

Quantile Regression with Monotone Dropout Missingness

May 11, 2013

Abstract

1 Introduction

Quantile regression is an attractive way to study the relationship between response and covariates when one (or several) quantiles are of interest as compared to mean regression. The dependence between upper or lower quantiles of the response variable and the covariates are expected to vary differentially relative to that of the average. This is often of interest in econometrics, educational studies, biomedical studies, and environment studies (Yu and Moyeed (2001), Buchinsky (1994), Buchinsky (1998), He et al. (1998), Koenker and Machado (1999), Wei et al. (2006), Yu et al. (2003)). More comprehensive review of applications of quantile regression was presented in Koenker (2005).

Unlike mean regression, quantile regression is more robust to outliers and provides more information about how covariates affect quantiles. For example, both as a summary statistics of data, median gives a measure more robust than arithmetic mean when outliers exist. In addition, mean regression only focus on the change of covariates on the mean, while quantile regression can offer a more complete description of the conditional distribution of the response. Different effects of covariates can be assumed for different quantiles.

The traditional frequentist approach was proposed by Koenker and Bassett (1978) for a single quantile (τ) with estimators derived by minimizing a loss function. The popularity of this approach is due to its computational efficiency by linear programming, well-developed asymptotic properties, and straightforward extensions to simultaneous quantile regression and random effect models. However, asymptotic inference may not be accurate for small sample sizes.

Bayesian approaches offer exact inference. Motivated by the loss (check) function, Yu and Moyeed (2001) proposed an asymmetric Laplace distribution for the error term, such that maximizing the posterior distribution is equivalent to minimizing the check function. Other than parametric Bayesian approaches, some semiparametric methods have been proposed for median regression. Walker and Mallick (1999) used a diffuse finite Pólya Tree prior for the error term. Kottas and Gelfand (2001) modeled the error by two families of median zero distribution using a mixture Dirichlet process priors, which is very useful for unimodal

error distributions. [Hanson and Johnson \(2002\)](#) adopted mixture of Pólya Tree prior in median regression, which is more robust in terms of multimodality and skewness. Other recent approaches include quantile pyramid priors, mixture of Dirichlet process priors of multivariate normal distributions and infinite mixture of Gaussian densities which put quantile constraints on the residuals ([Hjort and Petrone \(2007\)](#), [Hjort and Walker \(2009\)](#), [Kottas and Krnjajić \(2009\)](#), [Reich et al. \(2010\)](#)).

However, above methods deal with complete data without missingness. There are still few articles about quantile regression with missingness. [Yuan and Yin \(2010\)](#) introduced a Bayesian quantile regression approach for longitudinal study with nonignorable missing data. They used random effects to explain the within-subject correlation and applied a l_2 penalty in the traditional quantile regression check function to shrink toward the common population values. However, the quantile regression coefficients are still conditional on the random effects, which is not of interest if we are looking into the marginal relationship. [Roy and Daniels \(2008\)](#) proposed a pattern mixture model for data with nonignorable dropout, which borrowed idea from [Heagerty \(1999\)](#). But their methods are studying the marginal covariate effects on the mean. We still use the ideas from them, but we will fit them into quantile regression models, which will be illustrated below.

The structure of this article is as below: first, we introduce our quantile regression methods to deal with monotone dropout in general case in section 2, including sensitivity analysis and computational details. And we also used simulation studies to demonstrate the performance of our algorithm in section 3. We also applied our approach on tours data, and we demonstrated the process in section 4. Finally, discussion and conclusions are in section 5.

2 Model

In this section, we first introduce some notations on monotone dropout, then describe our proposed quantile regression model in section 2.1. And we will give more details on deploying MAR and MNAR and computation in section 2.2 and section 2.3.

Under monotone dropout scenario, without loss of generality, denote $S_i \in \{1, 2, \dots, J\}$ to be the follow up time, and $\mathbf{Y}_i = (Y_{i1}, Y_{i2}, \dots, Y_{iJ})^T$ to be the response vector for subject i , where J is the maximum follow up time. Y_{i1} is always observed. What we are interested in are the marginal τ -th quantile regression coefficients $\boldsymbol{\gamma}_j = (\gamma_{j0}, \gamma_{j2}, \dots, \gamma_{jp})^T$ such that:

$$\Pr(Y_{ij} \leq \mathbf{x}_i^T \boldsymbol{\gamma}_j) = \tau, \text{ for } j = 1, \dots, J, \quad (1)$$

where \mathbf{x}_i is a $p \times 1$ covariates vector for subject i .

Let

$$\begin{aligned} p_k(Y) &= p(Y|S = k), \\ p_{\geq k}(Y) &= p(Y|S \geq k) \end{aligned}$$

be the density of response Y given follow-up time $S = k$ and $S \geq k$. And \Pr_k be the corresponding probability given $S = k$.

2.1 Settings

Instead of selection model, we adopt mixture model to deal with missingness, because it is easier to deploy missing mechanism. Meanwhile, we want to capture the slope change of marginal quantile lines not only by mixture of distributions, but also by the heterogeneity of variance.

Without loss of clarity, we suppress the subscript i for subject i . Specify the conditional distribution as:

$$\begin{aligned} p_k(y_1) &= N(\Delta_1 + \mathbf{x}_1^T \boldsymbol{\beta}_1^{(k)}, \exp(\mathbf{x}_1^T \boldsymbol{\alpha}_1^{(k)})), k = 1, \dots, J, \\ p_k(y_j | y_1, \dots, y_{j-1}) &= \begin{cases} N(\Delta_j + \mathbf{x}_{j*}^T \boldsymbol{\beta}_{j*}^{(k)}, \exp(\mathbf{x}_j^T \boldsymbol{\alpha}_j^{(k)})) & k < j; \\ N(\Delta_j + \mathbf{x}_{j*}^T \boldsymbol{\beta}_{j*}^{(\geq j)}, \exp(\mathbf{x}_j^T \boldsymbol{\alpha}_j^{(\geq j)})) & k \geq j; \end{cases}, \text{ for } 2 \leq j \leq J, \\ \Pr(S = k) &= \pi_k, \\ \sum_{k=1}^J \pi_k &= 1, \end{aligned} \quad (2)$$

where π_k do not depend on covariates, $\mathbf{x}_{j*} = (\mathbf{x}_j^T, y_1, \dots, y_{j-1})^T$ is a $(p + j - 1) \times 1$ modified covariates vector, and $\boldsymbol{\alpha}_j^{(k)}$ is a $p \times 1$ vector controlling heterogeneity of response component j under follow up time $S = k$, $\boldsymbol{\beta}_{j*}^{(k)} = (\boldsymbol{\beta}_j^T, \beta_{y_1 j}, \dots, \beta_{y_{j-1} j})^T$, which can be regarded as coefficients of interaction of pattern k and modified covariates analog to mean regression with length $(p + j - 1) \times 1$.

We model the heterogeneity parameters $\boldsymbol{\alpha}_j$ inside the exponential because there is no restriction on those heterogeneity parameters, therefore it is computationally more stable under both frequentist and Bayesian framework.

In equations group for model settings (2), Δ_j are deterministic by

$$\tau = \Pr(Y_j \leq \mathbf{x}_j^T \boldsymbol{\gamma}_j) = \sum_{k=1}^J \pi_k \Pr_k(Y_j \leq \mathbf{x}_j^T \boldsymbol{\gamma}_j), \quad (3)$$

for $j = 1$ and

$$\begin{aligned} \tau &= \Pr(Y_j \leq \mathbf{x}_j^T \boldsymbol{\gamma}_j) = \sum_{k=1}^J \pi_k \Pr_k(Y_j \leq \mathbf{x}_j^T \boldsymbol{\gamma}_j) \\ &= \sum_{k=1}^J \pi_k \int \cdots \int \Pr_k(Y_j \leq \mathbf{x}_j^T \boldsymbol{\gamma}_j | y_1, \dots, y_{j-1}) p_k(y_{j-1} | y_1, \dots, y_{j-2}) \\ &\quad \cdots p_k(y_2 | y_1) p_k(y_1) dy_{j-1} \cdots dy_1. \end{aligned} \quad (4)$$

for $j = 2, \dots, J$. Thus Δ_j are a function of other identifiable parameters. More computational details will be given in section 2.3.

The whole idea is to model the marginal quantile regression coefficients directly, then to involve them in the likelihood through restrictions by iterated expectation from mixture model, and finally to estimate them in either frequentist or Bayesian way. Meanwhile we model the mixture patterns and heterogeneity between subjects to explain the change of

marginal quantile regression coefficients, either case can make the slope change over different quantiles. Otherwise, the quantile lines would just be parallel to each other for different quantiles.

Furthermore, we have to put another set of restrictions because of identifiability:

$$\begin{aligned} \sum_{k=1}^J \beta_{l1}^{(k)} &= 0, l = 1, \dots, p, \\ \sum_{k=1}^{j-1} \beta_{lj}^{(k)} + \beta_{lj}^{(\geq j)} &= 0, l = 1, \dots, p, 2 \leq j \leq J; \end{aligned}$$

where $\beta_1 = (\beta_{11}, \dots, \beta_{p1})^T$ and $\beta_j = (\beta_{1j}, \dots, \beta_{pj})^T$. We show the need to keep those restrictions in appendix A.

2.2 Sensitivity Analysis

Theorem of [Molenberghs et al. \(1998\)](#) (MAR for discrete-time pattern mixture models under monotone dropout) is used to specify models under different missing mechanism assumptions, MAR holds if and only if, for each $j \geq 2$ and $k < j$:

$$p_k(y_j | y_1, \dots, y_{j-1}) = p_{\geq j}(y_j | y_1, \dots, y_{j-1}). \quad (5)$$

We illustrate how to deploy MAR and MNAR assumption from both frequentist way and Bayesian framework.

In model settings equation (2), when $2 \leq j \leq J, k < j$, Y_j is not observed, thus $\beta_j^{(k)}, k = 1, \dots, J$ and $\alpha_j^{(k)}, \beta_{y,j}^{(k)} = (\beta_{y_1,j}^{(k)}, \dots, \beta_{y_{j-1},j}^{(k)}), k < j$ can not be identified. Denote

$$\begin{aligned} \alpha_j^{(k)} &= \alpha_j^{(\geq j)} + h_j^{(k)}, \\ \beta_{y,j}^{(k)} &= \beta_{y,j}^{(\geq j)} + \eta_j^{(k)}, \end{aligned}$$

where $h_j^{(k)} = (h_{1j}^{(k)}, \dots, h_{pj}^{(k)})$ and $\eta_j^{(k)} = (\eta_{y_1,j}^{(k)}, \dots, \eta_{y_{j-1},j}^{(k)})$ for $k < j$, then $\xi_s = (\xi_\beta, \xi_\alpha)$ could be a group of sensitivity parameters, where $\xi_\beta = (\beta_j^{(k)}, \beta_j^{(\geq j)}, \eta_j^{(k)}), k < j, 2 \leq j \leq J$, and $\xi_\alpha = (h_j^{(k)}), k < j, 2 \leq j \leq J$.

- **Frequentist way:**

When $\xi_s = \xi_{s0} = \mathbf{0}$, it yields Molenberghs condition (5), therefore MAR condition satisfies. If ξ_s is fixed at $\xi_s \neq \xi_{s0}$, then Molenberghs condition fails, thus the missing mechanism is missing not at random.

- **Bayesian Framework:**

We put priors on (ξ_s, ξ_m) (ξ_m are identifiable parameters) as :

$$p(\xi_s, \xi_m) = p(\xi_s)p(\xi_m).$$

If we assume MAR with no uncertainty, the prior of ξ_s is $p(\xi_s = \mathbf{0}) = 1$. Sensitivity analysis can be executed through putting a set of priors on ξ_s to check the effect of priors on the posterior inference about quantile regression coefficients γ_{ij}^τ . For example, if MAR is assumed with uncertainty, priors can be assigned as $E(\xi_s) = \xi_{s0} = \mathbf{0}$ with $\text{Var}(\xi_s) \neq \mathbf{0}$. If we assume MNAR with no uncertainty, we can put priors satisfying $E(\xi_s) = \Delta_\xi$, where $\Delta_\xi \neq \mathbf{0}$ and $\text{Var}(\xi_s) = \mathbf{0}$. If MNAR is assumed with uncertainty, then priors could be $E(\xi_s) = \Delta_\xi$, where $\Delta_\xi \neq \mathbf{0}$ and $\text{Var}(\xi_s) \neq \mathbf{0}$.

2.3 Computation

In section 2.3.1, we give details on how to calculate Δ_{ij} in model (2) for $j = 1, \dots, J$. Then we illustrate how to get maximum likelihood estimator using adaptive gradient descent algorithm in section 2.3.2. Last, we present the Bayesian sampling procedure in section 2.3.3.

2.3.1 Δ Calculation

Seen from equation (3) and (4), Δ_{ij} depends on subject covariates \mathbf{x}_i , thus Δ_{ij} needs to be calculated for each subject generally. We illustrate how to calculate Δ_{ij} given all the other parameters $\xi = (\xi_m, \xi_s)$.

- Δ_{i1} : Expand equation (3):

$$\tau = \sum_{k=1}^J \pi_k \Phi \left(\frac{\mathbf{x}_{i1}^T \gamma_1 - \Delta_{i1} - \mathbf{x}_{i1}^T \beta_1^{(k)}}{\exp(\mathbf{x}_{i1}^T \alpha_1^{(k)})} \right),$$

where Φ is the standard normal CDF. Because above equation is continuous and monotone on Δ_{i1} , it can be solved by standard numerical root-find method without much difficulty, for example, the bisection method.

- $\Delta_{ij}, 2 \leq j \leq J$: First we introduce a lemma:

Lemma 2.1. *An integral of a normal CDF over a non-standard normal distribution can be simplified to a closed form in terms of another normal CDF:*

$$\int \Phi \left(\frac{x - b}{a} \right) d\Phi(x; \mu, \sigma) = \begin{cases} 1 - \Phi \left(\frac{b - \mu}{\sigma} / \sqrt{\frac{a^2}{\sigma^2} + 1} \right) & a > 0, \\ \Phi \left(\frac{b - \mu}{\sigma} / \sqrt{\frac{a^2}{\sigma^2} + 1} \right) & a < 0, \end{cases} \quad (6)$$

where $\Phi(x; \mu, \sigma)$ stands for a CDF of normal distribution with mean μ and standard deviation σ .

Proof of lemma can be seen in appendix B.

For computation of Δ_{ij} when $j \geq 2$, to solve equation (4), we propose two approaches: (for convenience, we use Δ_j to save typing for subject i .)

1. **Assume first order relationship:** We assume

$$p(Y_j|S, x, Y_{j-1}, \dots, Y_1) = p(Y_j|S, x, Y_{j-1}).$$

After obtaining Δ_{j-1} , for each component in equation (4):

$$\begin{aligned} p(Y_j \leq \mathbf{x}^T \boldsymbol{\gamma}_j | S = k) &= \int \dots \int p(Y_j \leq \mathbf{x}^T \boldsymbol{\gamma}_j | S = k, x, Y_{j-1}, \dots, Y_1) \\ &\quad dF(Y_{j-1} | S = k, Y_{j-2}, \dots, Y_1) \dots dF(Y_2 | S = s, Y_1), \\ &= \int p(Y_j \leq \mathbf{x}^T \boldsymbol{\gamma}_j | S = s, x, Y_{j-1}) dF(Y_{j-1} | S = k, Y_{j-2}). \end{aligned}$$

Thus, only one integral is needed. Furthermore, by lemma 2.1, we can simplify above integral by closed form in terms of normal CDF.

2. **Recursive Computation:**

From equation (6), we can find after single integral, the kernel part is still a normal CDF, but with other coefficients. So recursive simplification can be applied. After recursively applying lemma 2.1 $j - 1$ times, equation (4) becomes a closed form in terms of normal CDF analytically without calculating integral numerically, thus it can be solved again using standard numerical root-find method for Δ_{ij} .

Recursive computation is exact and put less restrictions than assuming first order relationship. In the later simulation and real data analysis, we adopt the recursive computation method. However, it is computationally more complicated than assuming first order relationship.

2.3.2 Maximum Likelihood Estimation

The contributed observed likelihood for $\mathbf{y}_i = (y_1, \dots, y_k)$ with follow-up time $S = k$ is

$$\begin{aligned} L_i(\boldsymbol{\xi} | \mathbf{y}_i, S_i = k) &= \pi_k p_k(y_k | y_1, \dots, y_{k-1}) p_k(y_{k-1} | y_1, \dots, y_{k-2}) \dots p_k(y_1) \\ &= \pi_k p_{\geq k}(y_k | y_1, \dots, y_{k-1}) p_{\geq k-1}(y_{k-1} | y_1, \dots, y_{k-2}) \dots p_k(y_1) \end{aligned} \quad (7)$$

We use adaptive gradient descent algorithm to get the maximum likelihood estimates [Riedmiller and Braun \(1993\)](#). Denote $J(\boldsymbol{\xi}) = -\log L = -\log \sum_{i=1}^n L_i$. Then to maximize likelihood is equivalent to minimize the target function $J(\boldsymbol{\xi})$. Under MAR assumption, we fix $\boldsymbol{\xi}_s = \mathbf{0}$, while under MNAR assumption, $\boldsymbol{\xi}_s$ can be chosen to assume there is an intercept shift between the conditional distributions of $Y_j | Y_1, \dots, Y_{j-1}, S$, or there is heterogeneity between those distributions.

During each step in adaptive gradient descent algorithm, Δ_{ij} has to be calculated for each subject and component, as well as partial derivatives for each parameter. Because it is computationally expensive, we use compiled fortran language within R to make it fast.

Details about the maximization algorithm are presented in the appendix [ref].

2.3.3 Bayesian Framework

Under Bayesian framework, we are going to put priors on the parameters ξ and make exact inference through posterior samples.

We can use the block Gibbs sampling method to draw sample from posterior distribution. Denote all the parameters (including sensitivity parameters) to sample as : (TODO: specifying priors)

$$\xi = \left(\gamma_1, \gamma_2, \dots, \gamma_J, \beta_{j*}^{(k)}, \alpha_j^{(k)} \text{ for } k = 1, \dots, J, j = 1, \dots, J \right).$$

Comma separated parameters are marked to sample as a block in block Gibbs sampling. All the procedures require Metropolis-Hasting algorithm to update samples because the likelihood is intractable.

As mentioned in section 2.2, MAR or MNAR assumptions are adopted using specific priors. For example, if MAR is assumed with no uncertainty, then $\xi_s = \mathbf{0}$ with probability 1. Details for updating parameters are:

- γ_1 : Use Metropolis-Hasting algorithm.
 1. Draw (γ_1^c) candidates from candidate distribution;
 2. Based on the new candidate parameter ξ^c , calculate candidate Δ_{i1}^c for each subject i as we described in section 2.3.1. If $S > 1$ for corresponding subject i , update candidate $\Delta_{ij}^c, j \geq 2$ as well because $\Delta_{ij}, j \geq 2$ depend on Δ_{i1} . (For $S = 1$, we only need to update Δ_{i1}^c);
 3. Plug in Δ_{i1}^c or $(\Delta_{i1}^c, \Delta_{ij}^c, j \geq 2)$ in likelihood (7) to get candidate likelihood;
 4. Obtain Metropolis-Hasting ratio, move the chain or keep the previous one.
- $(\gamma_j, j \geq 2)$: similar algorithm but only update Δ_{ij} for subjects with $S \geq j$.
- $(\beta_{j*}^{(k)})$: similar to γ_j .
- For the rest of the parameters, algorithms for updating the samples are all similar to γ_j .

Computation is still expensive due to need to calculate Δ and likelihood in each iteration. Compiled language is recommended to implement the algorithm.

3 Simulation Study

In this section, we test performance of our proposed heterogeneous model in section 2.1 combined with adaptive gradient descent algorithm, compared with *rq* function in *quantreg* R package.

3.1 MAR

The simulation study included 1000 data sets. Each data set consists 200 bivariate observations $Y_i = (Y_{i1}, Y_{i2})$ for $i = 1, \dots, 200$. Y_{i1} was always observed, while some of Y_{i2} were missing with probability 0.5. Covariates x were sampled from uniform (0,2). We sampled Y_i from:

$$\begin{aligned} Y_{i1}|R = 1 &\sim N(2 + x, 1 + 0.5x) \\ Y_{i2}|R = 1, y_{i1} &\sim N(1 - x - 1/2y_{i1}, 1) \\ Y_{i1}|R = 0 &\sim N(-2 - x, 1 + 0.5x) \\ Y_{i2}|R = 0, y_{i1} &\sim N(1 - x - 1/2y_{i1}, 1) \\ p(R = 1) &= 0.5. \end{aligned}$$

$Y_{i2}|R = 0, y_{i1}$ is not observed. And in this setting, we assume distribution of $Y_{i2}|R = 0, y_{i1}$ is equal to $Y_{i2}|R = 1, Y_{i1}$, thus MAR condition satisfies.

By integrating $Y_{i1}|R$ out of $Y_{i2}|R, y_{i1}$, we have

$$\begin{aligned} Y_{i2}|R = 1 &\sim N(-3x/2, 5/4 + x/8) \\ Y_{i2}|R = 0 &\sim N(2 - x/2, 5/4 + x/8) \end{aligned}$$

To apply MAR assumption, we fix sensitivity parameter $\xi_s = (0, 0, 0, 0, 0)$ as discussed in section 2.2 for our proposed model. For *rq* function from *quantreg* R package, because only $Y_{i2}|R = 1$ is observed, quantile regression for Y_{i2} can only be fit from the information of $Y_{i2}|R = 1$ vs x .

For each dataset in both scenario, we fit quantile regression for quantiles $\tau = 0.1, 0.3, 0.5, 0.7, 0.9$.

Algorithms were evaluated by mean squared error (MSE), which is defined by :

$$\text{MSE}(\gamma_{ij}) = \frac{1}{1000} \sum_{k=1}^{1000} \left(\hat{\gamma}_{ij}^{(k)} - \gamma_{ij} \right)^2,$$

where γ_{ij} is the true value for quantile regression coefficient, $\hat{\gamma}_{ij}^{(k)}$ is the estimates in k -th simulated dataset ($(\gamma_{01}, \gamma_{11})$ for Y_{i1} , $(\gamma_{02}, \gamma_{12})$ for Y_{i2}).

Mean squared errors are shown in table 1. Results show that our proposed method has smaller MSE than *rq* function in all cases. Furthermore, when Y_{i2} are missing at random, our method shows tremendous advantage over *rq* method, because *quantreg* does not consider the missing mechanism. The difference in MSE becomes larger for upper quantile is because $Y_2|R = 0$ tends to be larger than $Y_2|R = 1$, therefore, the *rq* method only using the observed Y_2 yields larger bias for the marginal quantile estimates.

Table 1: Simulation result: MSE for coefficients estimates of quantiles 0.1, 0.3, 0.5, 0.7, 0.9 under MAR assumptions. $(\gamma_{01}, \gamma_{11})$ are quantile regression coefficients for Y_{i1} , and $(\gamma_{02}, \gamma_{12})$ are ones for Y_{i2} . MM stands for our proposed method, and RQ stands for the 'rq' function in R package 'quantreg'.

	MAR									
	0.1		0.3		0.5		0.7		0.9	
	MM	RQ	MM	RQ	MM	RQ	MM	RQ	MM	RQ
γ_{01}	0.09	0.15	0.12	0.19	0.11	1.08	0.16	0.19	0.10	0.15
γ_{11}	0.09	0.15	0.07	0.19	0.14	1.19	0.08	0.20	0.10	0.15
γ_{02}	0.08	0.27	0.07	0.59	0.06	1.08	0.12	1.75	0.24	2.92
γ_{12}	0.06	0.17	0.05	0.13	0.06	0.33	0.07	0.75	0.09	0.96

3.2 MCAR and MNAR

Another group of dataset were simulated to test our algorithm performance under MCAR and MNAR situation, comparing to *rq* function in *quantreg* R package. We sampled Y_i from:

$$\begin{aligned} \begin{pmatrix} Y_{i1} \\ Y_{i2} \end{pmatrix} \Big| R = 1 &\sim N \left(\begin{pmatrix} 1+x \\ 1-x \end{pmatrix}, \begin{pmatrix} 1 & 0 \\ 0 & 1 \end{pmatrix} \right), \\ Y_{i1} \Big| R = 0 &\sim N(-1-x, 1), \\ p(R = 1) &= 0.5. \end{aligned}$$

We conducted simulation study under two different situations: MCAR and MNAR. Under MCAR scenario, we still fixed sensitivity parameter ξ_s in section 2.2 at 0 (assuming it is MAR), while under MNAR scenario, we fixed ξ_s at the true value $(1, 0, 0, 0)$, assuming there was an intercept shift between distribution of $Y_{i2} \mid Y_{i1}, R = 1$ and $Y_{i2} \mid Y_{i1}, R = 0$.

For each dataset in both scenario, we fit quantile regression for quantiles $\tau = 0.1, 0.3, 0.5, 0.7, 0.9$.

Algorithms were evaluated by mean squared error (MSE) as we described above.

Simulation results show estimates from our algorithm are closer to the true value for all quantiles from 0.1 to 0.9. Table 2 and 3 list the MSE for coefficients estimates of quantile 0.1, 0.3, 0.5, 0.7, 0.9 under MAR and MNAR assumptions. Even for extreme quantiles ($\tau = 0.1$ and $\tau = 0.9$), our algorithm behave as good as for non-extreme quantile ($\tau = 0.3, 0.5, 0.7$) in terms of MSE. Furthermore, 'rq' function did not consider the missing mechanism, so under MNAR assumption, 'quantreg' method led to tremendous MSE and our proposed method were much closer to the true value.

4 Real Data Analysis

Here is the analysis for *tours* data. *Weight2* stands for weight at 6th month after the baseline measure, and *weight3* stands for the one at 18th month after the baseline. There were three

Table 2: Simulation result: MSE for coefficients estimates of quantiles 0.1, 0.3, 0.5, 0.7, 0.9 under MCAR scenario. $(\gamma_{01}, \gamma_{11})$ are quantile regression coefficients for Y_{i1} , and $(\gamma_{02}, \gamma_{12})$ are ones for Y_{i2} . MM stands for our proposed method, and RQ stands for the 'rq' function in R package 'quantreg'.

	MAR									
	0.1		0.3		0.5		0.7		0.9	
	MM	RQ	MM	RQ	MM	RQ	MM	RQ	MM	RQ
γ_{01}	0.05	0.09	0.04	0.10	0.03	0.24	0.04	0.10	0.05	0.10
γ_{11}	0.03	0.07	0.02	0.08	0.58	0.74	0.03	0.08	0.03	0.07
γ_{02}	0.04	0.12	0.05	0.07	0.04	0.06	0.05	0.07	0.05	0.11
γ_{12}	0.03	0.09	0.03	0.05	0.03	0.05	0.03	0.05	0.03	0.09

Table 3: Simulation result: MSE for coefficients estimates of quantiles 0.1, 0.3, 0.5, 0.7, 0.9 under MNAR scenario. $(\gamma_{01}, \gamma_{11})$ are quantile regression coefficients for Y_{i1} , and $(\gamma_{02}, \gamma_{12})$ are ones for Y_{i2} . MM stands for our proposed method, and RQ stands for the 'rq' function in R package 'quantreg'.

	MNAR									
	0.1		0.3		0.5		0.7		0.9	
	MM	RQ	MM	RQ	MM	RQ	MM	RQ	MM	RQ
γ_{01}	0.04	0.09	0.04	0.10	0.03	0.24	0.04	0.10	0.04	0.10
γ_{11}	0.03	0.07	0.02	0.08	0.64	0.74	0.03	0.08	0.03	0.07
γ_{02}	0.04	0.30	0.05	0.52	0.07	1.06	0.05	1.79	0.05	2.59
γ_{12}	0.03	0.09	0.03	0.05	0.03	0.05	0.03	0.05	0.03	0.09

treatments and two main races in this study (Treatment M, Treatment O and Treatment T; Race 1(black) and Race 3(white)). Weights at 6th month were always observed and some weights at 18th month were missing (211 observed out of 224 , 94%). All weights are scaled by 1/100.

Here is the boxplot of weights vs treatments and races.

Quantile regression models were fitted for responses *weight2* and *weight3* together ($Y_i = (Y_{i1}, Y_{i2})$) for quantiles (10%, 30%, 50%, 70%, 90%). Covariates are treatments and races, and we assume their effects are additive. Treatment M and Race 1 are baseline references. We fit 1,000 bootstraps to obtain the 95% confidence intervals.

Estimates are presented in table 4. For *weight2*, quantile estimates show there is no significant difference for three treatment group through all the quantiles, because all 95 % confidence intervals include 0. However, when comparing weights from two races, weights at 6th month from race 3 are generally lower than ones from race 1. Estimates of race 3 effect to weights quantiles 10% up to 70% are all significantly away from zero (negative). However,

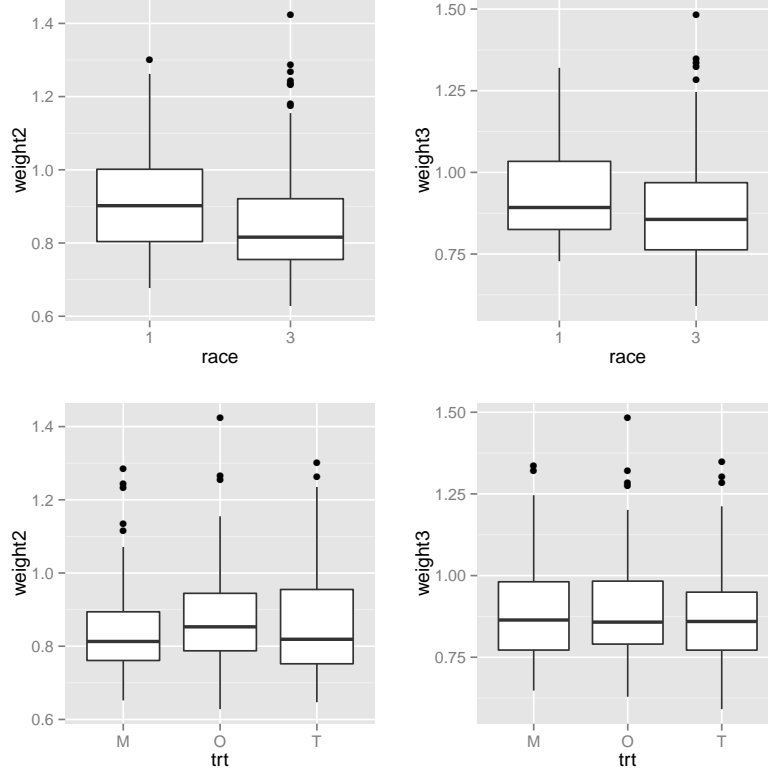


Figure 1: Boxplot of weights at 6th month (weight2) and 18th months (weight3) vs treatments (M, O, T) and race (black:1, white:3)

for top weights of two races (90% quantile), the difference is not significant.

For weights at 18th month (weight3), we have similar conclusions. Confidence intervals of treatment effects on weight3 for all quantiles (10% up to 90%) include zero. But after 18 months, weights of patients from race 3 are significantly lower than ones from race 1 only for lower quantiles (10% to 30%). They are not significantly different for quantiles (50% to 90%).

5 Discussion

In this paper, we developed a marginal quantile regression model for data with monotone dropout missingness. We use pattern mixture model to explain the missing mechanism instead of selection model. Meanwhile, marginal quantile regression coefficients are of interest instead of ones conditional on random effects as in [Yuan and Yin \(2010\)](#). In addition, our approach capture nontrivial change of quantile lines over different quantiles by modeling the mixture distribution and heterogeneity of variance.

It is also easy for our method to adopt missing mechanism assumptions for nonignorable missing data. We illustrated how to put informative priors in Bayesian way and how to assign values for sensitivity parameters in frequentist way. to apply different missing mechanism assumptions. We also recommend recursive integration to simplify the computation and it is also easy to implement even in high dimension. Simulation study demonstrates that our approach has smaller MSE than the traditional frequentist method and it behaves even

Table 4: Marginal quantile regression coefficients for Weight2

	Intercept	Trt.O	Trt.T	Race.3
Weight2				
$\tau = 0.1$	0.80 (0.70, 0.86)	0.01 (-0.04, 0.07)	-0.01 (-0.06, 0.06)	-0.13 (-0.19, -0.04)
$\tau = 0.3$	0.83 (0.79, 0.92)	0.04 (-0.02, 0.07)	0.02 (-0.04, 0.05)	-0.07 (-0.16, -0.03)
$\tau = 0.5$	0.85 (0.82, 0.98)	0.05 (-0.03, 0.09)	0.04 (-0.06, 0.07)	-0.03 (-0.14, 0.00)
$\tau = 0.7$	0.95 (0.89, 1.03)	0.03 (-0.02, 0.10)	0.02 (-0.04, 0.08)	-0.04 (-0.12, 0.00)
$\tau = 0.9$	0.98 (0.94, 1.11)	0.07 (-0.02, 0.14)	0.06 (-0.02, 0.14)	-0.01 (-0.10, 0.05)
Weight3				
$\tau = 0.1$	0.78 (0.38, 0.84)	-0.01 (-0.07, 0.06)	-0.04 (-0.10, 0.04)	-0.13 (-0.18, -0.02)
$\tau = 0.3$	0.82 (0.78, 0.93)	0.01 (-0.04, 0.06)	-0.01 (-0.07, 0.05)	-0.06 (-0.16, -0.01)
$\tau = 0.5$	0.88 (0.84, 1.00)	0.02 (-0.06, 0.06)	0.02 (-0.08, 0.06)	-0.03 (-0.13, 0.02)
$\tau = 0.7$	0.99 (0.92, 1.07)	0.00 (-0.05, 0.08)	-0.00 (-0.07, 0.06)	-0.04 (-0.12, 0.01)
$\tau = 0.9$	1.02 (0.98, 1.16)	0.05 (-0.06, 0.11)	0.04 (-0.05, 0.12)	0.00 (-0.10, 0.06)

better in MAR and MNAR case.

Our model assumes heterogeneous normal distribution for each component, which might be too restrictive. It is possible to replace them with non or semi parametric model, for example Dirichlet process or Pólya tree. Correlation within subject is not explicit in our method. It is hard to look into marginal correlations in a pattern mixture model. It would be also interesting to assume that component probabilities depend on covariates. Our future work also include to propose a goodness of fit test to check the model fit.

6 Acknowledgments

References

- Moshe Buchinsky. Changes in the u.s. wage structure 1963-1987: Application of quantile regression. *Econometrica*, 62(2):pp. 405–458, 1994. ISSN 00129682. URL <http://www.jstor.org/stable/2951618>.
- Moshe Buchinsky. Recent advances in quantile regression models: A practical guideline for empirical research. *The Journal of Human Resources*, 33(1):pp. 88–126, 1998. ISSN 0022166X. URL <http://www.jstor.org/stable/146316>.
- T. Hanson and W.O. Johnson. Modeling regression error with a mixture of polya trees. *Journal of the American Statistical Association*, 97(460):1020–1033, 2002.

- X. He, P. Ng, and S. Portnoy. Bivariate quantile smoothing splines. *Journal of the Royal Statistical Society: Series B (Statistical Methodology)*, 60(3):537–550, 1998. ISSN 1467-9868. doi: 10.1111/1467-9868.00138. URL <http://dx.doi.org/10.1111/1467-9868.00138>.
- Patrick J Heagerty. Marginally specified logistic-normal models for longitudinal binary data. *Biometrics*, 55(3):688–698, 1999.
- N.L. Hjort and S. Petrone. Nonparametric quantile inference using dirichlet processes. *Advances in statistical modeling and inference*, pages 463–492, 2007.
- N.L. Hjort and S.G. Walker. Quantile pyramids for bayesian nonparametrics. *The Annals of Statistics*, 37(1):105–131, 2009.
- R. Koenker. *Quantile regression*, volume 38. Cambridge Univ Pr, 2005.
- Roger Koenker and Jr. Bassett, Gilbert. Regression quantiles. *Econometrica*, 46(1):pp. 33–50, 1978. ISSN 00129682. URL <http://www.jstor.org/stable/1913643>.
- Roger Koenker and Jose A. F. Machado. Goodness of fit and related inference processes for quantile regression. *Journal of the American Statistical Association*, 94(448):pp. 1296–1310, 1999. ISSN 01621459. URL <http://www.jstor.org/stable/2669943>.
- A. Kottas and A.E. Gelfand. Bayesian semiparametric median regression modeling. *Journal of the American Statistical Association*, 96(456):1458–1468, 2001.
- A. Kottas and M. Krnjajić. Bayesian semiparametric modelling in quantile regression. *Scandinavian Journal of Statistics*, 36(2):297–319, 2009.
- Geert Molenberghs, Bart Michiels, MG Kenward, and Peter J Diggle. Monotone missing data and pattern-mixture models. *Statistica Neerlandica*, 52(2):153–161, 1998.
- B.J. Reich, H.D. Bondell, and H.J. Wang. Flexible bayesian quantile regression for independent and clustered data. *Biostatistics*, 11(2):337–352, 2010.
- Martin Riedmiller and Heinrich Braun. A direct adaptive method for faster backpropagation learning: The rprop algorithm. In *Neural Networks, 1993., IEEE International Conference on*, pages 586–591. IEEE, 1993.
- Jason Roy and Michael J Daniels. A general class of pattern mixture models for nonignorable dropout with many possible dropout times. *Biometrics*, 64(2):538–545, 2008.
- Stephen Walker and Bani K. Mallick. A bayesian semiparametric accelerated failure time model. *Biometrics*, 55(2):477–483, 1999. ISSN 1541-0420. doi: 10.1111/j.0006-341X.1999.00477.x. URL <http://dx.doi.org/10.1111/j.0006-341X.1999.00477.x>.
- Ying Wei, Anneli Pere, Roger Koenker, and Xuming He. Quantile regression methods for reference growth charts. *Statistics in Medicine*, 25(8):1369–1382, 2006. ISSN 1097-0258. doi: 10.1002/sim.2271. URL <http://dx.doi.org/10.1002/sim.2271>.
- Keming Yu and Rana A. Moyeed. Bayesian quantile regression. *Statistics & Probability Letters*, 54(4):437 – 447, 2001. ISSN 0167-7152. doi: 10.1016/S0167-7152(01)00124-9. URL <http://www.sciencedirect.com/science/article/pii/S0167715201001249>.

Keming Yu, Zudi Lu, and Julian Stander. Quantile regression: applications and current research areas. *Journal of the Royal Statistical Society: Series D (The Statistician)*, 52(3):331–350, 2003. ISSN 1467-9884. doi: 10.1111/1467-9884.00363. URL <http://dx.doi.org/10.1111/1467-9884.00363>.

Ying Yuan and Guosheng Yin. Bayesian quantile regression for longitudinal studies with nonignorable missing data. *Biometrics*, 66(1):105–114, 2010.

A Identifiability

First consider univariate case with two patterns. Suppose y is univariate and there are two patterns $R = 1$ and $R = 0$.

Before going forward to quantile regression, first we consider identifiability problem in mean regression.

Consider a pattern mixture model:

$$\begin{aligned} y|R = 1 &\sim N(\Delta + R_1, \sigma_1) \\ y|R = 0 &\sim N(\Delta + R_0, \sigma_0) \\ \Pr(R = 1) &= \pi \\ E(y) &= \theta \end{aligned}$$

Thus by iterated expectation, we have

$$\begin{aligned} \theta &= \Delta + R_1\pi + R_0(1 - \pi) \\ \Delta &= \theta - \pi R_1 - (1 - \pi)R_0 \end{aligned}$$

We can see Δ is deterministic by θ, R_1, R_0 . If plugged in likelihood, we have

$$\begin{aligned} y|R = 1 &\sim N(\theta + (1 - \pi)R_1 - (1 - \pi)R_0, \sigma_1) \\ y|R = 0 &\sim N(\theta - \pi R_1 + \pi R_0, \sigma_0) \end{aligned}$$

Denote $\xi_1 = (\theta, R_1, R_0)$, and if $\xi_2 = (\theta, R_1 + 1, R_0 + 1)$, both parameters lead to the same distribution of $p(y, R) = p(y|R)p(R)$. Therefore, ξ is not identifiable. If we put constraints on R_1 and R_0 , for example $R_0 = -R_1$, then

$$\begin{aligned} y|R = 1 &\sim N(\theta + 2(1 - \pi)R_1, \sigma_1) \\ y|R = 0 &\sim N(\theta - 2\pi R_1, \sigma_0) \end{aligned}$$

thus it is identifiable. If $\xi_2 \neq \xi_1$, then $p_2(y, R) \neq p_1(y, R)$.

Secondly, we consider quantile regression for pattern mixture model:

$$\begin{aligned} y|R = 1 &\sim N(\Delta + R_1, \sigma_1) \\ y|R = 0 &\sim N(\Delta + R_0, \sigma_0) \\ \Pr(R = 1) &= \pi \\ p(y \leq \theta) &= \tau \end{aligned}$$

where θ is the quantile estimate of interest and we does not include covariates so far. We will show $\xi = (\theta, R_1, R_0)$ is not identifiable.

Again by iterated expectation, we have

$$\tau = \pi \Phi \left(\frac{\theta - \Delta - R_1}{\sigma_1} \right) + (1 - \pi) \Phi \left(\frac{\theta - \Delta - R_0}{\sigma_0} \right)$$

thus Δ is again deterministic by other parameters:

$$\Delta = h(\theta, R_1, R_0, \sigma_1, \sigma_0, \pi, \tau)$$

To show $\xi = (\theta, R_1, R_0, \sigma_1, \sigma_0)$ is not identifiable, we need to find $\xi' \neq \xi$, such that $p(y|R) = p'(y|R)$. If the last equation holds, then we must have $\sigma'_1 = \sigma_1, \sigma'_0 = \sigma_0$, thus we still need to find θ', R'_1, R'_0 such that

$$\begin{aligned} h(\xi) + R_1 &= h(\xi') + R'_1 \\ h(\xi) + R_0 &= h(\xi') + R'_0 \end{aligned}$$

By substracting previous equations, we have $R'_1 - R'_0 = R_1 - R_0$, thus denote $R'_1 = R_1 + \delta$ and $R'_0 = R_0 + \delta$, and let $\theta' = \theta$ such that

$$\Delta' = h(\theta', R_1, R_0, \sigma_1, \sigma_0, \delta) = h(\xi) - \delta = \Delta - \delta$$

then the new parameter ξ' yields the same distribution with one from ξ . Therefore ξ is not identifiable.

Instead, if we put constraint, for example $R_1 = -R_0$, then by calculation, $p(y|R; \xi) = p(y|R; \xi')$ yields $\xi = \xi'$.

Now consider the case with covariates. Suppose the model is

$$\begin{aligned} y|R=1, x &\sim N(\Delta + R_1 + \beta_{x1}x, \sigma_1) \\ y|R=0, x &\sim N(\Delta - R_1 + \beta_{x0}x, \sigma_0) \\ \Pr(R=1) &= \pi \\ p(y \leq \gamma_0 + \gamma_1 x) &= \tau \end{aligned}$$

Still Δ can be determined by

$$\Delta = h(x, \gamma_0, \gamma_1, R_1, \beta_{x1}, \beta_{x0}, \sigma_1, \sigma_0, \pi, \tau)$$

We want to show parameter $\xi = (\gamma_0, \gamma_1, R_1, \beta_{x1}, \beta_{x0}, \sigma_1, \sigma_0, \pi)$ is not identifiable by finding $\xi' \neq \xi$, but $p(y|R; \xi) = p(y|R; \xi')$. Still if the last equation holds, first we have $\sigma'_1 = \sigma_1, \sigma'_0 = \sigma_0$, then to equate the two means, we have

$$\begin{aligned} \Delta + R_1 + \beta_{x1}x &= \Delta' + R'_1 + \beta'_{x1}x \\ \Delta - R_1 + \beta_{x0}x &= \Delta' - R'_1 + \beta'_{x0}x \end{aligned}$$

By substracting the two equations, we have

$$2R_1 + (\beta_{x1} - \beta_{x0})x = 2R'_1 + (\beta'_{x1} - \beta'_{x0})x$$

which holds for all x . Thus $R_1 = R'_1$ and $(\beta_{x1} - \beta_{x0}) = (\beta'_{x1} - \beta'_{x0})$. Then let

$$\begin{aligned}\beta'_{x1} &= \beta_{x1} + \delta \\ \beta'_{x0} &= \beta_{x0} + \delta\end{aligned}$$

and all the other parameters in ξ' keep the same, we can still have the same distribution of $y|R; \xi$ but with different ξ . Therefore, ξ is not identifiable, especially for β_{x1} and β_{x0} . One solution is to restrict $\beta_{x1} = -\beta_{x0}$ to make all the parameters identifiable.

Now consider the bivariate (y_1, y_2) case, and we focus on the identifiability issue especially on $y_2|y_1$. Suppose the model is

$$\begin{aligned}y_2|y_1, x, R = 1 &\sim N(\Delta + R_1 + x\beta_{x1} + \beta_{11}y_1, \sigma_1) \\ y_2|y_1, x, R = 0 &\sim N(\Delta - R_1 - x\beta_{x1} + \beta_{10}y_1, \sigma_0)\end{aligned}$$

Here R stands for two different patterns, and missingness is not considered.

we are wondering if β_{11} and β_{10} are identifiable, say if there exists β'_{11} and β'_{10} , such that

$$\begin{aligned}\Delta + R_1 + x\beta_x + \beta_{11}y_1 &= \Delta' + R'_1 + x\beta'_x + \beta'_{11}y_1 \\ \Delta - R_1 - x\beta_x + \beta_{10}y_1 &= \Delta' - R'_1 - x\beta'_x + \beta'_{10}y_1\end{aligned}$$

still by subtracting two equations, we have $R_1 = R'_1$ and $\beta_x = \beta'_x$. Considering Δ is determined by integrating out y_1 , such that matching the two sides of the above equation for coefficient of y_1 , we must have $\beta_{11} = \beta'_{11}$ and $\beta_{10} = \beta'_{10}$, therefore, ξ is identifiable.

For identifiability issue in heterogeneous model described in section [ref], it is easy to show there is no trouble with heterogeneity parameters α , analog to the linear model case. For the other parameters, it can be found similar to the above discussion.

B Proof of Lemma 2.1

- Denote

$$I(a, b) = \int \Phi\left(\frac{x-b}{a}\right) \phi(x) dx$$

where Φ is the standard normal cdf and ϕ is the standard normal pdf and $a > 0$.

$$\begin{aligned}\frac{\partial I(a, b)}{\partial b} &= -\frac{1}{a} \int \phi\left(\frac{x-b}{a}\right) \phi(x) dx \\ &= -\frac{1}{\sqrt{2\pi}\sqrt{a^2+1}} \exp\left(-\frac{b^2}{2(a^2+1)}\right) \\ &= -\frac{1}{\sqrt{a^2+1}} \phi\left(\frac{b}{\sqrt{a^2+1}}\right)\end{aligned}$$

Since $I(a, \infty) = 0$,

$$\begin{aligned}
I(a, b) &= -\frac{1}{\sqrt{a^2 + 1}} \int_b^\infty \phi\left(\frac{s}{\sqrt{a^2 + 1}}\right) ds \\
&= \int_{b/\sqrt{a^2 + 1}}^\infty \phi(t) dt \\
&= 1 - \Phi(b/\sqrt{a^2 + 1})
\end{aligned} \tag{8}$$

For $a < 0$,

$$\begin{aligned}
\frac{\partial I(a, b)}{\partial b} &= -\frac{1}{a} \int \phi\left(\frac{x - b}{a}\right) \phi(x) dx \\
&= -\frac{\text{sgn}(a)}{\sqrt{2\pi}\sqrt{a^2 + 1}} \exp\left(-\frac{b^2}{2(a^2 + 1)}\right) \\
&= -\frac{\text{sgn}(a)}{\sqrt{a^2 + 1}} \phi\left(\frac{b}{\sqrt{a^2 + 1}}\right)
\end{aligned}$$

Since $I(a, -\infty) = 0$:

$$\begin{aligned}
I(a, b) &= \int_{-\infty}^{b/\sqrt{a^2 + 1}} \phi(t) dt \\
&= \Phi(b/\sqrt{a^2 + 1})
\end{aligned} \tag{9}$$

- If integrating over non-standard normal distribution:

$$\begin{aligned}
\int \Phi(x) d\Phi(x; \mu, \sigma) &= \int \Phi(x) \frac{1}{\sigma} \phi\left(\frac{x - \mu}{\sigma}\right) dx \\
&= \int \Phi(\sigma t + \mu) \phi(t) dt \\
\text{by equation (8):} \quad &= 1 - \Phi(-\mu/\sigma / \sqrt{1/\sigma^2 + 1})
\end{aligned}$$

- If integrating a non standard cdf over a non-standard normal distribution:

$$\begin{aligned}
\int \Phi\left(\frac{x - b}{a}\right) d\Phi(x; \mu, \sigma) &= \int \Phi\left(\frac{x - b}{a}\right) \frac{1}{\sigma} \phi\left(\frac{x - \mu}{\sigma}\right) dx \\
&= \int \Phi\left(\frac{\sigma y + \mu - b}{a}\right) \phi(y) dy \\
&= 1 - \Phi\left(\frac{b - \mu}{\sigma} / \sqrt{\frac{a^2}{\sigma^2} + 1}\right)
\end{aligned} \tag{10}$$

If $a < 0$,

$$\int \Phi\left(\frac{x - b}{a}\right) d\Phi(x; \mu, \sigma) = \Phi\left(\frac{b - \mu}{\sigma} / \sqrt{\frac{a^2}{\sigma^2} + 1}\right) \tag{11}$$