ (presented in section 2.3) with the objective function in Step 2 modified to

$$l^*_{k+1}(h) = \sum_{i=1}^n \left| I(\Delta_i = 1)\log(X'_i) - I(\Delta_i = 1)Z_i^\top h \right| + \left| Y^* - \sum_{l=1}^n \left( -I(\Delta_l = 1)Z_l^\top h \right) \right|$$
$$+ \left\| Y^* - \sum_{r=1}^n \left[ \left( 2Z_r^\top h \right) \sum_{l=0}^k I \left\{ \log(X'_r) \geqslant Z_r^\top \hat{\beta}_{PH, *}(\tau_l) \geqslant \log(L'_r) \right\} \{H(\tau_{l+1}) - H(\tau_l)\} \right] \right\|.$$

Theoretical properties, such as uniform consistency and weak convergence to a Gaussian process, can be established for $\hat{\beta}_{PH, *}(\tau)$ with similar lines of Peng & Huang (2008).

### 3.2. Quantile regression with competing risks data

Competing risks data arise in scientific studies involving multiple types of failures that are mutually exclusive. For example, a cancer patient may die from tumor recurrence or nonrecurrence-related reasons. This gives rise to a competing risks scenario, where death from tumor recurrence and death from nonrecurrence-related reasons are two competing failure types.

We adopt standard formulation of competing risks data. Let $T_k$ denote the latent time to failure of type $k$ ($k = 1, \ldots, K$). Define $T = \min(T_1, \ldots, T_K)$. Let $\epsilon$ denote the failure type corresponding to $T$ (i.e. $T = T_\epsilon$), $C$ denote independent censoring to $T$, and $\widetilde{Z}$ denote a $p{\times}1$ vector of covariates. Define $X = T \wedge C$, $\delta = I(T \leqslant C)\epsilon$, and $Z = \left(1, \widetilde{Z}^{\top}\right)^{\top}$. Here $\wedge$ is the minimum operator and $I(\cdot)$ is the indicator function. The observed competing risks data consist of $n$ iid replicates of $(X, \delta, Z)$, denoted by $\{(X_i, \delta_i, Z_i), i = 1, \ldots, n\}$.

Analysis of competing risks data generally follows two different perspectives. One perspective focuses on crude quantities, such as the cumulative incidence function or cause-specific hazard function. Studying crude quantities for a failure type naturally accounts for the presence of competing risks from the other types of failure. The other perspective concerns net quantities defined upon latent failure times $T_k$'s. Inference on the latent failure time for a failure type however implicitly hypothesizes a setting where the other types of failure do not exist. Such a setting may be controversial but can be meaningful in some situations. For example, patient dropouts can be a competing risk for time to death but may be avoided by diligent follow-up efforts. When the elimination of other types of failures is not possible, competing risks analysis oriented to crude quantities would be more appropriate. In the following, we discuss quantile regression methods for competing risks data developed under these two different perspectives.

### 3.2.1. Competing risks quantile regression based on cumulative incidence functions.—Peng & Fine (2009) proposed to formulate competing risks quantile regression using cumulative incidence function, which is the cause-specific analog of the usual survival function for an event time. Specifically, the type-$k$ cumulative incidence conditional quantile function is defined as $Q_k(\tau|Z) \doteq \inf\{t : F_k(t|Z) \geqslant \tau\}$, where $F_k(t|Z) \doteq \Pr(T \leqslant t, \epsilon = k|Z)$ denotes the type-$k$ cumulative incidence function ($k = 1, \ldots, K$). This quantity can be interpreted as the first time given covariate $Z$ at which the probability of type-$k$ failure having occurred exceeds $\tau$, in the presence of other types of failures.

A competing risks quantile regression model based on type-$k$ cumulative incidence function takes the form,

$$Q_k(\tau \mid Z) = \exp\left\{Z^T \beta_0(\tau)\right\}, \ \tau \in [\tau_L, \tau_U], \tag{13}$$

where $\beta_0(\tau)$ is a $(p{+}1){\times}1$ vector of unknown regression coefficients, and $0 \leqslant \tau_L \leqslant \tau_U < 1$. Under model (13), the non-intercept coefficients in $\beta_0(\tau)$ represent covariate effects on the $\tau$th-cumulative incidence quantile, $Q_k(\tau|Z)$, which may change with $\tau$. The $\exp(\cdot)$ function in (13) can be replaced by any other monotone link function.

To estimate $\boldsymbol{\beta}_0(\tau)$ in model (13), Peng & Fine (2009) proposed the following estimating equation,

$$n^{-1/2} \sum_{i=1}^{n} \boldsymbol{Z}_i \left( \frac{I\left\{ X_i \leqslant \exp\left( \boldsymbol{Z}_i^T \boldsymbol{b} \right) \right\} I(\delta_i = 1)}{\widehat{G}(X_i \mid \boldsymbol{Z}_i)} - \tau \right) = 0, \tag{14}$$

where $\widehat{G}(\,\cdot\, \mid \boldsymbol{Z})$ is a reasonable estimate for $G(x|\boldsymbol{Z}) \doteq \Pr(C \geqslant x|\boldsymbol{Z})$, which can be obtained by following the discussions about $\widehat{G}_U(\,\cdot\, \mid \boldsymbol{Z})$ in Section 2.3.4.

Solving equation (14) can be reformulated as locating the minimizer of the convex $L_1$-type function,

$$\sum_{i=1}^{n} I(\delta_i = 1) \left| \frac{\log(X_i)}{\widehat{G}(X_i \mid \boldsymbol{Z}_i)} - \boldsymbol{b}^T \frac{\boldsymbol{Z}_i}{\widehat{G}(X_i \mid \boldsymbol{Z}_i)} \right| + \left| M^* - \boldsymbol{b}^T \sum_{l=1}^{n} \frac{-\boldsymbol{Z}_l I(\delta_l = 1)}{\widehat{G}(X_l \mid \boldsymbol{Z}_l)} \right| + \left| M^* - \boldsymbol{b}^T \sum_{k=1}^{n} (2\boldsymbol{Z}_k \tau) \right|,$$

where $M^*$ is an extremely large positive number.

Peng & Fine (2009) showed that the resulting estimator is uniformly consistent in $\tau \in [\tau_L, \tau_U]$, and converges weakly to a tight mean-zero Gaussian process. They developed inference procedures about $\boldsymbol{\beta}_0(\tau)$ in model (13), which follow similar lines to those presented in section 2.4 for randomly censored data with known or unconditional independent censoring. Following the same framework, Sun et al. (2012) studied model (13) for the competing risks setting with missing failure types, where IPCW technique was used to to deal with unobserved failure types under the missing at random assumption.

### 3.2.2. Quantile regression based on latent failure time distributions in the presence of competing risks.—The analysis of competing risks data based on net quantities, such as the marginal distributions of $T_k$'s ($k = 1, \ldots, K$), is complicated by their nonparametric nonidentifiability (Tsiatis 1975). Without loss of generality, we consider the situation with $K = 2$. This special case coincides with the typical dependent censoring scenario, where the dependent censoring event and the event of interest can be viewed as a pair of competing risks.

Concerning the latent failure times $T_1$ and $T_2$, one may consider the following quantile regression models:

$$Q_{T_k}(\tau \mid \boldsymbol{Z}) = \exp\left\{ \boldsymbol{Z}^\top \boldsymbol{\beta}_{0, k}(\tau) \right\}, \ \tau \in (0, 1), \ k = 1, 2, \tag{15}$$

where $\boldsymbol{\beta}_{0,k}(\tau)$ is an vector of unknown coefficients, representing covariate effects on $Q_{T_k}(\tau \mid \boldsymbol{Z})$. Here, $\exp(\cdot)$ can be replaced by another monotone link function, which may take different forms in the models for $Q_{T_1}(\tau \mid \boldsymbol{Z})$ and $Q_{T_2}(\tau \mid \boldsymbol{Z})$.

Ji et al. (2014) studied the estimation of the marginal quantile regression models (15) with competing risks data. To mitigate the identifiability issue, additional modeling is imposed for the dependence structure between $T_1$ and $T_2$. Specifically, it is assumed that

$$\Pr(T > t_1, D > t_2 \mid \mathbf{Z}) = H\{\Pr(T > t_1 \mid \mathbf{Z}), \Pr(D > t_2 \mid \mathbf{Z})\}, \tag{16}$$

where $H(\cdot, \cdot)$ is a known copula function, for example, the Clayton copula (Clayton 1978),

i.e. $H(u, v) = \{u^{-r} + v^{-r} - 1\}^{-\frac{1}{r}}$, $r > 0$, and the Frank copula (Genest 1987), i.e.,

$H(u, v) = \log_r\left\{1 + \frac{(r^u - 1)(r^v - 1)}{r - 1}\right\}$, $r > 0$ and $r \ne 1$. Here $r$ is a known copula parameter, which

may be specified based on prior knowledge on the strength of the association between $T_1$
and $T_2$. In practice, one may perform a sensitivity analysis to obtain bounds for $Q_T(\tau|\mathbf{Z})$ by
varying $r$ in a plausible range.

To estimate $\boldsymbol{\beta}_{0,k}(\tau)$ in (15), Ji et al. (2014) utilized the martingales associated with cause-specific hazard functions. Let $N_k(t) \doteq I(X \leqslant t, \epsilon = k)$ denote the counting process for $T_k$ and define $M_k(t) = N_k(t) - \int_0^t I(X \geqslant u)\lambda_k^*(u \mid \mathbf{Z})\mathrm{d}u$, where

$\lambda_k^*(t \mid \mathbf{Z}) = \lim_{h \to 0} \Pr(t \leqslant T_k < t + h, \epsilon = k \mid T_1 \geqslant t, T_2 \geqslant t; \mathbf{Z})/h$, which is the cause-specific

hazard function for type-$k$ failure. As shown by Kalbfleisch & Prentice (2002), $M_k(t)$ is a martingale with respect to the filtration, $\mathscr{F}_{t, k} = \{N_k(t), Y(t+), \mathbf{Z}\}$. This implies

$E\left\{N_k(t) - \int_0^t I(X \geqslant s)\lambda_k^*(s \mid \mathbf{Z})\mathrm{d}s\right\} = 0$ for all $t \geqslant 0$. Under models (15) and (16), it can be

shown with stochastic integral manipulations that

$$\int_0^t I(X \geqslant s)\lambda_k^*(s \mid \mathbf{Z})\mathrm{d}s = \int_0^{F_{T_k}(t \mid \mathbf{Z})} I\{X \geqslant Q_{T_k}(u \mid \mathbf{Z})\}\phi_k$$
$$\left(1 - u, 1 - \int_0^1 I\left[\exp\left\{\mathbf{Z}_i^T \boldsymbol{\beta}_{0, 3 - k}(v)\right\} \leqslant Q_{T_k}(u \mid \mathbf{Z})\right]\mathrm{d}v\right)\mathrm{d}u,$$

where $F_{T_k}(t \mid \mathbf{Z}) = \Pr\{T_k \leqslant t \mid \mathbf{Z}\}$, $\phi_k(v_1, v_2) = \partial \log\{H(v_1, v_2)\}/\partial v_k$, and $k = 1, 2$. These facts motivate the estimating equations,

$$n^{\frac{1}{2}}S_n^{(k)}(\boldsymbol{\beta}_1, \boldsymbol{\beta}_2, \tau) = 0, \ k = 1, 2, \tag{17}$$

where $S_n^{(k)}(\boldsymbol{\beta}_1, \boldsymbol{\beta}_2, \tau) = n^{-1}\sum_{i=1}^n \mathbf{Z}_i\left\{N_{ki}\left[\exp\left\{\mathbf{Z}_i^\top\boldsymbol{\beta}_k(\tau)\right\}\right] - \int_0^\tau I\left(X_i \geqslant \exp\left\{\mathbf{Z}_i^\top\boldsymbol{\beta}_{3-k}(u)\right\}\right).\right.$
$\left. \times \phi_k\left(1 - u, 1 - \int_0^1 I\left\{\mathbf{Z}_i^T\boldsymbol{\beta}_{3-k}(v) \leqslant \mathbf{Z}_i^\top\boldsymbol{\beta}_k(u)\right\}\mathrm{d}v\right)\mathrm{d}u\right\}$

Note that $\boldsymbol{\beta}_{0,k}(\tau)$ may not be identifiable for all $\tau \in (0, 1)$ due to censoring to $T_k$ ($k = 1, 2$). Truncating the time scale by an upper bound, $\min(\exp\{\mathbf{Z}^\top\boldsymbol{\beta}_{0,1}(\tau_{U,1})\}, \exp\{\mathbf{Z}^\top\boldsymbol{\beta}_{0,2}(\tau_{U,2})\})$, leading to a modified estimating equation,

$$n^{\frac{1}{2}}S_n^{*(k)}(\boldsymbol{\beta}_1, \boldsymbol{\beta}_2, \tau) = 0, \ k = 1, 2, \tag{18}$$

where $S_n^{*(k)}(\boldsymbol{\beta}_1, \boldsymbol{\beta}_2, \tau) = n^{-1}\sum_{i=1}^n \mathbf{Z}_i\left\{N_{ki}\left[\exp\left\{\mathbf{Z}_i^\top\boldsymbol{\beta}_k(\tau)\right\}\right]I\left\{\log(X_i) \leqslant \mathbf{Z}_i^\top\boldsymbol{\beta}_{3-k}(\tau_{U, 3-k})\right\}.\right.$
$\left. - \int_0^\tau I\left(X_i \geqslant \exp\left\{\mathbf{Z}_i^\top\boldsymbol{\beta}_k(u)\right\}\right)I\left\{\mathbf{Z}_i^\top\boldsymbol{\beta}_k(u) \leqslant \mathbf{Z}_i^\top\boldsymbol{\beta}_{3-k}(\tau_{U, 3-k})\right\}\right.$
$\left. \times \phi_k\left(1 - u, 1 - \int_0^{\tau_{U, 3-k}} I\left\{\mathbf{Z}_i^\top\boldsymbol{\beta}_{3-k}(v) \leqslant \mathbf{Z}_i^\top\boldsymbol{\beta}_k(u)\right\}\mathrm{d}v\right)\mathrm{d}u\right\}$

Equations in (18) may be solved via an iterative algorithm:

**Step B.1**     Set $m = 0$. Choose the initial value $\widehat{\boldsymbol{\beta}}_2^{[m]}(\tau), \tau \in (0, \tau_{U,2}]$.

**Step B.2**     Solve $S_n^{*(1)}\left(\boldsymbol{\beta}_1, \widehat{\boldsymbol{\beta}}_2^{[m]}, \tau\right) = 0$ for $\widehat{\boldsymbol{\beta}}_1^{[m+1]}(\tau), \tau \in \left(0, \tau_{U,1}^{[m+1]}\right]$. Update $\tau_{U,1}$ with $\tau_{U,1}^{[m+1]}$.

**Step B.3**     Solve $S_n^{*(2)}\left(\widehat{\boldsymbol{\beta}}_1^{[m+1]}, \boldsymbol{\beta}_2, \tau\right) = 0$ for $\widehat{\boldsymbol{\beta}}_2^{[m+1]}(\tau), \tau \in \left(0, \tau_{U,2}^{[m+1]}\right]$. Update $\tau_{U,2}$ with $\tau_{U,2}^{[m+1]}$.

**Step B.4**     Let $m = m + 1$. Repeat Steps B.2 and B.3 until convergence criteria are met.

Here, the initial value in Step B.1 can be set as Peng & Huang (2008)'s estimator which treats $T_1$ and $T_2$ are independent. The equations in Steps B.2–B.3 can be solved by $L_1$ minimization problems along similar lines of Peng & Huang (2008).

Asymptotic properties were established for the resulting estimators of $\boldsymbol{\beta}_{0,k}$ ($k = 1, 2$), including uniform consistency and weak convergence to a Gaussian process. Inference can be conducted through a standard bootstrapping procedure.

### 3.3.   Quantile regression with semi-competing risks data

Semi-competing risks (Fine et al. 2001) refers to as a situation where time to a nonterminal event (e.g. a non-fetal disease landmark event) can be censored by time to a terminal event (e.g. death or dropout) but not vice versa. Let $T_1$, $T_2$, and $C$ denote time to the nonterminal event, time to the terminal event, and time to random censoring, respectively. Let $\widetilde{Z}$ be a $p \times 1$ vector of covariates and $Z = \left(1, \widetilde{Z}^{\top}\right)^{\top}$. Define $X = T_1 \wedge T_2 \wedge C$, $Y = T_2 \wedge C$, $\delta = I(T_1 < Y)$, and $\eta = I(T_2 < C)$. The standard semi-competing risks data consist of $n$ replicates of ($X$, $Y$, $\delta$, $\eta$, $\widetilde{Z}$), denoted by $\{(X_i, Y_i, \delta_i, \eta_i, \widetilde{Z}_i), i = 1, \dots, n\}$. In a standard semi-competing risks setting, $T_2$ is only subject to random censoring by $C$; thus quantile regression for $T_2$ can follow the approaches developed for randomly censored data; see section 2.

Semi-competing risks methods are usually focused on the inference about $T_1$, which is complicated by the dependent censoring by $T_2$. Intuitively, one may first coerce semi-competing risks data into classic competing risk data by ignoring the extra information on $T_2$ when $\delta = 1$, and then apply quantile regression approaches developed for competing risks data. For example, targeting crude quantities for the nonterminal event, one can directly perform competing risks cumulative incidence quantile regression presented in section 3.2.1. Of note, this approach does not incur information loss from only using the competing risks portion of the data. This is because when $\delta = 1$, the cumulative incidence function for the non-terminal event, by definition, does not involve the extra information on the terminal event after the occurrence of the non-terminal event. An exception arises when left truncation is present. In that case, the semi-competing risks data are observable only when $Y > L$, where $L$ is a known left truncation time. Coercing semi-competing risks data into competing risks data induces artificial left truncation defined as $X > L$, thereby leading to information loss.

Li & Peng (2011) developed an extension of Peng & Fine (2009)'s method for competing risks cumulative incidence quantile regression tailored to semi-competing risks data subject to left truncation. In this case, the observed data include $n$ i.i.d. replicates of $(X^*, Y^*, \delta^*, \eta^*, L^*, \mathbf{Z}^*)$, which follow the conditional distribution of $(X, Y, \delta, \eta, L, \mathbf{Z})$ given $L < Y$. Assume the cumulative incidence quantile regression model for $T_1$, which is model (13) with $k = 1$. The basic estimation idea is to employ the IPCW technique with an inverse weight derived to properly account for both censoring by $C$ and left truncation $Y$. Under the assumption of $(L, C)$ is independent of $(T_1, T_2, \mathbf{Z})$, an estimating equation is given by

$$n^{-1/2} \sum_{i=1}^{n} \mathbf{Z}_i^* \left[ \frac{I\left\{ X_i^* \leqslant \exp\left( \mathbf{Z}_i^{*T} \mathbf{b} \right), \delta_i^* = 1, \eta_i^* = 1 \right\}}{\widehat{W}(Y_i^*, \mathbf{Z}_i^*)} - \tau \right] = 0, \tag{19}$$

where $\widehat{W}(y, z) = \widehat{G}(y)/\widehat{G}(z)$ with

$$\hat{\alpha}(z) = \int_0^v \hat{S}_{T_2}(u \mid z) \widehat{F}_L(\mathrm{d}u), \quad \widehat{G}(y) = \frac{1}{n} \sum_{i=1}^{n} \frac{I\left( L_i^* < y \leqslant Y_i^* \right) \hat{\alpha}\left( \mathbf{Z}_i^* \right)}{\hat{S}_{T_2}\left( y - \mid \mathbf{Z}_i^* \right)}.$$

Here $\widehat{F}_L(y)$ represents the Lynden-Bell estimator of $F_L(y) \doteq P(L \leqslant y)$, and $\hat{S}_{T_2}(u \mid z)$ is an adequate estimator of $P(T_2 > u|z)$. In practice, given $T_2$ is only subject to random right censoring by $C$ and random left truncation by $L$, $\hat{S}_{T_2 \mid Z = z}(t)$ may be obtained by using any existing regression method for left truncated and right censored data, such as the Cox proportional hazards model. After obtaining $\widehat{W}(Y_i^*, \mathbf{Z}_i^*)$, equation (19) can be solved by an algorithm similar to that presented for equation (14). Desirable theoretical properties, including uniform consistency and weak convergence to a Gaussian process, can also be established for the resulting estimator.

When interests lie in net quantities related to the latent time to nonterminal event $T_1$, utilizing the extra information in semi-competing risk data (beyond its competing risks portion) generally leads to better identifiability as well as improved statistical efficiency. Along this line, Li & Peng (2015) developed a quantile regression method tailored to study the conditional quantile of $T_1$ in the semi-competing risks setting. Specifically, Li & Peng (2015) assumed the following models:

$$\Pr(T_1 > s, T_2 > t \mid \mathbf{Z}) = C\left\{ 1 - F_{T_1}(s \mid \mathbf{Z}), 1 - F_{T_2}(t \mid \mathbf{Z}); g\left( \overline{\mathbf{Z}}^T \mathbf{r}_0 \right) \right\}, \tag{20}$$

$$Q_{T_1}(\tau \mid \mathbf{Z}) = \exp\left\{ \mathbf{Z}^T \boldsymbol{\beta}_0(\tau) \right\}, \quad Q_{T_2}(\tau \mid \mathbf{Z}) = \exp\left\{ \mathbf{Z}^T \boldsymbol{\alpha}_0(\tau) \right\}, \quad 0 < \tau < 1, \tag{21}$$

where $\overline{\mathbf{Z}}$ is a subvector of $\mathbf{Z}$ or $\mathbf{Z}$ itself, $C(\cdot, \cdot; a)$ is a known copula function with a given copula parameter $a$, and $g(\cdot)$ is a known function. In copula model (20), the unknown parameter $\mathbf{r}_0$ depicts how covariates may influence the copula parameter, which is often closely linked to the association between $T_1$ and $T_2$. In (21), the non-intercept coefficients in $\boldsymbol{\beta}_0(\tau)$ and $\boldsymbol{\alpha}_0(\tau)$ represent covariate effects on the $\tau$-th quantile of $T_1$ and $T_2$ respectively,

which are permitted to change with $\tau$. When these coefficients are constant over $\tau$, the models in (21) reduce to AFT models for $T_1$ and $T_2$.

To estimate models (20) and (21), a useful fact is that (20) implies

$$\Pr(X > t \mid Y > t, \mathbf{Z}) = K_A\Big\{\Pr(T_1 > s \mid \mathbf{Z}), \Pr(T_2 > t \mid \mathbf{Z}), g\big(\bar{\mathbf{Z}}^\top \mathbf{r}_0\big)\Big\}, \; t > 0, \qquad (22)$$

$$\Pr(X \leqslant s \mid Y > t, \mathbf{Z}) = K_B\Big\{\Pr(T_1 > s \mid \mathbf{Z}), \Pr(T_2 > t \mid \mathbf{Z}), g\big(\bar{\mathbf{Z}}^\top \mathbf{r}_0\big)\Big\}, \; s \leqslant t, \qquad (23)$$

where $K_A(u, v, \theta) = C(u, v; a)/v$ and $K_B(u, v, a) = \{v - C(u, v; a)\}/v$. In addition, the model assumptions in (21) imply $\Pr(T_1 \leqslant \exp\{\mathbf{Z}^T\boldsymbol{\beta}_0(\tau)\} \mid \mathbf{Z}) = \tau$ and $\Pr(T_2 \leqslant t \mid \mathbf{Z})) = \int_0^1 I\big[t \geqslant \exp\big\{\mathbf{Z}^T\boldsymbol{\alpha}_0(u)\big\}\big]du$. Li & Peng (2015) utilized these results to construct the following estimating equations,

$$n^{-1/2} \sum_{i=1}^{n} \widetilde{\mathbf{Z}}_i I\Big\{ \widetilde{\mathbf{Z}}_i^T \boldsymbol{\beta}(\tau) \leqslant \widetilde{\mathbf{Z}}_i^T \widehat{\boldsymbol{\alpha}}(\tau_{U,2}) \Big\} P_i(\boldsymbol{\beta}, \widehat{\boldsymbol{\alpha}}, \mathbf{r}, \tau) = 0, \; n^{-1/2} \sum_{i=1}^{n} \int_{\tau_a}^{\tau_b} \overline{\overline{\mathbf{Z}}}_i Q_i(\boldsymbol{\beta}, \widehat{\boldsymbol{\alpha}}, \mathbf{r}, \tau) = 0,$$

where $\widehat{\boldsymbol{\alpha}}(\cdot)$ is Peng & Huang (2008)'s estimator of $\boldsymbol{\alpha}_0(\cdot)$ given $T_2$ is only subject to random censoring by $C$, $\tau_{U,2}$ is an upper bound of a $\tau$-range where $a_0(\tau)$ is identifiable, and

$$P_i(\boldsymbol{\beta}, \boldsymbol{\alpha}, \mathbf{r}, \tau) = I\Big\{\log X_i > \widetilde{\mathbf{Z}}_i^T \boldsymbol{\beta}(\tau)\Big\} - I\Big\{\log Y_i > \widetilde{\mathbf{Z}}_i^T \boldsymbol{\beta}(\tau)\Big\} \times K_A\Big\{\tau, \int_0^{\tau_{U,2}} I\Big\{\widetilde{\mathbf{Z}}_i^T \boldsymbol{\beta}(\tau) \geqslant \widetilde{\mathbf{Z}}_i^T \boldsymbol{\alpha}(u)\Big\}du, g\big(\overline{\overline{\mathbf{Z}}}_i^T \mathbf{r}\big)\Big\},$$

$$Q_i(\boldsymbol{\beta}, \boldsymbol{\alpha}, \mathbf{r}, \tau) = \int_{t \in (0, \infty)} I\Big\{\widetilde{\mathbf{Z}}_i^T \boldsymbol{\beta}(\tau) \leqslant \log t \leqslant \widetilde{\mathbf{Z}}_i^T \boldsymbol{\alpha}(\tau_{U2}) \wedge \log Y_i\Big\}$$
$$\times \Big(I\Big\{\log X_i \leqslant \widetilde{\mathbf{Z}}_i^T \boldsymbol{\beta}(\tau)\Big\} - K_B\Big[\tau, \int_0^{\tau_{U,2}} I\Big\{\log t \geqslant \widetilde{\mathbf{Z}}_i^T \boldsymbol{\alpha}(u)\Big\}du, g\big(\overline{\overline{\mathbf{Z}}}_i^T \mathbf{r}\big)\Big]\Big)dt.$$

To compute $Q_i(\boldsymbol{\beta}, \widehat{\boldsymbol{\alpha}}, \mathbf{r}, \tau)$, one only needs to evaluate the integration over $t \in \big(0, \max_{i=1}^n Y_i \wedge \exp\big\{\mathbf{Z}_i^\top \widehat{\boldsymbol{\alpha}}(\tau_{U,2})\big\}\big]$. Confining $\boldsymbol{\beta}(\cdot)$ to be a cadlag step function, the integrand in $Q_i(\boldsymbol{\beta}, \widehat{\boldsymbol{\alpha}}, \mathbf{r}, \tau)$ is a piecewise constant function of $\tau$, and hence $Q_i(\boldsymbol{\beta}, \widehat{\boldsymbol{\alpha}}, \mathbf{r}, \tau)$ can be calculated as a finite sum. Li & Peng (2015) presented an iterative algorithm to solve these estimating equations. Li & Peng (2015) showed that the resulting estimator of $\mathbf{r}_0$ is consistent and asymptotic normal. Desirable theoretical properties, including uniform consistency and weak convergence to a Gaussian process, were established for the resulting estimator of $\boldsymbol{\beta}_0(\tau)$ for $\tau \in [\nu_1, \tau_{U,1}]$, where $0 < \nu_1 < \tau_{U,1} < 1$.

## 4. Quantile Regression and Its Adaptations for Recurrent events data

Recurrent events data are frequently encountered in clinical or epidemiological studies when the event of interest, such as infection and hospitalization, can occur repeatedly over time. Consider a general recurrent events data setting, where the observation of recurrent events is

subject to an observation window specified as a time interval $(L, R]$ (Nelson 2003). The counting process for the observed recurrent events is given by $N^{re}(t) = \sum_{j=1}^{\infty} I\left(L \leq T_j^{(i)} \leq t \wedge R\right)$, where $T^{(j)}$ denotes time to the $j$th recurrent event ($j = 1, 2,$ …), and the at-risk process is given by $Y^{re}(t) = I(L < t \leq R)$. Let $\widetilde{Z}$ be a $p \times 1$ vector of covariates and $Z = \left(1, \widetilde{Z}^{\top}\right)^{\top}$. The observed recurrent events data include $n$ i.i.d. replicates of $N^{re}(\cdot)$, $Z$, $L$, and $R$, denoted by $\left\{N_i^{re}(\cdot), Z_i, L_i, R_i\right\}_{i=1}^{n}$. In this section, we introduce three different ways to apply or adapt quantile regression to recurrent events data.

### 4.1. Quantile regression of recurrent event gap time.

Luo et al. (2013) proposed to model the gap time between recurrent events, namely, $G_{i,j} \doteq T_i^{(j)} - T_i^{(j-1)}$. By this approach, it is assumed that conditioning on $Z_i$ and a nonnegative subject-specific frailty variable $\gamma_i$, $N_i^{re}(\cdot)$ is a renewal process, and furthermore,

$$Q_{G_{i,j}}(\tau \mid Z_i) = \exp\left\{Z_i^{\top} \beta_0(\tau)\right\}, \ \tau \in (0, \tau_U]. \tag{24}$$

Consider the case where $L_i = 0$ and $R_i$ is independent of $\gamma_i$ and $\left\{T_i^{(j)}\right\}_{j=1}^{\infty}$ given $Z_i$. Let $m_i = N_i^{re}(R_i)$, $m_i^* = \max(m_i - 1, 1)$ and $\quad_i = I(m_i > 1)$. Define $X_{i,j} = G_{i,j}$ if $j < m_i$ and $X_{i,j} = R_i - T_i^{(m_i-1)}$ if $j = m_i$. Define $N_{i,j}(t) = I(G_{i,j} \leq t, \quad_i = 1)$, $R_{i,j}(t) = I(G_{i,j} \geq t)$, $H(x) = -\log(1-x)$. Note that uncensored gap times, $\{X_{i,j}, j = 1, \ldots, m_i - 1\}$, combined with the censored first gap time, $X_{i,1}$ with $\quad_i = 0$, can be viewed as clustered event times subject to random censoring. Under this view and by adapting the estimation framework of Peng & Huang (2008), Luo et al. (2013) proposed the following estimating equation for model (24):

$$n^{-1/2} \sum_{i=1}^{n} Z_i\left[N_i^*\left(\exp\left\{Z_i^{\top}\beta(\tau)\right\}\right) - \int_0^{\tau} R_i^*\left(\exp\left\{Z_i^{\top}\beta(u)\right\}\right)dH(u)\right] = 0, \tag{25}$$

where $N_i^*(t) = (m_i^*)^{-1}\sum_{j=1}^{m_i^*} N_{i,j}(t)$ and $R_i^*(t) = (m_i^*)^{-1}\sum_{j=1}^{m_i^*} R_{i,j}(t)$. An efficient algorithm to solve equation (25) can be developed along the lines of Peng & Huang (2008).

### 4.2. Generalized accelerated recurrence time model.

Huang & Peng (2009) and Sun et al. (2016) adopted a different strategy to adapt quantile regression modeling to recurrence events data. The main idea is to utilize the concept of time to expected frequency, which is a generalized version of conditional quantile that fits the recurrent events setting. Specifically, time to expected frequency is defined as $\tau_Z(u) \doteq \inf\{t \geq 0 : \mu_Z(t) \geq u\}$ for $u > 0$, where $\widetilde{N}(t) = \sum_{j=1}^{\infty} I\left(T^{(j)} \leq t\right)$ and $\mu_Z(t) = E\left\{\widetilde{N}(t) \mid Z\right\}$. It is easy to see that when the event of interest is not recurrent (i.e. $T^{(j)} = \infty$ for $j \geq 2$), $\tau_Z(u)$ becomes the conditional quantile $Q_{T^{(1)}}(\tau|Z)$. With recurrent events data, an adaptation of quantile regression modeling is to formulate covariate effects on $\tau_Z(u)$. This leads to the generalized accelerated recurrence time model, which is given by

$$\tau_Z(G(u)) = \exp\left\{Z^\top \beta_0(u)\right\}, \ u \in (0, U], \tag{26}$$

where $G(\cdot)$ is a known positive increasing function, the non-intercept coefficients in $\beta_0(u)$ represent covariate effects on time to expected frequency $G(u)$, and $U > 0$ is a prespecified constant.

The estimation of model (26) is facilitated by the counting process representation of model (26) justified in Sun et al. (2016). That is, model (26) is equivalent to

$$E\left\{N^{re}\left(\exp\left\{Z^\top \beta_0(u)\right\}\right) \mid Z\right\} = E\left\{\int_0^u Y^{re}\left(\exp\left\{Z^\top \beta_0(s)\right\}g(s)ds \mid Z\right\}, \ u \in (0, U], \tag{27}$$

where $g(u) = dG(u)/du$. This motivates a stochastic integral equation taking the form,

$$n^{-1/2} \sum_{i=1}^n X_i\left\{N_i^{re}\left(\exp\left\{X_i^\top \beta(u)\right\}\right) - \int_0^u Y_i^{re}\left(\exp\left\{X_i^\top \beta(s)\right\}\right)g(s)ds\right\} = 0, \ u \in (0, U]. \tag{28}$$

As commented in Sun et al. (2016), the theoretical and computational framework of Peng & Huang (2008) can be readily extended to study the recurrent events model (26). The algorithm to solve equation (28) is very similar to that for Peng & Huang (2008)'s martingale-based estimator (see section 2). The key modifications include adopting a grid on the frequency scale (instead of the $\tau$-scale), $\{0 = u_0 < u_1 < \cdots < u_{L(n)} = U\}$, and replace the objective function in Step 2 by

$$l_k(h) = \sum_{i=1}^n \sum_{j=1}^\infty I\left(L_i \leqslant T_i^{(j)} \leqslant R_i\right)\left|\log T_i^{(j)} - X_i^\top h\right| + \left|R^* - \left\{\sum_{i=1}^n \sum_{j=1}^\infty I\left(L_i \leqslant T_i^{(j)} \leqslant R_i\right)(-X_i)^\top h\right\}\right|$$
$$+ \left|R^* - \left\{\sum_{i=1}^n 2X_i^\top h \sum_{m=0}^{k-1} Y_i\left(\exp\left\{X_i^\top \widehat{\beta}(u_m)\right\}\right)\int_{u_m}^{u_{m+1}} g(s)ds,\right\}\right|,$$

where $R^*$ is an extremely large number. Theoretical arguments for Peng & Huang (2008)' estimator can also be generalized to establish the asymptotic properties of the estimator derived based on equation (28), including the uniform consistency for $u \in [v, U]$, where $0 < v < U$, and weak convergence to a Gaussian process at the root $n$ rate.

### 4.3. Quantile regression of individual recurrent risk measure.

More recently, Ma et al. (2020) proposed quantile regression of a sensible individual risk measure formulated upon the intensity process of recurrent events. Let $\widetilde{N}(t) \doteq \sum_{j=1}^\infty I\left(T^{(j)} \leqslant t\right)$ denote the underlying recurrent event process. Ma et al. (2020) assumed that given a nonnegative random variable $\gamma_i$, $\widetilde{N}_i(t)$ is a nonstationary Poisson process with the intensity function,

$$\lambda(t \mid \gamma_i) = \gamma_i \cdot \lambda_0(t). \tag{29}$$

Here $\lambda_0(t)$ stands for an unknown baseline intensity function, which is nonnegative and continuous, and is subject to the constraint, $\int_0^{v^*} \lambda_0(t)dt = 1$, with a predetermined constant $v^*$. This constraint is necessary for the purpose of model identifiability.

Under model (29), $\gamma_i$ captures the scale shift of subject $i$'s intensity process from the unknown baseline intensity $\lambda_0(t)$. A special case of model (29) is Wang et al. (2001)'s semi-parametric multiplicative intensity model, where $\gamma_i = \xi_i \exp\left(\widetilde{Z}_i^\top b_0\right)$, and $\xi_i$ is an unobservable frailty. This connection suggests that $\gamma_i$ can serve as a sensible measure of the latent subject-specific risk of recurrent events, which may naturally account for both observed covariates and unobservable frailty.

Ma et al. (2020) proposed to use quantile regression to explore the heterogeneity in $\gamma_i$ which quantifies the subject-specific risk of recurrent events. Specifically, it is assumed that

$$Q_{\gamma_i}(\tau \mid Z_i) = \exp\left\{Z_i^\top \beta_0(\tau)\right\}. \tag{30}$$

The non-intercept coefficients in $\beta_0(\tau)$ represent covariate effects on the $\tau$-th quantiles of $\gamma_i$.

A main challenge to estimate model (30) is that $\gamma_i$'s are not observed. Considering the setting with $L_i = 0$ and assuming $R_i$ is independent of $\widetilde{N}_i(\cdot)$ given $\gamma_i$, and $R_i$ is independent of $\gamma_i$ given $Z_i$, Ma et al. (2020) employed the principle of conditional score (Stefanski & Carroll 1987) and proposed the estimating equation,

$$n^{1/2}S_n(\beta, \widehat{\mu}, \tau) = 0, \tag{31}$$

where $S_n(\beta, \mu, \tau) \doteq n^{-1}\sum_{i=1}^n \int_r Z_i \cdot \psi_\tau\left\{\log(r) - Z_i^\top\beta(\tau)\right\}f\{r \mid m_i, C_i, Z_i; \beta(\cdot), \mu(\cdot)\}dr$ and

$\widehat{\mu}(t) = \exp\left\{\widehat{H}(t)\right\}$ with $\widehat{H}(t) = -\int_t^{v^*}\dfrac{\sum_{i=1}^n dN_i^{re}(s)}{\sum_{i=1}^n I(R_i \geqslant s)N_i^{re}(s)}$. Here $\psi_\tau(v) = \tau - I(v < 0)$, and

$$f\{\gamma \mid m, C, X; \beta(\cdot), \mu(\cdot)\} = \frac{\rho\{m \mid \gamma, C; \mu(\cdot)\}g\{\gamma \mid X; \beta(\cdot)\}}{\int_r \rho\{m \mid r, C; \mu(\cdot)\}g\{r \mid X; \beta(\cdot)\}dr},$$

where $\rho\{m \mid \gamma, C; \mu(\cdot)\} = \dfrac{\{\gamma\mu(C)\}^m}{m!}\exp\{-\gamma\mu(C)\}$ and

$g\{\gamma \mid X; \beta(\cdot)\} = \lim_{\delta \to 0}\dfrac{\delta}{\exp\left\{X^\top\beta(\tau_\gamma + \delta)\right\} - \exp\left\{X^\top\beta(\tau_\gamma)\right\}}$, with $\tau_\gamma = \{\tau \in (0,1) : \exp\{X^\top\beta(\tau) = \gamma\}$. It can be shown that $f\{\gamma \mid m, C, X; \beta_0(\cdot), \mu_0(\cdot)\}$ denotes the conditional density of $\gamma$ given $m$, $C$ and $X$ under the assumed models, and hence $E[S_n(\beta_0, \mu_0, \tau)] = 0$.

To solve equation (31), Ma et al. (2020) approximated $\beta(\tau)$ by using splines with $K(n)$ knots, and developed an iterative algorithm to find an estimate for the $\beta_0(\tau)$ in model (30) based on equation (31). The details are omitted here. Under certain regularity conditions, the resulting

estimator was shown to be uniformly consistent for $\tau \in [\zeta_1, \zeta_2]$, where $0 < \zeta_1 < \zeta_2 < 1$. Weak convergence to a Gaussian process was also established.

## 5.    Illustrations of quantile regression for survival data

### 5.1.    An example of quantile regression analysis with randomly censored data

We use a dataset from a dialysis study that investigated predictors of mortality in a cohort of 191 incident dialysis patients with chronic renal failure, aged 20 years and older, who started on chronic hemodialysis or peritoneal dialysis therapy between July 1996 and August 1997, recruited from metro-Atlanta area (Kutner et al. 2002). Of particular interest is a risk factor on symptoms of restless legs syndrome (RLS), which negatively affect quality of life and mortality risk as evidenced by prior studies. In this study, baseline measures were collected between 1996 and 1997 and vital status was monitored to December, 2005. In this dataset, the survival time $T$ of 35% dialysis patients were censored due to renal transplant or end of study.

Figure 1 plots the Kaplan-Merier curves for survival time stratified by the binary variable indicating moderate to severe RLS symptoms versus mild RLS symptoms (denoted by BLEGS). It is noted that the 25th percentiles of survival time for the the severe RLS group and the mild RLS groups are 0.95 versus 2.45 years, which are statistically significantly different. The 75th survival time percentiles for these two groups are rather similar, both between 7 and 8 years. This observation suggests that BLEGS may have an inhomogeneous effect on the distribution or quantile function of $T$. We next consider BLEGS, along with other potential predictors including patient's age (AGE), the indicator of fish consumption over the first year of dialysis (FISHH), the indicator of baseline HD dialysis modality (BHDPD), the indicator of eduction equal or higher than college (HIEDU), and the indicator of being black (BLACK). We fit the data with the standard Cox PH model and AFT model. In Table 1, we present the estimation results including the estimated coefficients and the associated p values. It is shown that both Cox PH model and AFT model do not suggest a significant effect of BLEGS on dialysis survival, though Figure 1 demonstrates its potential influence on the lower part of the survival distribution.

We next conduct quantile regression based on model (4) using Peng & Huang (2008)'s method for the same dataset. Figure 2 displays Peng & Huang (2008)'s estimator of $\boldsymbol{\beta}_0(\tau)$ along with 95% pointwise confidence intervals. In Figure 2, we observe that the coefficient for BLEGS diminishes gradually with $\tau$ whereas estimates for the other coefficients seem to be fairly constant. We apply the second-stage inference to formally investigate the constancy of each coefficient. The results confirm our observation from Figure 1, suggesting a varying effect of BLEGS and constant effects of the other covariates. This may lead to an interesting scientific implication that BLEGS may affect the survival experience of dialysis patients with short survival times but may have little impact on that of long-term survivors. The confirmed nonconstancy of the BLEGS coefficients further indicates the lack-of-fit of an AFT model for this dialysis data.

We also estimate the average quantile effects defined as $\int_l^u \beta_0^{(i)}(u)du (i = 2, ..., 7)$. The results are given in Table 2. We observe that the estimated average effect of BLEGS based on quantile regression has a larger magnitude compared to that based on the AFT model. The associated $p$ value is less than 0.05, providing some evidence for the association between RLS and dialysis survival. This example suggests that naively treating varying effects as constant ones may lead to attenuated covariate effect estimates and consequently result in biased conclusions.

### 5.2. An example of quantile regression analysis with competing risks data

We use the dataset from the breast cancer trial E1178 by the Eastern Cooperative Oncology Group (Cummings et al. 1993). In this study, patients were followed-up until breast cancer recurrence (BCR) or non-recurrence related death (NRD), whichever occurred first. This dataset includes 82 patients assigned to placebo and 85 patients assigned to tamoxifen. In the tamoxifen group, 42 patients experienced breast cancer recurrence and 23 died without recurrence; in the placebo group, 59 patients had breast cancer recurrence and 19 died without recurrence.

We apply the quantile regression strategy to evaluate the difference between two-year tamoxifen therapy versus placebo, while adjusting for other potential risk factors, including age, tumor size, number of positive nodes. Since it is more clinically relevant to evaluate BCR in the presence of NRD than with the unrealistic exclusion of NRD, we choose to use the cumulative incidence quantile regression method (Peng & Fine 2009) to analyze this competing risks dataset.

In Figure 3, we plot the BCR and NRD cumulative incidence functions separately for patient groups stratified by treatment, or age, number of positive nodes, and tumor size dichotomized at their median values, which are 71 years, 3, and 25mm, respectively. From Figure 3, we observe that all BCR cumulative incidence curves exceed 0.45 in the right tails. In contrast, the cumulative incidence curves for NRD are below 0.20. A visual impression from Figure 3 is that tamoxifen, number of positive nodes and tumor size may impact the cumulative incidence of BCR but not NRD, and their effects on BCR may not be constant.

We apply Peng & Fine (2009)'s method to fit the data with the competing risks quantile regression model (13), where the failure type corresponds to BCR and the exponential link function is replaced by the identify function. The number of positive nodes is incorporated into the model after log-transformation. Based on the results in Figure 3, we let $[\tau_L, \tau_U] = [0.10, 0.45]$. The analysis results displayed in Figure 4 suggest that patients who received placebo tend to experience breast cancer recurrence sooner than those on tamoxifen. In this example, age does not show a significant effect on the timing of breast cancer recurrence in the presence of nonrecurrence death. The effects of tumor size and number of nodes demonstrate some interesting increasing trend. The coefficient estimates, coupled with the 95% confidence intervals, suggest that tumor size and node number may only have a significant influence on the BCR cumulative incidence quantiles with relatively larger $\tau$'s, such as $\tau = 0.35$ or 0.4. The changing trend of the effects of tumor size and node number over $\tau$ is confirmed by second-stage constancy tests. The clinical implication may be that

either tumor size or number of positive nodes may significantly shorten the time to BCR for patients with moderate or low risk of BCR (corresponding to large $\tau$'s), while such an impact may vanish when patients are subject to high risk of BCR (corresponding to small $\tau$'s) possibly due to worse pre-existing health condition or other unknown factors. The treatment coefficients are rather constant, and are significantly above zero for many $\tau$'s. This reflects the beneficial effect of Tamoxifen treatment in term of prolonging the progression to BCR.

## 6. Remarks

Applying quantile regression to analyze survival data can provide robust and dynamic insight about the association between covariates and survival outcomes, which may not be offered by traditional survival regression methods. There have been rich developments of quantile regression methods for survival data in the last two decades. In this paper, we provide a selective review of approaches available to handle various types of survival data, including randomly censored data, competing and semi-competing risks data, truncated data, recurrent events data. Most of these methods are easy and stable to implement. This feature can help foster the applications of quantile regression in survival analysis.

Due to space limit, we omit many important relevant method developments. These include, but are not limited to, cure rate quantile regression methods (Wu & Yin 2013, 2017b,a) and censored quantile regression methods attending to regression quantile monotonicity across quantile levels, such as semiparametric copula quantile regression (De Backer et al. 2017).

Some important problems not covered in this paper but worth attention are quantile regression for survival data with high-dimensional covariates, survival data with time-dependent covariates, and survival data with missing covariates. This paper also does not discuss scenarios where the collection of survival data is attached to a special epidemiologic design, such as case-cohort design, and nested case-cohort design. Another interesting direction for extending survival quantile regression is to integrate quantile regression with causal inference. Work has emerged along these directions and merits further research efforts.
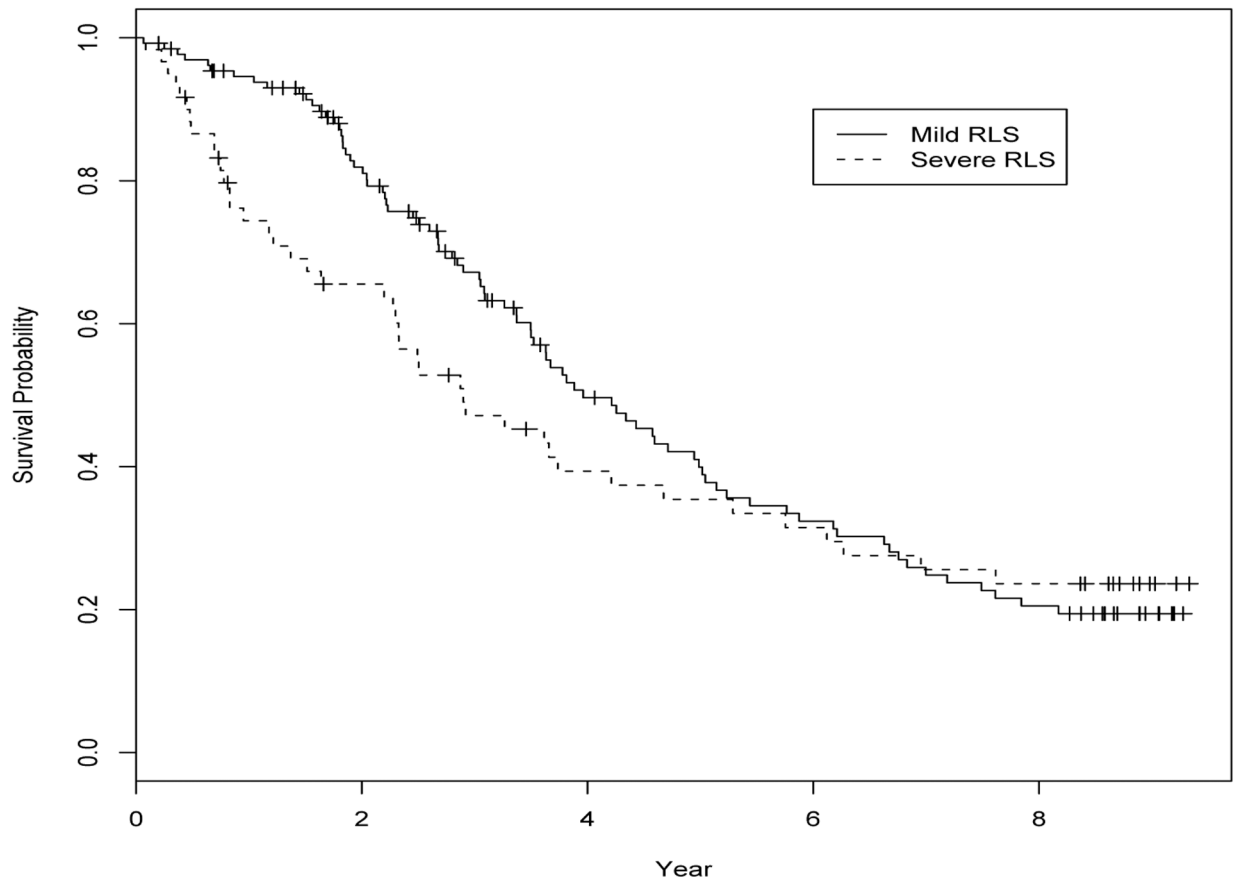
## ACKNOWLEDGMENTS

## LITERATURE CITED

Andersen PK, Gill RD. 1982. Cox's regression model for counting processes: a large sample study. The annals of statistics :1100–1120

Beran R. 1981. Nonparametric regression with randomly censored survival data

Buchinsky M, Hahn J. 1998. A alternative estimator for censored quantile regression. Econometrica 66:653–671

Buckley J, James I. 1979. Linear regression with censored data. Biometrika 66:429–436

Chernozhukov V, Hong H. 2001. Three-step censored quantile regression and extramarital affairs. Journal of the American Statistical Association :872–882

Clayton D. 1978. A model for association in bivariate life tables and its application in epidemiological studies of familial tendency in chronic disease incidence 65:141–151

Cox DR. 1972. Regression models and life-tables. Journal of the Royal Statistical Society: Series B (Methodological) 34:187–202

Cummings F, Gray R, Tormey D, Davis T, Volk H, et al. 1993. Adjuvant tamoxifen versus placebo in elderly women with node-positive breast cancer: long-term follow-up and causes of death. Journal of clinical oncology 11:29–35 [PubMed: 8418238]

De Backer M, El Ghouch A, Van Keilegom I. 2020. Linear censored quantile regression: A novel minimum-distance approach. Scandinavian Journal of Statistics

De Backer M, El Ghouch A, Van Keilegom I, et al. 2017. Semiparametric copula quantile regression for complete or censored data. Electronic Journal of Statistics 11:1660–1698

De Backer M, Ghouch AE, Van Keilegom I. 2019. An adapted loss function for censored quantile regression. Journal of the American Statistical Association 114:1126–1137

Efron B. 1967. The two-sample problem with censored data. Proc fifth berkley symposium in mathematical statistics, IV :831–553

Fine JP, Jiang H, Chappell R. 2001. On semi-competing risks data. Biometrika 88:907–919

Fitzenberger B. 1997. A guide to censored quantile regressions. Handbooks of Statistics: Robust Inference 15:405–437

Fygenson M, Ritov Y. 1994. Monotone estimating equations for censored data. The Annals of Statistics 22:732–746

Genest C. 1987. Frank's family of bivariate distributions 74:549–555

Huang Y. 2002. Calibration regression of censored lifetime medical cost. Journal of the American Statistical Association 98:318–327

Huang Y. 2010. Quantile Calculus and Censored Regression. The Annals of Statistics 38:1607–1637 [PubMed: 20592942]

Huang Y, Peng L. 2009. Accelerated recurrence time models. Scandinavian Journal of Statistics 36:636–648

Ji S, Peng L, Cheng Y, Lai H. 2012. Quantile regression for doubly censored data. Biometrics :101–112 [PubMed: 21950348]

Ji S, Peng L, Li R, Lynn MJ. 2014. Analysis of dependently censored data based on quantile regression. Statistica Sinica 24:1411 [PubMed: 25382953]

Kalbfleisch JD, Prentice RL. 2002. The statistical analysis of failure time data (2nd ed.). New York: Wiley Koenker R. 2005. Quantile regression, no. 9780521845731 in cambridge books

Koenker R. 2008. Censored quantile regression redux. Journal of Statistical Software 27:http://www.jstatsoft.com

Koenker R. 2017. Quantile regression: 40 years on. Annual Review of Economics 9:155–176

Koenker R, Bassett G. 1978. Regression quantiles. Econometrica 46:33–50

Koenker R, Portnoy S, Ng PT, Zeileis A, Grosjean P, Ripley BD. 2019. Package ?quantreg?

Kutner NG, Clow PW, Zhang R.and Aviles X. 2002. Association of fish intake and survival in a cohort of incident dialysis patients. American Journal of Kidney Diseases 39:1018–1024 [PubMed: 11979345]

Li KC, Wang JL, Chen CH, et al. 1999. Dimension reduction for censored regression data. The Annals of Statistics 27:1–23

Li R, Peng L. 2011. Quantile regression for left-truncated semi-competing risks data. Biometrics 67:701–710 [PubMed: 21133883]

Li R, Peng L. 2015. Quantile regression adjusting for dependent censoring from semicompeting risks. Journal of the Royal Statistical Society: Series B (Statistical Methodology) 77:107–130 [PubMed: 25574152]

Luo X, Huang CY, Wang L. 2013. Quantile regression for recurrent gap time data. Biometrics 69:375–385 [PubMed: 23489055]

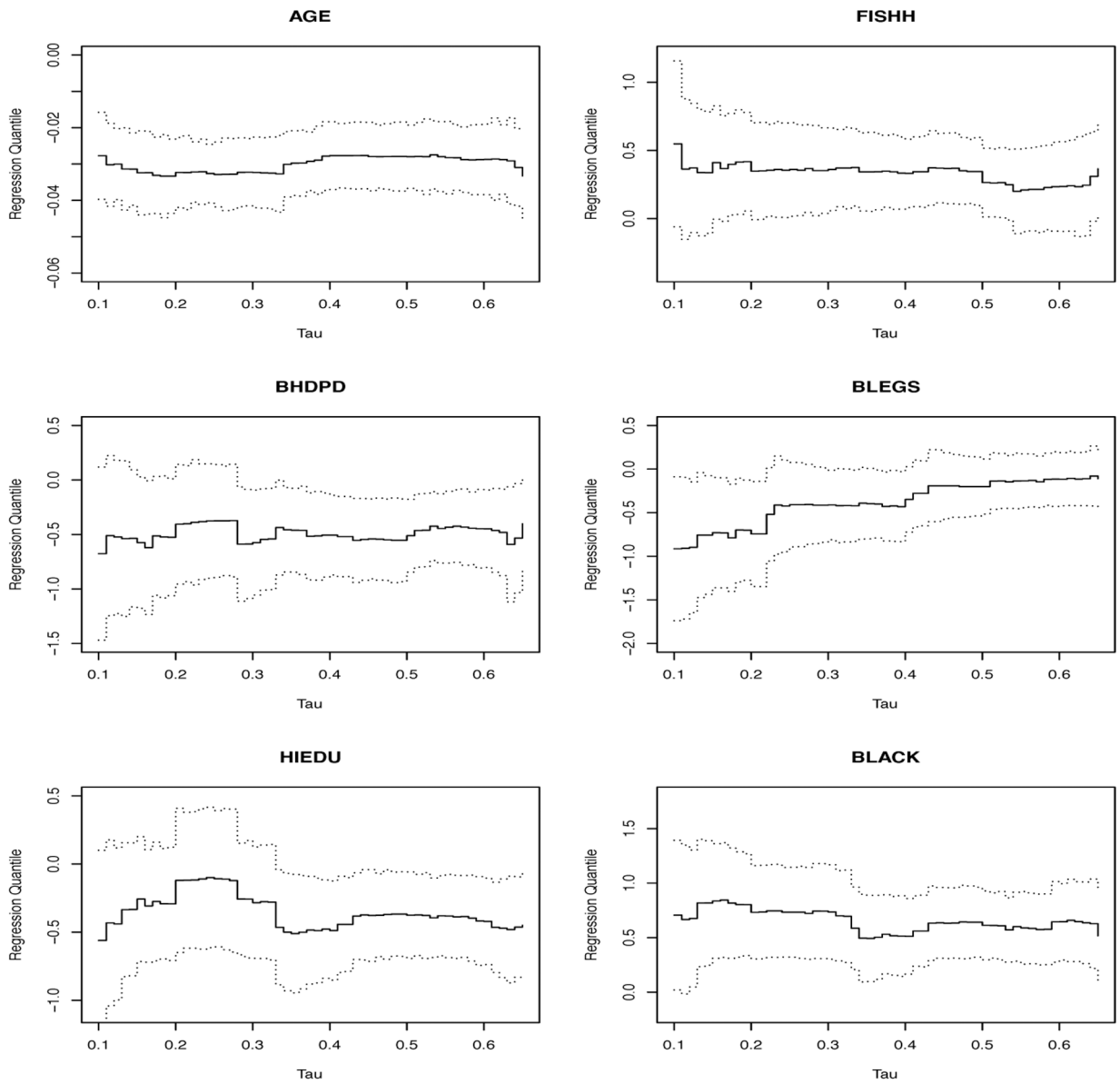Ma H, Peng L, Huang CY, Lai H. 2020. Heterogeneous individual risk modeling of recurrent events. Biometrika (Accepted)

Miller RG. 1976. Least squares regression with censored data. Biometrika 63:449–464

Nelson W. 2003. Recurrent events data analysis for product repairs, disease recurrences and other applications. Philadelphia: ASA-SIAM

Neocleous T, Vanden Branden K, Portnoy S. 2006. Correction to censored regression quantiles by s. portnoy, 98 (2003), 1001–1012. Journal of the American Statistical Association 101:860–861

Parzen MIWLJ, Ying Z. 1994. A resampling method based on pivotal estimating functions. Biometrika 81:341–350

Peng L. 2012. A note on self-consistent estimation of censored regression quantiles. Journal of Multivariate Analysis 105:368–379

Peng L, Fine J. 2009. Competing risks quantile regression. Journal of the American Statistical Association 104:1440–1453

Peng L, Huang Y. 2008. Survival analysis with quantile regression models. Journal of the American Statistical Association 103:637–649

Portnoy S. 2003. Censored regression quantiles. Journal of the American Statistical Association 98:1001–1012

Portnoy S, Lin G. 2010. Asymptotics for censored regression quantiles. Journal of Nonparametric Statistics 22:115–130

Powell J. 1984. Least absolute deviations estimation for the censored regression model. Journal of Econometrics 25:303–325

Powell J. 1986. Censored regression quantiles. Journal of Econometrics 32:143–155

Prentice RL. 1978. Linear rank tests with right censored data. Biometrika 65:167–179

Reid N. 1994. A conversation with sir david cox. Statistical Science 9:439–455

Ritov Y. 1990. Estimation in a linear regression model with censored data. The Annals of Statistics :303–328

Robins J, Rotnitzky A. 1992. Recovery of information and adjustment for dependent censoring using surrogate markers. In AIDS Epidemiology-Methodological Issues, eds. Jewell N, Dietz K, Farewell V. Boston: Birkhauser, 24–33

Stefanski L, Carroll R. 1987. Conditional scores and optimal scores for generalized linear measurement-error models. Biometrika 74:703–716

Sun X, Peng L, Huang Y, Lai HJ. 2016. Generalizing quantile regression for counting processes with applications to recurrent events. Journal of the American Statistical Association 111:145–156 [PubMed: 27212738]

Sun Y, Wang HJ, Gilbert J. 2012. Quantile regression for competing risks data with missing cause of failure. Statistica Sinica 22

Tanner MA, Wong WH. 1987. The calculation of posterior distributions by data augmentation. Journal of the American statistical Association 82:528–540

Tsiatis A. 1975. A nonidentifiability aspect of the problem of competing risks. Proceedings of the National Academy of Sciences 72:20–22

Tsiatis AA. 1990. Estimating regression parameters using linear rank tests for censored data. The Annals of Statistics :354–372

Wang H, Wang L. 2009. Locally weighted censored quantile regression. Journal of the American Statistical Association 104:1117–1128

Wang MC, Qin J, Chiang CT. 2001. Analyzing recurrent event data with informative censoring. Journal of the American Statistical Association 96:1057–1065

Wei L, Gail M. 1983. Nonparametric estimation for a scale-change with censored observations. Journal of the American Statistical Association 78:382–388

Wei LJ, Ying Z, Lin D. 1990. Linear regression analysis of censored survival data based on rank tests. Biometrika 77:845–851

Wu Y, Yin G. 2013. Cure rate quantile regression for censored data with a survival fraction. Journal of the American Statistical Association 108:1517–1531

Wu Y, Yin G. 2017a. Cure rate quantile regression accommodating both finite and infinite survival times. Canadian Journal of Statistics 45:29–43
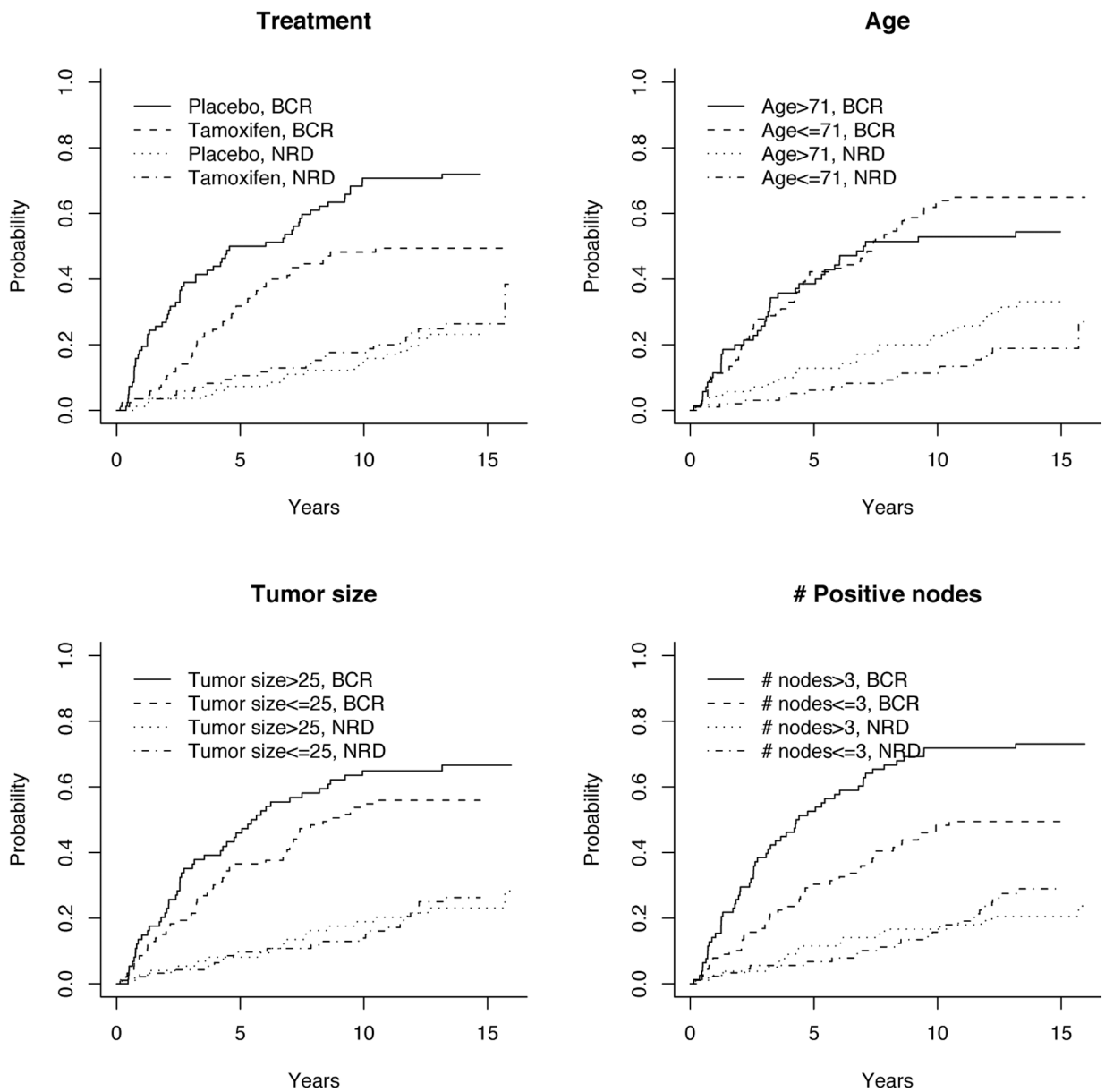
Wu Y, Yin G. 2017b. Multiple imputation for cure rate quantile regression with censored data. Biometrics 73:94–103 [PubMed: 27479513]

Xia Y, Zhang D, Xu J. 2010. Dimension reduction and semiparametric estimation of survival models. Journal of the American Statistical Association 105:278–290

Yang X, Narisetty NN, He X. 2018. A new approach to censored quantile regression estimation. Journal of Computational and Graphical Statistics 27:417–425

Ying Z, Jung SH, Wei LJ. 1995. Survival analysis with median regression models. Journal of the American Statistical Association 90:178–184

Zhou L. 2006. A simple censored median regression estimator. Statistica Sinica 16:1043–1058
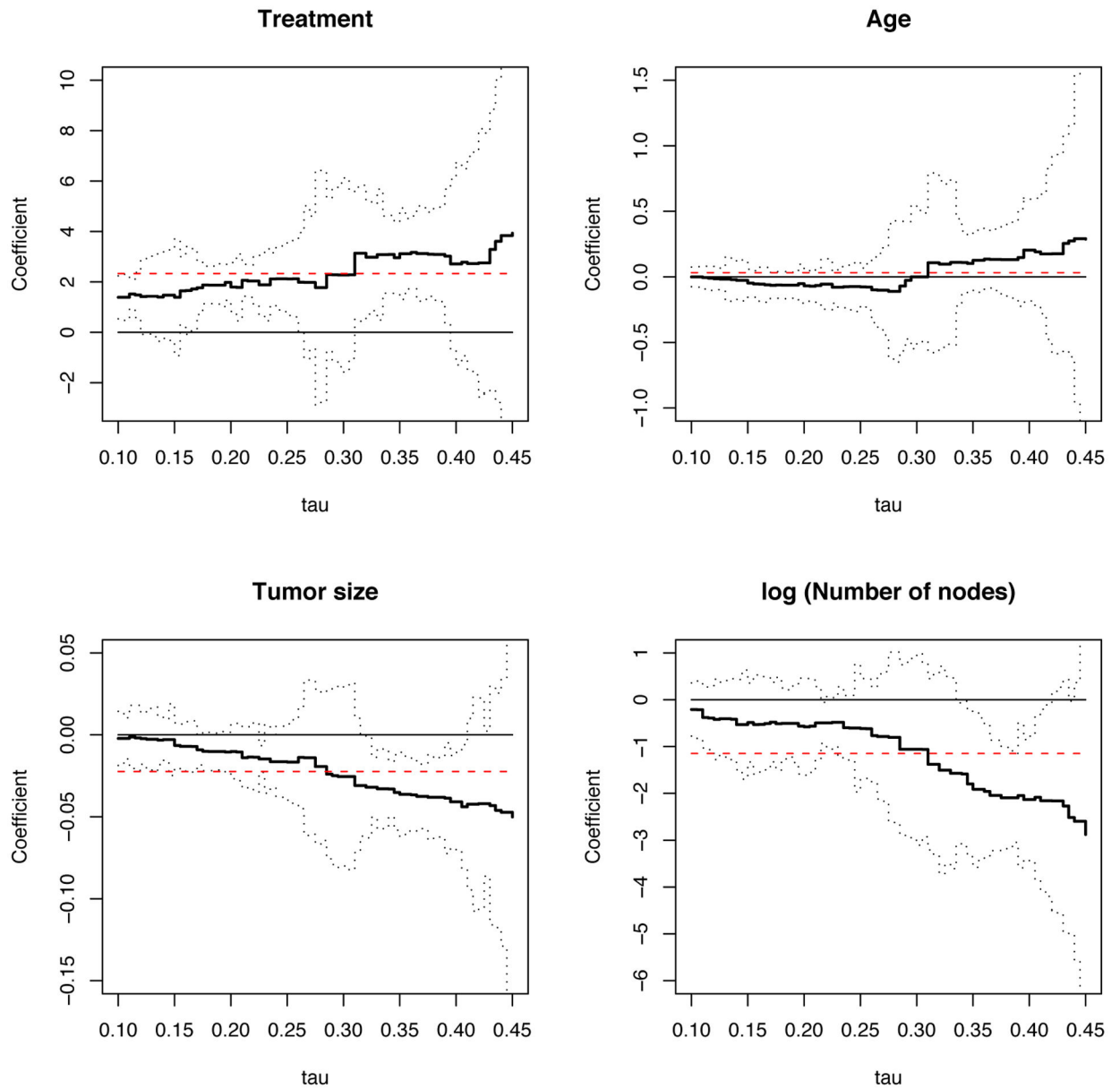
**Figure 1.**
The dialysis example: Kaplan Meier curves of survival time stratified by the status of RLS symptoms

**Figure 2.**
The dialysis example: Peng and Huang's estimator (solid lines) and 95% pointwise confidence intervals (dotted lines) of regression quantiles.

**Figure 3.**
E1178 Trial Example: Estimated Cumulative Incidence Functions.

**Figure 4.**
E1178 Trial Example: Estimated Regression Coefficients for the Breast Cancer Recurrence
Endpoint. Bold Solid Lines Represent Coefficient Estimates; Dotted Lines Represent 95%
Pointwise Confidence Intervals; Dashed Lines Represent Estimates for Trimmed Mean
Covariate Effects.

**Table 1**

Results from fitting the Cox PH model and AFT model to the dialysis dataset.

|  | Cox Model | | AFT Model | |
|---|---|---|---|---|
|  | Coef | *p* value | Coef | *p* value |
| AGE | 0.059 | <0.001 | −0.035 | <0.001 |
| FISHH | −0.831 | <0.001 | 0.485 | <0.001 |
| BHDPD | 0.837 | <0.001 | −0.473 | <0.001 |
| BLEGS | 0.264 | 0.197 | −0.173 | 0.232 |
| HIEDU | 0.625 | 0.009 | 0.364 | 0.024 |
| BLACK | −1.014 | <0.001 | 0.591 | <0.001 |

Coef: coefficient estimate; SE: standard error

**Table 2**

Estimation of average covariate effects based on quantile regression.

|        | AveEff  | SE    | *p* value |
|--------|---------|-------|-----------|
| AGE    | −0.030  | 0.003 | < 0.001   |
| FISHH  | 0.327   | 0.116 | 0.005     |
| BHDPD  | −0.489  | 0.162 | 0.003     |
| BLEGS  | −0.369  | 0.161 | 0.022     |
| HIEDU  | −0.350  | 0.137 | 0.011     |
| BLACK  | 0.654   | 0.144 | < 0.001   |

AveEff: Estimated average effect; SE: standard error