

Quantile Regression Varying-Coefficient Model in The Case of Longitudinal Data: A Study on The Success Rates of Secondary Schools in Al-Diwaniyah Governorate \ Iraq

Ali Mohammed Farhan⁽¹⁾

Mohammed Sadiq Abdul Razzaq⁽²⁾

ali.mohammed1101a@coadec.uobaghdad.edu.iq dr_alldouri@coadec.uobaghdad.edu.iq

College of Administration and Economics / University of Baghdad / Iraq

Abstract

In this paper, we study the quantile regression varying-coefficient model in the case of longitudinal data. two methods are used to estimate the model: the traditional method, which ignores within-subject correlations, and the weighted method, which accounts for within-subject correlations by incorporating weights. the local polynomial method is used to estimate the nonparametric functions. Five quantile levels are examined. these methods are applied to data on the success rates of the third intermediate grade in Al-Diwaniyah Governorate. the study investigates the impact of four variables on the success rates of 337 middle and high schools over five years. the results show that the weighted estimation method is more efficient than the traditional method across all quantile levels.

Keywords:

Longitudinal data analysis, Varying-coefficient, quantile regression, the traditional local polynomial, weighted local polynomial.

1. Introduction

Regression analysis is undeniably one of the most important topics in statistics, focusing on studying and analyzing the relationship between a response variable and one or more explanatory variables. It has wide applications in various fields such as medicine, economics, and sociology. Regression analysis comes in various forms, one of which is

quantile regression proposed by (Koenker and Bassett, 1978). On the contrary ordinary regression, quantile regression measures the relationship between variables by estimating conditional quantile functions of the response variable. this provides researchers with insights into the relationship between variables across different conditional distributions, especially those at the beginning or end of the data. In recent years, researchers have shown significant interest in quantile regression for longitudinal data, highlighting the importance of longitudinal data for containing more information than time series or cross-sectional data. Longitudinal data integrates both types of data, providing richer information about the phenomenon under study. to estimate model parameters and analyze the relationship between variables, there are parametric, semi-parametric, and non-parametric methods. Parametric methods often make strict assumptions that may be challenging to apply in practice, whereas non-parametric methods have fewer assumptions. to increase flexibility, we use varying coefficient models, which allow the effect of explanatory variables on the response variable to vary based on the values of other variables, such as time or any other variable.

The problem addressed in this research revolves around the fact that some data used in regression analysis do not meet the assumption of normal distribution, or they may not satisfy the assumptions regarding random error or suffer from skewness. One of the challenges researcher's faces is that regression models may not accurately represent the true relationship between variables due to neglecting the dynamic nature of the data. Additionally, the problem of multicollinearity exacerbates the issue, as regression models struggle to capture the multidimensional nature of the data. to address these challenges, this study focuses on analyzing quantile regression varying coefficient with longitudinal data as a solution to the aforementioned problems.

This research aims to estimate the parameters of varying coefficient regression when the data is longitudinal data.

Many previous studies have discussed this topic, among which we will review: **(Wu and Chiang, 2000)** Proposed two types of kernel estimators based on the local two-stage least squares method to estimate time-varying coefficients for the varying coefficient regression model in the case of longitudinal response data and cross-sectional explanatory variables. Utilized simulation to study the properties of the proposed methods and applied

them to real data in the health domain, specifically investigating the effect of smoking on CD4 cells before and after HIV infection.

(Kim, 2007) Studied the varying coefficient regression model and presented a methodology for estimating the model using a Polynomial spline function. the estimation calculation relied on the standard partitioning algorithm. the experimental side was employed to test the method, and it was applied to forced expiratory volume data.

(Cai and Xu, 2008) Proposed a method for estimating quantile varying coefficient time series models and smoothing the model using both local liner and local constant. Additionally, suggested a method for selecting the bandwidth based on Akaike information criterion and demonstrated the effectiveness of the approaches through simulation. the methods were applied to Boston housing price data.

(Tang and Leng, 2011) Proposed a novel approach for estimating quantile regression with longitudinal data using the empirical likelihood function while considering within-subject correlation. demonstrated that this approach is more efficient than neglecting these correlations using simulation. the method was applied to medical data on the effect of smoking on CD4 cells.

(Kim and Yang, 2011) Adopted the varying-coefficient regression model with random effects for analyzing clustered data. proposed a semi-parametric approach using empirical likelihood for estimating the random effects parameters of the model and resorted to simulation to demonstrate the method's efficiency. the method was applied to real data from two phenomena: first, data from alcohol and drug addiction treatment centers, and second, data on choking incidents while swallowing food.

(Saif-alddin, 2012) Reviewed some nonparametric techniques for estimating time-varying coefficient functions in the context of the nonparametric varying-coefficient model for balanced longitudinal data. the techniques employed included local linear boundary regression and cubic smoothing spline techniques. the two-stage method was utilized to estimate the coefficient functions using the aforementioned techniques. Furthermore, the researcher suggested the adoption of some robust methods and employed simulation to verify the performance of both traditional and robust methods. these methods were then applied to economic sectors in Iraq.

(Rashed, 2014) Studied the varying coefficient model as well as the partial varying coefficient model. estimated both the varying coefficient model (VCM) and the partial varying coefficient model (PVCM) using nonparametric and semi-parametric estimation methods, respectively. Comparisons between these methods were conducted. the thesis included a proposal comprising general formulations for kernel functions. Simulation was employed to compare the nonparametric methods, and practical applications were demonstrated using stock closing prices and trading volumes based on the varying coefficient model and the partial varying coefficient model.

(Badr, 2016) Employed the nonparametric regression method to diagnose and estimate the longitudinal data model in cases where certain assumptions regarding the random error vector are not met, particularly in the problem of heteroscedasticity and autocorrelation problem. these problems can render the estimation process inaccurate or sometimes infeasible. the researcher calculated nonparametric estimators, addressing each problem separately, and formulated simulation experiments for the models used in this study to assess the performance of both traditional and proposed methods. Furthermore, the methods were practically applied to gross domestic product data in the state's general budget.

(Kim and Cho, 2018) Proposed two types of weights for estimating the varying-coefficient regression model in order to address within-subject correlations. The first type is global weight, which incorporates all observations in its calculation, while the second type is local weight, which considers nearby observations. evaluated the estimation method using simulation and found that employing weights enhances the efficiency of the estimators. the method was applied to real data of patients with acquired immunodeficiency syndrome (AIDS).

(Lin, Tang, and Zhu, 2020) Developed a weighted approach to enhance the efficiency of the varying-coefficient autoregressive model for longitudinal data. obtained the weights via empirical likelihood method, utilized spline method for obtaining smoothers, and employed the quadratic inference function for modeling the inverse conditional correlation matrix. Simulation results indicated that the weighted methods More efficiency from conventional methods. the approaches were applied to nursing home data in New Jersey.

2. Methodology

Quantile regression varying-coefficient model in the case of longitudinal data in the case of longitudinal data, the quantile varying-coefficient regression model takes the following form (Lin, Tang, and Zhu, 2020):

$$Q_{\tau}(Y_{ij} | \mathbf{Z}_{ij}) = \beta_{1,\tau}(T_{ij})x_{ij1} + \cdots + \beta_{p,\tau}(T_{ij})x_{ijp} + e_{ij}(\tau) \quad (1)$$

$$i = 1, 2, \dots, n$$

$$j = 1, 2, \dots, m_i$$

Where:

Y_{ij} : is the response variable.

$\mathbf{Z}_{ij} = (x_{ij1}, \dots, x_{ijp}, T_{ij})^T$: represent the explanatory variables observed with order j in the subject with order i

$\tau \in (0,1)$: denotes the quantile level.

$e_{ij}(\tau)$: represents the random error.

$\beta_{1,\tau}(T_{ij}), \dots, \beta_{p,\tau}(T_{ij})$: are the coefficients of the explanatory variables in the regression equation, which are functions of the variable $T_{ij} \in [0,1]$ (Rashed, 2014). In other words, T_{ij} changes the coefficients of x_{ij1}, \dots, x_{ijp} using unknown functions $\beta_{1,\tau}(\cdot), \dots, \beta_{p,\tau}(\cdot)$, and the dependence of $\beta_{p,\tau}(\cdot)$ on T_{ij} involves a specific type of interaction between T_{ij} and x_{ijp} . In some cases, the variable T_{ij} cannot be distinguished from x_{ijp} , while in others, T_{ij} may be a special variable such as time, as in this study.

2.1 The traditional local polynomial method for estimating VCQR models

Equation (1) can be rewritten as follows:

$$Q_{\tau}(Y_{ij} | \mathbf{Z}_i) = \mathbf{X}_{ij}^T \boldsymbol{\beta}_{\tau}(T_{ij}) + e_{ij}(\tau) \quad (2)$$

Where:

$$\mathbf{Z}_i = (\mathbf{Z}_{i1}, \dots, \mathbf{Z}_{im_i})^T \text{ and}$$

$\beta_{\tau}(t) = (\beta_{1,\tau}(t), \dots, \beta_{p,\tau}(t))^T$ is a vector of unknown smoothing functions, and the error vector is

$$e_{ij}(\tau) = Y_{ij} - Q_{\tau}(Y_{ij} | Z_i)$$

Under the assumption that the conditional quantile τ of $e_{ij}(\tau)$ given Z_i is equal to zero (Tang, and Leng, 2011) and that the observations Y_{ij} are independent, the errors $e_{ij}(\tau)$ are therefore independent. The estimated varying-coefficient regression model for longitudinal data can be expressed as follows in Equation (2) (Lin, Tang, and Zhu, 2020):

$$Q_{\tau}(\hat{Y}_{ij}/Z_i) = X_{ij}^T \beta_{\tau}(T_{ij}) \quad (3)$$

This method heavily relies on local polynomial due to their favorable properties such as statistical efficiency and model adaptability. Assuming that the coefficient function $\beta(\cdot)$ has $(K + 1)$ derivatives, where (K) represents the number of variables $K \geq 1$, it can be approximated by local polynomial around the grid point t_0 (Cai and Xu 2008).

$$\beta_{\tau}(T_{ij}) \cong \sum_{k=0}^P \alpha_k (T_{ij} - t_0)^k \quad (4)$$

Where

$$\alpha_k = \beta^k(t_0)/k!$$

Where $\beta^k(t_0)$ represents the (K) th derivative of $\beta(t_0)$. Equation (1) can be approximately written as:

$$Q_{\tau}(Y_{ij}/Z_i) = \sum_{k=0}^P X_{ij}^T \alpha_k (T_{ij} - t_0)^k \quad (5)$$

This leads to the local loss weighting function, which can be used to obtain model parameter estimates as follows:

$$\sum_{i=1}^n \sum_{j=1}^m \rho_{\tau} \{Y_{ij} - \sum_{k=0}^P X_{ij}^T \alpha_k (T_{ij} - t_0)^k\} K_h(T_{ij} - t_0) \quad (6)$$

Where $K(0)$ is the kernel function and $K_h(X) = (\frac{k(x/h)}{h})$ when $h = h_{nm}$ is a sequence in positive numbers approaching zero, controlling the amount of smoothing used in the estimation. By minimizing Equation (6), we obtain:

$$\hat{\beta}(t_0) = \hat{\alpha}_0$$

Which represents the local polynomial estimation of $\beta(t_0)$, with the remaining parameters as follows:

$$\hat{\beta}^k(t_0) = K! \hat{\beta}_k \quad (K \geq 1) \quad (7)$$

Which represents the local polynomial estimation of K derivatives of $\beta^k(t_0)$ (Saif-alddin, 2012).

Local polynomial estimations for nonparametric quantile regression models can be approximate, especially for small or large values of τ , because a small number of data points may be available. However, they are crucial in describing the local properties of nonparametric functions (Cai and Xu, 2008).

2.2 smooth parameter h Selection

As is well-known, the smoothing parameter (h) plays a crucial role in balancing bias and variance, the components from which the mean squared error is composed. Thus, the value of the mean squared error is influenced by the choice of this parameter. To select the smoothing parameter(h), we follow the following approach, where we rely on the nonparametric Akaike Information Criterion (AIC) (Cai and Xu, 2008), representing the corrected criterion for the bias of nonparametric regression models:

$$AIC(h) = \{\hat{\sigma}_\tau^2\} + \frac{2(ph + 1)}{[nm - (ph + 2)]} \quad (8)$$

Where:

$\hat{\sigma}_\tau^2$: is the estimated variance of the errors for quantile regression, calculated using the following formula:

$$\hat{\sigma}_\tau^2 = \sum_{i=1}^n \sum_{j=1}^m P_\tau \{Y_{ij} - X_{ij}^T \alpha_k(t_0)\} / nm \quad (9)$$

ph : represents the nonparametric degrees of freedom, also known as the effective number of parameters, which depend on the trace of the hat matrix in the linear estimations in nonparametric quantile regression. There is no explicit expression for the hat matrix due to its non-linearity. However, we can use a first-order approximation (the local Bahadur

representation) to derive an explicit expression, which can be interpreted as an approximation of the hat matrix.

$$(T_{ij})_h = \frac{(T_{ij} - t_0)}{h} \quad (10)$$

$$X_{ij}^* \begin{pmatrix} X_{ij} \\ (T_{ij})_h X_{ij} \end{pmatrix} \quad (11)$$

$$Y_{ij}^* = Y_{ij} - X_{ij}^T [\beta(t_0) + \beta^1(T_{ij} - t_0)] \quad (12)$$

Where:

$$H = \text{diag}[I_p, hI_d]$$

Where:

I_p is the $p * p$ identity matrix.

$$\theta = \sqrt{nm}H \begin{pmatrix} \alpha_0 - \beta(t_0) \\ \alpha_1 - \beta^1(t_0) \end{pmatrix} \quad (13)$$

$$S_n = S_n(t_0) = a_n \sum_{i=1}^n \sum_{j=1}^m \epsilon_{ij} X_{ij}^* X_{ij}^{*T} K(T_{ij})_h \quad (14)$$

Where:

$$a_n = (nmh)^{-\frac{1}{2}}$$

$$\epsilon_{ij} = I(Y_{ij} \leq X_{ij}^T \beta(t_0) + a_n) - I(Y_{ij} \leq X_{ij}^T \beta(t_0)) \quad (15)$$

Equation 15 can be rewritten as follows:

$$S_n(t_0) = f_t(t_0)\Omega_1^*(t_0) + o_p(1) \quad (16)$$

Where:

$f_t(t_0)$: represents the density function of (T) , and that

$$\Omega_1^*(t_0) = \text{diag}\{1, M_2\} \otimes \Omega^*(t_0)$$

Where:

$$M_2 = \int t^2 K(t) dt$$

$$\Omega^*(t_0) \equiv E \left[X_{ij} X_{ij}^T X \times f\left(\frac{y}{tx}\right)(y) \right] \quad (17)$$

Where:

$f\left(\frac{y}{tx}\right)(y)$ is the conditional density function of $f\left(\frac{y}{tx}\right)(y)$ given X, T .

From the above, $\hat{\theta}$ can be calculated through:

$$\hat{\theta} = a_n S_n^{-1} \sum_{i=1}^n \sum_{j=1}^m \delta_{\tau}(Y_{ij}^*) X_{ij} K(T_{ij})_h + o_p(1) \quad (18)$$

Where:

$$\delta_{\tau}(\cdot) = \tau - I(\cdot < 0)$$

This leads to:

$$\hat{Q}_{\tau}(Y_{ij}/Z_i) - Q_{\tau}(Y_{ij}/Z_i) =$$

$$1/nm \sum_{i=1}^n \sum_{j=1}^m \delta_{\tau}(Y_{ij}^*(T_{ij})) K_h((T_{ij} - T)/h) X_{ij}^{0T} S_n^{-1}(T_{ij}) X_{ij}^* + o_p(an) \quad (19)$$

Where:

$$X_{ij}^0 = \begin{pmatrix} X_{ij} \\ 0 \end{pmatrix}$$

And the coefficient $\delta_{\tau}(Y_{ij}^*(T_{ij}))$ on the right-hand side is:

$$\eta_{ij} = a_n^2 K(0) X_{ij}^{0T} S_n^{-1}(T_{ij}) X_{ij}^0 \quad (20)$$

Now we have:

$$Ph = \sum_{i=1}^n \sum_{j=1}^m \eta_{ij} \quad (21)$$

Which represents an approximation of the trace of the hat matrix. In practical applications to achieve this, we first need to estimate $\beta(t_0)$ because the estimation $S_n(t_0)$ involves $\beta(t_0)$. Therefore, a trial smoothing parameter is used in estimating $\beta(t_0)$.

2.3 The weighted local polynomial method for estimating VCQR models

The previous estimation method ignores within-subject correlations, while this weighted method takes these correlations into account. The fundamental idea behind this method is to utilize weights in estimating the VCQR model, which can be summarized as follows: the coefficients vector is estimated according to the formula (Lin, Tang, and Zhu, 2020):

$$\hat{\alpha}_\tau^{wp} = \underset{\alpha_\tau \in \mathbb{R}^{pN}}{\operatorname{argmin}} \sum_{i=1}^n \omega_i(\hat{\alpha}) \sum_{j=1}^{m_i} \rho_\tau(Y_{ij} - \Lambda_{ij}^T \alpha_\tau) \quad (22)$$

Where:

$$\Lambda_{ij} = \sum_{k=0}^P X_{ij}^T (T_{ij} - t_0)^k K_h(T_{ij} - t_0)$$

The weights are calculated through:

$$\omega_i(\hat{\alpha}) = n^{-1} \{1 + \lambda_{\hat{\alpha}}^T g_i(\hat{\alpha})\}^{-1} \quad (23)$$

$\lambda_{\hat{\alpha}}$ and $\hat{\alpha}$ are obtained by solving the following equations:

$$n^{-1} \sum_{i=1}^n \{\partial g_i(\hat{\alpha})^T / \partial \hat{\alpha}\} \lambda_{\hat{\alpha}} \{1 + \lambda_{\hat{\alpha}}^T g_i(\hat{\alpha})\}^{-1} = \mathbf{0} \quad (24)$$

and

$$n^{-1} \sum_{i=1}^n g_i(\hat{\alpha}) \{1 + \lambda_{\hat{\alpha}}^T g_i(\hat{\alpha})\}^{-1} = \mathbf{0} \quad (25)$$

Calculated through $(\alpha) g_i(\alpha)$ (Qu and Li, 2006)

$$g_i(\alpha) = \begin{pmatrix} \Lambda_i^T A_i^{-1/2} M_{i1} A_i^{-1/2} (Y_i^k - \Lambda_i \alpha) \\ \vdots \\ \Lambda_i^T A_i^{-1/2} M_{is} A_i^{-1/2} (Y_i^k - \Lambda_i \alpha) \end{pmatrix}$$

Where:

$$Y_i^k = Y_i k_h (T_{ij} - t_0)$$

$$A_i = \text{diag}(\sigma^2(t_{i1}), \dots, \sigma^2(t_{im}))$$

$$\Lambda_i = (\Lambda_{i1}, \dots, \Lambda_{im})^T$$

3. Real data

In this practical aspect of the research, the methods will be applied to the data of schools from the Directorate General of Education in Al-Diwaniyah Governorate. The dataset consists of information from 337 intermediate and secondary schools for the years 2018 to 2023, excluding the year 2020 where the ministry exams for the third intermediate grade were not conducted due to the COVID-19 pandemic.

The response variable Y_{ij} , representing the pass rate from the first round for the third intermediate grade for each school over five years, will be examined. As for the explanatory variables:

- 1- The first explanatory variable $x_{ij,1}$ is nominal, represents the type of school, whether governmental or private (0 for private schools and 1 for governmental schools).
- 2- The second explanatory variable $x_{ij,2}$, represents the ratio of the number of teachers to the number of students in the same school over five years.
- 3- The third explanatory variable $x_{ij,3}$, represents the average years of service for the school's teachers over five years.
- 4- The fourth variable $x_{ij,4}$ is nominal, indicating the gender of the students in the school (0 for boys' schools, 1 for mixed-gender schools, and 2 for girls' schools).

The data were collected and obtained from the Directorate General of Education in Al-Qadisiyah Governorate, and the model will take the following form

$$Q_\tau(Y_{ij} | \mathbf{Z}_{ij}) = \alpha_{1,\tau}(T_{ij})x_{ij,1} + \alpha_{2,\tau}(T_{ij})x_{ij,2} + \alpha_{3,\tau}(T_{ij})x_{ij,3} + \alpha_{4,\tau}(T_{ij})x_{ij,4}$$

Estimation was performed using two methods: the traditional method and the weighted method. Five different quantile levels were used for the estimation: (0.1, 0.3, 0.5, 0.7, 0.9). The smooth parameter h was selected using the roll in Section (2.2). Table (1) shows the

mean squared error for the model and for both estimation methods used in this study across the five different quantile levels.

Table 1 represents the values of MSE for different quantile level τ for the traditional local polynomial method (TLP) and weighted local polynomial method (WLP).

τ	TLP	WLP
0.1	30.7497	23.6098
0.3	21.4255	20.1122
0.5	19.8148	18.9135
0.7	22.3005	21.6039
0.9	29.2595	28.2496

From Table 1, we observe that the weighted method is more efficient than the traditional method across all quantile levels, with its efficiency increasing at both higher and lower quantile levels.

Figure 1 show the estimated varying coefficient curves over time for both methods and for five quantile levels. The horizontal axis in this figure 1 represents time. From these figures, we notice that all explanatory variables have a positive effect on the response variable. Additionally, we can observe that the effect of school type and the teacher-student ratio is at its highest at the 0.5 quantile level, and its impact is lower at both the higher and lower quantile levels when using the weighted method. However, with the traditional method, the effect of school type and the teacher-student ratio is at its highest at the 0.9 quantile level. On the other hand, the influence of the average years of service and school gender is similar for both methods across all quantile levels.

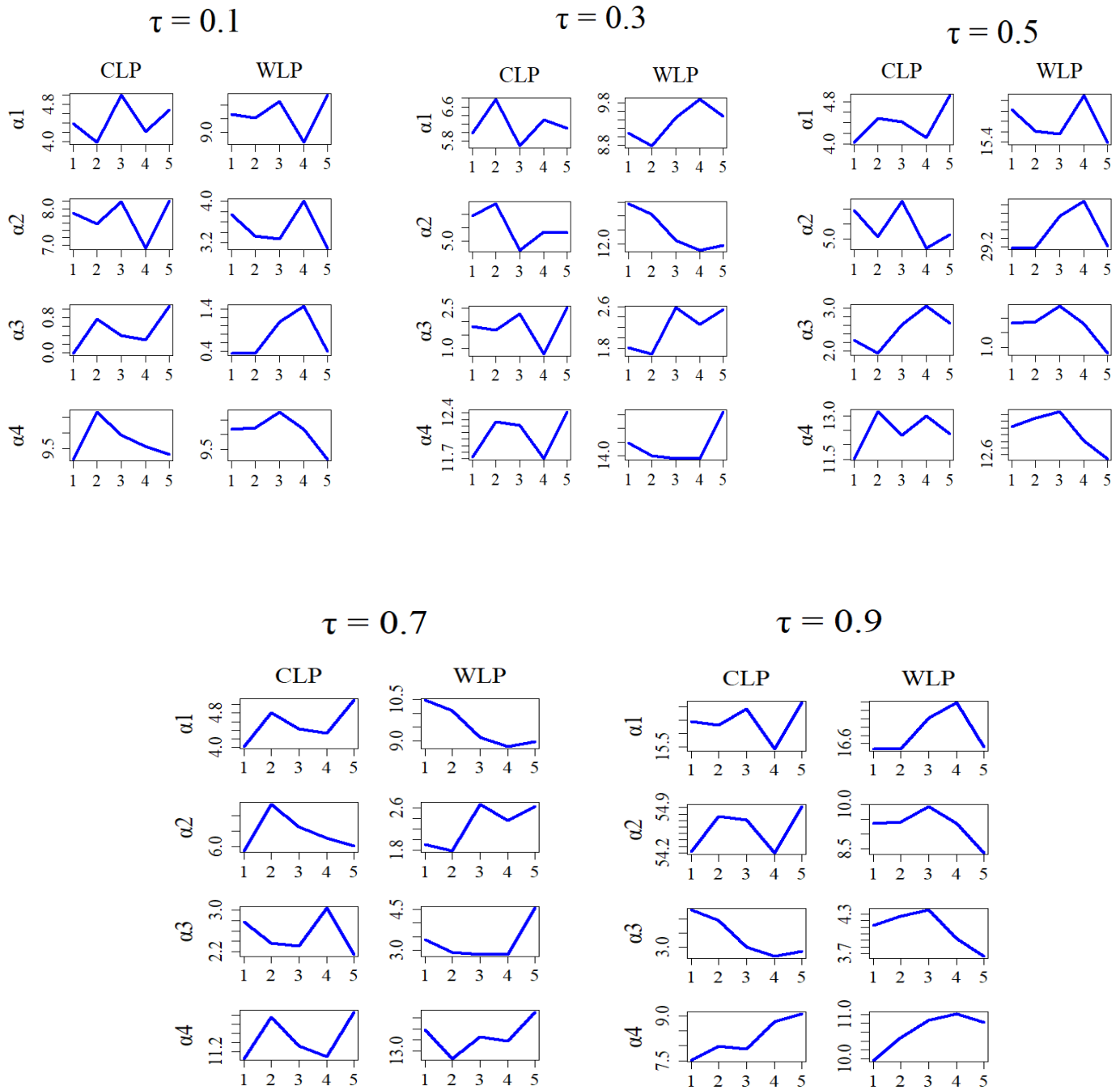


Figure 1: Estimated varying coefficient curves over time for the variables type of school, ratio of teachers to students, average years of service for teachers and sex of school at different quantile levels (0.1, 0.3, 0.5, 0.7, and 0.9)

4. Conclusions

At the end of this paper, the data analysis showed that the weighted local polynomial method for estimating the Varying Coefficient Quantile Regression (VCQR) model in the case of longitudinal data is more efficient than the traditional local polynomial method for

estimating the same model because the weighted method takes into account within-subject correlations. The results also indicated that the use of quantiles reveals the effect of explanatory variables on the response variable at different points of the distribution, providing a more comprehensive view of the impact of explanatory variables on success rates. Therefore, we hope that the findings of this paper will be used to increase success rates by improving the ratio of teachers to students, considering factors like the average years of service.

Acknowledgements

I would like to extend my gratitude to everyone who contributed to the completion of this research paper, and I also thank those who provided feedback to ensure that this paper is presented in its current form.

References

1. Kim, M. O., & Yang, Y. 2011. Semiparametric approach to a random effects quantile regression model. *Journal of the American Statistical Association*, 106(496), 1405-1417.
2. Badr, Duraïd Hussein 2016. The diagnosis estimation of the nonparametric regression function of the panel data in Case some of its hypotheses are not verified. the Ph. D to the College of Administration and Economics, University of Baghdad.
3. Cai, Z., & Xu, X. 2008. Nonparametric quantile estimations for dynamic smooth coefficient models. *Journal of the American Statistical Association*, 103(484), 1595-1608.
4. Kim, M. O. (2007). Quantile regression with varying coefficients.
5. Kim, S., & Cho, H. R. 2018. Efficient estimation in the partially linear quantile regression model for longitudinal data.
6. Koenker, R., & Bassett Jr, G. 1978. Regression quantiles. *Econometrica: journal of the Econometric Society*, 33-50.
7. Lin, F., Tang, Y., & Zhu, Z. 2020. Weighted quantile regression in varying-coefficient model with longitudinal data. *Computational Statistics & Data Analysis*, 145, 106915.
8. Qu, A., & Li, R. 2006. Quadratic inference functions for varying-coefficient models with longitudinal data. *Biometrics*, 62(2), 379-391.
9. Rashed, Husam A. 2014. Nonparametric smoothers for varying coefficient model and partial varying coefficient model. the Ph. D to the College of Administration and Economics, University of Baghdad.

10. Saifalddin, Ali. 2012. Nonparametric model estimation for longitudinal data of economic activities in Iraq. The Ph. D. to the College of Administration and Economics, University of Baghdad.
11. Tang, C. Y., & Leng, C. 2011. Empirical likelihood and quantile regression in longitudinal data analysis. *Biometrika*, 98(4), 1001-1006.
12. Wu, C. O., & Chiang, C. T. 2000. Kernel smoothing on varying coefficient models with longitudinal dependent variable. *Statistica Sinica*, 433-456.