

Estimation of Quantile Regressions with Multiple Fixed Effects

Fernando Rios-Avila
Levy Economics Institute
Annandale-on-Hudson, NY
friosavi@levy.org

Leonardo Siles
Universidad de Chile
Santiago, Chile
lsiles@fen.uchile.cl

Gustavo Canavire-Bacarreza
The World Bank
Washington, DC
gcanavire@worldbank.org

Abstract.

This paper introduces two new Stata commands, `qregfe` and `qregplot`, designed for estimating and visualizing quantile regression models with fixed effects. `qregfe` implements several methodologies including Correlated Random Effects, Canay (2011), a modified Canay estimator, and Method of Moments Quantile Regression (Machado and Santos Silva, 2019). This allows researchers to control for unobserved heterogeneity in panel data settings while examining heterogeneous effects across the conditional distribution. `qregplot` facilitates visualization of quantile regression results, enabling graphical examination of covariate effects across different quantiles.

Keywords: st0001, Quantile Regression, Fixed Effects, Panel Data

1 Introduction

Quantile regression, introduced by Koenker and Bassett (1978), has become an important tool in economic analysis, allowing researchers to examine how the relationship between the dependent and independent variables varies across different points of the conditional distribution of the outcome. While ordinary least squares focuses on analyzing the conditional mean, quantile regression provides a more comprehensive view of how covariates impact the entire conditional distribution of the dependent variable. This can reveal heterogeneous effects that may be otherwise overlooked when analyzing the conditional mean.

A relatively recent development in the literature has focused on extending quantile regression analysis in a panel data setting to account for unobserved, but time-fixed heterogeneity. This is particularly important in empirical research, where unobserved heterogeneity can bias estimates of the effects of interest. However, as is common in the estimation of non-linear models with fixed effects, introducing fixed effects in quantile regression models poses several challenges. On the one hand, the simple inclusion of fixed effects can lead to an incidental parameter problem, which can bias estimates of the quantile coefficients (Neyman and Scott 1948; Lancaster 2000). On the other hand, the computational complexity of estimating quantile regression models with fixed effects can be prohibitive, particularly for large datasets with multiple high-dimensional fixed effects. While many strategies have been proposed for estimating this type of model (see Galvao and Kengo (2017) for a review), none has become standard due to restrictive

assumptions regarding the inclusion of fixed effects and the computational complexity.

In spite of the growing interest in estimating quantile regression models with fixed effects in applied research, particularly in the fields of labor economics, health economics, and public policy, among others, there are few commands that allow the estimation of such models. In Stata, there are three main built-in commands available for estimating quantile regression: `qreg`, `ivqregress`, and `bayes: qreg`, and none of them allow for the inclusion of fixed effects, other than using the dummy variable approach. From the community-contributed commands, there is `xtqreg`, which implements a quantile regression model with fixed effects based on the method of moments proposed by Machado and Santos Silva (2019), and more recently `xtmdqr` which implements a minimum distance estimation of quantile regression models with fixed effects described in Melly and Pons (2023). In both cases, these commands are constrained to a single set of fixed effects.¹

To address this, in this paper we introduce a new Stata command for estimating quantile regressions with multiple fixed effects: `qregfe`. This command implements: the extended estimator of quantile regression via moments (`mmqreg`) proposed by Machado and Santos Silva (2019) and Rios-Avila et al. (2024); an implementation of a correlated random effects estimator based on Abrevaya and Dahl (2008), Wooldridge (2019) and Wooldridge (2010, Ch12.10.3); the estimator proposed by Canay (2011), and a proposed modification of this approach. In addition, we also present an auxiliary command `qregplot` for the visualization of the quantile regression models.

This command offers the advantage of allowing for the estimation of conditional quantile regressions while controlling for multiple fixed effects. First, they leverage existing Stata commands, as well as other community-contributed commands, to allow users to estimate quantile regression models and their standard errors under different assumptions. Second, they reduce the impact of the incidental parameters problem depending on the assumptions underlying the data generating process. In terms of standard errors, `mmqreg` allows for the estimation of analytical standard errors (see Machado and Santos Silva (2019) and Rios-Avila et al. (2024)), whereas `qregfe` emphasizes the use of bootstrap standard errors. Finally, the command is designed to be user-friendly, allowing for the estimation of quantile regression models with fixed effects in a single line of code.

The remainder of the paper is organized as follows. Section 2 reviews the methodological framework for quantile regression, along with the methods and formulas behind the estimators implemented by `qregfe` command. Section 3 introduces the commands, along with a brief description of their syntax and options. Section 4 introduces an auxiliary command for the visualization of quantile regression models. Section 6 provides an empirical application demonstrating their use. Section 7 concludes.

1. There are other community-contributed commands like `xtrifreg`, `rifhdfe`, `qregpd`, `rqr` among others that allow for the estimation of quantile regression models, but do not estimate conditional quantile regressions, but instead focus on unconditional quantile regressions, or quantile treatment effects.

2 The Basics

Quantile regressions allow researchers to identify the heterogeneous effect covariates could have over the entire conditional distribution of the dependent variable. Let y_i be the dependent variable, x_i the vector of covariates excluding a constant, and $0 < \tau < 1$ is a parameter such that $q_\tau(y_i|X)$ identifies the τ th quantile of the conditional distribution of $y_i|X$. Under the assumption that conditional quantiles are linear functions of the parameters, the quantile regression model can be written as:

$$q_\tau(y_i|X) = \beta_0(\tau) + x_i'\beta(\tau) \quad (1)$$

Where $\beta(\tau)$ is the vector of coefficients that may vary across τ and needs to be estimated, and x_i is a vector of exogenous covariates that may include nonlinear functions of underlying variables. This expression indicates that, conditional on τ , the τ -th quantile of y can be approximated by a linear function of X .

Under the assumption that the conditional quantile function is linear and correctly specified, a useful way to think about the data generating process is to consider the following model:

$$y_i = \beta_0(U_i) + x_i'\beta(U_i) \quad (2)$$

where U_i is a random variable that follows a uniform distribution. It can be seen as the rank an individual belongs to among all individuals with the same characteristics. In addition, β_0 and $\beta(U_i)$ are smooth functions that depend on U_i .²

As explained in Wooldridge (2010), the coefficient of quantile regression models can be identified by minimizing the following loss function, with respect to $\beta(\tau)$:

$$\hat{\beta}_0(\tau), \hat{\beta}(\tau) = \min_{\beta(\tau)} \sum_{i=1}^n \rho_\tau(y_i - \beta_0(\tau) - x_i'\beta(\tau)) \quad (3)$$

Where $\rho_\tau(u) = u(\tau - I(u < 0))$ is the check function, and $I(\cdot)$ is the indicator function. In essence, quantile regressions are estimating the parameters locally around the τ -th quantile, although other approaches are possible.³

Most commands for estimating quantile regression models focus on estimating the above loss function, using linear programming techniques, while others like Kaplan and Sun (2017) (`sivqr`) and Chernozhukov et al. (2022) (`qrprocess`) use other optimization techniques.

2. This way of thinking about quantile regression coefficients is similar to the use of Smooth varying coefficient models, except that the running variable is not observed.

3. Kaplan and Sun (2017) for example uses a nonparametric approach that produces a smooth set of beta coefficients. And Bottai and Orsini (2019), proposes methods for estimating parametric quantile regression models, imposing parametric restrictions on the quantile coefficients across the distribution.

When no unobserved heterogeneity is present, quantile regression models can be easily implemented in a panel setting (see Wooldridge (2010)), using a pooled version of the model. However, when unobserved heterogeneity is explicitly present, the estimation of quantile regressions is more challenging. Consider the case of panel data and the following data generating process:

$$y_{it} = \beta_0(U_{it}) + x'_{it}\beta(U_{it}) + \alpha_i(U_{it}) \quad (4)$$

Where U_{it} is a random variable that follows a uniform distribution, and $\alpha_i(U_{it})$ is the unobserved effect that varies across individuals. In this case, the conditional quantile regression model can be written as:

$$q_\tau(y_{it}|x_{it}, \alpha_i(\tau)) = \beta_0(\tau) + x'_{it}\beta(\tau) + \alpha_i(\tau) \quad (5)$$

This specification explicitly considers that the unobserved effect is identified for each i th observation, and that it varies across quantiles ($\alpha_i(\tau)$). A common approach used, yet incorrect due to the incidental parameter problem, is to estimate this model by adding dummy variables for each individual in the quantile regression model (as in Budig and England (2001)), or by demeaning the explanatory variables (as in Budig and Hodges (2010)). In contrast with standard linear models, there is no transformation of the data that can eliminate the individual fixed effects for non-linear models like quantile regressions.

In this framework, the problem of the incidental parameter problem occurs because the unobserved factors cannot be consistently identified. However, because the number of available observations per individual fixed effect is limited, they cannot be estimated with precision. In turn, the cumulative errors in the estimation of the fixed effects will also affect the identification of the conditional distribution of the outcome, which quantile regressions leverage, leading to inconsistent estimates of all parameters.⁴

In the next section, we present a few solutions and implementations for the estimation of quantile regression models with multiple fixed effects.

2.1 Correlated Random Effects: CRE

The first approach we discuss is the use of Correlated Random Effects (CRE) models for the estimation of quantile regression models. The CRE model is an alternative methodology for the estimation of fixed effects models that was proposed by Mundlak (1978) and generalized by Chamberlain (1982). In contrast with standard fixed effects, the approach allows users to control for time-fixed covariates in addition to time-varying covariates. And, in contrast with the random effects model, it does not make the assumption that the unobserved effect is uncorrelated with the observed covariates. Interestingly, in the context of linear models, the CRE model is equivalent to the fixed effects model (Wooldridge 2010). Consider the following model:

4. This is similar to the measuring error problem of dependent variables in quantile regression models discussed in Hausman et al. (2021).

$$y_{it} = \beta_0 + x_{it}\beta + \alpha_i + u_{it} \quad (6)$$

It is well known that if α_i is correlated with x_{it} , the Random Effects (RE) estimator will be inconsistent, due to the omitted variable bias. The solution proposed by Mundlak (1978) and Chamberlain (1982) was to explicitly account for that correlation in the model, by assuming the unobserved effect α_i is a linear projection of the observed time-varying variables plus an uncorrelated disturbance. Specifically:

$$\begin{aligned} \text{Mundlak : } \alpha_i &= \gamma_0 + \bar{x}_i\gamma + v_i \\ \text{Chamberlain : } \alpha_i &= \gamma_0 + x_{i1}\gamma_1 + x_{i2}\gamma_2 + \cdots + x_{iT}\gamma_T + v_i \end{aligned} \quad (7)$$

The main difference between both approaches was that Chamberlain (1982) proposes a more flexible specification allowing all realizations of the time-varying variables to explain the unobserved effect. In contrast, Mundlak's approach only considers the average of the time-varying variables, which is a more restrictive specification. Using either model specification, if we substitute Equation 7 into Equation 6, the final model can be written as:

$$y_{it} = \beta_0 + x_{it}\beta + \gamma_0 + f(x_{it})\Gamma + v_i + u_{it} \quad (8)$$

where $f(x_{it})$ can be the full set of time-varying variables or just the average of them. Notice that in this specification, β_0 and γ_0 cannot be independently identified, and that the new model now has a compound error $v_i + u_{it} = \mu_{it}$, which is uncorrelated with x_{it} . To account for the within-individual correlation driven by v_i , the CRE model should be estimated using either random effects, or clustering standard errors at the individual level (see Wooldridge (2010) for a discussion). Interestingly, either method provides the same results if the panel data is balanced, and all covariates are strictly exogenous. However, this identity breaks down in other cases (see Abrevaya (2013)).

The strategy proposed by Abrevaya and Dahl (2008) was to extend the CRE model (Chamberlain (1982) style) for the estimation of quantile regression models. This, however, has some limitations. First, when the number of periods is large, the number of additional regressors grows quickly, which can lead to other problems during estimation. Second, while the application of Chamberlain (1982) projection approach for unbalanced data is possible (see Abrevaya (2013)), it is not straightforward to implement in practice, especially for the framework of quantile regressions. Instead, we follow Wooldridge (2010) and Wooldridge (2019), and use the Mundlak representation of the CRE model for the estimation of quantile regression models. Wooldridge (2019) has shown that this can be easily applied for cases with unbalanced panels, and the estimation of non-linear models.

Specifically, Wooldridge (2010) suggests that we could use a local projection of the quantile-specific unobserved effect. If we concentrate on $\alpha(U_{it})$, where U_{it} is a random variable that follows a uniform distribution, we could write the unobserved effect as:

$$\alpha_i(U_{it}) = \gamma_0(U_{it}) + \bar{x}'_i \gamma(U_{it}) + v_i^{U_{it}} \quad (9)$$

Then, we can use Equation 9 to write the new Data Generating Process (DGP):

$$y_{it} = \beta_0(U_{it}) + x_{it}\beta(U_{it}) + \gamma_0(U_{it}) + \bar{x}'_i \gamma(U_{it}) + v_i^{U_{it}} \quad (10)$$

Two important points to note here. First, as before, $\beta_0(\cdot)$ and $\gamma_0(\cdot)$ cannot be independently identified, which makes the interpretation of the constant term difficult. Second, $v_i^{U_{it}}$ is not a smooth function of U_{it} , but rather an unrelated disturbance that is left after modeling the unobserved effect, and remains unobserved. If we assume that $v_i^{U_{it}}$ is small enough compared to the overall variation driven by U_{it} , we could identify the quantile regression coefficients as follows:

$$q_\tau(y_{it}|x_{it}, \bar{x}_i) = b_0(\tau) + x_{it}\beta(\tau) + \bar{x}'_i \gamma(\tau) \quad (11)$$

Which can be estimated using any standard quantile regression method. However, if $v_i^{U_{it}}$ is large, standard estimators will leverage the distribution of the compound error $v_i^{U_{it}}$ and U_{it} , which may lead to inconsistent estimates of the quantile coefficients.⁵

Nevertheless, assuming that the residual $v_i^{U_{it}}$ is small, the CRE-quantile regression approach has a few other benefits that may be of interest. First, as discussed in Wooldridge (2019), it can be easily used in the presence of unbalanced panels. Second, it may also provide an approach to control for multiple fixed effects.⁶ For example, let us expand on Equation 4, and consider the case of a two-way fixed effects model:

$$y_{it} = \beta_0(U_{it}) + x'_{it}\beta(U_{it}) + \alpha_i(U_{it}) + \alpha_t(U_{it})$$

To apply the two-way CRE model, we could use the following representation of the unobserved effects:

$$\alpha_i(U_{it}) + \alpha_t(U_{it}) = \gamma_0(U_{it}) + \lambda_i^x \gamma_i(U_{it}) + \lambda_t^x \gamma_t(U_{it}) + v_{it}^{U_{it}} \quad (12)$$

where λ_i^x and λ_t^x are obtained by estimating the following model for each explanatory variable x_{it} :

$$x_{it} - \bar{x} = \lambda_i^x + \lambda_t^x + \epsilon_{it} \quad (13)$$

We use the centered transformation of the explanatory variable, that is $x_{it} - \bar{x}$, so that all λ 's have an expected value of zero. In contrast with Baltagi (2023), we

5. This is the main critique raised by Canay (2011) to the estimator proposed by Abrevaya and Dahl (2008). In fact, the more dominant $v_i^{U_{it}}$ becomes, the more the estimates will resemble the OLS estimates.

6. Baltagi (2023) and Wooldridge (2021) discuss this for the two-way Mundlak estimator

suggest that rather than modeling each individual component separately, it is easier to think of the problem of modeling the combination of the two (or many) unobserved components as a function of λ_i^x and λ_t^x , which are the equivalent to \bar{x}_i in the Mundlak one-way fixed effect model. Additionally, different from Wooldridge (2021) and Baltagi (2023), we emphasize that the estimation of λ_i^x and λ_t^x should be done simultaneously (Equation 13), rather than estimating the conditional means separately. This is more general and applicable to any number of fixed effects. This can be done using an iterative process similar to Rios-Avila (2015) or Correia (2016).⁷

With these considerations, the conditional quantile regression can be written as:

$$q_\tau(y_{it}|x_{it}, \lambda_i^x, \lambda_t^x) = x'_{it}\beta(\tau) + \lambda_i^{x'}\gamma_i(\tau) + \lambda_t^{x'}\gamma_t(\tau)$$

Which could be extended to any number of fixed effects. As before, this approach is valid if the residual $v_{it}^{U_{it}}$ from the time and individual fixed effects (or all fixed effects considered) are small enough compared to the variation driven by the latent rank variable U_{it} .

In terms of the standard errors, for the linear CRE model, it is suggested to use the random effects estimator, or clustering standard errors at the individual level. For the quantile regression model, clustering standard errors at the individual level is also suggested by Wooldridge (2010), and some routines already implement this feature. When multiple fixed effects are considered, it is suggested to use the bootstrap methods for the estimation of the standard errors.

2.2 Canay (2011) Estimator

The second approach under consideration is the estimator proposed by Canay (2011). As mentioned before, this paper argues that the estimator proposed by Abrevaya and Dahl (2008), and thus the implementation described above, may not provide consistent estimates of the quantile regression coefficients, as long as there is a disturbance $v_i^{U_{it}}$ left after modeling Equation 9. Instead, under the assumption that the unobserved effect is a pure location shift, they propose an alternative estimator that can be used to consistently estimate the quantile regression coefficients.

Before presenting the estimator, it is convenient to review a second approach that has been used to understand quantile regression models: The location-scale model. Under this specification, consider the following data generating process:

$$\begin{aligned} y_i &= \beta_0 + \beta_1 x_i + \gamma_0(U_i) + \gamma_1(U_i)x_i \\ \text{or} \\ y_i &= \beta_0 + \beta_1 x_i + \mu_i \end{aligned} \tag{14}$$

7. If the panel is perfectly balanced, estimating λ_i^x and λ_t^x separately will provide the same results as estimating them simultaneously. Internally, we use Correia (2016) `reghdfe` to obtain the predicted fixed effects

In this specification, we assume that β_0 and β_1 are the location parameters that capture how the whole distribution of y_i is affected by x_i . In contrast, $\gamma_0(U_i)$ and $\gamma_1(U_i)$ are the scale parameters that capture the heterogeneous effect of x_i on y_i that deviates from the location effect. When using a simple linear regression model, estimated via OLS, we could assume that the compound component $\gamma_0(U_i) + \gamma_1(U_i)x_i$ is fully captured by the error term μ_i .

This provides three insights. First, that OLS could be used to identify the location effect of x_i on y_i , which we know as the average or conditional mean effect. Second, if $\gamma_1(U_i)$ is different from zero, the model is heteroskedastic, and the quantile regression could be used to identify this type of unobserved heterogeneity. Lastly, if a covariate has no scale effect, all quantile coefficients will be the same as the OLS coefficient, except for the constant.

Although of little use, this location-scale model can be easily estimated using a two-step approach. First, estimate the location effect of x_i on y_i using OLS:

$$y_i = \beta_0 + x_i' \beta + \mu_i \quad (15)$$

Then, using the predicted residuals $\hat{\mu}_i$, estimate the quantile regression model:

$$q_\tau(\hat{\mu}_i | x_i) = \gamma_0(\tau) + x_i' \gamma(\tau) \quad (16)$$

It is a simple exercise to show that adding $\beta + \gamma(\tau)$ provides the same point estimates as estimating the full quantile regression model. However, it does provide a simple connection between OLS and quantile regression models.

Now, let's reconsider the data generating process in Equation 4. Canay (2011) imposes the assumption that α_i is a pure location shift that should be constant across quantiles. More explicitly, the data generating process in a panel data setting can be written as follows:

$$y_{it} = \beta_0(U_{it}) + x_{it} \beta(U_{it}) + \alpha_i \quad (17)$$

which imply that the conditional quantile regression model can be written as:

$$q_\tau(y_{it} | x_{it}, \alpha_i) = \beta_0(\tau) + x_{it} \beta(\tau) + \alpha_i \quad (18)$$

This assumption has important implications for the identification of the quantile regression coefficients. First, by assuming that α_i is a pure location shift, it reduces the number of parameters that need to be estimated in the model, because α_i can now be estimated globally, while the quantile regression coefficients can be estimated locally. Second, based on the previous insights of the location scale model, it suggests that we can estimate the quantile regression model using a two-step approach. First, estimate all location effects using OLS, and then estimate the scale effects using the predicted

residuals, but excluding variables we assume have a pure location shift effect. More formally, the estimator proposed by Canay (2011) can be described as follows:

1. Estimate the location effect of x_{it} and the unobserved heterogeneity α_i on y_{it} using OLS:

$$y_{it} = \beta_0 + x'_{it}\beta_1 + \alpha_i + \varepsilon_{it}$$

2. Use the predicted fixed effects $\hat{\alpha}_i$ ⁸ to transform the dependent variable as $\tilde{y}_{it} = y_{it} - \hat{\alpha}_i$, and estimate the quantile regression model:

$$q_\tau(\tilde{y}_{it}|x_{it}) = \beta_0(\tau) + x'_{it}\beta_1(\tau)$$

This simple approach allows for the identification of the quantile coefficients by imposing the assumption that the unobserved characteristics only have a location shift effect on the outcome. In addition, like the CRE model, it can be extended to multiple fixed effects, as long as one is willing to assume that the unobserved effects are pure location shifts. For example, consider a case with two fixed effects dimensions (individual and time). Under the assumption that the unobserved effects are pure location shifts, the data generating process can be written as:

$$y_{it} = \beta_0(U_{it}) + x'_{it}\beta(U_{it}) + \alpha_i + \alpha_t \quad (19)$$

As before, if we assume that α_i and α_t are constant across quantiles, we could use the same two-step approach to estimate the quantile regression coefficients. First, estimate the location effects using OLS, and then estimate the quantile regression model using the transformed dependent variable, after absorbing the predicted fixed effects.

$$q_\tau(y_{it} - \hat{\alpha}_i - \hat{\alpha}_t|x_{it}) = q_\tau(\tilde{y}_{it}|x_{it}) = x_{it}\beta(\tau) \quad (20)$$

This can again be easily estimated using standard quantile regression methods and extended to any number of fixed effects.

2.3 Modified Canay(2011) Estimator

Perhaps one of the main limitations of Canay (2011) is that it assumes that the unobserved effect is a pure location shift. In fact, this is one of the criticisms raised by Machado and Santos Silva (2019) to the estimator. While this makes sense intuitively, because an individual will only be assigned to a single rank at a given point in time, it

8. Empirically, this can be done using the `reghdfe` command, as part of the `abs()` suboptions.

is not consistent with the idea that the unobserved effect is in fact a proxy for an unobserved characteristic of the individual, and that characteristic could have a different impact on the dependent variable across quantiles. In this case, if this assumption is violated, it may lead to inconsistent estimates of the quantile coefficients. To address this limitation, we propose a small modification to the Canay estimator.

We start by assuming that the unobserved effect represents some characteristics of the individual that are constant across quantiles, and that can be compared across individuals. Under this consideration, the data generating process can be written as:

$$y_{it} = \beta_0 + x'_{it}\beta + \beta_\alpha\alpha_i + \gamma_0(U_{it}) + x'_{it}\gamma(U_{it}) + \gamma_\alpha(U_{it})\alpha_i \quad (21)$$

Where α_i is the unobserved effect, β_α is the location coefficient of the unobserved heterogeneity, and $\gamma(U_{it})$ is a smooth function that varies across quantiles. For the identification of α_i , we start with the same approach as Canay (2011), imposing the assumption that $\beta_\alpha = 1$. In other words, the first step from Canay (2011) is the same as the first step presented before.

The second step, however, suggests that rather than transforming the dependent variable using the predicted fixed effects, we should estimate the quantile regression model using the predicted unobserved effects as an additional explanatory variable. This can be done by estimating the following model:

$$q_\tau(y_{it}|x_{it}, \hat{\alpha}_i) = x_{it}\beta(\tau) + \beta(\tau)\hat{\alpha}_i$$

As before, this model can be extended to multiple fixed effects by simply estimating the unobserved effects using OLS, and then estimating the quantile regression model using the predicted unobserved effects. The main advantage over Canay (2011) is that this estimator allows for the unobserved effect to have a different impact on the dependent variable across quantiles, which may be more realistic in many applications. However, it assumes the OLS estimator does allow for the consistent estimation of an unobserved effect that is comparable across individuals, which may not always be the case.

In terms of standard errors, Canay (2011) provides some guidance for the derivation of analytical standard errors for their estimator. Recently, however, Besstremyannaya and Golovan (2019) has shown that the analytical standard errors derivations are incorrect. Instead, based on their recommendations, we suggest that the bootstrap method should be used for the estimation of the standard errors for both the Canay and Modified Canay estimators.

2.4 Method of Moments Quantile Regression Machado and Santos Silva (2019)

The last methodology we consider is the Method of Moments Quantile Regression (MMQREG) estimator proposed by Machado and Santos Silva (2019), and extended by Rios-Avila et al. (2024). The methodology was proposed as a feasible approach to incorporate fixed effects in a quantile regression model, allowing for unobserved effects to have a different impact on the dependent variable across quantiles. This is done by separating the identification of quantile coefficients into a location, scale, and quantile effect, using a method of moments approach.

To understand this approach, let's consider the data generating process from Equation 14. As mentioned earlier, this approach suggests that a quantile regression model can be identified using a location-scale model, where the location effect shows the average effect of the covariates on the dependent variable, and the scale effect shows the heterogeneous effect of the covariates, as a deviation from the location effect. Machado and Santos Silva (2019) extends this idea by suggesting that the scale component can be further decomposed into a pure scale effect and a mediating factor. Specifically, the author considers the case where the data generating process can be written as:

$$\begin{aligned} Y_i &= \beta_0 + x'_i\beta + (\delta_0 + X'_i\delta) * \mu_i \text{ or} \\ Y_i &= \beta_0 + x'_i\beta + (\delta_0 + X'_i\delta) * F^{-1}(U_i) \text{ or} \end{aligned} \quad (22)$$

This specification assumes that μ_i is an identically and independently distributed random variable with any arbitrary distribution. $\delta_0 + X'_i\delta$ denotes the multiplicative scale component (heteroskedasticity generating component), and $\beta_0 + x'_i\beta$ denotes the location coefficients. The second line in Equation 22 represents the same model, but using the inverse of the distribution function of U_i as the mediating factor. As before, U_i is a random variable that follows a uniform distribution and captures the rank of the individual among all individuals with the same characteristics.

What is interesting about this specification is that it simplifies the identification of the quantile coefficients by imposing a strict parametric relationship across the quantile-specific coefficients. Specifically, the location and scale coefficients can be identified globally, requiring only the local identification of the distribution of μ_i to identify all quantile coefficients.

For this simple case, the MMQREG estimator can be described as follows:⁹

1. Estimate the location effect of x_i on y_i using OLS:

$$y_i = \beta_0 + x'_i\beta + R_i$$

9. Further details on the implementation of the methodology can be found in Machado and Santos Silva (2019), with additional extensions in Rios-Avila et al. (2024).

2. Use the predicted residuals \hat{R}_i to estimate the scale effect using OLS, where the dependent variable is defined as $|\hat{R}_i|$:

$$|\hat{R}_i| = \delta_0 + x'_i \delta$$

3. Obtain a standardized residual by dividing the residuals of (1) by the predicted scale effect from (2), and estimate the τ th unconditional quantile of this distribution:

$$\hat{\mu}_i = \frac{\hat{R}_i}{\delta_0 + x'_i \delta}$$

$$q_\tau(\hat{\mu}_i) = q_0(\tau)$$

4. Finally, the quantile regression coefficients can be estimated by the following equation:

$$\beta_\tau = \beta + q_0(\tau)\delta$$

It can be seen that because Steps 1 and 2 are estimated globally using OLS, they can easily be extended to multiple fixed effects without major difficulties. One can simply assume, for example, that the fixed effects are estimated using a dummy variable approach, and all the steps described above follow.¹⁰ The derivation of standard errors follows from the use of the method of moments approach and the identification of the empirical influence functions, based on the following moment conditions:

$$\begin{aligned} E[y_i - \beta_0 - x'_i \beta] &= 0 \\ E[|y_i - \beta_0 - x'_i \beta| - \delta_0 + x'_i \delta] &= 0 \\ E\left[I\left(\frac{y_i - \beta_0 - x'_i \beta}{\delta_0 + x'_i \delta} < q_0(\tau)\right)\right] &= \tau \end{aligned} \tag{23}$$

10. In Rios-Avila et al. (2024), additional steps are presented because of the use of partial-out covariates. However, the same results can be obtained using the simpler approach described here. Despite the flexibility of this approach, it is important to note that the MMQR estimator is based on the assumption that all quantile coefficients are constructed as a simple combination of the location and scale coefficients. They only differ in the value of the mediating factor, the unconditional quantile of the standardized residual. This may be a strong assumption and may not be applicable in all cases. However, it does provide a few advantages over the other approaches. First, as suggested by Machado and Santos Silva (2019), with a sufficiently flexible specification of the scale component, the MMQR estimator can capture most of the important features related to the heterogeneous effect of the covariates on the dependent variable. Second, because of the identification assumptions, this estimator avoids the problem of quantile crossing¹¹, which is common in other quantile regression models. Nevertheless, it is important to note that the MMQR assumptions may not be applicable in all cases.

3 Implementation: `qregfe`

While the estimation of the methods described above can be easily implemented using standard Stata commands, the process can be cumbersome and prone to errors. To address this, we introduce the Stata command `qregfe`. This command allows users to estimate quantile regression models with fixed effects using the CRE, Canay, and Modified Canay estimators, as well as the MMQR estimator.¹² The command is designed to be user-friendly, allowing for the estimation of quantile regression models with fixed effects in a single line of code, with only minor changes for model estimation.

The command uses the standard Stata syntax:

```
qregfe depvar [indepvars] [if] [in] [pw], quantile(#) [options]
```

Where `depvar` is the dependent variable, `indepvars` are the independent variables, `if` and `in` are the standard Stata options. Because most of the estimators require the use of bootstrapping, Canay and CRE estimators do not allow for the use of weights. Only Machado and Santos Silva (2019) estimator allows for it.

The main options for the command are:

- Fixed effect estimator:
 - `cre`: Correlated Random Effects estimator
 - `canay`: Canay (2011) estimator
 - `canay(modified)`: Modified Canay (2011) estimator
 - `mmqreg`: Machado and Santos Silva (2019) estimator. This is the only estimator that also has its own command `mmqreg`.
- `qmethod(qmethod_options)`: Specifies the method used to estimate the quantile regression component. The default is `qreg`, but any other quantile regression command can be used. This will affect how `quantile(#)` is specified.
 - If necessary, one can also request `qmethod_options` that are specific to the quantile regression command used.
- `q(#)`: specifies the specific quantile for which coefficients will be obtained. The default is the median `q(50)`. The values this could take depend on the `qmethod()` used. For example, if using `qmethod(sqreg)`, one could specify `q(10 25 90)` to obtain the 10th, 25th, and 90th quantiles. However, if one uses `qmethod(qrprocess)`¹³, one could specify `q(0.1 0.25 0.9)` to obtain the same quantiles. If using `mmqreg`, multiple quantiles are always allowed.

12. The MMQR estimator is also implemented by the `xtqreg` command, which was provided by the authors. There is also `mmqreg`, which was originally implemented as a stand-alone command, but it is now integrated within this command.

13. `qrprocess` is a user-written command that uses a different algorithm for the estimation of quantile regression models.

- **abs(varlist)**: specifies the variables that would be used to absorb the fixed effects. This is necessary for CRE and Canay estimators. For the **mmqreg** estimator, this is not necessary, as the estimator still works for the case of no fixed effects.
- **boot[(bootstrap_options)]**: request bootstrap standard errors to be computed.
 - if not specified, standard errors correspond to the default from **qmethod()**, except for **mmqreg**, which uses the GLS standard errors, as proposed by Machado and Santos Silva (2019) and Rios-Avila et al. (2024).
 - **bootstrap options** are the standard options for the **bootstrap** in Stata, using the same results.
- **parallel**: request the estimation of bootstrap standard errors to be performed using **parallel** package (Vega Yon and Quistorff (2019)), based on the specifications given in **bootstrap_options**.
 - **parallel_cluster(#)** specifies the number of clusters to be used for the parallelization. The default is 2.
- **seed(#)**: specifies the seed for the random number generator. The default is to use none.
- **other_options**: Other options that are specific to the **qmethod()** used can also be passed to the command directly.

mmqreg specific options are:

- **robust** and **cluster(varname)**: These options are available when using **mmqreg** estimator. They follow the derivations from Rios-Avila et al. (2024), based on robust and cluster standard errors for GMM estimators. The default are the GLS standard errors.
- **denopt()**: This option requests alternative options for the estimator of the point density for the residual quantile. See **help qreg**, specifically standard errors **vceopts** section.
- **dfadj**: Request to adjust the degrees of freedom for the standard errors. The default is to use no adjustment.
- **ls**: Request displaying the location and scale coefficients. The default is not to display them.

4 Visualization

In addition to the introduction of the **qregfe** command, we also introduce the **qregplot** command. This command allows users to visualize the results of the quantile regression coefficients obtained using the **qregfe** command, as well as other quantile regression commands. It works similarly to the command **grqreg** (Azevedo 2004), but allows for the visualization of coefficients across many quantile regression estimators. The command uses the standard Stata syntax:

```
qregplot [indep varlist] , quantiles(numlist) cons options
```

This command does not work independently. Instead, it requires users to first estimate a quantile regression model using any of the standard quantile regression commands (`qreg`, `bsqreg`, `sqreg`, `qreg2`, `qrprocess`, `qregfe`, etc). Then, calling on the `qregplot` command will identify the general syntax used by the previous command, and estimate the same model across the specified quantiles, storing the point estimates and confidence intervals, which are then used to plot the coefficients across the quantiles. The confidence intervals are estimated at the same level of confidence as the original model. The main options for the command are:

- **indep varlist**: specifies the independent variables to be plotted. The default is to plot all the independent variables' coefficients except for the constant. One should use the same names as identified in the coefficients matrix.
- **quantiles(numlist)**: is the list of quantiles to be plotted. The default is to use `quantiles(10(5)90)`. If the underlying model was estimated using `sqreg`, this is superseded by the quantiles used in the simultaneous quantile regression model.
- **cons**: request the plotting of the constant term. The default is not to plot it.
- **ols**: request estimating and plotting of the OLS coefficients along with the quantile regression coefficients. The default is not to plot them. When requested, the OLS specification will follow the same specification as the quantile regression model, but will not include the absorbed fixed effects, unless specified using `olsopt()`
- **olsopt()**: specifies the options to be used for the OLS estimation. All options available for `regress` can be specified here. If one requests absorbing fixed effects, the command will use `reghdfe` to estimate the OLS model.
- **seed(#)**: specifies the seed for the random number generator, so the same seed is used across quantiles. The default is to use none.
- **label**: Request using variable labels as titles for each coefficient plot. The default is to use variable names.
- **labelopt(options)**: Specifies options for labels. Two options are available: `lines(#L)` which requests breaking the labels into `#L` number of lines, default is 1; `maxlength(#K)` requests breaking a label into lines of a max length `#K`.
- **mtitles(titles)**: Provides titles for each plot. If not enough titles are provided, the command will use the variable names, or variable labels.

Figure design options:

- **raopt(rarea options)**: Specifies options to be used for plotting the Confidence Intervals. The default options are `pstyle(p1) fintenity(30) lwidth(none)`.
- **lnopt(line options)**: Specifies options to be used for plotting the line of point estimates. The default options are `pstyle(p1) lwidth(0.3)`.
- **twopt(twoway options)**: Specifies options to be used for the two-way graph that combines the `rarea` and `line` plots. The default option is to set graph and plot region margins to `vsmall`.

- **graph_combine options:** One can specify any **graph combine** options directly to the command. This will be added on the last step that combines the graphs, and will not affect the individual plots. However, if only one coefficient is plotted, options that are specified directly will affect the two-way plot.

Storing, saving, and recycling options:

Because the estimation of multiple quantile regression models can be time-consuming, the **qregplot** command allows for storing or saving the results for later use.

- **estore(name):** Request storing the results into memory. Similar to **estimates store**.
- **esave(filename):** Request saving the results into a **ster** file. Similar to **estimates save**.
- **from(name):** Request using the stored results from memory to plot the quantile coefficients. Similar to **estimates restore**. If results were saved to a file, one should first use **est use filename** to load the results into memory, and then store them into memory with a name.

5 Illustration

To illustrate the use of the **qregfe** and **qregplot** commands, we use the **nlswork** dataset, which is part of the Stata distribution. This data is a subset of the National Longitudinal Survey of Young Women who were 14-24 years old in 1968, and follows individuals from 1968 to 1988. Given the nature of the data, the data is unbalanced, with some observations being observed only once or twice, while few others observed for every year. For this illustration we consider only cases where individuals were observed for at least 5 periods. This leaves us with a sample of 2771 observations over 15 years.

For the illustration of the model, we consider a simple mincer-type wage equation, where the dependent variable is the log of wages, and the independent variables are age, years of experience, years of education, and if he/she did not live in an MSA. We only consider individual fixed effects. The code to estimate the model is as follows:

```
webuse nlswork, clear
bysort idcode:gen nobs=_N
drop if nobs<5
* Linear Regression model with fixed effects
reghdfe ln_wage age tenure ttl_exp not_smsa , abs(idcode)
est sto ols
* Quantile Regression CRE
qregfe ln_wage age tenure ttl_exp not_smsa, abs(idcode) cre q(50)
est sto cre_q50
* Quantile Regression Canay
qregfe ln_wage age tenure ttl_exp not_smsa, abs(idcode) canay q(50)
```



```

est sto canay_q50
* Quantile Regression mCanay
qregfe ln_wage age tenure ttl_exp not_smsa, abs(idcode) canay(modified) q(50)
est sto mcanay_q50
* Quantile Regression MMQREG
qregfe ln_wage age tenure ttl_exp not_smsa, abs(idcode) mmqreg q(50)
est sto mmqreg_q50
esttab ols cre_q50 canay_q50 mcanay_q50 mmqreg_q50, ///
       se b(3) tex mtitle(OLS CRE Canay MCanay MMqreg)

```

Which produces the results shown in Table 1.

Table 1: Mincer regression: OLS, CRE, Canay, Modified Canay, and MMQREG estimators

	(1) OLS	(2) CRE	(3) Canay	(4) MCanay	(5) MMqreg
main					
age	-0.004*** (0.001)	-0.009*** (0.001)	-0.005*** (0.000)	-0.005*** (0.000)	-0.004*** (0.001)
tenure	0.011*** (0.001)	0.016*** (0.001)	0.011*** (0.001)	0.011*** (0.001)	0.010*** (0.001)
ttl_exp	0.030*** (0.002)	0.031*** (0.002)	0.029*** (0.001)	0.029*** (0.001)	0.030*** (0.002)
not_smsa	-0.098*** (0.010)	-0.088*** (0.015)	-0.098*** (0.003)	-0.099*** (0.003)	-0.099*** (0.010)
m1_age		0.009*** (0.002)			
m1_tenure		0.016*** (0.002)			
m1_ttl_exp		0.010*** (0.003)			
m1_not_smsa		-0.150*** (0.016)			
__f1__				0.991*** (0.005)	
__cons	1.586*** (0.019)	1.696*** (0.028)	1.654*** (0.009)	1.653*** (0.009)	1.606*** (0.021)
N	23534	23534	23534	23534	23534

Standard errors in parentheses

* $p < 0.05$, ** $p < 0.01$, *** $p < 0.001$

All models were estimated using default options. Thus, standard errors may be biased, even for the simple OLS model. Because the quantile regression models are estimated at the median, its not surprising that the coefficients are almost identical

to the OLS model. However, the standard errors differ, as they are estimated using different methods. Nevertheless, the exercise allow us to further describe the small differences in the output compared to standard quantile regression models.

1st. For the CRE model, there are 4 more variables added to the specification. This corresponds to the individual level averages of all variables used in the main specification. If more fixed effects were used, more variables would be added with the prefix `m?_`, where `?` indicates with respect to which fixed effect the variable is “averaged”. They do not required additional interpretation, but in linear models could be used to test between random and fixed effects estimators.

2nd. For the Canay model, the output is identical to the OLS specification. This is because the model transformation is done internally.

3rd. For the Modified Canay model, there is also an additional variable added to the specification `__f1__`. This represents the predicted fixed effects. If more Fixed effects were used, more variables would be added with the following form `__f?__`. Could be potentially used to test the validity of the Canay estimator.

4th. For the MMQREG model, the output is identical to the OLS specification. However, `ls` were specified, we could also add the location and scale coefficients to the output.

To visualize the results, we will use the `qregplot` command. We will store the estimates into memory, so it is easier to recall them later, and construct the plots for the coefficients.

```
qregfe ln_wage age tenure ttl_exp not_smsa, abs(idcode) cre q(50)
qregplot, estore(cre)
qregfe ln_wage age tenure ttl_exp not_smsa, abs(idcode) canay q(50)
qregplot, estore(canay)
qregfe ln_wage age tenure ttl_exp not_smsa, abs(idcode) canay(modified) q(50)
qregplot, estore(mcanay)
qregfe ln_wage age tenure ttl_exp not_smsa, abs(idcode) mmqreg q(50)
qregplot, estore(mqr)
```

Once all models have been estimated and stored, we can build the plots using the following code, with only minor modifications for the creation of the plots for other variables:

```
qregplot age, from(cre) name(cre_age) title("CRE")
qregplot age, from(canay) name(canay_age) title("Canay")
qregplot age, from(mcanay) name(mcanay_age) title("Modified Canay")
qregplot age, from(mqr) name(mqr_age) title("MMQREG")
graph combine cre_age canay_age mcanay_age mqr_age , ycommon
graph export "qregplot_age.pdf", replace
```

The code above will produce the plots shown in Figure 1, Figure 2 and Figure 3.

The plots show the considerable differences in the coefficients across the estimators, due to differences in the underlying assumptions.

When considering age, Figure 1, the CRE estimate suggests that, once individual fixed effects have been accounted for, the effect of age on wages is mostly negative, except on the top and bottom of the conditional distribution. Canay and modified Canay, on the other hand suggest a more stable but negative effect on wages, although that effect is less pronounced at the top of the distribution. The MMQREG estimator also suggest a negative effect, somewhat smaller than the other estimators, but with a slight decrease across the quantiles.

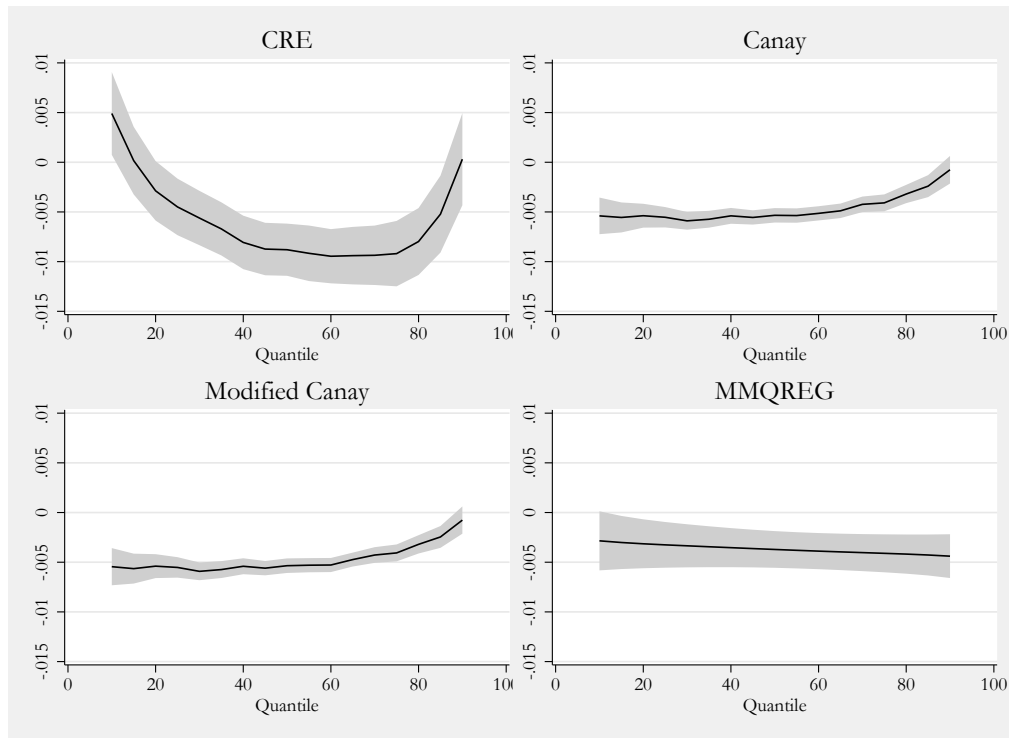


Figure 1: QR Coefficients for Age

When considering Tenure, Figure 2, all estimators suggest very similar results. As tenure increases, one would also expect higher wages, with largest effects at the bottom of the distribution, but decreasing as we move to the top. CRE results decline the least, while Canay and Modified Canay show a steeper decline.

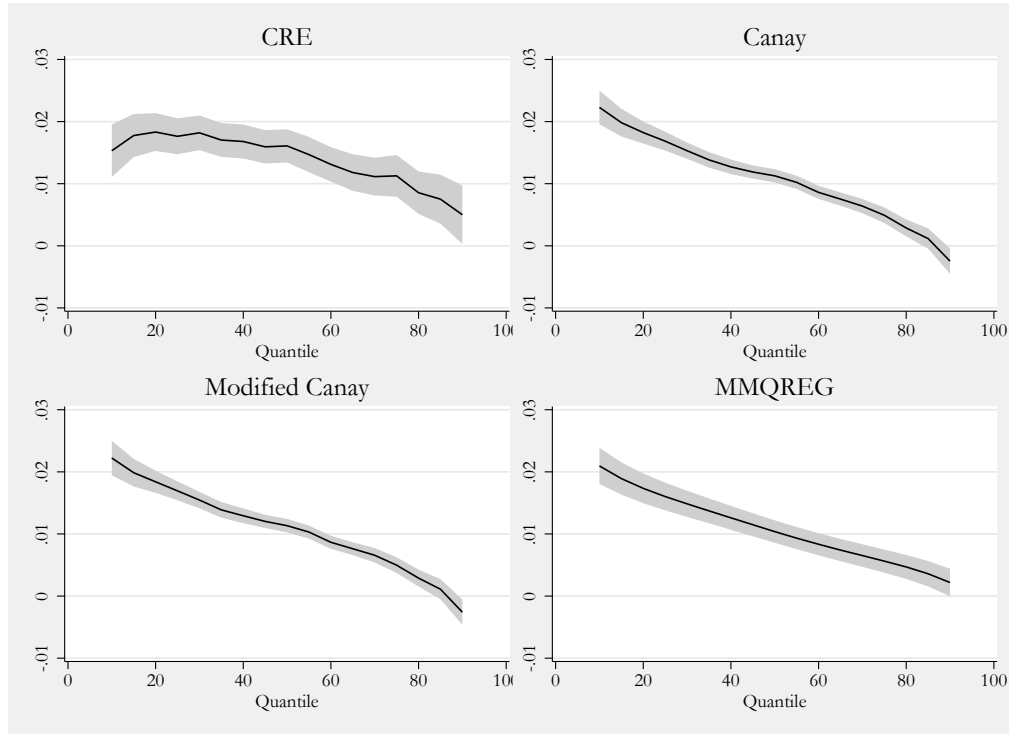


Figure 2: QR Coefficients for Tenure

Finally, when considering experience, Figure 3, the results across estimators is mixed. They all suggest that, people with more experience, also tend to earn higher wages, with the effect being more pronounced the further up of the distribution we analyze. However, the magnitude of the effect varies across estimators. CRE estimates suggest the steepest change in this relationship, with Canay and Modified Canay suggesting a much flatter trend. MMQREG, on the other hand shows results that are in between the other estimators.

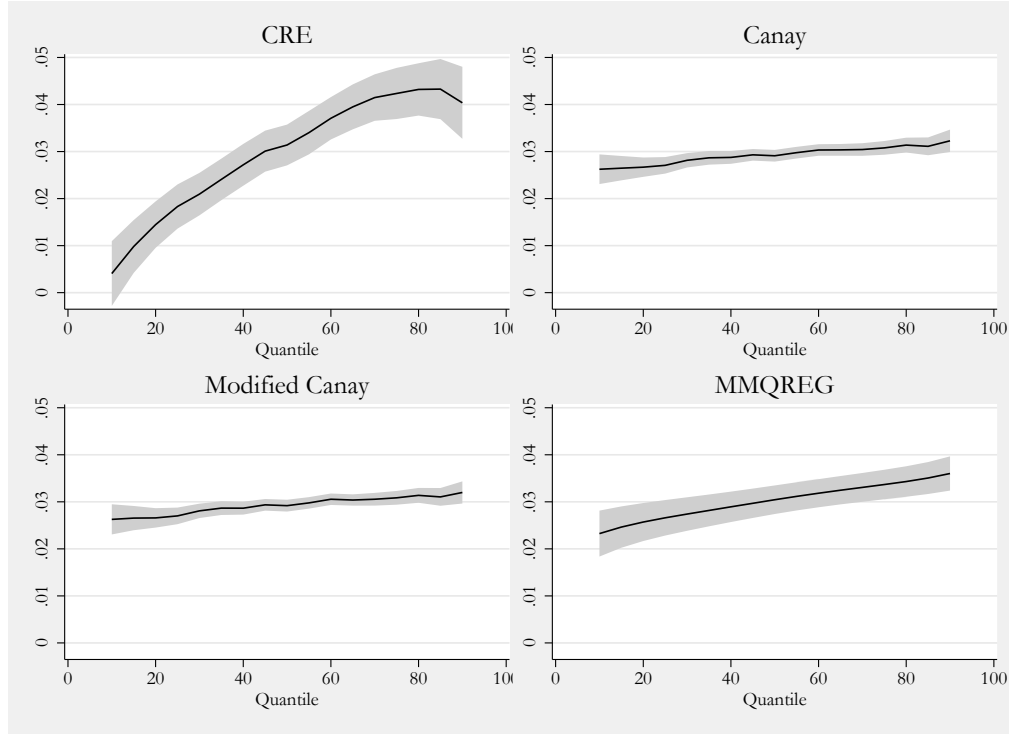


Figure 3: QR Coefficients for Experience

6 Conclusions

This paper has introduced two new Stata commands, `qregfe` and `qregplot`, designed to estimate and visualize quantile regression models with fixed effects. These commands address the growing need for tools that can handle unobserved heterogeneity in quantile regression frameworks, particularly in panel data settings.

The `qregfe` command offers an implementation of several methodologies for estimating quantile regressions with fixed effects, including the Correlated Random Effects (CRE) approach, Canay (2011) estimator, a modified version of Canay's estimator, and the Method of Moments Quantile Regression (MMQREG) proposed by Machado and Santos Silva (2019). By providing these different estimation methods within a single command, `qregfe` allows researchers to easily compare results across methods and select the most appropriate approach for their specific research context.

The companion command `qregplot` facilitates the visualization of quantile regression results, enabling researchers to graphically examine how the effects of covariates vary across different points of the conditional distribution of the outcome variable. It can be used not only with `qregfe`, but also with a large set of other quantile regression

commands. This type of visual representation is valuable for identifying and interpreting heterogeneous effects that are harder to observed simply using regression tables.

While these commands aim to make this methodologies more accessible to applied researchers, it should be noted that each method relies on specific assumptions that may not hold in all empirical settings. Thus, one should carefully consider the underlying assumptions of each approach and potentially compare results across methods to ensure robust conclusions.

7 References

- Abrevaya, J. 2013. The projection approach for unbalanced panel data. *The Econometrics Journal* 16(2): 161–178. <http://dx.doi.org/10.1111/j.1368-423X.2012.00389.x>.
- Abrevaya, J., and C. M. Dahl. 2008. The Effects of Birth Inputs on Birthweight. *Journal of Business & Economic Statistics* 26(4): 379–397.
- Azevedo, J. P. 2004. GRQREG: Stata module to graph the coefficients of a quantile regression. Statistical Software Components, Boston College Department of Economics. <https://ideas.repec.org/c/boc/bocode/s437001.html>.
- Baltagi, B. H. 2023. The two-way Mundlak estimator. *Econometric Reviews* 42(2): 240–246. <https://www.tandfonline.com/doi/full/10.1080/07474938.2023.2178139>.
- Besstremyannaya, G., and S. Golovan. 2019. Reconsideration of a simple approach to quantile regression for panel data. *The Econometrics Journal* 22(3): 292–308. <http://dx.doi.org/10.1093/ectj/utz012>.
- Bottai, M., and N. Orsini. 2019. qmodel: A command for fitting parametric quantile models. *The Stata Journal: Promoting communications on statistics and Stata* 19(2): 261–293. <http://dx.doi.org/10.1177/1536867X19854002>.
- Budig, M. J., and P. England. 2001. The Wage Penalty for Motherhood. *American Sociological Review* 66(2): 204–225. <http://dx.doi.org/10.1177/000312240106600203>.
- Budig, M. J., and M. J. Hodges. 2010. Differences in Disadvantage. *American Sociological Review* 75(5): 705–728. <http://dx.doi.org/10.1177/0003122410381593>.
- Canay, I. A. 2011. A simple approach to quantile regression for panel data. *The Econometrics Journal* 14(3): 368–386.
- Chamberlain, G. 1982. Multivariate regression models for panel data. *Journal of Econometrics* 18(1): 5–46.
- Chernozhukov, V., I. Fernández-Val, and B. Melly. 2022. Fast algorithms for the quantile regression process. *Empirical Economics* 62(1): 7–33. <https://doi.org/10.1007/s00181-020-01898-0>.
- Correia, S. 2016. A Feasible Estimator for Linear Models with Multi-Way Fixed Effects. *Unpublished Manuscript*.
- Galvao, A. F., and K. Kengo. 2017. Quantile regression methods for longitudinal data. In *Handbook of quantile regression*, 363–380. Chapman and Hall/CRC.
- Hausman, J., H. Liu, Y. Luo, and C. Palmer. 2021. Errors in the Dependent Variable of Quantile Regression Models. *Econometrica* 89(2): 849–873. <http://dx.doi.org/10.3982/ECTA14667>.

- Kaplan, D. M., and Y. Sun. 2017. Smoothed estimating equations for instrumental variables quantile regression. *Econometric Theory* 33(1): 105–157. Publisher: Cambridge University Press.
- Koenker, R., and G. Bassett. 1978. Regression Quantiles. *Econometrica* 46(1): 33–50.
- Lancaster, T. 2000. The incidental parameter problem since 1948. *Journal of Econometrics* 95(2): 391–413.
- Machado, J. A., and J. Santos Silva. 2019. Quantiles via moments. *Journal of Econometrics* 213(1): 145–173.
- Melly, B., and M. Pons. 2023. Minimum Distance Estimation of Quantile Panel Data Models. Unpublished working paper.
- Mundlak, Y. 1978. On the Pooling of Time Series and Cross Section Data. *Econometrica* 46(1): 69. <http://dx.doi.org/10.2307/1913646>.
- Neyman, J., and E. L. Scott. 1948. Consistent Estimates Based on Partially Consistent Observations. *Econometrica* 16(1): 1–32.
- Rios-Avila, F. 2015. Feasible Fitting of Linear Models with N Fixed Effects. *The Stata Journal* 15(3): 881–898. Publisher: SAGE Publications.
- Rios-Avila, F., L. Siles, and G. Canavire-Bacarreza. 2024. Estimating Quantile Regressions with Multiple Fixed Effects through Method of Moments. *IZA discussion Paper* .
- Vega Yon, G. G., and B. Quistorff. 2019. parallel: A command for parallel computing. *The Stata Journal* 19(3): 667–684. <https://doi.org/10.1177/1536867X19874242>.
- Wooldridge, J. M. 2010. *Econometric analysis of cross section and panel data*. MIT press.
- . 2019. Correlated Random Effects Models with Unbalanced Panels. *Journal of Econometrics* 211(1): 137–150.
- . 2021. Two-Way Fixed Effects, the Two-Way Mundlak Regression, and Difference-in-Differences Estimators. Working paper. Available at SSRN: <https://ssrn.com/abstract=3906345>.

About the authors

Fernando Rios-Avila is a Research Scholar at the Levy Economics Institute of Bard College. Leonardo Siles is a master student at the Universidad de Chile. Gustavo Canavire-Bacarreza is a Senior Economist at the World Bank.