

# Composite Quantile Regression for Nonparametric Model with Random Censored Data

Rong Jiang, Weimin Qian

Department of Mathematics, Tongji University, Shanghai, China  
Email: jrtrying@126.com, wmqian2003@yahoo.com.cn

Received May 29, 2012; revised June 30, 2012; accepted July 15, 2012

Copyright © 2013 Rong Jiang, Weimin Qian. This is an open access article distributed under the Creative Commons Attribution License, which permits unrestricted use, distribution, and reproduction in any medium, provided the original work is properly cited.

## ABSTRACT

The composite quantile regression should provide estimation efficiency gain over a single quantile regression. In this paper, we extend composite quantile regression to nonparametric model with random censored data. The asymptotic normality of the proposed estimator is established. The proposed methods are applied to the lung cancer data. Extensive simulations are reported, showing that the proposed method works well in practical settings.

**Keywords:** Kaplan-Meier Estimator; Censored Data; Composite Quantile Regression; Kernel Estimator; Nonparametric Model

## 1. Introduction

Consider the following nonparametric regression model with random censored data:

$$T = m(U) + \sigma(U)\varepsilon, \quad (1)$$

where  $m(\cdot)$  is an unknown smoothing function,  $\sigma(\cdot)$  is a positive function representing the standard deviation and  $\varepsilon$  is the random error with mean 0 and variance 1. Let  $C$  denote the censoring variable, whose distribution may depend on  $U$ , where  $U$  is vector of observed covariates. In this paper, we focus on random right censoring, we only observe the triples  $(U, Y, \delta)$ , where  $Y = \min(T, C)$  and  $\delta = I(T \leq C)$  are the observed response variable and the censoring indicator respectively, where  $T$  is the survival time.

Censored quantile regression was first studied by [1] for fixed censoring. [2] proposed an estimator for a conditional quantile assuming that the regression models at lower quantiles are all linear. A recursively weighted estimation procedure that can be regarded as a generalization of the Kaplan-Meier estimator to conditional quantiles was described in their paper. Afterward, [3] presented an alternative approach that is based on the Nelson-Aalen estimator of the cumulative hazard function but still requires the same global-linearity assumption as Portnoy's. Their method provides a more direct approach to the asymptotic theory and a simpler computation algorithm. More recent studies by [4], proposed to overcome the global-linearity assumption by directly

estimating the conditional censoring distribution nonparametrically using the local Kaplan-Meier method. Their computational algorithm is more stable and simpler to implement than Portnoy's or Peng and Huang's. Moreover, the local nonparametric estimator on which the model is based performs best when the covariates can be assumed independent.

Intuitively, the composite quantile regression (CQR) should provide estimation efficiency gain over a single quantile regression; see [5]. A composite quantile regression model assumes that there exist common covariate effects in a range of quantiles such that the quantile levels only differ in terms of the intercept. From a more general regression perspective, composite quantile regression seeks to model a set of parallel regression curves, and thus it can be viewed as a compromise between a set of quantile regression curves with different intercepts and slopes and a single summary regression curve. [6] proposed the local polynomial CQR estimators (LCQR) for estimating the nonparametric regression function and its derivative. It is shown that the local CQR method can significantly improve the estimation efficiency of the local least squares estimator for commonly-used non-normal error distributions. Furthermore, [7] studied semiparametric CQR estimates for semiparametric varying-coefficient partially linear model. They compared CQR with least squares and quantile regression, and the results showed that CQR outperformed both least squares and quantile regression. [8] considered CQR

estimates for single-index models. Recently, [9] extended the CQR method to linear model with randomly censored data. This motivates us to extend the CQR method to nonparametric model with censored data (LCQRC).

The paper is organized as follows. In Section 2, local composite quantile regression for nonparametric model with censored data is introduced, and the main theoretical results are also given in this section. Both simulation examples and a real data application are given in Section 3 to illustrate the proposed procedures. Final remarks are given in Section 4. The technical proofs are deferred to the Appendix.

## 2. Methodology

### 2.1. Local Composite Quantile Regression with Censored Data

We first consider an ideal situation where  $F_0(t|U_i)$ , the conditional cumulative distribution function of the survival time  $T_i$  given  $U_i$ , is assumed to be known. In this case, we define the following weight function:

$$\omega_i(F_0, \tau) = \begin{cases} 1, & \delta_i = 1 \text{ or } F_0(C_i|U_i) > \tau, \\ \frac{\tau - F_0(C_i|U_i)}{1 - F_0(C_i|U_i)}, & \delta_i = 0 \text{ and } F_0(C_i|U_i) < \tau, \end{cases} \quad (2)$$

$i = 1, \dots, n$ . In reality,  $F_0(t|U_i)$  is unknown and has to be estimated. We propose to estimate  $F_0(\cdot|U)$  nonparametrically using the local Kaplan-Meier estimator

$$\hat{F}(t|U) = 1 - \prod_{j=1}^n \left\{ 1 - \frac{B_{nj}(U)}{\sum_{s=1}^n I(Y_s \geq Y_j) B_{ns}(U)} \right\}^{\xi_j(t)},$$

where  $\xi_j(t) = I(Y_j \geq t, \delta_j = 1)$  and

$B_{nk}(U) = K([U - U_k]/h^*) / \sum_{i=1}^n K([U - U_i]/h^*)$ , where

$K(\cdot)$  is a smooth kernel function,  $h^* \in R^+$  is the bandwidth converging to zero as  $n \rightarrow \infty$ . By plugging  $\hat{F}(\cdot|U)$  into (2), we obtain the estimated local weights  $\omega_i(\hat{F}, \tau)$ . Consider estimating the value of  $m(U)$  at  $u_0$ . The LCQRC procedure estimates  $m(u_0)$ , defined by  $\hat{m}(u_0) = \frac{1}{q} \sum_{k=1}^q \hat{a}_k$ , via minimizing the locally weighted objective function

$$\sum_{k=1}^q \sum_{i=1}^n \left[ \omega_i(\hat{F}, \tau_k) \rho_{\tau_k} \{Y_i - a_k - b(U_i - u_0)\} + (1 - \omega_i(\hat{F}, \tau_k)) \rho_{\tau_k} \{Y_i^{+\infty} - a_k - b(U_i - u_0)\} \right] K\left(\frac{U_i - u_0}{h}\right)$$

where  $\rho_{\tau_k}(r) = \tau_k r - rI(r < 0)$ ,  $k = 1, 2, \dots, q$ , be  $q$  check loss functions at  $q$  quantile positions:  $\tau_k = k/(q+1)$  and  $Y^{+\infty}$  is any value sufficiently large to exceed all  $m(U_i)$ .

**Remark 1.** The detail explant of  $\omega_i(F_0, \tau)$  can see Remark 1 of [4].

### 2.2. Asymptotic Properties

Denote by  $f_U(\cdot)$  the marginal density function of the covariate  $U$ ,  $\mu_j = \int u^j K(u) du$  and

$\nu_j = \int u^j K^2(u) du$ . To prove main results in this paper, the following technical conditions are imposed.

**A1.** The functions  $F_0(t|U)$  and  $G(t|U)$  have first derivatives with respect to  $t$ , denoted as  $f_0(t|U)$  and  $g(t|U)$ , which are uniformly bounded away from infinity. In addition,  $F_0(t|U)$  and  $G(t|U)$  have bounded second order partial derivatives with respect to  $U$ .

**A2.**  $\Sigma$  is positive definite matrix.

**A3.**  $m(U)$  has a continuous second derivative in the neighborhood of  $u_0$ .

**A4.**  $f_U(\cdot)$  is differentiable and positive in the neighborhood of  $u_0$ .

**A5.** The conditional variance  $\sigma^2(u)$  is continuous in the neighborhood of  $u_0$ .

**A6.** Assume that the error has a symmetric distribution with a positive density  $f_0(\cdot)$ .

**Remark 2.** Assumption A1 is needed for the local Kaplan-Meier estimator. It allows us to obtain the local expansions of  $F_0(t|U)$  and  $G(t|U)$  in the neighborhood of  $m(U)$ , and to obtain the uniform consistency and the linear representation of  $\hat{F}(t|U)$ , which are needed for deriving the asymptotic normality result. Assumption A2 ensures that the expectation of the estimating function has a unique zero, and it is needed to establish the asymptotic distribution. Assumptions A3-A6 are the same conditions for establishing the asymptotic normality of local composite quantile regression ([6]).

We state the asymptotic normality for  $\hat{m}(u_0)$  in the following theorem.

**Theorem 1.** Assume that the triples  $(U_i, Y_i, \delta_i)$  constitute and i.i.d. multivariate random sample, and that the censoring variable  $C_i$  is independent of  $T_i$  conditional on the covariate  $U_i$ . Suppose that  $u_0$  is an interior of the support of  $f_U(\cdot)$ . Under the regularity conditions A1-A6, if  $h \rightarrow 0$  and  $nh \rightarrow \infty$ , then

$$\sqrt{nh} \left[ \hat{m}(u_0) - m(u_0) - \frac{1}{2} m''(u_0) \mu_2 h^2 \right] \xrightarrow{L} N(0, \Sigma),$$

where  $\xrightarrow{L}$  stands for convergence in distribution and

$$\Sigma = \frac{\nu_0}{f_U(u_0)} \frac{1}{q^2} \sum_{k, k'=1}^q \Gamma_k^{-1} \Lambda_{kk'} \Gamma_{k'}^{-1}, \text{ where}$$

$$\Gamma_k = E \cdot \left[ (1 - G(m(U) + \sigma(U)c_k | U)) f_0(m(U) + \sigma(U)c_k | U) \right],$$

$$\Lambda_{kk'} = E_{(U,C)} \left[ I(C \leq m(U) + \sigma(U)c_k \wedge c_{k'}) \cdot \frac{F_0(C|U)(1 - \tau_k)(1 - \tau_{k'})}{1 - F_0(C|U)} + I(C > m(U) + \sigma(U)c_k \wedge c_{k'}) \cdot (\tau_k \wedge \tau_{k'} - \tau_k \tau_{k'}) \right]$$

and  $c_k = F^{-1}(\tau_k)$ .

### 3. Numerical Studies

In this section, we conduct simulation studies to assess the finite sample performance of the proposed procedures and illustrate the proposed methodology on a lung cancer data set. Moreover, we compare the performance of the newly proposed method with LCQR ([6]) and nonparametric quantile regression with censored data (NQRC) that was proposed by [10].

In the proposed compute process, we take  $Y^{+\infty} = 100 \max\{Y_1, \dots, Y_n\}$  and  $K = \frac{15}{16}(1 - U^2)^2 I(|U| \leq 1)$ . The bandwidth  $h^*$  can be obtained by 10-fold cross-validation method (see [4]), and we use the short-cut strategy method to select  $h$  (see [6]).

#### 3.1. Example 1

The data are generated from the following model

$$T_i = 10U_i \sin(2\pi U_i) + \varepsilon_i, i = 1, \dots, n,$$

where  $U_i$  is uniformly distributed on  $[0, 1]$  and  $\varepsilon$  is i.i.d. standard normal random variables. The censoring variable  $C_i \sim U(0, c)$  and  $Y_i = \min(T_i, C_i)$ . The value of the constant  $c$  in the model determines the censoring proportion. In our simulations, we consider three

censoring rates (CR): 20%, 30% and 40%. For each censoring rate, the sample sizes are taken to be 100 and 200. To evaluate the finite sample performance of our estimator. Two distance measures are approximated, the first one the mean absolute deviation error (MADE) is given by  $n^{-1} \sum_{i=1}^n |\hat{m}(U_i) - m(U_i)|$ , and the second one the mean squared error (MSE) defines as

$$n^{-1} \sum_{i=1}^n [\hat{m}(U_i) - m(U_i)]^2.$$

Furthermore, we define the rate of MADE and MSE which are  $RMAD = MADE/MADE_{LCQR}$  and  $RMSE = MSE/MSE_{LCQR}$ .

For right censored data, quantile functions with  $\tau$  close to 1 may not be identifiable due to censorship. In our simulations, we consider  $q = 5$  for LCQR and LCQRC estimators. The means and standard deviations of MADE, MSE, RMAD and RMSE are respectively reported in **Table 1** and **Table 2**. From **Tables 1** and **2**, we can make the following observations: the performance of proposal method is better than that of LCRQ and NQRC. Moreover, LCQRC estimators are much more accurate when sample sizes increase. **Figure 1** summarize the Curve estimates for three censoring rates of 20%, 30% and 40% with different sample sizes. It shows that the performance of LCQRC is very close to the true value.

#### 3.2. Example 2

It is necessary to investigate the effect of heteroscedastic errors. The observations  $(U_i, T_i, C_i), i = 1, \dots, n$ , are generated from following model

$$T_i = 2 \exp(-16U_i^2) + \sin(\pi U_i) + 0.2U_i \varepsilon_i, i = 1, \dots, n,$$

where  $U_i, \varepsilon_i$ , and  $C_i$  are generated following the same way as in Example 1. The means and standard deviations of MADE, MSE, RMAD and RMSE are respectively reported in **Table 3** and **Table 4**. The

**Table 1. Simulation results of  $\hat{m}(\cdot)$  with  $n = 100$  for Example 1.**

CR	Method	RMAD	MADE	RMSE	MSE
20%	LCQR5	-	0.7169 (0.5047)	-	0.7804 (0.9798)
	LCQRC5	0.7946 (0.1195)	0.5643 (0.4679)	0.7297 (0.2452)	0.5481 (0.8718)
	NQRC0.5	0.9655 (0.1837)	0.5919 (0.4155)	0.9338 (0.3355)	0.5476 (0.6646)
30%	LCQR5	-	0.8284 (0.6031)	-	1.0578 (1.2602)
	LCQRC5	0.6607 (0.1187)	0.5431 (0.4468)	0.4903 (0.1639)	0.5060 (0.7765)
	NQRC0.5	0.8378 (0.3501)	0.6960 (0.4994)	0.7920 (0.6298)	0.8650 (1.0371)
40%	LCQR5	-	1.0102 (0.7520)	-	1.6054 (1.8037)
	LCQRC5	0.5501 (0.1348)	0.5507 (0.5448)	0.4332 (0.3251)	0.6810 (1.9258)
	NQRC0.5	0.7222 (0.1882)	0.7169 (0.5424)	0.5548 (0.2795)	0.8530 (1.1186)

**Table 2. Simulation results of  $\hat{m}(\cdot)$  with  $n = 200$  for Example 1.**

CR	Method	RMADE	MADE	RMSE	MSE
20%	LCQR5	-	0.5374 (0.3750)	-	0.4386 (0.5481)
	LCQRC5	0.7327 (0.1089)	0.3893 (0.3443)	0.6612 (0.2327)	0.2737 (0.5262)
	NQRC0.5	0.8356 (0.2114)	0.5443 (0.3769)	0.8254 (0.2910)	0.4732 (0.5346)
30%	LCQR5	-	0.7342 (0.5154)	-	0.8121 (0.9382)
	LCQRC5	0.5768 (0.1191)	0.4191 (0.3379)	0.3815 (0.1658)	0.2988 (0.4677)
	NQRC0.5	0.7539 (0.1344)	0.5484 (0.4015)	0.5985 (0.1898)	0.4729 (0.5870)
40%	LCQR5	-	0.9755 (0.7164)	-	1.4825 (1.5772)
	LCQRC5	0.4231 (0.0771)	0.4084 (0.3504)	0.2084 (0.0839)	0.2957 (0.5196)
	NQRC0.5	0.7689 (0.2634)	0.7355 (0.4962)	0.6059 (0.4493)	0.8477 (0.9701)

**Table 3. Simulation results of  $\hat{m}(\cdot)$  with  $n = 100$  for Example 2.**

CR	Method	RMADE	MADE	RMSE	MSE
20%	LCQR5	-	0.0639 (0.0538)	-	0.0074 (0.0175)
	LCQRC5	0.9620 (0.1600)	0.0603 (0.0451)	0.9230 (0.4398)	0.0059 (0.0123)
	NQRC0.5	1.2231 (0.3356)	0.0750 (0.0500)	1.5544 (1.0808)	0.0085 (0.0110)
30%	LCQR5	-	0.0936 (0.0856)	-	0.0191 (0.0337)
	LCQRC5	0.8380 (0.2094)	0.0791 (0.1318)	0.6470 (0.4510)	0.0065 (0.0028)
	NQRC0.5	0.8502 (0.3245)	0.0714 (0.0556)	0.7788 (0.6206)	0.0089 (0.0137)
40%	LCQR5	-	0.1564 (0.1549)	-	0.0560 (0.0957)
	LCQRC5	0.5269 (0.1720)	0.0676 (0.0763)	0.3816 (0.5250)	0.0167 (0.0962)
	NQRC0.5	0.5836 (0.3207)	0.0796 (0.0613)	0.3874 (0.5181)	0.0117 (0.0181)

**Table 4. Simulation results of  $\hat{m}(\cdot)$  with  $n = 200$  for Example 2.**

CR	Method	RMADE	MADE	RMSE	MSE
20%	LCQR5	-	0.0540 (0.0439)	-	0.0051 (0.0103)
	LCQRC5	0.9244 (0.2066)	0.0488 (0.0339)	0.8220 (0.3890)	0.0036 (0.0068)
	NQRC0.5	1.1134 (0.2889)	0.0599 (0.0433)	1.2714 (0.7178)	0.0059 (0.0082)
30%	LCQR5	-	0.0811 (0.0812)	-	0.0134 (0.0248)
	LCQRC5	0.6377 (0.1167)	0.0506 (0.0378)	0.3222 (0.1032)	0.0040 (0.0074)
	NQRC0.5	0.8157 (0.2516)	0.0638 (0.0493)	0.5476 (0.2898)	0.0067 (0.0095)
40%	LCQR5	-	0.1492 (0.1297)	-	0.0432 (0.0678)
	LCQRC5	0.4230 (0.1894)	0.0554 (0.0451)	0.1769 (0.1246)	0.0054 (0.0163)
	NQRC0.5	0.5235 (0.1757)	0.0714 (0.0474)	0.2591 (0.1744)	0.0077 (0.0104)

performance of LCQRC is presented in **Figure 2**. The results of Example 1 and Example 2 show very similar messages.

### 3.3. Example 3

As an illustration, we now apply the proposed LCQRC to the lung cancer data. The data contain 228 observations on ten variables. The censoring percentage is 27%, so the estimators are expected to perform well. More details about the study can be found in [11], and the dataset is included in the R package *lung*. We are interested in estimating the conditional of survival time (in days) given age (in years). Here, we use model (1) to fit the lung cancer data, where  $Y$  is the  $\log_{10}$  (survival time) and  $U$  is the age/100. To evaluate the performance of our estimator. Two distance measures were approximated,

the first one the mean absolute deviation error ( $\text{MADE}_y$ ) given by  $n^{-1} \sum_{i=1}^n |\hat{y}_i - y_i|$ , and the second one

the mean squared error ( $\text{MSE}_y$ ) defined as

$$n^{-1} \sum_{i=1}^n [\hat{y}_i - y_i]^2, \text{ where } n = 228. \text{ Furthermore, we define the rate of } \text{MADE}_y \text{ and } \text{MSE}_y \text{ which are}$$

$$\text{RMADE}_y = \text{MADE}_y / \text{MADE}_{y_{\text{LCQR5}}} \text{ and}$$

$$\text{RMSE}_y = \text{MSE}_y / \text{MSE}_{y_{\text{LCQR5}}}. \text{ Next, we report and compare results with LCQR and NQRC for estimating the survival time. The simulation results for the LCQR, LCQRC and NQRC are given in Table 5. It shows that LCQRC is better than that of LCRQ and NQRC. Figure 3 summarize the simulation results for LCQRC5. It$$

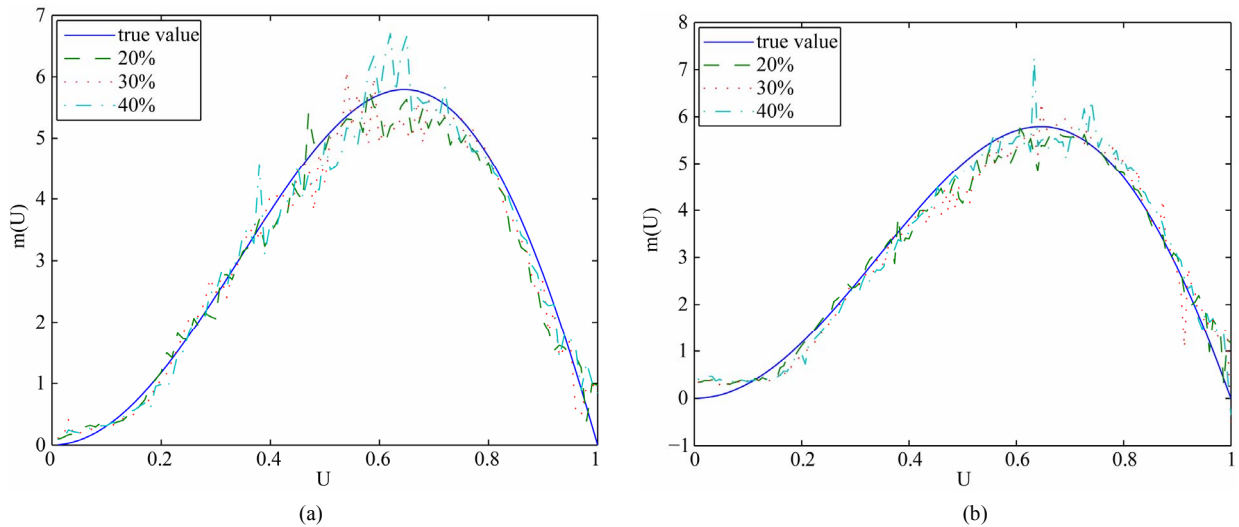


Figure 1. Curve estimates of  $\hat{m}(\cdot)$  for Example 1. (a)  $n = 100$ ; (b)  $n = 200$ .

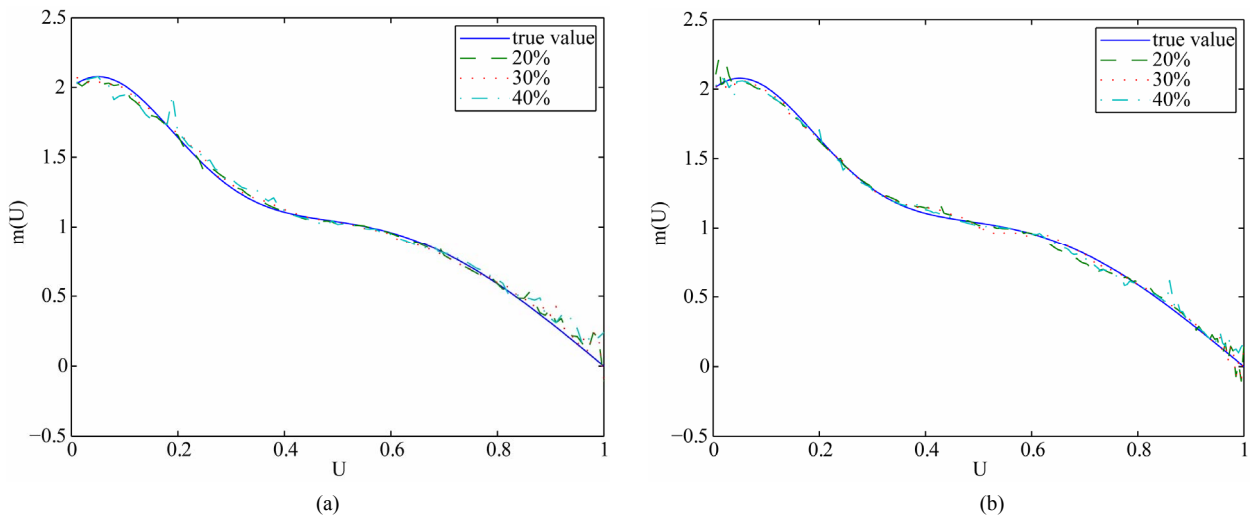
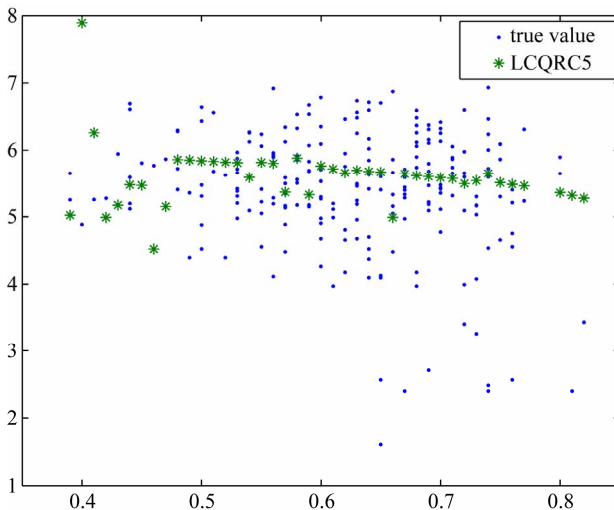


Figure 2. Curve estimates of  $\hat{m}(\cdot)$  for Example 2. (a)  $n = 100$ ; (b)  $n = 200$ .

**Table 5. Simulation results of  $\hat{y}$  for lung cancer data.**

Method	RMADE <sub>y</sub>	MADE <sub>y</sub>	RMSE <sub>y</sub>	MSE <sub>y</sub>
LCQR5	-	0.6877 (0.7074)	-	0.9711 (2.2640)
LCQRC5	0.9941	0.6836 (0.6788)	0.9536	0.9261 (2.0939)
NQRC0.5	1.1172	0.7683 (0.5738)	0.9454	0.9181 (1.3586)

**Figure 3. Curve estimates for lung cancer data.**

shows that the proposal is valid.

#### 4. Conclusion

In this work, we have focused on the LCQR for non-parametric model with censored data and its nice theoretical properties have been proven. The proposed approaches are demonstrated by simulation examples and real data applications. In addition, we believe the method can be extended to varying coefficient model (see [7]).

#### REFERENCES

- [1] J. L. Powell, "Least Absolute Deviations Estimation for the Censored Regression Model," *Journal of Econometrics*, Vol. 25, No. 3, 1984, pp. 303-325. [doi:10.1016/0304-4076\(84\)90004-6](https://doi.org/10.1016/0304-4076(84)90004-6)
- [2] S. Portnoy, "Censored Regression Quantiles," *Journal of the American Statistical Association*, Vol. 98, No. 464, 2003, pp. 1001-1012. [doi:10.1198/016214503000000954](https://doi.org/10.1198/016214503000000954)
- [3] L. Peng and Y. Huang, "Survival Analysis with Quantile Regression Models," *Journal of the American Statistical Association*, Vol. 103, No. 482, 2008, pp. 637-649. [doi:10.1198/016214508000000355](https://doi.org/10.1198/016214508000000355)
- [4] H. J. Wang and L. Wang, "Locally Weighted Censored Quantile Regression," *Journal of the American Statistical Association*, Vol. 104, No. 478, 2009, pp. 1117-1128. [doi:10.1198/jasa.2009.tm08230](https://doi.org/10.1198/jasa.2009.tm08230)
- [5] H. Zou and M. Yuan, "Composite Quantile Regression and the Oracle Model Selection Theory," *Annals of Statistics*, Vol. 36, No. 3, 2008, pp. 1108-1126. [doi:10.1214/07-AOS507](https://doi.org/10.1214/07-AOS507)
- [6] B. Kai, R. Li and H. Zou, "Local Composite Quantile Regression Smoothing: An Efficient and Safe Alternative to Local Polynomial Regression," *Journal of the Royal Statistical Society, Series B*, Vol. 72, No. 1, 2010, pp. 49-69. [doi:10.1111/j.1467-9868.2009.00725.x](https://doi.org/10.1111/j.1467-9868.2009.00725.x)
- [7] B. Kai, R. Li and H. Zou, "New Efficient Estimation and Variable Selection Methods for Semiparametric Varying-Coefficient Partially Linear Models," *Annals of Statistics*, Vol. 39, No. 1, 2011, pp. 305-332. [doi:10.1214/10-AOS842](https://doi.org/10.1214/10-AOS842)
- [8] R. Jiang, Z. G. Zhou, W. M. Qian and W. Q. Shao, "Single-Index Composite Quantile Regression," *Journal of the Korean Statistical Society*, Vol. 3, No. 3, 2012, pp. 323-332. [doi:10.1016/j.jkss.2011.11.001](https://doi.org/10.1016/j.jkss.2011.11.001)
- [9] R. Jiang, W. M. Qian and Z. G. Zhou, "Variable Selection and Coefficient Estimation via Composite Quantile Regression with Randomly Censored Data," *Statistics & Probability Letters*, Vol. 2, No. 2, 2012, pp. 308-317. [doi:10.1016/j.spl.2011.10.017](https://doi.org/10.1016/j.spl.2011.10.017)
- [10] A. Gannoun, J. Saracco, A. Yuan and G. Bonney, "Non-Parametric Quantile Regression with Censored Data," *Scandinavian Journal of Statistics*, Vol. 32, No. 4, 2005, pp. 527-550. [doi:10.1111/j.1467-9469.2005.00456.x](https://doi.org/10.1111/j.1467-9469.2005.00456.x)
- [11] C. L. Loprinzi, *et al.*, "Prospective Evaluation of Prognostic Variables from Patient-Completed Questionnaires. North Central Cancer Treatment Group," *Journal of Clinical Oncology*, Vol. 12, No. 3, 1994, pp. 601-607.
- [12] W. Gonzalez-Manteiga and C. Cadarso-Suarez, "Asymptotic Properties of a Generalized Kaplan-Meier Estimator with Some Applications," *Journal of Nonparametric Statistics*, Vol. 4, No. 1, 1994, pp. 65-78. [doi:10.1080/10485259408832601](https://doi.org/10.1080/10485259408832601)
- [13] K. Knight, "Limiting Distributions for  $L_1$  Regression Estimators under General Conditions," *Annals of Statistics*, Vol. 26, No. 2, 1998, pp. 755-770. [doi:10.1214/aos/1028144858](https://doi.org/10.1214/aos/1028144858)

### Appendix

**Lemma 1.** Assume assumption A1 hold. Then

$$\begin{aligned} \|\hat{F} - F_0\| &= \sup_i \sup_U |\hat{F}(t|U) - F_0(t|U)| \\ &= O_p\left((\log n)^{1/2} n^{-1/4-\gamma_0/2}\right), \end{aligned}$$

where  $0 < \gamma_0 < 1/4$ .

**Proof.** This follows directly from theorem 2.1 of [12].

**Proof of Theorem 1** Let

$$\begin{aligned} \sqrt{nh}(a_k - m(u_0) - \sigma(u_0)c_k) &= \eta_k, \\ h\sqrt{nh}(b - m'(u_0)) &= v, \quad \Delta_{i,k} = (\eta_k + vx_i)/\sqrt{nh}, \\ x_i &= (U_i - u_0)/h, \\ d_i(u_0) &= m(U_i) - m(u_0) - m'(u_0)(U_i - u_0) \\ &\quad + c_k(\sigma(U_i) - \sigma(u_0)) \end{aligned}$$

Then  $(\eta_1, \dots, \eta_q, v)$  is the minimizer of the following criterion:

$$\begin{aligned} L_n &= \sum_{k=1}^q \sum_{i=1}^n K_i \left[ \omega_i(\hat{F}, \tau_k) \left\{ \rho_{\tau_k}(A_i - \sigma(U_i)c_k + d_i(u_0) - \Delta_{i,k}) - \rho_{\tau_k}(A_i - \sigma(U_i)c_k + d_i(u_0)) \right\} \right. \\ &\quad \left. + (1 - \omega_i(\hat{F}, \tau_k)) \left\{ \rho_{\tau_k}(B_i - \sigma(U_i)c_k + d_i(u_0) - \Delta_{i,k}) - \rho_{\tau_k}(B_i - \sigma(U_i)c_k + d_i(u_0)) \right\} \right], \end{aligned}$$

where  $A_i = Y_i - m(U_i)$ ,  $B_i = Y_i^{+\infty} - m(U_i)$  and  $K_i = K\{(U_i - u_0)/h\}$ . To apply the identity ([13])

$$\rho_{\tau}(x - y) - \rho_{\tau}(x) = y \left\{ I(x \leq 0) - \tau \right\} + \int_0^y \left\{ I(x \leq z) - I(x \leq 0) \right\} dz,$$

we have

$$\begin{aligned} L_n &= \sum_{k=1}^q \sum_{i=1}^n K_i \Delta_{i,k} \omega_i(\hat{F}, \tau_k) \left[ I(A_i < \sigma(U_i)c_k - d_i(u_0)) - \tau_k \right] \\ &\quad + \sum_{k=1}^q \sum_{i=1}^n K_i \int_0^{\Delta_{i,k}} \omega_i(\hat{F}, \tau_k) \left[ I(A_i < \sigma(U_i)c_k - d_i(u_0) + t) - I(A_i < \sigma(U_i)c_k - d_i(u_0)) \right] dt \\ &\quad + \sum_{k=1}^q \sum_{i=1}^n K_i \Delta_{i,k} (1 - \omega_i(\hat{F}, \tau_k)) \left[ I(B_i < \sigma(U_i)c_k - d_i(u_0)) - \tau_k \right] \\ &\quad + \sum_{k=1}^q \sum_{i=1}^n K_i \int_0^{\Delta_{i,k}} (1 - \omega_i(\hat{F}, \tau_k)) \left[ I(B_i < \sigma(U_i)c_k - d_i(u_0) + t) - I(B_i < \sigma(U_i)c_k - d_i(u_0)) \right] dt \\ &\triangleq L_{1n} + L_{2n} + L_{3n} + L_{4n}. \end{aligned}$$

Since  $Y^{+\infty}$  is any value sufficiently large to exceed all  $m(U_i)$ ,  $L_{3n} = -\sum_{k=1}^q \sum_{i=1}^n K_i \Delta_{i,k} (1 - \omega_i(\hat{F}, \tau_k)) \tau_k$  and  $L_{4n} = 0$ ,

then  $L_{1n} + L_{3n} = \sum_{k=1}^q \sum_{i=1}^n K_i \Delta_{i,k} \left( \omega_i(\hat{F}, \tau_k) I(A_i < \sigma(U_i)c_k - d_i(u_0)) - \tau_k \right)$ .

Denote  $L_{2n} = \sum_{k=1}^q L_{2n}^{(k)}$ , where

$$L_{2n}^{(k)} = \sum_{i=1}^n K_i \int_0^{\Delta_{i,k}} \omega_i(\hat{F}, \tau_k) \left[ I(A_i < \sigma(U_i)c_k - d_i(u_0) + t) - I(A_i < \sigma(U_i)c_k - d_i(u_0)) \right] dt.$$

$$\begin{aligned} &\omega_i(F, \tau_k) I(A_i < \sigma(U_i)c_k - d_i(u_0)) \\ &= I(C_i > m(U_i) + \sigma(U_i)c_k - d_i(u_0), T_i \leq m(U_i) + \sigma(U_i)c_k - d_i(u_0)) + I(C_i \leq m(U_i) + \sigma(U_i)c_k - d_i(u_0), T_i \leq C_i) \\ &\quad + I(C_i \leq m(U_i) + \sigma(U_i)c_k - d_i(u_0), T_i > C_i) \times \left( 1 - \frac{1 - \tau_k}{1 - F_i(C_i)} I(F_i(C_i) < \tau_k) \right). \end{aligned}$$

By the conditional independence of  $T$  and  $C$  given  $U$ , we have

$$E[I(C > t, T < t)|U] = P(C > t|U)P(T < t|U) = [1 - G(t|U)]F_0(t|U),$$

$$E[I(C < t, T < C)|U] = E_{C|U}[I(C < t)P(T < C|C, U)] = \int_{-\infty}^t F_0(s|U)g(s|U)ds.$$

Therefore,

$$E[\omega_i(F_0, \tau_k)I(A_i \leq \sigma(U_i)c_k - d_i(u_0)|U)]$$

$$= \{1 - G(m(U) + \sigma(U_i)c_k - d|U)\}F_0(m(U) + \sigma(U_i)c_k - d|U) + \int_{-\infty}^{m(U) + \sigma(U_i)c_k - d} F_0(s|U)g(s|U)ds$$

$$+ \int_{-\infty}^{m(U) + \sigma(U_i)c_k - d} \{1 - F_0(s|U)\} \times \left[1 - \frac{1 - \tau_k}{1 - F_0(s|U)} I(F_0(s|U) < \tau_k)\right] g(s|U)ds.$$

By Lemma 1, we have

$$\left| \sum_{i=1}^n K_i \int_0^{\Delta_{i,k}} E\left\{[\omega_i(\hat{F}, \tau_k) - \omega_i(F_0, \tau_k)][I(A_i \leq \sigma(U_i)c_k + t) - I(A_i \leq \sigma(U_i)c_k)]|U\right\} dt \right|$$

$$\leq \sum_{i=1}^n K_i \int_0^{\Delta_{i,k}} E\left\{|\omega_i(\hat{F}, \tau_k) - \omega_i(F_0, \tau_k)| |I(A_i \leq \sigma(U_i)c_k + t) - I(A_i \leq \sigma(U_i)c_k)| |U\right\} dt$$

$$= \sum_{i=1}^n K_i \int_0^{\Delta_{i,k}} E\left\{O(|\hat{F} - F_0|) |I(A_i \leq \sigma(U_i)c_k + t) - I(A_i \leq \sigma(U_i)c_k)| |U\right\} dt$$

$$= \sum_{i=1}^n K_i \int_0^{\Delta_{i,k}} E\left\{|I(A_i \leq \sigma(U_i)c_k + t) - I(A_i \leq \sigma(U_i)c_k)| |U\right\} dt \cdot o_p(1) \rightarrow 0.$$

Then, we can obtain

$$E[L_{2n}^{(k)}|U]$$

$$= \sum_{i=1}^n K_i \int_0^{\Delta_{i,k}} E\left\{(\omega_i(\hat{F}, \tau_k)) [I(A_i < \sigma(U_i)c_k - d_i(u_0) + t) - I(A_i < \sigma(U_i)c_k - d_i(u_0))] |U\right\} dt$$

$$= \sum_{i=1}^n K_i \int_0^{\Delta_{i,k}} E(\omega_i(F_0, \tau_k) [I(A_i < \sigma(U_i)c_k - d_i(u_0) + t) - I(A_i < \sigma(U_i)c_k - d_i(u_0))] |U) dt + o_p(1)$$

$$\rightarrow \frac{1}{2} f_U(u_0)(\eta_k, v) \begin{pmatrix} \Gamma_k & 0 \\ 0 & \mu_2 \Gamma_k \end{pmatrix} (\eta_k, v)^T.$$

$$\text{Var}[L_{2n}^{(k)}|U]$$

$$= \sum_{i=1}^n \text{Var}\left\{\left[K_i \int_0^{\Delta_{i,k}} \omega_i(\hat{F}, \tau_k) [I(A_i < \sigma(U_i)c_k - d_i(u_0) + t) - I(A_i < \sigma(U_i)c_k - d_i(u_0))] dt\right] |U\right\}$$

$$\leq \sum_{i=1}^n E\left\{\left[K_i \int_0^{\Delta_{i,k}} (\omega_i(\hat{F}, \tau_k) [I(A_i < \sigma(U_i)c_k - d_i(u_0) + t) - I(A_i < \sigma(U_i)c_k - d_i(u_0))]) dt\right]^2 |U\right\}$$

$$= o\left(\sum_{i=1}^n K_i^2 \Delta_{i,k}^2\right) = o_p(1).$$

So, we can obtain  $L_{2n}^{(k)} \rightarrow \frac{1}{2} f_U(u_0)(\eta_k, v) \begin{pmatrix} \Gamma_k & 0 \\ 0 & \mu_2 \Gamma_k \end{pmatrix} (\eta_k, v)^T$ , then

$$\hat{\eta}_k = f_U^{-1}(u_0) \Gamma_k^{-1} M_k,$$

where



$$M_k(u_0) = \frac{1}{\sqrt{nh}} \sum_{i=1}^n K_i \left[ \tau_k - \omega_i(\hat{F}, \tau_k) I(Y_i < m(U_i) + \sigma(U_i)c_k - d_i(u_0)) \right].$$

Note that the error is symmetric, thus  $\sum_{k=1}^q c_k = 0$ , then it follows that

$$\sqrt{nh} [\hat{m}(U_0) - m(U_0)] = \frac{1}{q} \sum_{k=1}^q \hat{\eta}_k.$$

Since  $E[\tau_k - \omega_i(F_0, \tau_k) I(Y_i < m(U_i) + \sigma(U_i)c_k)] = 0$ , then

$$\begin{aligned} & \frac{1}{\sqrt{nh}} E[M_k(u_0)|U] \\ &= \frac{1}{nh} \sum_{i=1}^n K_i E\left[\left[\tau_k - \omega_i(\hat{F}, \tau_k) I(Y_i < m(U_i) + \sigma(U_i)c_k - d_i(u_0))\right] | U\right) \\ &= \frac{1}{nh} \sum_{i=1}^n K_i E\left[\left[\omega_i(F_0, \tau_k) \{I(Y_i < m(U_i) + \sigma(U_i)c_k) - I(Y_i < m(U_i) + \sigma(U_i)c_k - d_i(u_0))\}\right] | U\right) + o_p(1) \\ &= \frac{1}{nh} \sum_{i=1}^n K_i d_i(u_0) \Gamma_k (1 + o_p(1)) \\ &= \frac{1}{2} f_U(u_0) \mu_2 h^2 m''(u_0) \Gamma_k + o_p(h^2). \end{aligned}$$

So, we can obtain

$$\begin{aligned} \text{bias}[\hat{m}(u_0)|U] &= \frac{f_U^{-1}(u_0)}{q\sqrt{nh}} \sum_{k=1}^q \Gamma_k^{-1} E[M_k|U] = \frac{1}{2} \mu_2 h^2 m''(u_0), \\ \text{var}[\hat{m}(u_0)|U] &= \frac{1}{nh} \Sigma + o_p\left(\frac{1}{nh}\right). \end{aligned}$$

This completes the proof.