

# Nonparametric Estimation of a Truncated Regression Model with an Application to GATE Program

Hariyanto<sup>1</sup>, Ira Madanisa<sup>2</sup>, Julham Efendi Nasution<sup>3</sup>

<sup>1</sup>Department of Mathematics, University of North Sumatera, Indonesia, PP.AR-DHLATULHASANAH, Medan, Indonesia

<sup>2,3</sup>Department of Mathematics, University of North Sumatera, Indonesia

**Abstract:** In this paper we consider identification and estimation of a nonparametric location scale model. We first use the truncated data. Then we use truncated regression model. Truncated regression is used to model dependent variables for which some of the observations are not included in the analysis because of the value of the dependent variable. In the latter case we propose a simple estimation procedure based on combining conditional quantile estimators for three distinct quantiles. The new estimator is shown to converge at the optimal nonparametric rate with a limiting normal distribution. A small scale simulation study indicates that the proposed estimation procedure performs well in finite samples. We also present an empirical application on GATE Program using example data test.

**Keywords:** GATE Program, Language Score, Truncated Regression.

## 1. Introduction

The nonparametric location-scale model is usually of the form:

$$y_i = \mu(x_i) + \sigma(x_i)\epsilon_i$$

where  $x_i$  is an observed  $d$ -dimensional random vector and  $\epsilon_i$  is an unobserved random variable, distributed independently of  $x_i$ , and assumed to be centered around zero in some sense. The functions  $\mu(\cdot)$  and  $\sigma(\cdot)$  are unknown. In this paper, we consider extending the nonparametric location-scale model to accommodate censored data. The advantage of our nonparametric approach here is that economic theory rarely provides any guidance on functional forms in relationships between variables. To allow for censoring, we work within the latent dependent variable framework, as is typically done for parametric and semiparametric models. We thus consider a model of the form:

$$y_i^* = \mu(x_i) + \sigma(x_i)\epsilon_i$$

$$y_i = \max(y_i^*, 0)$$

where  $y_i^*$  is a latent dependent variable, which is only observed if it exceeds the fixed censoring point, which we assume without loss of generality is 0. We consider identification and estimation of  $\mu(x_i)$  after imposing the location restriction that the median of  $\epsilon_i = 0$ . We emphasize that our results allow for identification of  $\mu(x_i)$  on the entire support of  $x_i$ . This is in contrast to identifying and estimating  $\mu(x_i)$  only in the region where it exceeds the censoring point, which could be easily done by extending Powell's (1984) CLAD estimator to a nonparametric setting. One situation is when the data set is heavily censored. In this case,  $\mu(x_i)$  will be less than the censoring point for a large portion of the support of  $x_i$ , requiring estimation at these points necessary to draw meaningful inference regarding its shape. Our approach is based on a structural relationship between the conditional median and upper quantiles which

holds for observations where  $\mu(x_i) \geq 0$ . This relationship can be used to motivate an estimator for  $\mu(x_i)$  in the region where it is negative. Our results are thus based on the condition

$$P_X(x_i : \mu(x_i) \geq 0) > 0$$

where  $P_X(\cdot)$  denotes the probability measure of the random variable  $x_i$ .

The paper is organized as follows. The next section explains the key identification condition, and motivates a way to estimate the function  $\mu(\cdot)$  at each point in the support of  $x_i$ . Section 3 introduces the new estimation procedure and establishes the asymptotic properties of this estimator when the identification condition is satisfied. Section 4 considers an extension of the estimation procedure to estimate the distribution of the disturbance term. Section 5 explores the finite sample properties of the estimator through the results of a simulation study. Section 6 presents an empirical application STIFIN test, in which we estimate the survivor function in the region beyond the censoring point. Section 7 concludes by summarizing results.

## 2. Censored and Truncated Data : Comparison Definitions

•  $Y$  is *censored* when we observe  $X$  for all observations, but we only know the true value of  $Y$  for a restricted range of observations. Values of  $Y$  in a certain range are reported as a single value or there is significant clustering around a value, say 0.

- If  $Y = k$  or  $Y > k$  for all  $Y \Rightarrow Y$  is censored from below or left-censored.

- If  $Y = k$  or  $Y < k$  for all  $Y \Rightarrow Y$  is censored from above or right-censored.

We usually think of an uncensored  $Y$ ,  $Y^*$ , the true value of  $Y$  when the censoring mechanism is not applied. We typically have all the observations for  $\{Y, X\}$ , but not  $\{Y^*, X\}$ .

•  $Y$  is truncated when we only observe  $X$  for observations where  $Y$  would not be censored. We do not have a full sample for  $\{Y, X\}$ , we exclude observations based on characteristics of  $Y$ .

### 3. Estimation Procedure and Asymptotic Properties

#### 3.1 Estimation Procedure

In this section we consider estimation of the function  $\mu(\cdot)$ . Our procedure will be based on our identification results in the previous section, and involves nonparametric quantile regression at different quantiles and different points in the support of the regressors. Our asymptotic arguments are based on the local polynomial estimator for conditional quantile functions introduced in Chaudhuri(1991a,b). For expositional ease, we only describe this nonparametric estimator for a polynomial of degree 0, and refer readers to Chaudhuri(1991a,b), Chaudhuri et al.(1997), Chen and Khan(2000,2001), and Khan(2001) for the additional notation involved for polynomials of arbitrary degree.

First, we assume the regressor vector  $x_i$  can be partitioned as  $(x_i^{ds}, x_i^c)$  where the  $d_{ds}$ -dimensional vector  $x_i^{ds}$  is discretely distributed, and the  $d_c$ -dimensional vector  $x_i^c$  is continuously distributed. We let  $C_n(x_i)$  denote the cell of observation  $x_i$  and let  $h_n$  denote the sequence of bandwidths which govern the size of the cell. For some observation  $x_j, j \neq i$ , we let  $x_j \in C_n(x_i)$  denote that  $x_j^{(ds)} = x_i^{(ds)}$  and  $x_j^c$  lies in the  $d_c$ -dimensional cube centered at  $x_i^c$  with side length  $2h_n$ .

Let  $I[\cdot]$  be an indicator function, taking the value 1 if its argument is true, and 0 otherwise. Our estimator of the conditional  $\alpha^{\text{th}}$  quantile function at a point  $x_i$  for any  $\alpha \in (0, 1)$  involves  $\alpha$ -quantile regression (see Koenker and Bassett (1978)) on observations which lie in the defined cells of  $x_i$ . Specifically, let  $\theta$  minimize:

$$\sum_{j=1}^n I[x_j \in C_n(x_i)] \rho_{\alpha}(y_j - \theta)$$

where  $\rho_{\alpha}(\cdot) \equiv \alpha|\cdot| + (2\alpha - 1)(\cdot)I[\cdot < 0]$ .

Our estimation procedure will be based on a random sample of  $n$  observations of the vector  $(y_i, x_i)$  and involves applying the local polynomial estimator at three stages. Throughout our description,  $\hat{\cdot}$  will denote estimated values.

**1) Local Constant Estimation of the Conditional Median Function.** In the first stage, we estimate the conditional median at each point in the sample, using a polynomial of degree 0. We will let  $h_{1n}$  denote the bandwidth sequence used in this stage. Following the terminology of Fan(1992), we refer to this as a local constant estimator, and denote the estimated values by  $\hat{q}_{0.5}(x_i)$ . Recalling that our identification result is based on observations for which the median function is positive, we assigns weights to these estimated values using a weighting function, denoted by  $w(\cdot)$ . Essentially,  $w(\cdot)$  assigns 0 weight to observations in the sample for which the estimated value of the median function is 0, and assigns positive weight for estimated values which are positive.

**2) Weighted Average Estimation of the Disturbance Quantiles** In the second stage, the unknown quantiles  $c_{\alpha 1}, c_{\alpha 2}$  are estimated (up to the scalar constant  $c$ ) by a weighted average of local polynomial estimators of the quantile functions for the higher quantiles  $\alpha 1, \alpha 2$ . In this stage, we use a polynomial of degree  $k$ , and denote the second stage bandwidth sequence by  $h_{2n}$ .

We let  $\hat{c}_1, \hat{c}_2$  denote the estimators of the unknown

constants  $\frac{c_{\alpha 1}}{\Delta c}, \frac{c_{\alpha 2}}{\Delta c}$ , and define them as:

$$\hat{c}_1 = \frac{\frac{1}{n} \sum_{i=1}^n \tau(x_i) w(\hat{q}_{0.5}(x_i)) \cdot \frac{(\hat{q}_{\alpha 1}(x_i) - \hat{q}_{0.5}^{(p)}(x_i))}{(\hat{q}_{\alpha 2}(x_i) - \hat{q}_{\alpha 1}(x_i))}}{\frac{1}{n} \sum_{i=1}^n \tau(x_i) w(\hat{q}_{0.5}(x_i))}$$

$$\hat{c}_2 = \frac{\frac{1}{n} \sum_{i=1}^n \tau(x_i) w(\hat{q}_{0.5}(x_i)) \cdot \frac{(\hat{q}_{\alpha 2}(x_i) - \hat{q}_{0.5}^{(p)}(x_i))}{(\hat{q}_{\alpha 2}(x_i) - \hat{q}_{\alpha 1}(x_i))}}{\frac{1}{n} \sum_{i=1}^n \tau(x_i) w(\hat{q}_{0.5}(x_i))}$$

where  $\tau(x_i)$  is a trimming function, whose support, denoted by  $X_{\tau}$ , is a compact set which lies strictly in the interior of  $X$ . The trimming function serves to eliminate “boundary effects” that arise in nonparametric estimation. We use the superscript (p) to distinguish the estimator of the median function in this stage from that in the first stage.

#### 3) Local Polynomial Estimation at the Point of Interest

Letting  $x$  denote the point at which the function  $\mu(\cdot)$  is to be estimated at, we combine the local polynomial estimator, with polynomial order  $k$  and bandwidth sequence  $h_{3n}$ , of the conditional quantile function at  $x$  using quantiles  $\alpha 1, \alpha 2$ , with the estimator of the unknown disturbance quantiles, to yield the estimator of  $\mu(x)$ :

$$\hat{\mu}(x) = \hat{c}_2 \hat{q}_{\alpha 1}(x) - \hat{c}_1 \hat{q}_{\alpha 2}(x)$$

### 4. Estimating the Distribution of $\epsilon_i$

As mentioned in Section 2, the distribution of the random variable  $\epsilon_i$  is identified for all quantiles exceeding  $\alpha_0 \equiv \inf\{\alpha: \sup_{x \in X} q_{\alpha}(x) > 0\}$ . In this section we consider estimation of these quantiles, and the asymptotic properties of the estimator. Estimating the distribution of  $\epsilon_i$  is of interest for two reasons. First, the econometrician may be interested in estimating the entire model, which would require estimators of  $\sigma(x_i)$  and the distribution of  $\epsilon_i$  as well as of  $\mu(x_i)$ . Second, the estimator can be used to construct tests of various parametric forms of the distribution of  $\epsilon_i$ , and the results of these tests could then be used to adopt a (local) likelihood approach to estimating the function  $\mu(x_i)$ .

Before proceeding, we note that the distribution of  $\epsilon_i$  is only identified up to scale, and we impose the scale normalization that  $c_{0.75} - c_{0.25} \equiv 1$ . We also assume without loss of generality that  $\alpha 0 \leq 0.25$ . To estimate  $c_{\alpha}$  for any  $\alpha \geq \alpha_0$ , we let  $\alpha = \min(\alpha, 0.5)$  and define our estimator as

$$\hat{c}_{\alpha} = \frac{\frac{1}{n} \sum_{i=1}^n \tau(x_i) w(\hat{q}_{\alpha-}(x_i)) \cdot (\hat{q}_{\alpha}(x_i) - \hat{q}_{0.75}^{(p)}(x_i))}{\frac{1}{n} \sum_{i=1}^n \tau(x_i) w(\hat{q}_{\alpha-}(x_i)) \cdot (\hat{q}_{0.75}(x_i) - \hat{q}_{0.25}(x_i))}$$

The proposed estimator, which involves averaging nonparametric estimators, will converge at the parametric ( $\sqrt{n}$ ) rate and have a limiting normal distribution, as can be

rigorously shown using similar arguments found in Chen and Khan(1999b).

• The estimated  $\beta_k$  shows the effect of  $x_{ki}$  on  $y_i$ . Thus,

$$\frac{\partial E(y_i | X_i, y_i^* > 0)}{\partial X_{k,i}} = \beta_k + \frac{\partial E(\varepsilon_i | y_i^* > 0)}{\partial X_{k,i}}$$

$$= \beta_k + \sigma \frac{\partial \lambda}{\partial X_{k,i}} = \beta_k + \sigma (\lambda_i^2 - \alpha_i \lambda_i) \left( -\frac{\beta_k}{\sigma} \right)$$

$$= \beta_k (1 - \lambda_i^2 - \alpha_i \lambda_i) = \beta_k (1 - \delta_i)$$

where  $\delta_i = \lambda(\alpha_i) [\lambda(\alpha_i) + \alpha_i]$ ,  $0 < \delta_i < 1$

### 5. Truncated Regression

• Data truncation is (B-1): the truncation is based on the y-variable.

• We have the following regression satisfies all CLM assumptions:

$$y_i = x_i' \beta + \varepsilon_i, \quad \varepsilon_i \sim N(0, \sigma^2)$$

- We sample only if  $y_i < c_i$  - Observations dropped if  $y_i \geq c_i$  by design.

- We know the exact value of  $c_i$  for each person.

• Given the normality assumption for  $\varepsilon_i$ , ML is easy to apply.

- For each,  $\varepsilon_i = y_i - x_i' \beta$ , the likelihood contribution is  $f(\varepsilon_i)$ .

- But, we select sample only if  $y_i < c_i$

=> we have to use the density function of  $\varepsilon_i$  conditional on  $y_i < c_i$ :

$$f(\varepsilon_i | y_i < c) = \frac{f(\varepsilon_i | \varepsilon_i < c_i - x_i' \beta)}{P(u_i < c_i - x_i' \beta)}$$

$$= \frac{f(\varepsilon_i)}{P\left(\frac{\varepsilon_i}{\sigma} < \frac{c_i - x_i' \beta}{\sigma}\right)} = \frac{f(\varepsilon_i)}{\Phi\left(\frac{c_i - x_i' \beta}{\sigma}\right)}$$

$$= \frac{1}{\Phi\left(\frac{c_i - x_i' \beta}{\sigma}\right)} \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{\varepsilon_i^2}{2\sigma^2}\right)$$

$$= \frac{\frac{1}{\sigma} \phi\left(\frac{\varepsilon_i}{\sigma}\right)}{\Phi\left(\frac{c_i - x_i' \beta}{\sigma}\right)}$$

• The likelihood contribution for  $i^{th}$  observation is given by

$$L_i = \frac{\frac{1}{\sigma} \phi\left(\frac{y_i - x_i' \beta}{\sigma}\right)}{\Phi\left(\frac{x_i' \beta}{\sigma}\right)}$$

ln(joint density of N values of  $y^*$ )

• The likelihood function is given by

$$\text{Log } L(\beta, \sigma) = \sum_{i=1}^N \log L_i = -\frac{N}{2} [\log(2\pi) + \log(\sigma^2)] - \frac{1}{2\sigma^2} \sum_{i=1}^N \varepsilon_i^2$$

$$- \sum_{i=1}^N \log \left[ \Phi\left(\frac{x_i' \beta}{\sigma}\right) \right]$$

log(joint probability of  $y^* > 0$ )

• The values of  $(\beta, \sigma)$  that maximizes Log L are the ML estimators

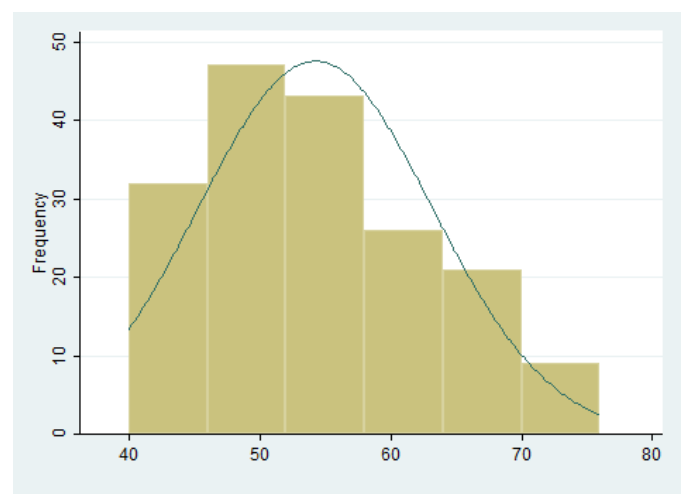
### 6. Application to GATE Program

A study of students in a special GATE (gifted and talented education) program wishes to model achievement as a function of language skills and the type of program in which the student is currently enrolled. A major concern is that students are required to have a minimum achievement score of 40 to enter the special program. Thus, the sample is truncated at an achievement score of 40.

Variable	Obs	Mean	Std. Dev.	Min	Max
achiv	178	54.23596	8.96323	41	76
langscore	178	54.01124	8.944896	31	67

Summary for variables: achiv  
by categories of: prog (type of program)

prog	N	mean	sd
general	40	51.575	7.97074
academic	101	56.89109	9.018759
vocation	37	49.86486	7.276912
<b>Total</b>	<b>178</b>	<b>54.23596</b>	<b>8.96323</b>



type of program	Freq.	Percent	Cum.
general	40	22.47	22.47

academic | 101 56.74 79.21  
 vocation | 37 20.79 100.00

-----+-----  
 Total | 178 100.00  
 Fitting full model:

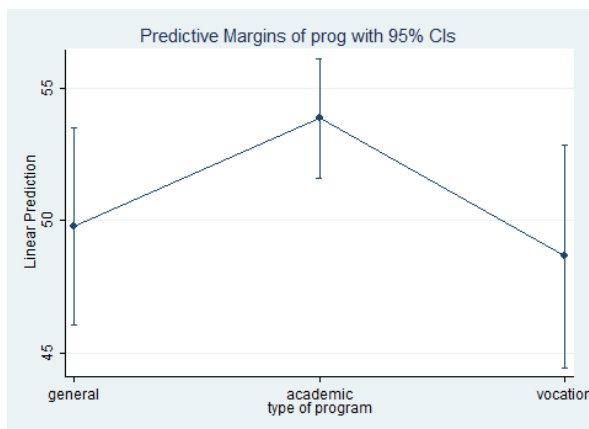
Iteration 0: log likelihood = -598.11669  
 Iteration 1: log likelihood = -591.68374  
 Iteration 2: log likelihood = -591.31208  
 Iteration 3: log likelihood = -591.30981  
 Iteration 4: log likelihood = -591.30981

```

Truncated regression
Limit: lower = 40          Number of obs = 178
      upper = +inf        Wald chi2(3) = 54.76
Log likelihood = -591.30981      Prob > chi2 = 0.0000
-----+-----
      achiev |   Coef.   Std. Err.   z   P>|z|   [95% Conf. Interval]
-----+-----
      langscore | .7125775   .1144719   6.22   0.000   .4882168   .9369383
      |
      prog |
      2 | 4.065219   2.054938   1.98   0.048   .0376131   8.092824
      3 | -1.135863   2.669961   -0.43   0.671   -6.368891   4.097165
      |
      _cons | 11.30152   6.772731   1.67   0.095   -1.97279   24.57583
-----+-----
      /sigma | 8.755315   .666803   13.13   0.000   7.448405   10.06222
-----+-----
Predictive margins      Number of obs = 178
Model VCE      : OIM
Expression      : Linear prediction, predict()
-----+-----
      |      Delta-method
      |      Margin   Std. Err.   z   P>|z|   [95% Conf. Interval]
-----+-----
      prog |
      1 | 49.78871   1.897166   26.24   0.000   46.07034   53.50709
      2 | 53.85393   1.150041   46.83   0.000   51.59989   56.10797
      3 | 48.65285   2.140489   22.73   0.000   44.45757   52.84813
    
```

In the table above, we can see that the expected mean of **avchiv** for the first level of **prog** is approximately 49.79; the expected mean for level 2 of **prog** is 53.85; the expected mean for the third level of **prog** is 48.65.

**Marginsplot**



**7. Conclusion**

In the output above, we can see that the expected mean of **avchiv** for the first level of **prog** is approximately 49.79; the expected mean for level 2 of **prog** is 53.85; the expected mean for the third level of **prog** is 48.65.

**References**

- [1] Chaudhuri, P. (1991a) "Nonparametric Quantile Regression", *Annals of Statistics*, 19, 760-777.
- [2] Chaudhuri, P. (1991b) "Global Nonparametric Estimation of Conditional Quantiles and their Derivatives", *Journal of Multivariate Analysis*, 39, 246-269.
- [3] Chaudhuri, P., K. Doksum, and A. Samarov (1997) "On Average Derivative Quantile Regression", *Annals of Statistics*, 25, 715-744.
- [4] Chen, S., Dahl. G. B., and S. Khan (2002) "Nonparametric Identification and Estimation of a Censored Regression Model with an Application to Unemployment Insurance Receipt", Center For Labor Economics University of California, Berkeley Working Paper no.54
- [5] Chen, S. and S. Khan (2000) "Estimation of Censored Regression Models in the Presence of Nonparametric Multiplicative Heteroskedasticity", *Journal of Econometrics*, 98, 283-316.
- [6] Chen, S. and S. Khan (2001) "Semiparametric Estimation of a Partially Linear Censored Regression Model", *Econometric Theory*, 17, 567-590.
- [7] Fan, J. (1992) "Design-adaptive Nonparametric Regression", *Journal of the American Statistical Association*, 87, 998-1004.
- [8] Fan, J. and I. Gijbels (1996) *Local Polynomial Modelling and its Applications*, New York: Chapman and Hall.
- [9] Khan, S. (2001) "Two Stage Rank Estimation of Quantile Index Models", *Journal of Econometrics*, 100, 319-355.
- [10] Koenker, R. and G.S. Bassett Jr. (1978) "Regression Quantiles", *Econometrica*, 46, 33-50.
- [11] Powell, J.L. (1984) "Least Absolute Deviations Estimation for the Censored Regression Model", *Journal of Econometrics*, 25, 303-325.
- [12] Rosenbaum, P.R. and D.B. Rubin (1983) "The Central Role of the Propensity Score in Observational Studies for Causal Effects", *Biometrika*, 70, 41-55