

王小刚,陈姜猛.基于光滑化方法的分段线性删失分位数回归模型估计[J].江西师范大学学报(自然科学版),2022,46(3): 268-276.

WANG Xiaogang, CHEN Jiangmeng. The piecewise linear censored quantile regression model estimation based on smoothing technique [J]. Journal of Jiangxi Normal University(Natural Science) 2022, 46(3): 268-276.

文章编号: 1000-5862(2022)03-0268-09

基于光滑化方法的分段线性删失分位数回归模型估计

王小刚,陈姜猛

(北方民族大学数学与信息科学学院,宁夏银川 750021)

摘要: 针对在分段线性删失分位数回归模型中的变点问题,该文通过引入光滑化方法得到了变点位置及模型系数的估计,推导了参数估计的大样本性质.光滑化方法解决了在变点估计方法中常用的格点搜索法存在计算烦琐、解释意义不强的问题,弥补了线性化技术无法证明渐近性的不足,提高了估计的有效性和稳健性.蒙特卡罗模拟结果验证了在同方差和异方差、固定和随机删失下在不同分位点时的估计效果都具有有效性和稳健性.药物滥用数据的实证分析表明:复发时间间隔与治疗时间存在正向影响,且复发时间在 0.498 处存在变点(0.5 分位数),治疗时间在 0.498 之前的复发时间间隔比在 0.498 之后的更长,即大约前一半时间的治疗更加有效.

关键词: 光滑化方法; 分段线性; 删失分位数回归模型; 变点

中图分类号: O 212.2 文献标志码: A DOI: 10.16357/j.cnki.issn1000-5862.2022.03.09

0 引言

相较于传统的回归模型,分位数回归模型能够根据不同分位数得到不同的估计结果,既能衡量在平均意义下的估计结果,又能反映在极端情况下的估计效果,具有灵活性和稳健性以及单调同变性等优点,从而它受到越来越多学者的关注.在管理、金融和医学等领域的研究中,由于研究提前结束、研究对象无法继续参与实验(如研究对象工作变迁、意外身亡等事件不能继续研究)会导致无法获取个体准确的生存时间,即数据存在删失,所以直接使用分位数回归会损失有用的信息,使得模型估计效果产生偏差.文献[1-2]提出了基于删失数据的分位数回归模型估计方法,得到有效的估计; Peng Liming 等^[3]使用与删失数据相关联的鞅结构,估计了删失分位数回归模型的参数; Wang Huixia Judy 等^[4]在假设生存时间和删失变量条件独立的条件下,提出用局部加权分位数回归方法估计删失分位数回归模

型参数; Tang Yanlin 等^[5]为了自动地估计和检测分位数之间的共性,提出了一种新的带有 2 种分位数惩罚变量的删失分位数回归算法; 张倩倩等^[6]对存在删失的医疗费用建立了分位数回归模型,提出了简单加权估计和模拟退火算法,并验证了模型具有稳健性; 李忠桂等^[7]基于 EL 法和 SEL 法对右删失数据构造了更高精度的检验统计量; 孙桂萍等^[8]对长度偏差右删失数据剩余寿命分位数模型提出了参数估计,推导了估计的相合性和渐近正态性; 张立文等^[9]研究了删失分位数回归模型的变点检验问题; 王江峰等^[10]研究了在删失指标随机缺失下的非参数删失回归模型,构造了回归函数的复合分位数回归估计,得到估计的大样本性质.

在实际应用中,删失数据会受到诸如新政策实施、新技术变革、经济金融危机、全球新冠肺炎疫情等重大事件冲击,可能会呈现出非线性的结构特征和趋势,若不考虑变点的存在性则会产生有偏的估计.在此种情况下,由于分段线性删失分位数回归模

收稿日期: 2022-01-10

基金项目: 宁夏自然科学基金(2021AAC03186),宁夏高等教育一流学科建设基金(NXYLXK2017B09)和北方民族大学服务宁夏九大产业(FWXX36)资助项目.

作者简介: 王小刚(1980—),男,宁夏银川人,教授,博士,主要从事变点理论研究.E-mail: wxg@nun.edu.cn

型既能解释模型的非线性结构特征,又能保持删失分位数回归模型原有优点,所以该模型就成为研究的热点,也是使用较为广泛的模型。

在分段线性删失分位数回归模型中,困难之处在于需要对未知的变点位置进行估计。常用的变点位置估计方法可以分为2类。一类方法是2步法,即事先给定某个变点初值,将模型按照变点的位置分成前后2个部分,利用已知方法分别估计参数,然后通过最优化目标函数得到变点估计。如常用的格点搜索法^[11-13],该方法的优点在于不需要对模型分布做假设,具有一定的灵活性,然而其缺点在于变点估计的精度与运算效率成反向变动关系,且变点估计较难给出符合实际意义的解释,因此常用来做变点位置的初步估计。另一类方法是同时估计方法,即通过某种方法将变点参数转化为模型参数,再利用已有方法得到模型参数及变点位置的估计。如常用的线性化技术^[14-15],线性化技术简便易行,但难以推导大样本性质,且变点估计效果依赖于迭代快慢与质量,常会遇到变点收敛速率过慢状况。线性化方法的缺陷可以通过光滑化的核函数方法加以弥补^[16-18],光滑化的核函数方法的思想是将目标函数在变点处展开,利用光滑的核函数替代在目标函数中的不可导项,从而得到变点位置和模型系数的估计。光滑化的核函数方法估计更加有效,收敛速度更快,且能够得到估计的大样本性质,因此光滑化的核函数方法被广泛应用于变点的位置估计研究。

本文首先给出了分段线性删失分位数回归模型及估计方法,推导估计量的大样本性质,然后通过数值模拟验证了在分段线性删失分位数回归模型中的变点位置及模型系数估计效果,并将该方法对线性化技术进行了比较。利用药物滥用数据进行的实证分析发现了存在的变点,并给出了合理解释。

1 分段线性分位数模型及其估计

1.1 模型定义与估计

给定观察值 (x_i, y_i, δ_i) , 考虑删失分位数回归模型如下:

$$Y_i = \max\{C_i, T_i\}, \quad \delta_i = I(T_i \geq C_i), \quad Q_{T_i}(\tau, \xi | X_i, Z_i) = \beta_0 + \beta_1 X_i + Z_i^T \gamma, \quad i = 1, 2, \dots, n,$$

其中 Y_i 是响应变量, β_0, β_1 是模型系数, $\xi = (\beta_0, \beta_1, \gamma^T)^T$ 为 $p+2$ 维参数向量, X_i 和 Z_i 是可观测的协变

量, T_i 是未删失的独立变量, C_i 和 δ_i 分别表示随机删失变量和删失指标, $Q_{T_i}(\tau, \xi | X_i, Z_i)$ 表示响应变量 Y_i 在协变量 (X_i, Z_i^T) 的条件下的 τ 分位数, $\tau \in (0, 1)$ 。

假设未知的变点位置为 ψ , 分段线性删失分位数回归模型定义如下:

$$Q_{T_i}(\tau, \xi | X_i, Z_i) = \beta_0 + \beta_1 X_i + \beta_2 (X_i - \psi)_+ + Z_i^T \gamma, \quad (1)$$

其中 $(X_i - \psi)_+ = (X_i - \psi) I(X_i > \psi)$, $I(\cdot)$ 是示性函数, 且 $\{(Y_i, X_i, Z_i^T) : i = 1, 2, \dots, n\}$ 是来自 (Y, X, Z) 的独立样本。在模型(1)中, 协变量 X_i 在变点 ψ 处的分位数存在分段效应, 即在变点发生之前系数为 β_1 , 而在变点发生之后系数为 $\beta_1 + \beta_2$ 。为了识别变点, 要求 $\beta_2 \neq 0$ 。定义 $\theta = (\beta_0, \beta_1, \beta_2, \gamma^T, \psi)^T$, 则参数 θ 的估计可通过最小化目标函数得到, 即

$$l_{n,\tau}(\theta) = \sum_{i=1}^n \rho_{\tau}(Y_i - \max\{C_i, \beta_0 + \beta_1 X_i + \beta_2 (X_i - \psi)_+ + Z_i^T \gamma\}), \quad (2)$$

其中 $\rho_{\tau}(k) = k(\tau - I(k < 0))$ 是检验函数。由于变点 ψ 的存在导致目标函数(2)不可导, 因此, 借鉴文献[17-18]的方法采用光滑的核函数 $K((X - \psi)/h)$ 代替示性函数 $I(X > \psi)$, $h(>0)$ 表示窗宽, 将核函数 $K((X - \psi)/h)$ 代入目标函数(2), 得到新的目标函数:

$$S_{n,\tau}(\theta) = \sum_{i=1}^n \rho_{\tau}(Y_i - \max\{C_i, \beta_0 + \beta_1 X_i + \beta_2 (X_i - \psi) K((X_i - \psi)/h) + Z_i^T \gamma\}). \quad (3)$$

最小化目标函数(3)可得到模型参数 $\xi = (\beta_0, \beta_1, \beta_2, \beta_3, \gamma^T)^T$ 和变点 ψ 的估计。为了寻求目标函数(3)的最优化解, 将非线性项 $\beta_2 (X_i - \psi) K((X_i - \psi)/h)$ 在给定初值 ψ^* 处进行泰勒展开, 即

$$\beta_2 (X_i - \psi) K((X_i - \psi)/h) \approx \beta_2 U_i + \beta_3 V_i, \quad \text{其中 } U_i = (X_i - \psi^*) K((X_i - \psi^*)/h), \quad \beta_3 = \beta_2 (\psi - \psi^*), \quad V_i = -\{K((X_i - \psi^*)/h) + (X_i - \psi^*) K'((X_i - \psi^*)/h)\}.$$

将 U_i, V_i 代入目标函数(3), 有

$$S_{n,\tau}(\theta) = \sum_{i=1}^n \rho_{\tau}(Y_i - \max\{C_i, \beta_0 + \beta_1 X_i + \beta_2 U_i + \beta_3 V_i + Z_i^T \gamma\}). \quad (4)$$

通过最小化目标函数(4)即可更新参数 $\xi = (\beta_0, \beta_1, \beta_2, \beta_3, \gamma^T)^T$, 估计算法如下。

算法1 基于光滑化方法的分段线性删失分位数模型估计。

Step 1 设置初始变点 $\psi^{(0)}$ 及参数值 $\xi^{(0)} =$

$(\hat{\beta}_0^{(0)} \hat{\beta}_1^{(0)} \hat{\beta}_2^{(0)} \hat{\beta}_3^{(0)} \hat{\gamma}^{(0)T})^T$, 为保证算法收敛 $\hat{\beta}_3^{(0)}$ 要足够小, 如 $\hat{\beta}_3^{(0)} = 0.02$.

Step 2 对于给定变点 $\psi^{(k)}$, 更新参数 $\hat{\xi}^{(k)}$, 即

$$\hat{\xi}^{(k)} = \arg \min_{\xi} \sum_{i=1}^n \rho_{\tau}(Y_i - \max\{C_i \beta_0 + \beta_1 X_i + \beta_2 U_i^{(k)} + \beta_3 V_i^{(k)} + Z_i^T \gamma\}) ,$$

其中 $U_i^{(k)} = (X_i - \psi^{(k)}) K((X_i - \psi^{(k)})/h)$, $V_i^{(k)} = -\{K((X_i - \psi^{(k)})/h) + (X_i - \psi^{(k)}) K'((X_i - \psi^{(k)})/h)/h\}$.

Step 3 利用 $\psi^{(k+1)} = \psi^{(k)} + \beta_3^{(k+1)}/\beta_2^{(k+1)}$ 更新变点 $\psi^{(k+1)}$.

Step 4 重复步骤 2 和步骤 3 直到收敛, 收敛条件为

$$\|\hat{\theta}^{(k+1)} - \hat{\theta}^{(k)}\|_{\infty} \leq 10^{-5} ,$$

其中 $\forall u \in \mathbf{R}^q$, $\|u\|_{\infty} = \max_j |u_j|$.

参数估计依赖窗宽的选择, 本文通过最小化交叉验证函数选择最优窗宽, 即

$$\hat{h} = \arg \min C_V(h) = \arg \min \sum_{i=1}^n \rho_{\tau}\{Y_i - \hat{Q}_{T_i}(\tau, \theta | X_i, Z_i^T)\} ,$$

其中 $\hat{Q}_{T_i}(\tau, \theta | X, Z)$ 是在观测值 (Y_i, X_i, Z_i^T) 中去掉第 i 个个体后的估计. 用符号 $\hat{\theta}_n = (\hat{\xi} \hat{\psi})^T$ 表示算法 1 得到的估计.

1.2 渐近性质

令 $g(\omega_i, \theta) = \beta_0 + \beta_1 X_i + \beta_2 (X_i - \psi) K((X_i - \psi)/h) + Z_i^T \gamma$ 其中 $\omega_i = (1, X_i, Z_i^T)$ 代入式(3) 可得

$$S_{n,\tau}(\theta) = \sum_{i=1}^n \rho_{\tau}(Y_i - \max\{C_i g(\omega_i, \theta)\}) . \quad (5)$$

定义

$$q(\omega_i, \theta) = \partial \max\{C_i g(\omega_i, \theta)\} / \partial \theta = \text{diag}(I(g(\omega_i, \theta) > C_i)) (1, X_i, (X_i - \psi) K((X_i - \psi)/h), Z_i^T, -\beta_1 K((X_i - \psi)/h) - \beta_1 (X_i - \psi) K'((X_i - \psi)/h)/h)^T ,$$

$$C_n = \sum_{i=1}^n \tau(1 - \tau) E(q(\omega_i, \theta_0) q^T(\omega_i, \theta_0)) / n ,$$

$$D_n = \partial E \sum_{i=1}^n \rho'_{\tau}(Y_i - \max\{C_i g(\omega_i, \theta)\}) \cdot$$

$$q(\omega_i, \theta) / \partial \theta |_{\theta=\theta_0} / n ,$$

其中 θ_0 为真值 $\rho'_{\tau}(k) = \tau - I(k < 0)$ 表示 ρ_{τ} 的 1 阶导数. 极小化目标函数(5) 等价于求解方程

$$\sum_{i=1}^n \rho'_{\tau}(Y_i - \max\{C_i g(\omega_i, \theta)\}) q(\omega_i, \theta) = 0 .$$

为了证明其渐近性, 给出如下正则性条件.

(A₁) 对于给定的 $\varepsilon > 0$, 有 $\lim_{n \rightarrow \infty} P(Y_i - \max\{C_i g(\omega_i, \theta)\} < \varepsilon) = 0$;

(A₂) $F_{\tau}(\hat{e}_i)$ 是误差项 \hat{e}_i 的条件分布函数, 并且有连续有界的条件密度函数 $f_{\tau}(\hat{e}_i)$, 其中 $i = 1, 2, \dots, n$;

(A₃) X_i 的密度函数 $P_i(x)$ 在有界紧集合 $[M_1, M_2]$ 上连续;

(A₄) $E(\|Z_i\|^3) < \infty$ 其中 $\|\cdot\|$ 表示 Euclidean 范数;

(A₅) 存在 2 个正定矩阵 C_0 和 D_0 , 分别满足

$$\lim_{n \rightarrow \infty} C_n \equiv C_0, \lim_{n \rightarrow \infty} D_n \equiv D_0 .$$

引理 1 假设正则条件(A₁) ~ (A₅) 成立, 则有 $\hat{\theta}_n \xrightarrow{P} \theta_0$.

证 为证明 $\hat{\theta}_n$ 依概率收敛于 θ_0 , 首先证明 $\sup_{\theta} |S_{n,\tau}(\theta) - l_{n,\tau}(\theta)|/n \xrightarrow{P} 0$ 在条件(A₁) 下, 对于给定的 $\varepsilon > 0$, 有

$$\lim_{n \rightarrow \infty} P(Y_i - \max\{C_i g(\omega_i, \theta)\} < \varepsilon) = 0 ,$$

容易得到 $\forall \delta > 0, \exists N$, 当 $n > N$ 时, 有

$$P(Y_i - \max\{C_i g(\omega_i, \theta)\} < \varepsilon) < \delta ,$$

用核函数 $K((X - \psi)/h)$ 代替示性函数 $I(X > \psi)$, 有 $I(X_i > \psi) = K((X_i - \psi)/h) + o(h)$ $h \rightarrow 0$.

令 $u(\omega_i, \theta) = \beta_0 + \beta_1 X_i + \beta_2 (X_i - \psi) + Z_i^T \gamma$ 则

$$u(\omega_i, \theta) = \beta_0 + \beta_1 X_i + \beta_2 (X_i - \psi) (K((X_i - \psi)/h) + o(h)) + Z_i^T \gamma = g(\omega_i, \theta) + \beta_2 (X_i - \psi) \cdot o(h) ,$$

可得

$$D_n = \sum_{i=1}^n \{\rho_{\tau}(Y_i - \max\{C_i g(\omega_i, \theta)\}) - \rho_{\tau}(Y_i -$$

$$\max\{C_i u(\omega_i, \theta)\})\} / n = \sum_{i=1}^n (\tau(Y_i - \max\{C_i g(\omega_i, \theta)\}) - \tau(Y_i - \max\{C_i u(\omega_i, \theta)\}) - (Y_i - \max\{C_i g(\omega_i, \theta)\}) I(Y_i < \max\{C_i g(\omega_i, \theta)\}) + (Y_i - \max\{C_i u(\omega_i, \theta)\}) I(Y_i < \max\{C_i u(\omega_i, \theta)\})) / n \leq$$

$$\sum_{i=1}^n (|\tau \beta_2 (X_i - \psi)| o(h) - (Y_i - \max\{C_i g(\omega_i, \theta)\}) I(Y_i < \max\{C_i g(\omega_i, \theta)\}) + (Y_i - \max\{C_i g(\omega_i, \theta)\}) I(Y_i < \max\{C_i u(\omega_i, \theta)\}) - (Y_i - \max\{C_i g(\omega_i, \theta)\}) I(Y_i < \max\{C_i u(\omega_i, \theta)\}) + (Y_i - \max\{C_i u(\omega_i, \theta)\}) I(Y_i < \max\{C_i u(\omega_i, \theta)\})) / n \leq$$

$$u(\omega_i, \theta) \} I(Y_i < \max\{C_i, \mu(\omega_i, \theta)\}) / n \leq \sum_{i=1}^n (|\tau\beta_2 \cdot (X_i - \psi) + o(h) - \beta_2(X_i - \psi)| \cdot o(h) I(Y_i < \max\{C_i, u(\omega_i, \theta)\}) + (Y_i - \max\{C_i, g(\omega_i, \theta)\}) I(|Y_i - \max\{C_i, g(\omega_i, \theta)\}| < |\max\{C_i, \beta_2(X_i - \psi) + o(h)\}|)) / n.$$

由正则性假设(A₁) 知有

$$\sup_{\theta} |S_{n\tau}(\theta) - l_{n\tau}(\theta)| / n \leq \sup_{\theta} |\sum_{i=1}^n \tau\beta_2(X_i - \psi)| / n o(h) + \sup_{\theta} |-\sum_{i=1}^n \beta_2(X_i - \psi)| / n o(h) + \sup_{\theta} |\sum_{i=1}^N Y_i - \max\{C_i, g(\omega_i, \theta)\}| / n + \sup_{\theta} |\sum_{i=N+1}^n Y_i - \max\{C_i, g(\omega_i, \theta)\}| / n \cdot \delta,$$

对于有限的 N 和任意小的 δ , 当 $n \rightarrow \infty$ 时, 上述不等式右边 4 项都收敛于 0, 即

$$\sup_{\theta} |S_{n\tau}(\theta) - l_{n\tau}(\theta)| / n \xrightarrow{P} 0.$$

结合 Li Chenxi 等^[12] 证明的一致性, 得到 $l_{n\tau}(\theta)$ 的连续性. 引理 1 得证.

引理 2 在假设(A₂) ~ (A₄) 条件下, 对任意正项数列 $\{a_n\}$ 且 $\lim_{n \rightarrow \infty} a_n = 0$, 有

$$\sup_{\|\theta - \theta_0\| \leq a_n} \|\sum_{i=1}^n (v_i(\theta - \theta_0) - E(v_i(\theta - \theta_0)))\| / \sqrt{n} = o_p(1),$$

其中

$$v_i(\theta - \theta_0) = \rho_{\tau}'(Y_i - \max\{C_i, g(\omega_i, \theta)\}) q(\omega_i, \theta) - \rho_{\tau}'(Y_i - \max\{C_i, g(\omega_i, \theta_0)\}) q(\omega_i, \theta_0).$$

证 引理 2 的证明与文献[19] 中的引理 4.6 的证明相似. 为证明引理 2, 需验证文献[19] 中引理 4.6 的条件(B₁)、(B₃) 和(B₅) 满足即可. 条件(B₁) 显然成立.

对于条件(B₃) 有

$$\|v_i(\theta - \theta_0(\tau))\| \leq \|\rho_{\tau}'(Y_i - \max\{C_i, g(\omega_i, \theta)\}) (q(\omega_i, \theta) - q(\omega_i, \theta_0))\| + \|(\rho_{\tau}'(Y_i - \max\{C_i, g(\omega_i, \theta)\}) - \rho_{\tau}'(Y_i - \max\{C_i, g(\omega_i, \theta_0)\})) q(\omega_i, \theta_0)\| = \|I_{1i}\| + \|I_{2i}\|,$$

其中

$$\|I_{1i}\| = \|\rho_{\tau}'(Y_i - \max\{C_i, g(\omega_i, \theta)\}) (q(\omega_i, \theta) - q(\omega_i, \theta_0))\| = \|\rho_{\tau}'(Y_i - \max\{C_i, g(\omega_i, \theta)\}) q(\omega_i, \theta_0) (\theta - \theta_0)\| = o_p(1),$$

所以有 $E(\|I_{1i}\|^2 | \omega_i) = o_p(1)$. 同时,

$$\|I_{2i}\| = \|(\rho_{\tau}'(Y_i - \max\{C_i, g(\omega_i, \theta)\}) - \rho_{\tau}'(Y_i - \max\{C_i, g(\omega_i, \theta_0)\})) q(\omega_i, \theta_0)\| = \|(I(Y_i \leq \max\{C_i, g(\omega_i, \theta)\}) - I(Y_i \leq \max\{C_i, g(\omega_i, \theta_0)\})) q(\omega_i, \theta_0)\| \leq \|q(\omega_i, \theta_0)\| I(\max\{C_i, g_1(\omega_i, \theta, \theta_0)\} \leq Y_i \leq \max\{C_i, g_2(\omega_i, \theta, \theta_0)\}),$$

这里的 $g_1(\omega_i, \theta, \theta_0)$ 和 $g_2(\omega_i, \theta, \theta_0)$ 分别为

$$g_1(\omega_i, \theta, \theta_0) = \min\{\max\{C_i, g(\omega_i, \theta)\}, \max\{C_i, g(\omega_i, \theta_0)\}\},$$

$$g_2(\omega_i, \theta, \theta_0) = \max\{\max\{C_i, g(\omega_i, \theta)\}, \max\{C_i, g(\omega_i, \theta_0)\}\}.$$

又由于 $g_1(\omega_i, \theta, \theta_0) \leq g_2(\omega_i, \theta, \theta_0)$, 有

$$\|\max\{C_i, g(\omega_i, \theta)\} - \max\{C_i, g(\omega_i, \theta_0)\}\| = \|q^T(\omega_i, \theta_0)(\theta - \theta_0) + R_n\| \leq \|q(\omega_i, \theta_0)(\theta - \theta_0)\| \leq a_n \|q(\omega_i, \theta_0)\|,$$

即

$$E(\|I_{2i}\|^2 | \omega_i) \leq \|q(\omega_i, \theta_0)\|^2 E(I(g_1(\omega_i, \theta, \theta_0) \leq Y_i \leq g_2(\omega_i, \theta, \theta_0))) \leq \|q(\omega_i, \theta_0)\|^2 \cdot f_{\tau, \omega_i}(\xi_i) \|\max\{C_i, g(\omega_i, \theta)\} - \max\{C_i, g(\omega_i, \theta_0)\}\| \leq a_n f_{\tau, \omega_i}(\xi_i) \|q(\omega_i, \theta_0)\|^3,$$

其中 $g_1(\omega_i, \theta, \theta_0) < \xi_i < g_2(\omega_i, \theta, \theta_0)$, 所以

$$E(\|v_i(\theta - \theta_0(\tau))\|^2 | \omega_i) \leq d_i^2 / \sqrt{n}, \text{ 这里}$$

$$d_i = \sqrt{a_n f_{\tau, \omega_i} \sqrt{n} \|q(\omega_i, \theta_0)\|^3},$$

即条件(B₃) 满足.

下面验证条件(B₅). 对于任意的常数 C , 令

$$A_n = \sum_{i=1}^n d_i^2 = a_n \sqrt{n} \sum_{i=1}^n (f_{\tau, \omega_i}(\xi_i) \|q(\omega_i, \theta_0)\|^3),$$

有

$$P(\max_{1 \leq i \leq n} \|v_i(\theta - \theta_0)\| \geq CA_n^{1/2} a_n^{1/2} (\ln n)^{-2}) \leq$$

$$\sum_{i=1}^n P(\|v_i(\theta - \theta_0)\| \geq CA_n^{1/2} a_n^{1/2} (\ln n)^{-2}) \leq$$

$$\sum_{i=1}^n E(\|v_i(\theta - \theta_0)\|^2) / (C^2 A_n a_n (\ln n)^{-4}) \leq n^{-1/2} A_n \cdot (\ln n)^4 / (C^2 A_n a_n) = (\ln n)^4 / (C^2 a_n \sqrt{n}) = o_p(1),$$

即 $\max_{1 \leq i \leq n} \|v_i(\theta - \theta_0)\| = o_p(A_n^{1/2} a_n^{1/2} (\ln n)^{-2})$, 条件(B₅) 满足. 引理 2 得证.

定理 1 在正则条件(A₁) ~ (A₅) 下, 有

$$\sqrt{n}(\hat{\theta}_n - \theta_0) = o_p(1) \sqrt{n}(\hat{\theta}_n - \theta_0) \rightarrow N(0, D_0^{-1} C_0 D_0^{-T}),$$

其中协方差阵中的 C_0 和 D_0 的估计分别为

$$\hat{C}_0 = \sum_{i=1}^n (\tau(1 - \tau) q(\omega_i, \hat{\theta}_n) q^T(\omega_i, \hat{\theta}_n)) / n,$$

$$\hat{D}_0 = - \sum_{i=1}^n (\hat{f}_{\tau}(\hat{e}_i) q(\omega_i, \hat{\theta}_n) q^T(\omega_i, \hat{\theta}_n)) / n,$$

这里 $\hat{e}_i = Y_i - \max\{C_i, g(\omega_i, \hat{\theta}_n)\}$, $f_\tau(\hat{e}_i)$ 是 Y_i 的条件密度估计, 可用 W. Hendricks 等^[20] 所提出的离散导数来估计, 即

$$\hat{f}_\tau(\hat{e}_i) = 2\Delta_n / (\hat{Q}_{\tau+\Delta_n}(\hat{e}_i) - \hat{Q}_{\tau-\Delta_n}(\hat{e}_i)).$$

Δ_n 是平滑参数, 借鉴文献[21]的方法选择如下:

$\Delta_n = 1.57(1.5\varphi^2(\Phi^{-1}(\tau)) / (2(\Phi^{-1}(\tau))^2 + 1))^{1/3} / \sqrt[3]{n}$, 其中 $\varphi(\cdot)$ 和 $\Phi(\cdot)$ 分别表示标准正态分布的密度函数和分布函数.

证 在引理 1 和引理 2 的条件下, 可得

$$\begin{aligned} \sqrt{n} \parallel \sum_{i=1}^n \rho'_\tau(Y_i - \max\{C_i, g(\omega_i, \hat{\theta}_n)\}) q(\omega_i, \hat{\theta}_n) - \sum_{i=1}^n \rho'_\tau(Y_i - \max\{C_i, g(\omega_i, \theta_0)\}) q(\omega_i, \theta_0) - \sum_{i=1}^n E(\rho'_\tau(Y_i - \max\{C_i, g(\omega_i, \hat{\theta}_n)\}) q(\omega_i, \hat{\theta}_n)) \parallel = o_p(1). \end{aligned}$$

将 $\sum_{i=1}^n E(\rho'_\tau(Y_i - \max\{C_i, g(\omega_i, \hat{\theta}_n)\}) q(\omega_i, \hat{\theta}_n))$ 在 θ_0 处进行泰勒展开, 有

$$\begin{aligned} \sum_{i=1}^n E(\rho'_\tau(Y_i - \max\{C_i, g(\omega_i, \hat{\theta}_n)\}) q(\omega_i, \hat{\theta}_n)) = \\ (\partial E \sum_{i=1}^n \rho'_\tau(Y_i - \max\{C_i, g(\omega_i, \theta_0)\}) q(\omega_i, \theta_0)) / \partial \theta \big|_{\theta=\theta_0} (\hat{\theta}_n - \theta_0) + R_n = n D_{n\pi} (\hat{\theta}_n - \theta_0) + R_n, \end{aligned}$$

其中 $R_n = o_p(\sqrt{n})$, $D_{n\pi} = \sum_{i=1}^n \partial E(\rho'_\tau(Y_i - \max\{C_i, g(\omega_i, \theta)\}) q(\omega_i, \theta)) / \partial \theta \big|_{\theta=\theta_0} / n = \sum_{i=1}^n (-f_\tau(Y_i - \max\{C_i, g(\omega_i, \theta_0)\}) q(\omega_i, \theta_0) q^T(\omega_i, \theta_0) + \{ \tau - F_i(Y_i - \max\{C_i, g(\omega_i, \theta_0)\}) \} \partial q(\omega_i, \theta) / \partial \theta \big|_{\theta=\theta_0}) / n = - \sum_{i=1}^n (f_\tau(Y_i - \max\{C_i, g(\omega_i, \theta_0)\}) q(\omega_i, \theta_0) q^T(\omega_i, \theta_0)) / n.$

由于 $-\sum_{i=1}^n \rho'_\tau(Y_i - \max\{C_i, g(\omega_i, \hat{\theta}_n)\}) q(\omega_i, \hat{\theta}_n) / \sqrt{n} = o_p(1)$, 即

$$-\sum_{i=1}^n \rho'_\tau(Y_i - \max\{C_i, g(\omega_i, \theta_0)\}) q(\omega_i, \theta_0) / \sqrt{n} - n D_{n\pi} (\hat{\theta}_n - \theta_0) / \sqrt{n} - R_n = o_p(1),$$

因此有

$$\sqrt{n}(\hat{\theta}_n - \theta_0) = -D_{n\pi}^{-1} \sum_{i=1}^n \rho'_\tau(Y_i - \max\{C_i, g(\omega_i, \theta_0)\}) q(\omega_i, \theta_0) / \sqrt{n} + o_p(\sqrt{n}),$$

即 $\sqrt{n}(\hat{\theta}_n - \theta_0) \rightarrow N(0, D_0^{-1} C_0 D_0^{-1})$. 定理 1 得证.

2 蒙特卡罗模拟

为验证本文方法的有效性和稳健型, 考虑 4 种随机模拟, 数据生成过程如下: (i) 不同方差形态. 当同方差时 $y_i^* = \beta_0 + \beta_1 x_i + \beta_2(x_i - \psi)_+ + Z_i^T \gamma + \lambda_1 \varepsilon_i$; 当异方差时 $y_i^* = \beta_0 + \beta_1 x_i + \beta_2(x_i - \psi)_+ + Z_i^T \gamma + (\lambda_2 + \lambda_1 x_i) \varepsilon_i$. (ii) 不同删失形态. 当固定删失时 $C_i = k$; 当随机删失时 $C_i = \alpha_0 + \alpha_1 z_i + \alpha_2 x_i + \lambda_1 \varepsilon_i$, 其中 $x_i \sim U(-2, 4)$, $z_i \sim N(1, 4)$, 变点 $\psi = 1.5$, 模型系数设置为 $(\beta_0, \beta_1, \beta_2, \gamma) = (1, 2, -3, 4)$, 窗宽 $h = 0.2n^{-1/5}$, $\varepsilon_i \sim N(0, 1)$, λ_1 和 λ_2 为控制模型方差类型的参数, 真值设置为 $(\lambda_1, \lambda_2) = (0.2, 0.3)$. 设定删失模型的删失度为 30%, $\alpha^T = (0.3, 0.6, 0.3)$, $k = -5$. 样本量取 $n = 200$ 和 $n = 500$, 在分位数 $\tau = 0.3, 0.5, 0.7$ 的情况下分别重复模拟 1 000 次. 分段线性删失分位数回归模型为

$$Y = \max\{C, T\}, \delta = I(T \geq C), Q_\tau(\tau | X, Z) = \beta_0 + \beta_1 X + \beta_2(X - \psi)_+ + Z^T \gamma.$$

在固定删失和随机删失下, 在不同分位数下的同方差和异方差数值模拟结果分别如表 1 和表 2 所示. 表 1 和表 2 的上部分和下部分分别为在样本量 $n = 200$ 和 $n = 500$ 下估计量的评价指标, 即偏差 (B_s)、标准差 (S_D)、标准误 (E_s)、95% 的覆盖率 (C_p) 和置信区间长度 (A_w).

从表 1 可以看出: 在固定删失情况下, 当分位数为 0.5 且同方差时, 模型估计的 B_s 和 S_D 都很小, 且 E_s 和 S_D 非常接近, 这说明估计效果是有效的. A_w 较小, 且 C_p 都接近于 95%. 从横向来, 在 0.5 分位数下, 异方差的参数估计 B_s 和 S_D 均较小, E_s 和 S_D 很接近, 这说明即使在异方差情况下模型估计效果也较好. 另外, 异方差的置信区间长度 A_w 和覆盖率 C_p 的估计结果要比在同方差情况下的估计结果更差. 当分位数为 0.3 和 0.7 时, 模型估计的 B_s 、 S_D 要比在分位数为 0.5 时的估计结果更差, E_s 和 S_D 的接近程度较好. A_w 与 C_p 的效果与在分位数为 0.5 时的效果差异不大, 这说明估计结果具有稳健性.

当样本量从 200 增长到 500 时, 可以发现无论是估计的 B_s 、 S_D 、 E_s 还是 A_w 都呈现减小趋势, 但 C_p 呈现增大趋势, 这也验证了定理 1 的结果. 同样 E_s 与 S_D 非常接近, C_p 更接近于 95%, 这也验证了估计具有可靠性.

表 1 在固定删失下的模拟结果

$n \quad \tau$		同方差					异方差				
		β_0	β_1	β_2	γ	ψ	β_0	β_1	β_2	γ	ψ
$n = 200$	B_s	- 0.104	0.000	0.005	- 0.008	0.001	- 0.161	- 0.000	- 0.096	- 0.008	0.002
	S_D	0.074	0.007	0.082	0.091	0.023	0.142	0.023	0.194	0.248	0.075
	E_S	0.070	0.007	0.080	0.090	0.022	0.136	0.021	0.181	0.248	0.067
	A_w	0.274	0.028	0.316	0.355	0.088	0.535	0.085	0.709	0.972	0.263
	C_P	0.679	0.934	0.927	0.939	0.930	0.733	0.924	0.899	0.945	0.899
$\tau = 0.3$	B_s	- 0.000	0.000	0.005	- 0.009	0.001	- 0.004	0.000	0.003	0.005	0.003
	S_D	0.066	0.007	0.079	0.087	0.022	0.134	0.021	0.185	0.231	0.070
	E_S	0.067	0.007	0.077	0.086	0.021	0.128	0.021	0.171	0.235	0.064
	A_w	0.264	0.027	0.304	0.340	0.085	0.503	0.082	0.673	0.924	0.252
	C_P	0.934	0.942	0.938	0.948	0.932	0.908	0.945	0.912	0.947	0.912
$n = 200$	B_s	0.098	- 0.000	0.012	- 0.015	0.001	0.160	- 0.000	0.098	0.006	- 0.001
	S_D	0.069	0.007	0.083	0.092	0.023	0.140	0.021	0.189	0.236	0.074
	E_S	0.072	0.007	0.081	0.091	0.022	0.136	0.021	0.182	0.249	0.067
	A_w	0.279	0.029	0.321	0.359	0.089	0.534	0.085	0.714	0.978	0.265
	C_P	0.700	0.939	0.926	0.938	0.925	0.741	0.945	0.884	0.954	0.916
$\tau = 0.7$	B_s	0.106	0.000	0.005	0.008	0.000	- 0.157	- 0.000	- 0.102	- 0.003	0.001
	S_D	0.042	0.004	0.052	0.059	0.014	0.083	0.013	0.115	0.145	0.047
	E_S	0.044	0.004	0.051	0.057	0.014	0.084	0.013	0.113	0.155	0.042
	A_w	0.174	0.018	0.201	0.223	0.056	0.330	0.054	0.446	0.610	0.167
	C_P	0.316	0.943	0.945	0.942	0.942	0.527	0.947	0.844	0.959	0.912
$n = 500$	B_s	- 0.001	0.000	0.004	- 0.006	0.001	- 0.003	0.001	0.004	0.000	0.002
	S_D	0.043	0.004	0.049	0.054	0.014	0.080	0.013	0.108	0.139	0.044
	E_S	0.041	0.004	0.048	0.054	0.013	0.080	0.013	0.108	0.147	0.040
	A_w	0.163	0.017	0.189	0.211	0.053	0.316	0.052	0.424	0.577	0.159
	C_P	0.936	0.947	0.935	0.928	0.937	0.942	0.942	0.949	0.952	0.919
$\tau = 0.5$	B_s	0.095	- 0.000	0.009	- 0.011	0.001	0.157	0.001	0.101	0.002	- 0.000
	S_D	0.044	0.004	0.051	0.057	0.015	0.087	0.014	0.120	0.152	0.046
	E_S	0.044	0.004	0.051	0.057	0.014	0.084	0.013	0.113	0.155	0.042
	A_w	0.175	0.018	0.201	0.224	0.056	0.329	0.054	0.446	0.609	0.168
	C_P	0.414	0.960	0.934	0.941	0.935	0.518	0.946	0.831	0.951	0.914
$n = 500$	B_s	- 0.001	0.000	0.004	- 0.006	0.001	- 0.003	0.001	0.004	0.000	0.002
	S_D	0.043	0.004	0.049	0.054	0.014	0.080	0.013	0.108	0.139	0.044
	E_S	0.041	0.004	0.048	0.054	0.013	0.080	0.013	0.108	0.147	0.040
	A_w	0.163	0.017	0.189	0.211	0.053	0.316	0.052	0.424	0.577	0.159
	C_P	0.936	0.947	0.935	0.928	0.937	0.942	0.942	0.949	0.952	0.919
$\tau = 0.7$	B_s	0.095	- 0.000	0.009	- 0.011	0.001	0.157	0.001	0.101	0.002	- 0.000
	S_D	0.044	0.004	0.051	0.057	0.015	0.087	0.014	0.120	0.152	0.046
	E_S	0.044	0.004	0.051	0.057	0.014	0.084	0.013	0.113	0.155	0.042
	A_w	0.175	0.018	0.201	0.224	0.056	0.329	0.054	0.446	0.609	0.168
	C_P	0.414	0.960	0.934	0.941	0.935	0.518	0.946	0.831	0.951	0.914

在随机删失下的模拟结果如表 2 所示.从表 2 可以看出: 在同方差、相同样本容量、相同分位数情况下 固定删失的估计效果要比随机删失的估计效果更

好. B_s 、 S_D 、 E_S 之间的关系以及 A_w 和 C_P 之间的特征与表 1 中所得的结论相同.

表 2 在随机删失下的模拟结果

$n \quad \tau$		同方差					异方差				
		β_0	β_1	β_2	γ	ψ	β_0	β_1	β_2	γ	ψ
$n = 200$	B_s	- 0.113	0.001	0.006	- 0.009	0.001	- 0.187	0.010	- 0.110	0.004	0.006
	S_D	0.076	0.008	0.087	0.095	0.024	0.144	0.024	0.190	0.243	0.075
	E_S	0.072	0.008	0.082	0.092	0.023	0.139	0.023	0.184	0.255	0.070
	A_w	0.283	0.030	0.324	0.364	0.091	0.547	0.091	0.725	1.003	0.274
	C_P	0.630	0.946	0.929	0.935	0.920	0.700	0.897	0.886	0.958	0.920
$\tau = 0.3$	B_s	- 0.008	0.001	0.008	- 0.009	0.002	- 0.024	0.009	- 0.017	0.005	- 0.001
	S_D	0.068	0.007	0.080	0.086	0.023	0.138	0.022	0.189	0.238	0.072
	E_S	0.069	0.007	0.079	0.088	0.022	0.135	0.022	0.177	0.242	0.066
	A_w	0.273	0.029	0.311	0.348	0.087	0.531	0.087	0.694	0.950	0.259
	C_P	0.941	0.951	0.941	0.951	0.924	0.925	0.910	0.921	0.947	0.921
$n = 200$	B_s	0.098	0.001	0.007	- 0.012	0.001	0.139	0.009	0.075	0.025	0.000
	S_D	0.071	0.008	0.083	0.093	0.023	0.140	0.024	0.193	0.248	0.075
	E_S	0.074	0.007	0.084	0.094	0.023	0.141	0.023	0.185	0.256	0.069
	A_w	0.291	0.031	0.330	0.370	0.092	0.554	0.092	0.728	1.005	0.272
	C_P	0.726	0.934	0.928	0.941	0.932	0.794	0.927	0.902	0.951	0.908
$\tau = 0.7$	B_s	0.098	0.001	0.007	- 0.012	0.001	0.139	0.009	0.075	0.025	0.000
	S_D	0.071	0.008	0.083	0.093	0.023	0.140	0.024	0.193	0.248	0.075
	E_S	0.074	0.007	0.084	0.094	0.023	0.141	0.023	0.185	0.256	0.069
	A_w	0.291	0.031	0.330	0.370	0.092	0.554	0.092	0.728	1.005	0.272
	C_P	0.726	0.934	0.928	0.941	0.932	0.794	0.927	0.902	0.951	0.908

表 2(续)

n	π	同方差					异方差				
		β_0	β_1	β_2	γ	ψ	β_0	β_1	β_2	γ	ψ
$n = 500$	B_s	-0.109	0.001	0.004	-0.007	0.001	-0.182	0.010	-0.118	-0.001	0.001
	S_D	0.046	0.004	0.053	0.059	0.015	0.086	0.014	0.117	0.155	0.048
	E_S	0.046	0.005	0.052	0.058	0.014	0.088	0.014	0.116	0.159	0.043
	A_W	0.181	0.019	0.206	0.231	0.057	0.348	0.058	0.456	0.626	0.172
	C_P	0.342	0.955	0.948	0.945	0.935	0.466	0.895	0.817	0.960	0.920
$n = 500$	B_s	-0.005	0.001	0.005	-0.006	0.001	-0.030	0.009	-0.013	0.004	0.001
	S_D	0.045	0.004	0.048	0.053	0.014	0.083	0.013	0.110	0.149	0.044
	E_S	0.043	0.004	0.049	0.055	0.013	0.084	0.014	0.111	0.151	0.041
	A_W	0.171	0.018	0.195	0.218	0.054	0.332	0.055	0.435	0.595	0.164
	C_P	0.938	0.950	0.947	0.950	0.933	0.936	0.898	0.943	0.955	0.933
$n = 500$	B_s	0.097	0.001	0.006	-0.007	0.001	0.132	0.009	0.090	0.000	0.000
	S_D	0.046	0.005	0.055	0.062	0.014	0.086	0.015	0.117	0.154	0.046
	E_S	0.046	0.005	0.052	0.058	0.014	0.089	0.014	0.116	0.159	0.043
	A_W	0.180	0.019	0.205	0.229	0.057	0.350	0.057	0.455	0.625	0.170
	C_P	0.415	0.944	0.921	0.919	0.944	0.681	0.880	0.848	0.942	0.929

总体上,在其他情况不变的条件下,同方差下的估计效果要优于异方差下的估计效果,更大的样本容量的估计效果更好,这验证了参数估计满足相合性,固定删失的估计效果优于随机删失的估计效果,0.5的分位数的估计效果优于其他分位数的估计效果。蒙

特卡罗模拟结果表明本文的光滑化方法具有有效性和稳健性。

为对比光滑化方法的可行性,对本文方法和线性化技术得到的模拟结果进行比较。当样本容量为1 000且分位数为0.5时的模拟结果如表3和表4所示。

表 3 在固定删失下线性化技术与本文方法的模拟结果对比

方差类型		$\tau = 0.5 \quad n = 1\,000$									
		线性化技术					本文方法				
		B_s	S_D	E_S	A_W	C_P	B_s	S_D	E_S	A_W	C_P
同方差	β_0	-0.001	0.029	0.029	0.116	0.950	-0.002	0.029	0.029	0.116	0.951
	β_1	0.000	0.003	0.003	0.012	0.954	0.000	0.003	0.003	0.012	0.936
	β_2	0.002	0.034	0.034	0.135	0.950	0.006	0.034	0.034	0.134	0.938
	γ	-0.002	0.038	0.038	0.150	0.946	-0.007	0.037	0.038	0.149	0.938
	ψ	0.000	0.009	0.009	0.038	0.956	0.000	0.009	0.009	0.037	0.940
异方差	β_0	0.002	0.053	0.057	0.224	0.946	-0.002	0.056	0.059	0.221	0.943
	β_1	-0.000	0.009	0.009	0.036	0.944	-0.000	0.009	0.009	0.036	0.952
	β_2	-0.001	0.079	0.076	0.301	0.942	0.003	0.076	0.076	0.298	0.935
	γ	0.004	0.101	0.104	0.407	0.962	-0.004	0.100	0.103	0.405	0.952
	ψ	0.000 4	0.031	0.029	0.113	0.942	0.001	0.028	0.029	0.112	0.923

表 4 在随机删失下线性化技术与本文方法的模拟结果对比

方差类型		$\tau = 0.5 \quad n = 1\,000$									
		线性化技术					本文方法				
		B_s	S_D	E_S	A_W	C_P	B_s	S_D	E_S	A_W	C_P
同方差	β_0	-0.001	0.030	0.031	0.122	0.950	-0.005	0.031	0.030	0.120	0.939
	β_1	0.001	0.003	0.003	0.013	0.940	0.001	0.003	0.003	0.013	0.933
	β_2	-0.005	0.035	0.035	0.139	0.952	0.004	0.035	0.035	0.137	0.946
	γ	0.005	0.038	0.039	0.155	0.952	-0.005	0.038	0.039	0.153	0.950
	ψ	-0.000	0.010	0.010	0.039	0.930	0.001	0.010	0.009	0.038	0.947
异方差	β_0	-0.023	0.060	0.059	0.233	0.924	-0.026	0.059	0.059	0.232	0.922
	β_1	0.010	0.009	0.010	0.039	0.810	0.009	0.010	0.010	0.039	0.827
	β_2	-0.022	0.081	0.079	0.309	0.936	-0.012	0.079	0.078	0.306	0.933
	γ	0.005	0.104	0.107	0.419	0.950	0.001	0.106	0.106	0.418	0.951
	ψ	-0.001	0.030	0.029	0.117	0.958	0.001	0.030	0.029	0.116	0.941

从表3和表4中可以看出所有的 B_s 都非常小,因此2种方法的估计是渐近一致的。2种方法的 E_S 与 S_D 非常接近,这表明估计具有相合性。此外,2种方

法的 C_P 均接近95%。但在相同显著性水平下,本文方法在同方差情况下略好于线性化技术,在异方差情况下优于线性化技术。本文方法估计的 A_W 比线性

化技术估计的 A_w 更小,因此,本文方法在估计变点位置及模型参数的有效性上要优于线性化技术,且本文方法的收敛速率较高,而线性化技术在模拟时可能会存在不收敛情形。

3 实证分析

在医学领域中,药物的滥用会给个人、家庭及社会带来巨大的危害性。若不采取有效的预防措施和控制,则与之有关的疾病将会很快在全球泛滥成灾,所有国家都将处于这种危险之中,因此研究药物滥用复发时间的影响因素和变化趋势有着重要的实际意义,也有利于对药物滥用复发时间的控制和采取有效的预防措施。但是在对药物滥用复发治疗时,由于病人的药物没有复发而研究结束了或在研究结束前有病人提前结束治疗,所以导致数据存在删失。另外,考虑到模型会受到药物复发治疗的影响因素以及药物复发的反应快慢的冲击而产生变点,因此可通过用含变点的分段线性删失分位数回归模型研究药物滥用复发时间的影响因素,并利用本文方法进行估计。

文献 [22] 给出了关于 UIS 药物治疗的研究数据,原始数据的样本容量为 628,包含 575 个无缺失的数据。响应变量为复发时间间隔,协变量为随机化治疗 Z_T ($Z_T = 1$ 表示接受了 180 d 的治疗, $Z_T = 0$ 表示接受了 90 d 的治疗), Z_A 为登记年龄, Z_B 为抑郁评分(取值为 0 ~ 54), Z_V 为最近静脉用药史, $Z_V = 1$ 表示用药, $Z_V = 0$ 表示没有用药, Z_N 表示既往药物治疗次数(取值 0 ~ 40),由于 Z_N 在一定情况下会受到治疗效应解释的影响,所以将 Z_N 分成 Z_{N_1} 和 Z_{N_2} 2 个指标, Z_R 表示种族, $Z_R = 0$ 为白人, $Z_R = 1$ 表示非白人, Z_S 表示治疗地点, $Z_S = 0$ 表示治疗点 A, $Z_S = 1$ 表示治疗点 B, X_F 表示治疗时间,用治疗天数

T_L 定义 $X_F = T_L/90$ 为短治疗, $X_F = T_L/180$ 为长治疗。在研究结束前有病人提前结束治疗,因此数据存在删失。选取协变量 Z_{N_1} 、 Z_{N_2} 、 Z_V 、 Z_T 、 X_F ^[23] 构建分段线性删失分位数回归模型。

从药物复发时间间隔(取对数)和 X_F 的散点图(见图 1)可以看出: X_F 值在 0.5 前后复发时间间隔呈现不同特点。在 0.5 之前复发时间间隔较大,但在 0.5 之后逐渐变小,呈现出非线性特征,这表明模型可能存在变点。因此建立分段线性删失分位数回归来分析药物滥用复发时间的治疗数据。模型如下:

$$Y_i = \max(C_i, T_i) \quad \delta_i = I(T_i \geq C_i) \quad i = 1, 2, \dots, 575,$$
$$Q_{T_i}(\tau | X_i, Z_i) = \beta_0 + \beta_1 X_{Fi} + \beta_2 (X_{Fi} - \psi)_+ + \gamma_1 Z_{N1i} + \gamma_2 Z_{N2i} + \gamma_3 Z_{Vi} + \gamma_4 Z_{Ti},$$

其中 Y_i 为药物复发时间间隔, δ_i 为删失指标, $\theta = (\beta_0, \beta_1, \beta_2, \gamma_1, \gamma_2, \gamma_3, \gamma_4, \psi)^T$ 为未知参数。在分位数 $\tau = 0.3, 0.5, 0.7$ 时利用本文方法估计了变点位置及模型参数。模型估计结果如表 5 所示。

由表 5 可知:在 0.5 的分位数下,所有参数估计的 P 值都小于 0.05,通过了显著性检验。变点系数 β 不为 0 表示存在变点,变点位置为 0.498。其他协变量的系数都为正,这表示其他协变量与复发时间间隔正相关。在 0.3 分位数下的估计与在 0.5 分位数下的估计相同,变点位置为 0.485,与在 0.5 分位数下的情况相似。在 0.7 的分位数下,随机化治疗 Z_T 没通过显著性检验,其他协变量都通过了检验,且也存在变点,变点位置为 0.814。在不同分位数下的分段线性删失回归模型的估计都表明:复发时间间隔与治疗时间之间存在变点,前一半治疗时间(在 0.5 的分位数时为 0.498)与后一半治疗时间的治疗效果有显著差异,即前一半治疗时间的复发时间间隔更长。

表 5 估计结果

τ		β_0	β_1	β_2	γ_1	γ_2	γ_3	γ_4	ψ	$\beta_1 + \beta_2$
0.3	估计值	1.793	4.451	- 3.588	0.132	0.051	0.690	0.244	0.485	0.863
	标准差	0.097	0.409	0.395	0.051	0.021	0.040	0.060	0.083	
	P 值	0.000	0.000	0.000	0.020	0.020	0.000	0.000		
0.5	估计值	2.026	3.527	- 2.590	0.382	0.140	0.761	0.403	0.498	0.937
	标准差	0.396	0.919	0.912	0.089	0.036	0.095	0.148	0.151	
	P 值	0.000	0.000	0.007	0.000	0.000	0.000	0.020		
0.7	估计值	2.683	2.318	- 1.921	0.701	0.262	0.457	0.473	0.814	0.397
	标准差	0.369	0.747	0.796	0.233	0.095	0.210	0.451	0.451	
	P 值	0.000	0.001	0.019	0.000	0.000	0.012	0.270		

为更清晰地得到复发时间间隔与治疗时间之间的变化特征,在复发时间和 X_F 散点图上画出不同

分位数下的折线图(见图1)。

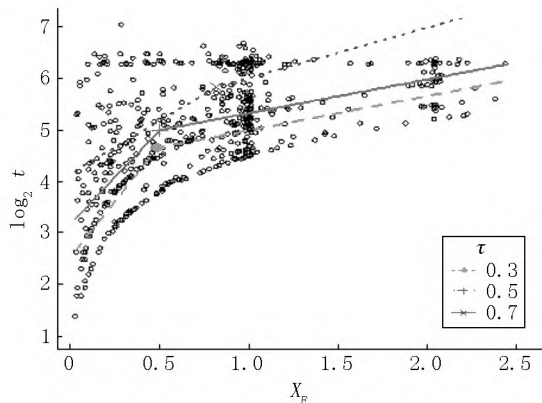


图1 药物复发时间间隔与治疗时间散点图

从图1可以看出:当分位数为0.3和0.5时,治疗时间在0.5附近(估计值分别为0.485和0.498)有显著差异.当分位数为0.7时,治疗时间在0.814处存在变点.变点前的复发时间间隔随着治疗时间的增加而增加,变点后的复发时间间隔明显比变点前的复发时间更小,也就是说,在前半段治疗时间里复发时间间隔更长,即变点之前的治疗比变点之后的治疗更有效。

4 结论

本文基于光滑化的核函数方法研究了分段线性删失分位数回归模型变点位置及模型系数的估计问题.利用核函数替代在目标函数中的非光滑项,通过迭代算法得到变点位置和模型系数估计,推导了估计量的渐近性质.蒙特卡罗模拟结果验证了估计效果的有效性和稳健性.同时利用该方法与线性化技术进行了比较,发现在置信区间较小的情况下,该方法的估计效果更好.针对药物滥用的数据实证分析表明:复发时间间隔与治疗时间存在正向影响,复发时间在0.5的分位数下存在变点,变点位置为0.498,即前一半时间的治疗比后一半时间的治疗更加有效.在0.3和0.7的分位数下也发现变点,且变点前后的治疗效果与在0.5分位数时的治疗效果相近。

5 参考文献

- [1] POWELL J L. Least absolute deviations estimation for the censored regression model [J]. *Journal of Econometrics*, 1984, 25(3): 303-325.
- [2] PORTNOY S. Censored regression quantiles [J]. *Journal of American Statistical Association*, 2003, 98(464): 1001-

- 1012.
- [3] PENG Limin, HUANG Yijian. Survival analysis with quantile regression models [J]. *Journal of the American Statistical Association*, 2008, 103(482): 637-649.
- [4] WANG Huixia Judy, WANG Lan. Locally weighted censored quantile regression [J]. *Journal of American Statistical Association*, 2009, 104(487): 1117-1128.
- [5] TANG Yanlin, WANG Huixia Judy. Penalized regression across multiple quantiles under random censoring [J]. *Journal of Multivariate Analysis*, 2015, 141: 132-146.
- [6] 张倩倩, 郑茜, 王纯杰, 等. 删失分位数回归在医疗费用中的应用 [J]. *数理统计与管理*, 2018, 37(6): 1050-1062.
- [7] 李忠桂, 何书元. 右删失数据下分位数回归的光滑经验似然检验 [J]. *应用概率统计*, 2019, 35(2): 153-164.
- [8] 孙桂萍, 厉诚博, 周勇. 长度偏差右删失数据剩余寿命的分位数回归 [J]. *数学学报(中文版)*, 2020, 63(1): 1-18.
- [9] 张立文, 倪中新, 何勇, 等. 删失分位数回归模型中的变点检测问题 [J]. *中国科学: 数学*, 2018, 48(9): 1159-1180.
- [10] 王江峰, 范国良, 温利民. 删失指标随机缺失下回归函数的复合分位数回归估计 [J]. *系统科学与数学*, 2018, 38(11): 1347-1362.
- [11] LERMAN P M. Fitting segmented regression models by grid search [J]. *Journal of the Royal Statistical Society: Series C(Applied Statistics)*, 1980, 29(1): 77-84.
- [12] LI Chenxi, WEI Ying, CHAPPEL R, et al. Bent line quantile regression with application to an allometric study of land mammals' speed and mass [J]. *Biometrics*, 2011, 67(1): 242-249.
- [13] LI Dong, TONG H. Nested sub-sample search algorithm for estimation of threshold models [J]. *Statistica Sinica*, 2016, 26(4): 1543-1554.
- [14] HOROWITZ J L. A smoothed maximum score estimator for the binary response model [J]. *Econometrica*, 1992, 60(3): 505-531.
- [15] MUGGIO V M R. Estimating regression models with unknown break-points [J]. *Statistics in Medicine*, 2003, 22(19): 3055-3071.
- [16] HIRUKAWA M. Asymmetric kernel smoothing: theory and applications in economics and finance [M]. Singapore: Springer, 2018: 59-71.
- [17] ZHOU Xiaoying, ZHANG Feipeng. Bent line quantile regression via a smoothing technique [J]. *Statistical Analysis and Data Mining*, 2020, 13(3): 216-228.
- [18] 王小刚, 李冰. 基于核函数方法的逐段线性 Tobit 回归模型估计 [J]. *山东大学学报(理学版)*, 2020, 55(6): 1-9.

(下转第290页)

The Measurement Structure of Meaning of Life and Its Computerize Adaptive Test: Based on the Bifactor Model

LIN Jingkai^{1,2}, TU Dongbo^{2*}

(1. Faculty of Education, Beijing Normal University, Beijing 100875, China;

2. School of Psychology, Jiangxi Normal University, Nanchang Jiangxi 330022, China)

Abstract: The main purpose of this study is to explore the measurement structure of meaning of life, then develop a Computerize Adaptive Test of Meaning of Life Measure (CAT-MLM) based on this structure, and discuss its improvement on the reliability and measurement accuracy of traditional paper-pencil test. The experimental results show that CAT-MLM not only has good measurement reliability and validity, but also reduces the number of questions and improves the efficiency of the test without losing the measurement accuracy. Under the same test item length, CAT-MLM can significantly improve the reliability and accuracy of traditional paper and pencil test.

Key words: computerized adaptive test; item response theory; meaning of life

(责任编辑: 冉小晓)

(上接第 276 页)

- [19] HE Xuming, SHAO Qiman. A general Bahadur representation of M -estimators and its application to linear regression [J]. The Annals of Statistics, 1996, 24(6): 2608-2630.
- [20] HENDRICKS W, KOEBKER R. Hierarchical spline models for conditional quantiles and the demand for electricity [J]. Journal of the American Statistical Association, 1992, 87(417): 58-68.
- [21] HALL P, SHEATHER S J. On the distribution of a studentized quantile [J]. Journal of the Royal Statistical Society: Series B (Methodological), 1988, 50(3): 381-391.
- [22] HOSMER DW Jr, LEMESHOW S. Applied survival analysis: regression modeling of time to event data [M]. New York: John Wiley & Sons, 1999.
- [23] KOENKER R. Censored quantile regression redux [J]. Journal of Statistical Software, 2008, 27(6): 1-25.

The Piecewise Linear Censored Quantile Regression Model Estimation Based on Smoothing Technique

WANG Xiaogang, CHEN Jiangmeng

(School of Mathematics and Information Science, North Minzu University, Yinchuan Ningxia 750021, China)

Abstract: The smoothing technique is proposed in the piecewise linear censored quantile regression model to solve the problem of change point, the estimator of change point and coefficients are obtained, and the large sample properties of the estimator is derived. The smoothing technique solves the cumbersome calculation and unreal meaning of the grid search method, and remedies the difficulty that linearization technology cannot prove the asymptotic properties. The validity and robustness of the estimation are verified by Monte Carlo simulation with homoscedasticity and heteroscedasticity, fixed and random censoring at different quantiles. The empirical analysis of drug abuse data shows that the recurrence interval and treatment time have a positive effect, and the recurrence time has a change point at 0.498 (0.5 quantile). The treatment time before 0.498 is longer than after 0.498, that is, the treatment in the first half of the time is more effective.

Key words: smoothing technique; piecewise linear; censored quantile regression model; change point

(责任编辑: 曾剑锋)