

分位数回归

相耐汀

2019.11.5

主要内容

- 为什么需要分位数回归
- 7.1分位数回归模型
 - 7.1.1删失分位数回归
 - 7.1.2分位数回归的近似性质
 - 7.1.3微妙之处
- 7.2对分位数处理效应的工具变量估计
 - 7.2.1分位数处理效应估计值

为什么需要分位数回归

为什么需要分位数回归

- 分位数回归是给定回归变量 X ，估计响应变量 Y 条件分位数的一个基本方法。
- 分位回归优点：
 - 能够更全面的描述被解释变量条件分布的全貌，而不仅仅分析被解释变量的均值。
 - 中位数回归与最小二乘法相比，估计结果对离群值表现得更加稳健。
- 举例：工资分布的变化

7.1 分位数回归模型

条件分位数函数

- 假定连续分布随机变量 Y_i 密度函数性质良好，那么给定回归元所在向量 X_i ，在分位数 τ 处的条件分位数函数可定义为：

$$Q_{\tau}(Y_i|X_i) = F_y^{-1}(\tau|X_i)$$

- 其中 $F_y(y|X_i)$ 是在给定 X_i 时 Y_i 在 Y 处的分布函数。

条件分位数函数

- 条件分布函数实际上是一类特殊的条件期望函数
- 条件期望函数：

$$E[Y_i|X_i] = \underset{m(X_i)}{\operatorname{argmin}} E[(Y_i - m(X_i))^2]$$

条件期望函数就是给定 X_i 后，对 Y_i 的最小均方误预测

- 条件分位函数也是最小化问题的解：

$$Q_\tau(Y_i|X_i) = \underset{q(X_i)}{\operatorname{argmin}} E[\rho_\tau(Y_i - q(X_i))]$$

其中 $\rho_\tau(u) = (\tau - 1(u \leq 0))u$ 被称为“校验函数”

条件分位数函数

■ 当 $\tau=0.5$ 时

- $\rho_{0.5}(u) = \frac{1}{2}(\text{sign } u)u = \frac{1}{2}|u|$
- 最小化问题转化为离差绝对值最小值

$$Q_{\tau}(Y_i|X_i) = \underset{q(X_i)}{\operatorname{argmin}} E\left[\frac{1}{2}|Y_i - q(X_i)|\right]$$

■ 其他情况下

- $\rho_{\tau}(u) = 1(u > 0) * \tau|u| + 1(u \leq 0) * (1 - \tau)|u|$
即

$$\rho_{\tau}(u) = \begin{cases} (\tau - 1)u, & u \leq 0 \\ \tau u, & u > 0 \end{cases}$$

- 通过不对称加权可以产生一个最小化元，从而计算出条件分位数

分位数回归

■ 分位数回归模型

$$Q_{\tau}(Y_i|X_i) = X_i' \beta_{\tau}$$

其中

$$\beta_{\tau} = \underset{b}{\operatorname{argmin}} E[\rho_{\tau}(Y_i - X_i' b)]$$

教育水平对工资分布影响

- 以工资分布的研究为例来探讨分位数回归的使用
- 劳动经济学家感兴趣的问题：给定类似教育水平和工作经验等协变量后，工资不平等是如何变化的。
- 在20世纪80和90年代，在不同教育水平组成的群体之间的总体收入差距有显著提升（比如接受高等教育的工资溢价），但是我们并不清楚在具有相同教育水平和工作经验的群体内部，工资水平是如何变化的。
- 以下用分位数回归的方法证明近年来群体内部工资不平等的扩大。

教育水平对工资分布影响

Table 7.1.1: Quantile regression coefficients for schooling in the 1970, 1980, and 2000 Censuses

Census	Obs.	Desc. Stats.		Quantile Regression Estimates					OLS Estimates	
		Mean	SD	0.1	0.25	0.5	0.75	0.9	Coeff.	Root MSE
1980	65023	6.4	0.67	.074 (.002)	.074 (.001)	.068 (.001)	.070 (.001)	.079 (.001)	.072 (.001)	0.63
1990	86785	6.46	0.06	.112 (.003)	.110 (.001)	.106 (.001)	.111 (.001)	.137 (.003)	.114 (.001)	0.64
2000	97397	6.5	0.75	.092 (.002)	.105 (.001)	.111 (.001)	.120 (.001)	.157 (.004)	.114 (.001)	0.69

Notes: Adapted from Angrist, Chernozhukov, and Fernandez-Val (2006). The tables reports quantile regression estimates of the returns to schooling, with OLS estimates shown at the right for comparison. The sample includes US-born white and black men aged 40-49. Standard errors are reported in parentheses. All models control for race and potential experience. Sampling weights were used for the 2000 Census estimates.

教育水平对工资分布影响

■ 如果教育水平对工资的影响类似以一种位移

□ 用经典线性回归模型来描述

$$Y_i \sim N(X_i' \beta, \sigma_\varepsilon^2)$$

□ 则如果不考虑变化的截距项，在每个分位数水平上，分位数回归系数应相同

$$P[Y_i - X_i' \beta < \sigma_\varepsilon \Phi^{-1}(\tau) | X_i] = \tau$$

■ 如果教育水平对工资的影响不是位移

□ 用经典线性回归模型来描述

$$Y_i \sim N(X_i' \beta, \sigma^2(X_i))$$

□ 令 $\sigma^2(X_i) = (\lambda' X_i)^2$

$$P[Y_i - X_i' \beta < (\lambda' X_i)^2 \Phi^{-1}(\tau) | X_i] = \tau$$

则在各个分位数水平上 $\beta_\tau = \beta + \lambda \Phi^{-1}(\tau)$ 分位数回归系数变化

7.1.1 删失分位数回归

删失数据定义

- 分位数回归允许我们在 Y_i 的分布一段被隐藏时仍可以考察 Y_i 的条件分布特征。
- 假设所得数据形式如下

$$Y_{i,obs} = Y_i * 1[Y_i < c] + c * 1[Y_i \geq c]$$

其中 $Y_{i,obs}$ 是看到的观测值， Y_i 是本应该看到的数据。

变量 $Y_{i,obs}$ 是 Y_i 的删失数据。

处理删失数据

- 由于我们不知道哪个条件分位数正处于删失点之下，所以（以高收入者收入被加密为例）考虑下式：

$$Q_{\tau}(Y_i|X_i) = \min(c, X_i'\beta_{\tau}^c)$$

- 参数 β_{τ}^c 为下式的解

$$\beta_{\tau}^c = \underset{b}{\operatorname{argmin}} E\{1[X_i'b < c] \cdot \rho_{\tau}(Y_i - X_i'b)\}$$

- 以上方法也就是针对 $X_i'\beta_{\tau}^c < c$ 求解分位数回归的最小化问题，只要有足够多的未删失数据，我们求得的估计值就给出了原本在不存在删失数据时才能求得的分位数回归函数。

7.1.2分位数回归的近似性质

分位数回归的近似性质

- 我们可以将分位数回归理解为一种最小均方误差意义下的对条件分位函数的线性近似。

- 定义分位数回归的设定偏误：

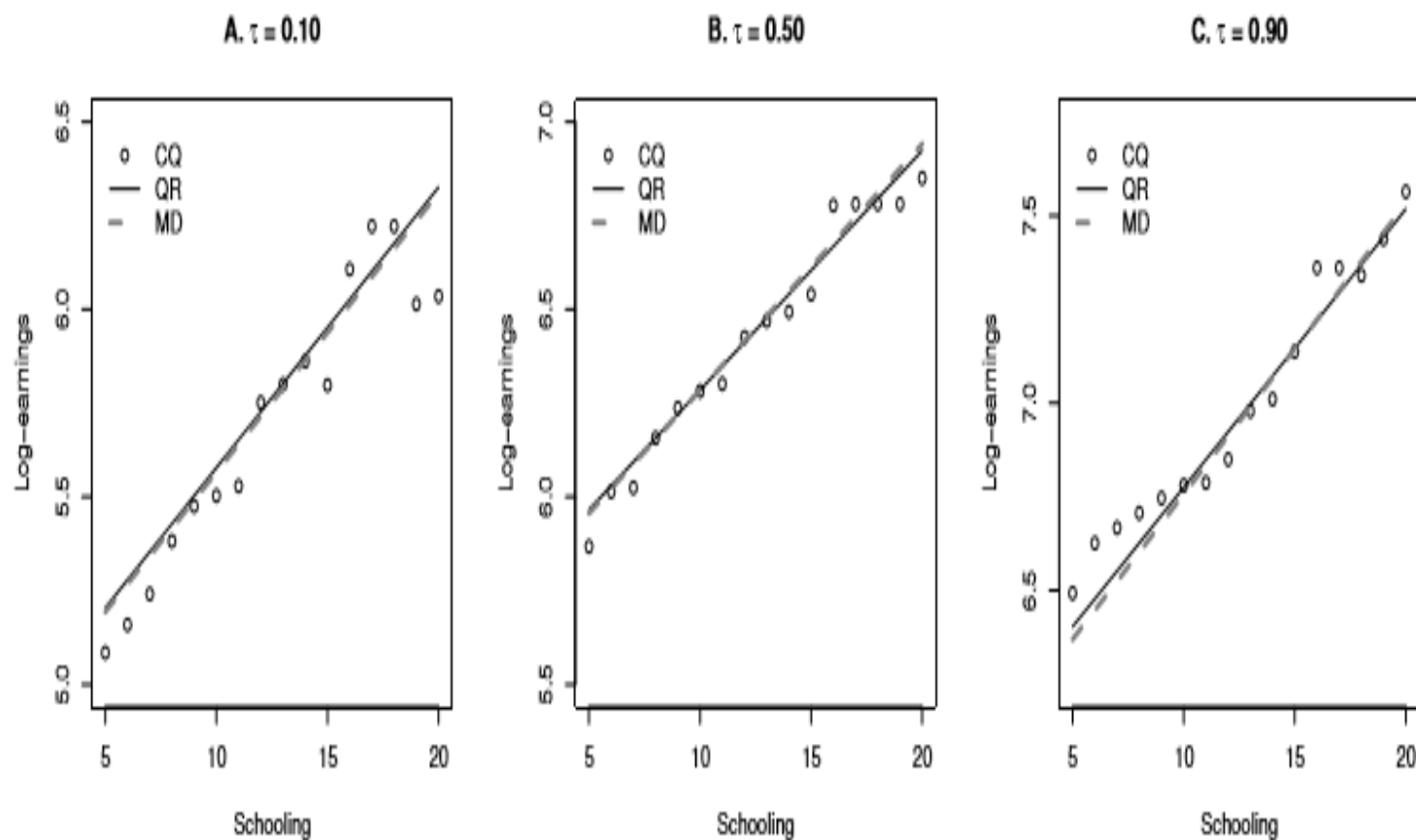
$$\Delta_{\tau}(X_i, \beta_{\tau}) = X_i' \beta - Q_{\tau}(Y_i | X_i)$$

- Angrist, Chernozhukov, and Fernandez-Val (2006) 指出：通过最小化模型设定偏误的加权平均值，我们可以得到总体分位数回归向量。

分位数回归的近似性质

- 可以用下述图片来阐明分位数回归的近似性质。
- 使用1980年的人口普查数据，给定最高完成学位水平下，绘出对数工资的条件分位数函数。
- 图A-C分别在0.1、0.5、0.9三个分位数水平下绘出了对 $Q_\tau(Y_i|X_i)$ 的线性分位数回归结果。
- 图中圆圈表示针对不同教育水平的个体分别估计出的条件分位数函数；分位数回归曲线则是实线。

分位数回归的近似性质



7.1.3微妙之处

微妙之处

- 分位数回归的系数告诉我们的是对分布的影响，而不是对个体的影响。
- 条件分位数和边际分位数的转变。

求解边际分位数

- 假设条件分位数函数是线性的, 即 $Q_\tau(Y_i|X_i) = X_i'\beta_\tau$
- 令 $F_y(y|X_i) = P[Y_i < y|X_i]$ 是在给定 X_i 时 Y_i 在 Y 处的条件累计密度函数, 其边际分布为 $F_y(y) = P[Y_i < y]$.
- 累计密度函数及其反函数关系为:

$$\int_0^1 1[F_y^{-1}(\tau|X_i) < y] d\tau = F_y(y|X_i)$$

-
- 用线性模型替换积分括号里面的条件分位函数，得

$$F_y(y|X_i) = \int_0^1 1[X_i'\beta_\tau < y]d\tau$$

- 用迭代期望率来求解边际分布函数 $F_y(y)$

$$F_y(y|X_i) = E\left[\int_0^1 1[X_i'\beta_\tau < y]d\tau\right]$$

- 对函数 $F_y(y)$ 进行转置，求得边际分位数 $Q_\tau(Y_i)$

$$Q_\tau(Y_i) = \inf\{y: F_y(y) \geq \tau\}$$

7.2对分位数处理效应的工具变量估计

-
- 用工具变量法来解决遗漏变量偏误问题
 - Abadie, Angrist, 和 Imbens (2002) 给出了使用工具变量估计分位数处理效应 (QTE) 的方法。
 - 该方法的假设和局部平均处理效应框架下的假设相同。
 - 对于 $\tau \in (0,1)$, 假设存在 α_τ 和 β_τ 满足:

$$Q_\tau(Y_i|X_i, D_i, D_{1i} > D_{0i}) = \alpha_\tau D_i + X_i' \beta_\tau$$

- 则参数 α_τ 给出了给定 X_i 后依从工具变量者在 Y_{1i} 和 Y_{0i} 上的分位数之差。

$$\alpha_\tau = Q_\tau(Y_{1i}|X_i, D_{1i} > D_{0i}) - Q_\tau(Y_{0i}|X_i, D_{1i} > D_{0i})$$

7.2.1 分位数处理效应估计值

分位数处理效应估计值

■ 通过对依从工具变量者总体进行分位数回归，我们可以估计出依从工具变量者的分位数回归系数。

■
$$(\alpha_\tau, \beta_\tau) = \underset{a,b}{\operatorname{argmin}} E\{\rho_\tau(Y_i - aD_i + X_i' b) | D_{1i} > D_{0i}\}$$
$$= \underset{a,b}{\operatorname{argmin}} E\{k_i \rho_\tau(Y_i - aD_i + X_i' b)\}$$

■ 其中,

$$k_i = 1 - \frac{D_i(1 - Z_i)}{1 - P(Z_i = 1 | X_i)} - \frac{(1 - D_i)Z_i}{P(Z_i = 1 | X_i)}$$

估计培训为受训者收入分位数影响

- 工作培训合法法案JTPA是一个大规模的联邦项目，他为20世纪80年代事业的美国工人提供资助性培训。
- 我们使用5102个成年按行来做培训对收入的影响的最小二乘回归、分位数回归、两阶段最小二乘回归和分位数处理效应估计值。
- JTPA实验中：
 - 控制组：没有人获得JTPA服务
 - 处理组：被分配到JTPA服务但只有60%的人接受。

估计培训为受训者收入分位数影响

A. OLS and Quantile Regression Estimates						
	OLS	Quantile				
		0.15	0.25	0.50	0.75	0.85
Training	3,754 (536)	1,187 (205)	2,510 (356)	4,420 (651)	4,678 (937)	4,806 (1,055)
% Impact of Training	21.20	135.56	75.20	34.50	17.24	13.43
High school or GED	4,015 (571)	339 (186)	1,280 (305)	3,665 (618)	6,045 (1,029)	6,224 (1,170)
Black	-2,354 (626)	-134 (194)	-500 (324)	-2,084 (684)	-3,576 (1,087)	-3,609 (1,331)
Hispanic	251 (883)	91 (315)	278 (512)	925 (1,066)	-877 (1,769)	-85 (2,047)
Married	6,546 (629)	587 (222)	1,964 (427)	7,113 (839)	10,073 (1,046)	11,062 (1,093)
Worked less than 13 weeks in past year	-6,582 (566)	-1,090 (190)	-3,097 (339)	-7,610 (665)	-9,834 (1,000)	-9,951 (1,099)
Constant	9,811 (1,541)	-216 (468)	365 (765)	6,110 (1,403)	14,874 (2,134)	21,527 (3,896)

估计培训为受训者收入分位数影响

B. 2SLS and QTE Estimates

	2SLS	Quantile				
		0.15	0.25	0.50	0.75	0.85
Training	1,593 (895)	121 (475)	702 (670)	1,544 (1,073)	3,131 (1,376)	3,378 (1,811)
% Impact of Training	8.55	5.19	11.99	9.64	10.69	9.02
High school or GED	4,075 (573)	714 (429)	1,752 (644)	4,024 (940)	5,392 (1,441)	5,954 (1,783)
Black	-2,349 (625)	-171 (439)	-377 (626)	-2,656 (1,136)	-4,182 (1,587)	-3,523 (1,867)
Hispanic	335 (888)	328 (757)	1,476 (1,128)	1,499 (1,390)	379 (2,294)	1,023 (2,427)
Married	6,647 (627)	1,564 (596)	3,190 (865)	7,683 (1,202)	9,509 (1,430)	10,185 (1,525)
Worked less than 13 weeks in past year	-6,575 (567)	-1,932 (442)	-4,195 (664)	-7,009 (1,040)	-9,289 (1,420)	-9,078 (1,596)
Constant	10,641 (1,569)	-134 (1,116)	1,049 (1,655)	7,689 (2,361)	14,901 (3,292)	22,412 (7,655)