

文章编号:1671-9352(2010)08-0123-04

核实数据下响应变量缺失的线性模型均值估计

宇世航

(齐齐哈尔大学理学院, 黑龙江 齐齐哈尔 161006)

摘要:借助核实数据,构造了回归系数的最小二乘估计,然后用借补方法和加权借补方法估计响应变量的均值,最后证明了估计量的渐近正态性。

关键词:随机缺失;线性 EV 模型;加权借补估计;渐近正态性

中图分类号:O212 **文献标志码:**A

Estimators of the mean of linear errors-in-variables models under validation data for missing response data

YU Shi-hang

(Department of Maths, Qiqihar University, Qiqihar 161006, Heilongjiang, China)

Abstract: The least squares estimate for the unknown regression coefficients are constructed with the help of validation data. Then input estimators are used to estimate the mean of the response. It is shown that the proposed estimators are asymptotically normal.

Key words: missing at random; linear errors-in-variables model; weighted imputation estimators; asymptotically normal

0 引言

考虑线性模型

$$Y = X^T \boldsymbol{\beta} + \varepsilon. \quad (0.1)$$

其中, $\boldsymbol{\beta} = (\beta_1, \beta_2, \dots, \beta_p)^T$ 是 $p \times 1$ 的回归系数向量, X 是 p 维随机向量, ε 是随机误差。

通常把带有变量误差的模型称为 EV (errors-variables) 模型,称之为测量误差模型。已有不少文献^[1-2]对线性测量误差的线性 EV 模型进行了研究,这种模型假定 $\tilde{X} = X + e$, 这里 X 是不能被直接观测的解释变量, \tilde{X} 是一个可观测变量。但是,在一般情况下 \tilde{X} 与 X 之间具有很复杂的关系,例如 $\tilde{X} = \varphi(X, e)$, 这里 e 是测量误差且与 (X, Y) 独立, φ 是一个任意已知的函数。在这种情况下要对未知参数 $\boldsymbol{\beta}$ 进行有效统计推断是相当困难的。基于替代数据与核实数据的统计推断已引起国内外学者的重视。例如: Sepanski 与 Lee 和 Xue 研究了基于核实数据的非线性 EV 模型^[3-4], Wang 在这方面做了许多工作^[5-6]。

事实上经常会遇到响应变量缺失的情况,这时通常的统计推断不能直接应用,通常处理缺失数据的方法是对每个缺失的响应变量值进行借补,然后利用标准方法和借补后的完全数据进行研究。对于线性模型在响应变量随机缺失和删失的情形, Wang 和 Qin 分别对这两种情况进行了研究^[7-8]。而对协变量带有测量误差且响应变量缺失的情况,目前的文献还较少。本文借助于核实数据,利用借补方法和加权借补方法估计响应变量的均值,并证明了其估计量的渐近正态性。

收稿日期:2009-10-26

基金项目:黑龙江省教育厅科学技术研究项目(11551543)

作者简介:宇世航(1972-),女,硕士,副教授,主要研究方向为数理统计、概率极限理论. Email: qqhrysh@163.com

1 回归系数的估计

假设主要数据 $\{(\tilde{X}_i, Y_i, \delta_i)_{i=1}^n\}$ 是独立同分布(i. i. d)的样本, 且与 i. i. d 的核实样本 $\{(\mathbf{X}_i, \tilde{X}_i)_{i=n+1}^{n+m}\}$ 相互独立, 其中当 Y_i 缺失时 $\delta_i = 0$, 否则 $\delta_i = 1$ 。进一步假定 $E(\varepsilon | \tilde{X}) = 0$, 令 $u(\tilde{X}) = E(\mathbf{X} | \tilde{X})$, 则基于完全观测数据, 模型(0.1)可改写为

$$\delta Y = \delta u^T(\tilde{X})\boldsymbol{\beta} + \delta\eta, \quad (1.1)$$

其中 $\eta = \varepsilon + \mathbf{X}^T\boldsymbol{\beta} - u^T(\tilde{X})\boldsymbol{\beta}$ 。

利用核实数据, 在 $\tilde{X} = \tilde{x}$ 的情况下用 \mathbf{X} 的核回归来估计模型(1.1)的回归系数 $u(\tilde{X})$, $u(\tilde{X})$ 的估计定义为

$$\hat{u}_m(\tilde{x}) = \frac{\sum_{i=n+1}^{n+m} \mathbf{X}_i K\left(\frac{\tilde{X}_i - \tilde{x}}{h_m}\right)}{\sum_{i=n+1}^{n+m} K\left(\frac{\tilde{X}_i - \tilde{x}}{h_m}\right)}, \quad (1.2)$$

其中 $K(\cdot)$ 为核函数, h_m 是收敛于 0 的窗宽。

定义 $\boldsymbol{\beta}$ 的估计量是

$$S_{m,n}(\boldsymbol{\beta}) = \frac{1}{n} \sum_{j=1}^n \delta_j (\hat{u}_m(\tilde{X}_j)\boldsymbol{\beta})^2 + \frac{1}{m} \sum_{i=n+1}^{n+m} \delta_i (Y_i - \mathbf{X}_i^T \boldsymbol{\beta})^2 \quad (1.3)$$

的最小值。由式(1.3)可得 $\boldsymbol{\beta}$ 的估计量为

$$\hat{\boldsymbol{\beta}}_{m,n} = \hat{\boldsymbol{\Sigma}}_{m,n}^{-1} \hat{\mathbf{A}}_{m,n}, \quad (1.4)$$

这里, $\hat{\boldsymbol{\Sigma}}_{m,n} = \frac{1}{n} \sum_{j=1}^n \delta_j \hat{u}_m(\tilde{X}_j) \hat{u}_m^T(\tilde{X}_j) + \frac{1}{m} \sum_{i=n+1}^{n+m} \delta_i \mathbf{X}_i \mathbf{X}_i^T$, $\hat{\mathbf{A}}_{m,n} = \frac{1}{n} \sum_{j=1}^n \delta_j \hat{u}_m(\tilde{X}_j) Y_j + \frac{1}{m} \sum_{i=n+1}^{n+m} \delta_i \mathbf{X}_i Y_i$ 。

令 D^k 是定义在 \mathbf{R}^p 上的连续函数族, 并且满足 $\frac{\partial^{i_1}}{\partial x_1^{i_1}} \frac{\partial^{i_2}}{\partial x_2^{i_2}} \cdots \frac{\partial^{i_p}}{\partial x_p^{i_p}} f(x_1, \dots, x_p)$ 一致有界。对任何向量 $\boldsymbol{\alpha}$, 用

$\|\boldsymbol{\alpha}\|$ 表示 Euclidean 模。给出正则条件(下列条件统称为条件 A):

$$(A \cdot X) E \|\mathbf{X}\|^2 < \infty。$$

$$(A \cdot Y) E Y^2 < \infty。$$

$$(A \cdot \varepsilon) \sup_{\tilde{x}} E[\varepsilon^2 | \tilde{X} = \tilde{x}] < \infty。$$

$$(A \cdot u) u(\cdot) \in D^k, k > p + 1。$$

(A · K) $K(\cdot)$ 是非负有界的 k 阶核函数, 且具有有界支撑。

(A · \tilde{X}) (1) 存在一正常数序列 η_n 满足 $nP(f_{\tilde{x}}(\tilde{X} < \eta_n) \rightarrow 0$; (2) $f_{\tilde{x}} \in D^k$ 。

(A · h_m) (1) $mh_m^{2p} \eta_n^2 \rightarrow \infty$; (2) $mh_m^{2k} \rightarrow 0$ 。

(A · nm) $\frac{n}{m} \rightarrow \lambda, \lambda \geq 0$ 为常数。

(A · $\boldsymbol{\Sigma}$) $\boldsymbol{\Sigma} = E[\delta u(\tilde{X}) u^T(\tilde{X})] + E[\delta \mathbf{X} \mathbf{X}^T]$ 为正定。

定理 1.1 假设条件 A 成立, 则有

$$\sqrt{n}(\hat{\boldsymbol{\beta}} - \boldsymbol{\beta}) \xrightarrow{d} N(0, \boldsymbol{\Sigma}^{-1} V \boldsymbol{\Sigma}^T),$$

其中 $V = E[\delta u(\tilde{X}) u^T(\tilde{X}) (Y - u^T(\tilde{X})\boldsymbol{\beta})^2] + \lambda \{E[\delta u(\tilde{X}) u^T(\tilde{X}) [\delta((\mathbf{X} - u^T(\tilde{X}))^T \boldsymbol{\beta})]^2] + E[\delta \mathbf{X} \mathbf{X}^T (Y - \mathbf{X}^T \boldsymbol{\beta})^2] + 2E[(Y - \mathbf{X}^T \boldsymbol{\beta}) ((\mathbf{X} - u^T(\tilde{X}))^T \boldsymbol{\beta}) \mathbf{X} u^T(\tilde{X})]\}$ 。

这里 $\hat{\boldsymbol{\beta}}_{m,n}$ 的渐近方差由 $\hat{V} = \hat{\boldsymbol{\Sigma}}_{n,m}^{-1} [\hat{V}_1 + \hat{V}_2] \hat{\boldsymbol{\Sigma}}_{n,m}^T$ 一致估计, $\hat{\boldsymbol{\Sigma}}_{n,m}^{-1}$ 同式(1.4)定义, 并且

$$\hat{V}_1 = \frac{1}{n} \sum_{j=1}^n [\delta_j \hat{u}_m(\tilde{X}_j) \hat{u}_m^T(\tilde{X}_j) (Y_j - \hat{u}_m^T(\tilde{X}_j) \hat{\boldsymbol{\beta}}_{m,n})^2],$$

$$\begin{aligned} \hat{V}_2 = & \frac{n}{m^2} \sum_{i=1}^m \{ [\delta_i \hat{u}_m(\tilde{X}_i) \hat{u}_m^T(\tilde{X}_i) [\delta_i ((\mathbf{X}_i - \hat{u}_i^T(\tilde{X}_i))^T \hat{\boldsymbol{\beta}}_{m,n})]^2] + [\delta_i \mathbf{X}_i \mathbf{X}_i^T (Y_i - \mathbf{X}_i^T \hat{\boldsymbol{\beta}}_{m,n})^2] + \\ & 2[(Y_i - \mathbf{X}_i^T \hat{\boldsymbol{\beta}}_{m,n}) ((\mathbf{X}_i - \hat{u}_i^T(\tilde{X}_i))^T \hat{\boldsymbol{\beta}}_{m,n}) \mathbf{X}_i \hat{u}_i^T(\tilde{X}_i)] \}. \end{aligned}$$

2 响应均值的估计

记响应均值 $EY = \theta$, 下面基于上述回归系数的估计量 $\hat{\boldsymbol{\beta}}_{m,n}$, 给出 θ 的估计量及其渐近性质。

定义 2.1 $\hat{\theta}_l = \frac{1}{n} \sum_{i=1}^n \{ \delta_i Y_i + (1 - \delta_i) \hat{u}_m^T(\tilde{X}_i) \hat{\boldsymbol{\beta}}_{m,n} \}$ 。

定义 2.2 $\hat{\theta}_p = \frac{1}{n} \sum_{i=1}^n \left\{ \frac{\delta_i}{\hat{P}(\tilde{X}_i)} Y_i + \left(1 - \frac{\delta_i}{\hat{P}(\tilde{X}_i)} \right) \hat{u}_m^T(\tilde{X}_i) \hat{\boldsymbol{\beta}}_{m,n} \right\}$,

这里 $\hat{P}(\tilde{X}_i) = \frac{\sum_{i=1}^n \delta_i K\left(\frac{\tilde{X}_i - \bar{x}}{b_n}\right)}{\sum_{i=1}^n K\left(\frac{\tilde{X}_i - \bar{x}}{b_n}\right)}$ 为 $P(\tilde{X}) = P(\delta_i = 1 | \tilde{X}_i = \bar{x})$ 的核估计。

定理 2.1 在条件 A 成立下, 有

$$\sqrt{n}(\hat{\theta}_l - \theta) \xrightarrow{d} N(0, \mathbf{V}_l),$$

其中 $\mathbf{V}_l = \boldsymbol{\Sigma}_1 + \boldsymbol{\beta}^T \boldsymbol{\Sigma}_3 \boldsymbol{\beta} - 2\boldsymbol{\Sigma}_4^T \boldsymbol{\beta} \theta + \theta^2 + \boldsymbol{\Sigma}_2^T (\boldsymbol{\Sigma}^{-1} \mathbf{V} \boldsymbol{\Sigma}^T) \boldsymbol{\Sigma}_2 + \boldsymbol{\Sigma}_2^T (\boldsymbol{\Sigma}^{-1} \mathbf{V} \boldsymbol{\Sigma}^T) \boldsymbol{\Sigma}_5$ 。记 $\boldsymbol{\Sigma}_1 = E[\delta(Y - u^T(\tilde{X})\boldsymbol{\beta})^2]$, $\boldsymbol{\Sigma}_2 = E[(1 - \delta)u(\tilde{X})]$, $\boldsymbol{\Sigma}_3 = E[u(\tilde{X})u^T(\tilde{X})]$, $\boldsymbol{\Sigma}_4 = E[u(\tilde{X})]$, $\boldsymbol{\Sigma}_5 = E[\delta u(\tilde{X})]$ 。这里 θ_l 的渐近方差 \mathbf{V}_l 的估计 $\hat{\mathbf{V}}_l$ 为

$$\hat{\mathbf{V}}_l = \hat{\boldsymbol{\Sigma}}_1 + \hat{\boldsymbol{\beta}}^T \hat{\boldsymbol{\Sigma}}_3 \hat{\boldsymbol{\beta}} - 2\hat{\boldsymbol{\Sigma}}_4^T \hat{\boldsymbol{\beta}} \hat{\theta} + \hat{\theta}^2 + \hat{\boldsymbol{\Sigma}}_2^T (\hat{\boldsymbol{\Sigma}}^{-1} \hat{\mathbf{V}} \hat{\boldsymbol{\Sigma}}^T) \hat{\boldsymbol{\Sigma}}_2 + \hat{\boldsymbol{\Sigma}}_2^T (\hat{\boldsymbol{\Sigma}}^{-1} \hat{\mathbf{V}} \hat{\boldsymbol{\Sigma}}^T) \hat{\boldsymbol{\Sigma}}_5,$$

且

$$\hat{\boldsymbol{\Sigma}}_1 = \frac{1}{n} \sum_{i=1}^n [\delta_i (Y_i - \hat{u}_m^T(\tilde{X}_i) \hat{\boldsymbol{\beta}})^2], \hat{\boldsymbol{\Sigma}}_2 = \frac{1}{n} \sum_{i=1}^n [(1 - \delta_i) \hat{u}_m(\tilde{X}_i)],$$

$$\hat{\boldsymbol{\Sigma}}_3 = \frac{1}{n} \sum_{i=1}^n [\hat{u}_m(\tilde{X}_i) \hat{u}_m^T(\tilde{X}_i)], \hat{\boldsymbol{\Sigma}}_4 = \frac{1}{n} \sum_{i=1}^n [\hat{u}_m(\tilde{X}_i)], \hat{\boldsymbol{\Sigma}}_5 = \frac{1}{n} \sum_{i=1}^n [\delta_i \hat{u}_m(\tilde{X}_i)].$$

定理 2.2 在条件 A 成立下, 当 $nb_n^4 \rightarrow 0, nb_n \rightarrow \infty$ 时, 有

$$\sqrt{n}(\hat{\theta}_p - \theta) \xrightarrow{d} N(0, \mathbf{V}_p),$$

其中, $\mathbf{V}_p = \boldsymbol{\Sigma}_{1p} + \boldsymbol{\beta}^T \boldsymbol{\Sigma}_3 \boldsymbol{\beta} - 2\boldsymbol{\Sigma}_4^T \boldsymbol{\beta} \theta + \theta^2 + \boldsymbol{\Sigma}_{2p}^T (\boldsymbol{\Sigma}^{-1} \mathbf{V} \boldsymbol{\Sigma}^T) \boldsymbol{\Sigma}_{2p} + \boldsymbol{\Sigma}_{2p}^T (\boldsymbol{\Sigma}^{-1} \mathbf{V} \boldsymbol{\Sigma}^T) \boldsymbol{\Sigma}_{5p}$, 且

$$\boldsymbol{\Sigma}_{1p} = E\left(\frac{\delta(Y - u^T(\tilde{X})\boldsymbol{\beta})^2}{P(\tilde{X})}\right) = \frac{\boldsymbol{\Sigma}_1}{E(\delta|\tilde{X})}, \boldsymbol{\Sigma}_{2p} = E\left(u(\tilde{X}) - \frac{E(\delta u|\tilde{X})}{E(\delta|\tilde{X})}\right), \boldsymbol{\Sigma}_{5p} = \frac{\boldsymbol{\Sigma}_5}{E(\delta|\tilde{X})}。$$

3 定理的证明

定理 1.1 的证明类似于文献[6]中定理 2.1 的证明。

定理 2.1 的证明

$$\begin{aligned} \sqrt{n}(\hat{\theta}_l - \theta) &= \frac{1}{\sqrt{n}} \sum_{i=1}^n [(\delta_i Y_i + (1 - \delta_i) \hat{u}_m^T(\tilde{X}_i) \hat{\boldsymbol{\beta}}_{m,n}) - \theta] = \\ &= \frac{1}{\sqrt{n}} \sum_{i=1}^n \delta_i (Y_i - u^T(\tilde{X}_i) \boldsymbol{\beta}) + \frac{1}{\sqrt{n}} \sum_{i=1}^n (1 - \delta_i) u^T(\tilde{X}_i) (\hat{\boldsymbol{\beta}}_{m,n} - \boldsymbol{\beta}) + \frac{1}{\sqrt{n}} \sum_{i=1}^n (u^T(\tilde{X}_i) \boldsymbol{\beta} - \theta) + \\ &= \frac{1}{\sqrt{n}} \sum_{i=1}^n (1 - \delta_i) (\hat{u}_m(\tilde{X}_i) - u(\tilde{X}_i))^T \hat{\boldsymbol{\beta}}_{m,n} = \\ &= B_{n1} + B_{n2} + B_{n3} + o_p(1), \end{aligned}$$

由中心极限定理, 显然有

$$\begin{aligned} B_{n1} &\xrightarrow{d} N(0, \boldsymbol{\Sigma}_1), \\ B_{n3} &\xrightarrow{d} N(0, \boldsymbol{\beta}^T \boldsymbol{\Sigma}_3 \boldsymbol{\beta} - 2\boldsymbol{\Sigma}_4^T \boldsymbol{\beta} \theta + \theta^2). \end{aligned}$$

再由定理 1.1, 有

$$B_{n_2} \xrightarrow{d} N(0, \Sigma_2^T (\Sigma^{-1} V \Sigma^{-T}) \Sigma_2)。$$

经简单推导得

$$\text{cov}(B_{n_1}, B_{n_3}) = 0, \text{cov}(B_{n_2}, B_{n_3}) = 0, \text{cov}(B_{n_1}, B_{n_2}) = \Sigma_2^T (\Sigma^{-1} V \Sigma^{-T}) \Sigma_5。$$

定理得证。

定理 2.2 的证明

$$\begin{aligned} \sqrt{n}(\hat{\theta}_p - \theta) &= \frac{1}{\sqrt{n} \sum_{i=1}^n} \left[\frac{\delta_i}{\hat{P}(\tilde{X}_i)} Y_i + \left(1 - \frac{\delta_i}{\hat{P}(\tilde{X}_i)} \right) \hat{u}_m^T(\tilde{X}_i) \hat{\beta}_{m,n} - \theta \right] = \\ &= \frac{1}{\sqrt{n} \sum_{i=1}^n} \left[\frac{\delta_i}{P(\tilde{X}_i)} Y_i + \left(1 - \frac{\delta_i}{P(\tilde{X}_i)} \right) u_m^T(\tilde{X}_i) \hat{\beta}_{m,n} + \right. \\ &\quad \left. \frac{\delta_i [\hat{P}(\tilde{X}_i) - P(\tilde{X}_i)]}{P^2(\tilde{X}_i)} (Y_i - \hat{u}_m^T(\tilde{X}_i) \hat{\beta}_{m,n}) - \theta \right] + o_p(1) = \\ &= \frac{1}{\sqrt{n} \sum_{i=1}^n} \frac{\delta_i}{P(\tilde{X}_i)} (Y_i - u^T(\tilde{X}_i) \beta) + \frac{1}{\sqrt{n} \sum_{i=1}^n} \left(1 - \frac{\delta_i}{P(\tilde{X}_i)} \right) u^T(\tilde{X}_i) (\hat{\beta}_{m,n} - \beta) + \\ &= \frac{1}{\sqrt{n} \sum_{i=1}^n} (u^T(\tilde{X}_i) \beta - \theta) + \frac{1}{\sqrt{n} \sum_{i=1}^n} \frac{\delta_i [\hat{P}(\tilde{X}_i) - P(\tilde{X}_i)]}{P^2(\tilde{X}_i)} (Y_i - u^T(\tilde{X}_i) \beta) + o_p(1) = \\ &= T_{n1} + T_{n2} + T_{n3} + T_{n4} + o_p(1)。 \end{aligned}$$

类似文献[7]的证明有 $T_{n4} = o_p(1)$,

由中心极限定理有

$$\begin{aligned} T_{n1} &\xrightarrow{d} N(0, \Sigma_{1p}), \\ T_{n3} &\xrightarrow{d} N(0, \beta^T \Sigma_3 \beta - 2 \Sigma_4^T \beta \theta + \theta^2)。 \end{aligned}$$

再由定理 1.1 有

$$T_{n2} \xrightarrow{d} N(0, \Sigma_{2p}^T (\Sigma^{-1} V \Sigma^{-T}) \Sigma_{2p})。$$

由
$$E \left(\frac{\delta(Y - u^T(\tilde{X})\beta)}{P(\tilde{X})} (u^T(\tilde{X})\beta - \theta) \right) = E \left\{ E \left[\left(\frac{\delta(Y - u^T(\tilde{X})\beta)}{P(\tilde{X})} (u^T(\tilde{X})\beta - \theta) \right) \middle| Y, \tilde{X} \right] \right\} = E [(Y - u^T(\tilde{X})\beta) (u^T(\tilde{X})\beta - \theta) | \tilde{X}] = 0$$

可知 $\text{cov}(T_{n1}, T_{n3}) = 0$, 同理可得 $\text{cov}(T_{n2}, T_{n3}) = 0$, 从而易知 $\text{cov}(T_{n1}, T_{n2}) = \Sigma_{2p}^T (\Sigma^{-1} V \Sigma^{-T}) \Sigma_{5p}$ 。定理得证。

参考文献:

- [1] FULLER W A. Measurement error models[M]. New York: Wiley, 1987.
- [2] CUI H J, CHEN S X. Empirical likelihood confidence region for parameter in the reeors-in-variables models[J]. Multivariate Anal, 2003, 84(1):101-115.
- [3] SEPANSKI J H, LEE L F. Semiparametric estimation of nonlinear error-in-variables models with validation study [J]. Nonparametric Statist, 1995, 4:365-394.
- [4] XUE L G. Empirical likelihood inference in nonlinear semiparametric EV models with validation data[J]. Acta Mathematica Sinica, Chinese Series, 2006, 49(1):145-154.
- [5] WANG Q H. Dimension reduction in partly linear error-in-response models with validation data[J]. Multivariate Anal, 2003, 85:234-252.
- [6] WANG Q H. Estimation of partial linear error-in-variables models with validation data[J]. Multivariate Anal, 1999, 69:30-64.
- [7] WANG Q H, LINTON Oliver. Semiparametric regression analysis with missing response at random[J]. Amer Statist Assoc, 2004, 99:334-345.
- [8] QIN G S, JING B Y. Empirical likelihood for censored linear regression[J]. Scan J Statist, 2001, 28:661-673.

(编辑:孙培芹)