

# Quantile Regression for Group Effect Analysis

Cristina Davino<sup>1</sup>    Domenico Vistocco<sup>2</sup>

<sup>1</sup>Dip.to di Studi sullo Sviluppo Economico  
Università di Macerata  
cdavino@unimc.it

<sup>2</sup>Dip.to di Scienze Economiche  
Università di Cassino  
vistocco@unicas.it

19<sup>th</sup> International Conference on Computational Statistics  
Paris, 22 – 27 August 2010



all computations and graphics were  
done in the R language using the  
packages *quantreg* and *ggplot2*

# Outline

- 1 Aim of the paper
- 2 QR for group effect analysis
  - Basic notation
  - The reference framework
  - The proposed approach
- 3 An empirical analysis
  - The dataset
  - Main results
- 4 Concluding remarks

# Outline

- 1 Aim of the paper
- 2 QR for group effect analysis
  - Basic notation
  - The reference framework
  - The proposed approach
- 3 An empirical analysis
  - The dataset
  - Main results
- 4 Concluding remarks

# Outline

- 1 Aim of the paper
- 2 QR for group effect analysis
  - Basic notation
  - The reference framework
  - The proposed approach
- 3 An empirical analysis
  - The dataset
  - Main results
- 4 Concluding remarks

# Outline

- 1 Aim of the paper
- 2 QR for group effect analysis
  - Basic notation
  - The reference framework
  - The proposed approach
- 3 An empirical analysis
  - The dataset
  - Main results
- 4 Concluding remarks

# Aim of the paper

## Identification of group effects in a quantile regression model

- 1 CONFIRMATIVE APPROACH
- 2 ROW-PARTITIONED DATA
  - Supervised approach
  - Unsupervised approach

# Aim of the paper

## Identification of group effects in a **quantile regression model**

- 1 **CONFIRMATIVE APPROACH**
- 2 **ROW-PARTITIONED DATA**
  - Supervised approach
  - Unsupervised approach

# Aim of the paper

## Identification of **group effects** in a quantile regression model

- 1 CONFIRMATIVE APPROACH
- 2 **ROW-PARTITIONED DATA**
  - Supervised approach
  - Unsupervised approach



# Aim of the paper

## Identification of group effects in a quantile regression model

- 1 CONFIRMATIVE APPROACH
- 2 ROW-PARTITIONED DATA
  - Supervised approach
  - Unsupervised approach

## Some solutions for group effect analysis

- Estimation of different models for each group
- Introduction of a dummy variable
- Multilevel modeling (Gelman and Hill, 2007)

# Basic notation

## The data structure

- $n$ : number of units
  - $p$ : number of regressors
  - $G$ : number of groups or levels
- 
- $\mathbf{X}_{[n \times p]}$ 
    - ${}_g X_{ij}$   
( $i=1, \dots, n; j=1, \dots, p; g=1, \dots, G$ )
  - $\mathbf{y}_{[n]}$ 
    - ${}_g y_i$   
( $i=1, \dots, n; g=1, \dots, G$ )
  - $n_g$ : number of units in group  $g$

# Classical vs quantile linear regression

## Classical linear regression (conditional expected value)

estimation of the conditional mean of a response variable ( $y$ ) distribution as a function of a set  $X$  of predictor variables

$$E(\mathbf{y} \mid \mathbf{X}) = \mathbf{X}\beta$$

$$\beta_i = \frac{\partial E(\mathbf{y})}{\partial x_i}$$

## Quantile regression (Koenker and Basset, 1978) (conditional quantiles)

estimation of the conditional quantiles of a response variable ( $y$ ) distribution as a function of a set  $X$  of predictor variables

$$Q_\theta(\mathbf{y} \mid \mathbf{X}) = \mathbf{X}\beta(\theta)$$

where:  $(0 < \theta < 1)$

$$\beta_i(\theta) = \frac{\partial Q_\theta(\mathbf{y})}{\partial x_i}$$

# Classical vs quantile linear regression

## Classical linear regression (conditional expected value)

estimation of the conditional mean of a response variable ( $y$ ) distribution as a function of a set  $X$  of predictor variables

$$E(\mathbf{y} \mid \mathbf{X}) = \mathbf{X}\beta$$

$$\beta_i = \frac{\partial E(\mathbf{y})}{\partial \mathbf{x}_i}$$

## Quantile regression (Koenker and Basset, 1978) (conditional quantiles)

estimation of the conditional quantiles of a response variable ( $y$ ) distribution as a function of a set  $X$  of predictor variables

$$Q_\theta(\mathbf{y} \mid \mathbf{X}) = \mathbf{X}\beta(\theta)$$

where: ( $0 < \theta < 1$ )

$$\beta_i(\theta) = \frac{\partial Q_\theta(\mathbf{y})}{\partial \mathbf{x}_i}$$

# Classical vs quantile linear regression

## Classical linear regression (conditional expected value)

estimation of the conditional mean of a response variable ( $y$ ) distribution as a function of a set  $X$  of predictor variables

$$E(\mathbf{y} \mid \mathbf{X}) = \mathbf{X}\beta$$

$$\beta_i = \frac{\partial E(\mathbf{y})}{\partial \mathbf{x}_i}$$

## Quantile regression (Koenker and Basset, 1978) (conditional quantiles)

estimation of the conditional quantiles of a response variable ( $y$ ) distribution as a function of a set  $X$  of predictor variables

$$Q_\theta(\mathbf{y} \mid \mathbf{X}) = \mathbf{X}\beta(\theta)$$

where: ( $0 < \theta < 1$ )

$$\beta_i(\theta) = \frac{\partial Q_\theta(\mathbf{y})}{\partial \mathbf{x}_i}$$

# Classical vs quantile linear regression

## Classical linear regression (conditional expected value)

estimation of the conditional mean of a response variable ( $y$ ) distribution as a function of a set  $X$  of predictor variables

$$E(\mathbf{y} \mid \mathbf{X}) = \mathbf{X}\beta$$

$${}_g\mathbf{y} = {}_g\mathbf{X}_g\beta + {}_g\mathbf{e}$$

## Quantile regression (conditional quantiles)

estimation of the conditional quantiles of a response variable ( $y$ ) distribution as a function of a set  $X$  of predictor variables

$$Q_\theta(\mathbf{y} \mid \mathbf{X}) = \mathbf{X}\beta(\theta)$$

where:  $(0 < \theta < 1)$

$$Q^\theta({}_g\mathbf{y} \mid {}_g\mathbf{X}) = {}_g\mathbf{X}_g\beta(\theta)$$

# Classical vs quantile linear regression

## Classical linear regression (conditional expected value)

estimation of the conditional mean of a response variable ( $y$ ) distribution as a function of a set  $X$  of predictor variables

$$E(\mathbf{y} | \mathbf{X}) = \mathbf{X}\beta$$

$${}_g\mathbf{y} = {}_g\mathbf{X}_g\beta + {}_g\mathbf{e}$$

## Quantile regression (conditional quantiles)

estimation of the conditional quantiles of a response variable ( $y$ ) distribution as a function of a set  $X$  of predictor variables

$$Q_\theta(\mathbf{y} | \mathbf{X}) = \mathbf{X}\beta(\theta)$$

where:  $(0 < \theta < 1)$

$$Q^\theta({}_g\mathbf{y}|{}_g\mathbf{X}) = {}_g\mathbf{X}_g\beta(\theta)$$

# The proposed approach

## 1 Global estimation

$$Q^\theta(\mathbf{y}|\mathbf{X}) = \mathbf{X}\hat{\mathbf{B}}(\theta)$$

## 2 Identification of the best model for each unit

### 1 density estimation

$$\hat{\mathbf{Y}} = \mathbf{X}\hat{\mathbf{B}}(\theta)$$

### 2 best model identification

$$\theta_i : \underset{\theta=1, \Theta}{\operatorname{argmin}} y_i - \hat{y}_i(\theta)$$

### 3 best density estimation vector

$$\hat{\mathbf{y}}_\theta^{best}$$

## 3 Identification of the best model for each group

$$g\theta^{best}, \text{ for } g = 1, G$$

## 4 Partial estimation

$$Q^\theta(\mathbf{y}|\mathbf{X}) = \mathbf{X}\hat{\mathbf{B}}(\theta)^{best}$$



# The proposed approach

## 1 Global estimation

$$Q^\theta(\mathbf{y}|\mathbf{X}) = \mathbf{X}\hat{\mathbf{B}}(\theta)$$

## 2 Identification of the best model for each unit

### 1 density estimation

$$\hat{\mathbf{Y}} = \mathbf{X}\hat{\mathbf{B}}(\theta)$$

### 2 best model identification

$$\theta_i : \underset{\theta=1, \Theta}{\operatorname{argmin}} y_i - \hat{y}_i(\theta)$$

### 3 best density estimation vector

$$\hat{\mathbf{y}}_\theta^{best}$$

## 3 Identification of the best model for each group

$$g\theta^{best}, \text{ for } g = 1, G$$

## 4 Partial estimation

$$Q^\theta(\mathbf{y}|\mathbf{X}) = \mathbf{X}\hat{\mathbf{B}}(\theta)^{best}$$

# The proposed approach

## 1 Global estimation

$$Q^\theta(\mathbf{y}|\mathbf{X}) = \mathbf{X}\hat{\mathbf{B}}(\theta)$$

## 2 Identification of the best model for each unit

### 1 density estimation

$$\hat{\mathbf{Y}} = \mathbf{X}\hat{\mathbf{B}}(\theta)$$

### 2 best model identification

$$\theta_i : \underset{\theta=1, \Theta}{\operatorname{argmin}} y_i - \hat{y}_i(\theta)$$

### 3 best density estimation vector

$$\hat{\mathbf{y}}_\theta^{best}$$

## 3 Identification of the best model for each group

$$g\theta^{best}, \text{ for } g = 1, G$$

## 4 Partial estimation

$$Q^\theta(\mathbf{y}|\mathbf{X}) = \mathbf{X}\hat{\mathbf{B}}(\theta)^{best}$$

# The proposed approach

## 1 Global estimation

$$Q^\theta(\mathbf{y}|\mathbf{X}) = \mathbf{X}\hat{\mathbf{B}}(\theta)$$

## 2 Identification of the best model for each unit

### 1 density estimation

$$\hat{\mathbf{Y}} = \mathbf{X}\hat{\mathbf{B}}(\theta)$$

### 2 best model identification

$$\theta_i : \underset{\theta=1, \Theta}{\operatorname{argmin}} y_i - \hat{y}_i(\theta)$$

### 3 best density estimation vector

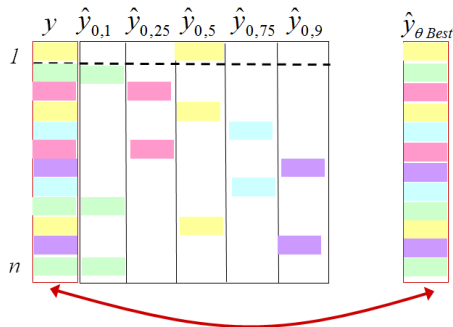
$$\hat{\mathbf{y}}_\theta^{\text{best}}$$

## 3 Identification of the best model for each group

$$g^{\theta^{\text{best}}}, \text{ for } g = 1, G$$

## 4 Partial estimation

$$Q^\theta(\mathbf{y}|\mathbf{X}) = \mathbf{X}\hat{\mathbf{B}}(\theta)^{\text{best}}$$



# The proposed approach

## 1 Global estimation

$$Q^\theta(\mathbf{y}|\mathbf{X}) = \mathbf{X}\hat{\mathbf{B}}(\theta)$$

## 2 Identification of the best model for each unit

### 1 density estimation

$$\hat{\mathbf{Y}} = \mathbf{X}\hat{\mathbf{B}}(\theta)$$

### 2 best model identification

$$\theta_i : \underset{\theta \in \Theta}{\operatorname{argmin}} y_i - \hat{y}_i(\theta)$$

### 3 best density estimation vector

$$\hat{\mathbf{y}}_\theta^{\text{best}}$$

## 3 Identification of the best model for each group

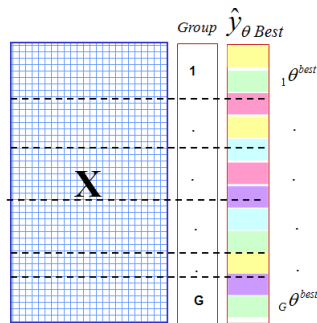
$$g\theta^{\text{best}}, \text{ for } g = 1, G$$

## 4 Partial estimation

$$Q^\theta(\mathbf{y}|\mathbf{X}) = \mathbf{X}\hat{\mathbf{B}}(\theta)^{\text{best}}$$

# The proposed approach

- 1 Global estimation  
 $Q^\theta(\mathbf{y}|\mathbf{X}) = \mathbf{X}\hat{\mathbf{B}}(\theta)$
- 2 Identification of the best model for each unit
  - 1 density estimation  
 $\hat{\mathbf{Y}} = \mathbf{X}\hat{\mathbf{B}}(\theta)$
  - 2 best model identification  
 $\theta_i : \underset{\theta=1, \Theta}{\operatorname{argmin}} y_i - \hat{y}_i(\theta)$
  - 3 best density estimation vector  
 $\hat{\mathbf{y}}_\theta^{best}$
- 3 Identification of the best model for each group  
 $g\theta^{best}$ , for  $g = 1, G$
- 4 Partial estimation  
 $Q^\theta(\mathbf{y}|\mathbf{X}) = \mathbf{X}\hat{\mathbf{B}}(\theta)^{best}$



# The proposed approach

## 1 Global estimation

$$Q^\theta(\mathbf{y}|\mathbf{X}) = \mathbf{X}\hat{\mathbf{B}}(\theta)$$

## 2 Identification of the best model for each unit

### 1 density estimation

$$\hat{\mathbf{Y}} = \mathbf{X}\hat{\mathbf{B}}(\theta)$$

### 2 best model identification

$$\theta_i : \underset{\theta \in \Theta}{\operatorname{argmin}} y_i - \hat{y}_i(\theta)$$

### 3 best density estimation vector

$$\hat{\mathbf{y}}_\theta^{\text{best}}$$

## 3 Identification of the best model for each group

$$g\theta^{\text{best}}, \text{ for } g = 1, G$$

## 4 Partial estimation

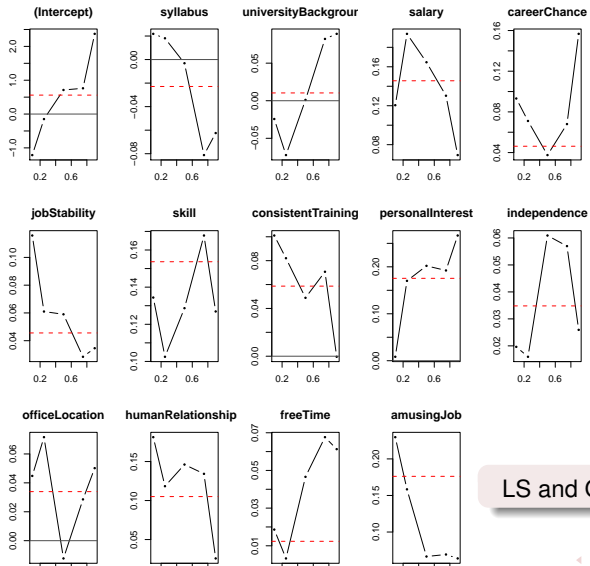
$$Q^\theta(\mathbf{y}|\mathbf{X}) = \mathbf{X}\hat{\mathbf{B}}(\theta)^{\text{best}}$$

# The dataset

## The evaluation of job satisfaction

- $n$ : random sample of 400 students graduated at University of Macerata and in a working condition at the time of the interview
- $p$ : 13 regressors (judgments of the different aspects related to the working experience)  
*syllabus, University background, consistent training, career chance, skill, personal interest, free time, salary, office location, job stability, human relationships, amusing job, independence*
- **dependent variable**: overall opinion on the job
- $G$ : 3 groups corresponding to the type of job  
*self-employed, private employee, public employee*

# Step 1: Global estimation



LS and QR coefficients

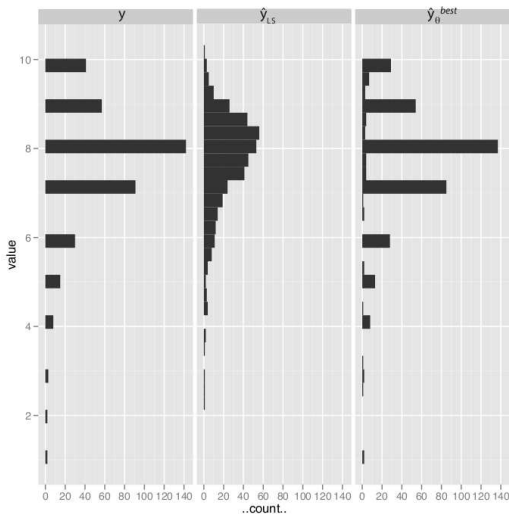


# Step 1: Global estimation

Variable	LS	$\theta=0.1$	$\theta=0.25$	$\theta=0.5$	$\theta=0.75$	$\theta=0.9$
Intercept	0.403	<b>-1.211</b>	-0.149	0.711	0.761	<b>2.370</b>
syllabus	-0.009	0.022	0.018	-0.003	-0.081	-0.062
University background	0.004	-0.024	-0.072	0.001	0.082	0.089
salary	<b>0.146</b>	<b>0.120</b>	<b>0.194</b>	<b>0.165</b>	<b>0.130</b>	0.069
career chance	<b>0.078</b>	0.093	0.071	0.037	0.068	<b>0.157</b>
job stability	<b>0.061</b>	<b>0.116</b>	0.061	0.059	0.028	0.035
skill	<b>0.117</b>	0.134	0.102	<b>0.129</b>	<b>0.168</b>	0.127
consistent training	0.043	0.101	0.082	0.049	0.070	-0.000
personal interest	<b>0.187</b>	0.008	<b>0.170</b>	<b>0.202</b>	<b>0.192</b>	<b>0.267</b>
independence	0.051	0.019	0.016	0.061	0.056	0.026
office location	0.031	0.044	0.072	-0.012	0.029	0.050
human relationships	0.126	<b>0.181</b>	0.118	<b>0.146</b>	<b>0.134</b>	0.026
free time	0.017	0.189	0.003	0.047	<b>0.067</b>	0.061
amusing job	<b>0.147</b>	<b>0.230</b>	<b>0.158</b>	0.066	0.069	0.064

LS and QR coefficients

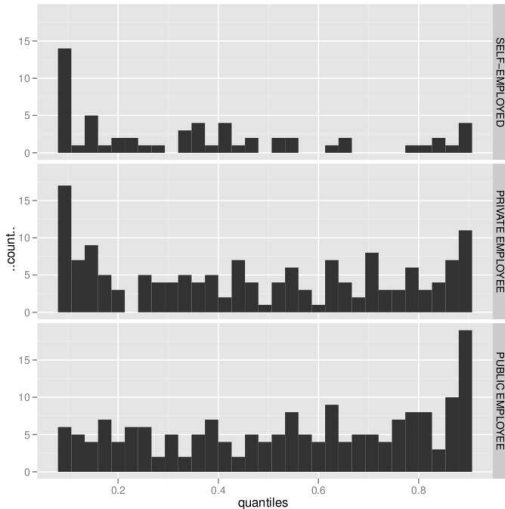
# Step 2: Identification of the best model for each unit



Distribution of the:

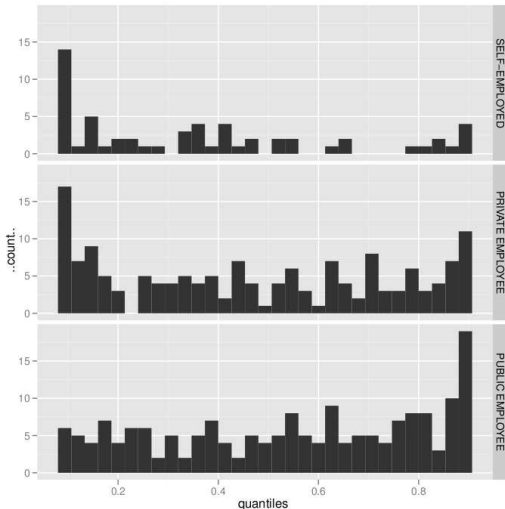
- dependent variable (*left panel*)
- LS estimated dependent variable (*middle panel*)
- best QR estimated dependent variable (*right panel*)

# Step 3: Identification of the best model for each group



Distribution of the “best” quantiles assigned to each unit grouped according to the type of job

# Step 3: Identification of the best model for each group



“Best” quantiles for each group:

Mean value of the “best” quantiles assigned to units belonging to the  $g^{th}$  group

- $\theta_1^{best} = 0.371$

- $\theta_2^{best} = 0.474$

- $\theta_3^{best} = 0.548$

## Step 4: Partial estimation

Variable	self-employed	private employee	public employee
intercept	0.646	0.683	0.694
syllabus	-0.007	-0.012	-0.035
University background	-0.030	0.006	0.026
salary	<b>0.201</b>	<b>0.152</b>	<b>0.160</b>
career chance	0.012	0.037	-0.008
job stability	0.049	0.034	0.054
skill	<b>0.118</b>	<b>0.156</b>	<b>0.184</b>
consistent training	0.065	0.066	0.064
personal interest	<b>0.200</b>	<b>0.175</b>	<b>0.202</b>
independence	0.022	0.035	0.035
office location	0.011	-0.006	0.007
human relationships	<b>0.114</b>	<b>0.152</b>	<b>0.107</b>
free time	0.018	0.032	0.026
amusing job	<b>0.148</b>	<b>0.124</b>	<b>0.141</b>

QR coefficients with group effects

# Concluding remarks and further issues

## The proposed approach

- Group effect analysis
- Impact of the regressors on the entire conditional distribution
- Semi-parametric approach

## Further developments

- Robust index for the identification of the “best” quantile
- Statistical significance of the differences among the “best” quantiles
- Time as grouping variable
- Unsupervised approach

# Concluding remarks and further issues

## The proposed approach

- Group effect analysis
- Impact of the regressors on the entire conditional distribution
- Semi-parametric approach

## Further developments

- Robust index for the identification of the “best” quantile
- Statistical significance of the differences among the “best” quantiles
- Time as grouping variable
- Unsupervised approach

# Main references



DAVINO, C., VISTOCCO, D. (2007): The evaluation of university educational processes: a quantile regression approach. *STATISTICA*, n.3, pp. 267-278.



DAVINO C., VISTOCCO D (2008): Quantile regression for the evaluation of student satisfaction. *STATISTICA APPLICATA*, vol. 20; p. 179-196.



EIDE, E. SHOWALTER, M.H. (1998): The effect of school quality on student performance: a quantile regression approach. *Economics Letters* 58, 345-350.



FURNO, M. (2010): Quantile regression analysis of the Italian school system. *Statistical Modelling*, vol. 4, 2010, in press.



GELMAN, A. HILL, J. (2006): *Data analysis using regression and multilevel/hierarchical models*. Cambridge University Press.



HAO, L. NAIMAN, D. Q. (2007): *Quantile Regression*, Series: Quantitative Applications in the Social Sciences, SAGE Publications.



LOCKHEED, M.E. HANUSHECK, E.R.(1994): Concepts of Educational Efficiency and Effectiveness, in Torsten Husén and T. Neville Postlethwaite (ed.), *International Encyclopedia of Education*, second edition, Volume 3 (Oxford: Pergamon, 1994), pp. 1779-1784.



KOENKER, R., BASSET, G.W. (1978): Regression Quantiles, *Econometrica* 46, 33-50.



KOENKER, R. (2005): *Quantile Regression*. Econometric Society Monographs.



KOENKER, R. (2009): *quantreg: Quantile Regression*. R package version 4.44. <http://CRAN.R-project.org/package=quantreg>.