

因变量缺失下线性回归模型的随机约束估计

李 静¹, 安佰玲²

¹中国劳动关系学院应用技术学院, 北京

²淮北师范大学数学科学学院, 安徽 淮北

收稿日期: 2022年11月27日; 录用日期: 2022年12月23日; 发布日期: 2022年12月30日

摘 要

本文研究了线性回归模型在因变量存在缺失的同时回归系数附加有随机约束条件时的估计问题, 基于完整数据方法和单点插补方法, 给出了模型系数的两种约束估计, 并研究了这两类估计量的渐近正态性。最后通过数值模拟验证了所提估计方法的有效性。

关键词

线性回归模型, 缺失数据, 插补方法, 随机约束估计

Stochastic Restricted Estimation of Linear Regression Models with Missing Responses

Jing Li¹, Bailing An²

¹School of Applied Technology, China University of Labor Relations, Beijing

²School of Mathematical Sciences, Huaibei Normal University, Huaibei Anhui

Received: Nov. 27th, 2022; accepted: Dec. 23rd, 2022; published: Dec. 30th, 2022

Abstract

This paper discusses estimation of linear regression models in the presence of multicollinearity and there are stochastic linear restrictions on the regression coefficients. Based on the complete-case method and single imputation technique, two stochastic restricted estimators are proposed. Asymptotically, properties of the proposed estimators are shown. Finally, some simulations are conducted to illustrate the proposed methods.

Keywords

Linear Regression, Missing Data, Imputation Method, Stochastic Restricted Estimation

Copyright © 2022 by author(s) and Hans Publishers Inc.

This work is licensed under the Creative Commons Attribution International License (CC BY 4.0).

<http://creativecommons.org/licenses/by/4.0/>



Open Access

1. 引言

线性回归模型作为最为常用的统计分析之一, 已经得到了深入的理论研究和广泛的应用。众所周知, 对于线性回归模型, 可以采用最小二乘估计方法或极大似然估计方法等。然而, 这些估计方法一般只利用了样本总体假定和样本信息。在分析一些实际领域的问题时, 常常可以通过该问题的专业背景以及前期研究等渠道获得一些额外信息, 这些额外信息经常被转化为模型中回归系数的约束条件, 比如某个回归系数按照实际问题来说应该是正值, 或者两个或多个回归系数之间存在某种关系, 详细介绍可参考 Toutenburg (1982) [1]。那么, 在模型的估计中如果能够将这些约束条件充分利用, 可以提高估计量的有效性。对于一般线性回归模型, 如果约束条件为针对系数的线性约束, 那么此时采用约束最小二乘估计是比普通最小二乘法更有效。除了线性约束之外, 还有随机线性约束以及不等式约束等情况, 具体内容可参考 Rao 和 Toutenburg (1999) [2]。

另一方面, 实际数据分析中, 数据的缺失是非常普遍的, 关于缺失数据的详细介绍可参考著作 Little 和 Rubin (1987) [3]。缺失数据下回归模型的研究得到了关注, 相关研究可参考 Cheng (1994) [4]、Chu & Cheng (1995) [5]、Wang & Rao (2002) [6]、Wang & Rao (2004) [7]和 Qin 等(2009) [8]等有关文献。

目前关于线性回归模型的约束估计, 大都是基于无缺失的情形, 关于缺失数据下约束估计的研究还很少。杨徐佳等(2011) [9]构造了因变量缺失下线性回归模型的估计。安佰玲等(2013) [10]在此基础上讨论了模型的精确约束估计问题。本文将在这两篇论文的基础上考虑因变量缺失下线性回归模型的随机约束估计。

考虑如下的线性回归模型

$$Y_i = X_i^T \beta + \varepsilon_i, \quad i = 1, 2, \dots, n \quad (1)$$

其中, Y_i 为因变量观测值, $X_i = (X_{i1}, X_{i2}, \dots, X_{iq})^T$, 为对应的自变量的观测值, β 为 $q \times 1$ 的未知参数向量, 模型误差 ε_i 为独立同分布的随机变量, 且有 $E(\varepsilon_i | X_i) = 0$ 和 $\text{Var}(\varepsilon_i | X_i) = \sigma^2$ 。为了方便介绍, 我们引入一个新的变量 δ 作为缺失的标志, $\delta_i = 1$ 表示 Y_i 的值可以被观测到, 而 $\delta_i = 0$ 表示 Y_i 值缺失。这里, 我们假定 Y_i 满足随机缺失的机制, 即

$$P(\delta = 1 | X, Y) = P(\delta = 1 | X). \quad (2)$$

该缺失机制是缺失数据分析中常用的假设条件。

考虑如下的随机约束条件

$$b = A\beta + \eta \quad (3)$$

其中 A 是 $k \times p$ 维的已知矩阵, 且 $\text{rank}(A) = k$, b 是 $k \times 1$ 维的已知向量。 η 为均值为 0, 协方差为 $\sigma^2 \Omega$ 的随机向量, 其中 Ω 为一已知正定矩阵。为了在线性回归模型的估计中考虑随机约束条件(3), Durbin (1953) [11], Theil 和 Goldberger (1961) [12]以及 Theil (1963) [13]提出了混合估计方法。

本文重点研究线性回归模型(1)在因变量缺失机制(2)和随机约束条件(3)下的估计问题, 从而将杨徐佳等(2011) [9]和安佰玲等(2013) [10]的结果推广到了随机约束情形。

第2节和第3节将分别基于完整数据方法和单点插补方法构造模型系数的随机约束估计, 并给出估计量的渐近性质。第4节将通过数值模拟验证所提方法的有效性, 定理的证明将放在第5节。

2. 基于完整数据方法的随机约束估计

完整数据方法就是只利用因变量和自变量都观测完整的数据, 将因变量存在缺失的那些观测值舍弃不用。假设独立同分布样本数据 $\{Y_i, \delta_i, X_i\}_{i=1}^n$ 来自模型(1), 则有

$$\delta_i Y_i = \delta_i X_i^T \beta + \delta_i \varepsilon_i, \quad i = 1, 2, \dots, n. \quad (4)$$

显然模型(4)可记为如下矩阵形式

$$\bar{Y} = \bar{X} \beta + \bar{\varepsilon}, \quad (5)$$

其中 $\bar{Y} = (\delta_1 Y_1, \delta_2 Y_2, \dots, \delta_n Y_n)^T$, $\bar{X} = (\delta_1 X_1, \delta_2 X_2, \dots, \delta_n X_n)^T$, $\bar{\varepsilon} = (\delta_1 \varepsilon_1, \delta_2 \varepsilon_2, \dots, \delta_n \varepsilon_n)^T$ 。

基于最小二乘法, 可得 β 基于完整数据方法的估计为

$$\hat{\beta}_C = \left(\sum_{i=1}^n \delta_i X_i X_i^T \right)^{-1} \left(\sum_{i=1}^n \delta_i X_i Y_i \right) = (\bar{X}^T \bar{X})^{-1} \bar{X}^T \bar{Y}. \quad (6)$$

下面基于文[11] [12] [13]中的混合估计方法来估计随机约束条件(3)下的线性回归模型(5)。构造的如下的辅助函数:

$$F(\beta) = (\bar{Y} - \bar{X} \beta)^T (\bar{Y} - \bar{X} \beta) + (b - A\beta)^T \Omega^{-1} (b - A\beta) \quad (7)$$

用辅助函数对于 β 求导数并令其为0, 可得:

$$\frac{\partial F(\beta)}{\partial \beta} = -2\bar{X}^T (\bar{Y} - \bar{X} \beta) - 2A^T \Omega^{-1} (b - A\beta)^T = 0. \quad (8)$$

简单整理可得 β 基于完整数据方法的随机约束估计

$$\hat{\beta}_C^{SR} = (\bar{X}^T \bar{X} + A^T \Omega^{-1} A)^{-1} (\bar{X}^T \bar{Y} + A^T \Omega^{-1} b). \quad (9)$$

根据文[9]中的定理 A.18, 可得:

$$(\bar{X}^T \bar{X} + A^T \Omega^{-1} A)^{-1} = (\bar{X}^T \bar{X})^{-1} - (\bar{X}^T \bar{X})^{-1} A^T \left[\Omega + A (\bar{X}^T \bar{X})^{-1} A^T \right]^{-1} A (\bar{X}^T \bar{X})^{-1} \quad (10)$$

则由(6)、(9)和(10), $\hat{\beta}_C^{SR}$ 可以等价表示为如下的形式

$$\hat{\beta}_C^{RS} = \hat{\beta}_C - (\bar{X}^T \bar{X})^{-1} A^T \left[\Omega + A (\bar{X}^T \bar{X})^{-1} A^T \right]^{-1} (A \hat{\beta}_C - b) \quad (11)$$

下面给出关于 $\hat{\beta}_C^{RS}$ 的渐近性质。

定理 1. 如果第5节的假设条件成立, $\hat{\beta}_C^{RS}$ 是渐近正态的, 有

$$\sqrt{n} (\hat{\beta}_C^{RS} - \beta) \xrightarrow{d} N(0, \sigma^2 \Sigma^{-1}),$$

其中 $\Sigma = E(\delta X_1 X_1^T)$ 。

显然, 我们所构造的因变量缺失下的随机约束估计 $\hat{\beta}_C^{RS}$ 的渐近性质与杨徐佳等(2011) [9]中构造的精确约束最小二乘估计的渐近性质相同。

3. 基于单点插补方法的约束估计及其性质

为了弥补完整数据方法丢弃数据从而损失信息的不足, 下面基于单点插补方法构造模型系数的随机约束估计。基于上一节得到的最小二乘估计 $\hat{\beta}_C$, 针对因变量存在缺失这一情况, 定义如下的新的因变量

$$Y_i^* = \delta_i Y_i + (1 - \delta_i) X_i^T \hat{\beta}_C, \quad (12)$$

显然, 当 Y_i 没有缺失时, $Y_i^* = Y_i$, 而当 Y_i 存在缺失时, 用插补值 $Y_i^* = X_i^T \hat{\beta}_C$ 代替其观测值。

基于构造的数据集 $(Y_i^*, X_i)_{i=1}^n$, 有如下的线性模型

$$Y_i^* = X_i^T \beta + e_i, \quad i = 1, 2, \dots, n. \quad (13)$$

其中 $Y_i^* = \delta_i Y_i + (1 - \delta_i) X_i^T \hat{\beta}_C$ 。

对模型(13)使用最小二乘方法, 得到 β 的单点插补估计

$$\hat{\beta}_I = \left(\sum_{i=1}^n X_i X_i^T \right)^{-1} \left(\sum_{i=1}^n X_i Y_i^* \right). \quad (14)$$

考虑随机约束条件(3), 构造如下的辅助函数

$$F^*(\beta, \lambda) = \sum_{i=1}^n (Y_i^* - X_i^T \beta)^2 + (b - A\beta)^T \Omega^{-1} (b - A\beta). \quad (15)$$

同第2节类似, 基于上面的辅助函数, 可得 β 单点插补的随机约束估计为

$$\hat{\beta}_I^{SR} = (X^T X + A^T \Omega^{-1} A)^{-1} (X^T Y^* + A^T \Omega^{-1} b) \quad (16)$$

同样, 该估计可以等价表达为

$$\hat{\beta}_I^{RS} = \hat{\beta}_I - (X^T X)^{-1} A^T \left[\Omega + A (X^T X)^{-1} A^T \right]^{-1} (A \hat{\beta}_I - b). \quad (17)$$

下面给出 $\hat{\beta}_I^{RS}$ 的渐近性质。

定理 2. 如果第5节的假设条件成立, $\hat{\beta}_I^{RS}$ 是渐近正态的, 满足

$$\sqrt{n} (\hat{\beta}_I^{RS} - \beta) \xrightarrow{d} N(0, \sigma^2 \Sigma^{-1}).$$

需要注意的是 $\hat{\beta}_I^{RS}$ 的渐近性质也与杨徐佳等(2013) [9]中构造的完整数据最小二乘估计的渐近性质相同。

4. 数值模拟

本节我们将通过数值模拟考察前面所提出估计方法的有效性。假设数据产生于如下的模型

$$y_i = x_{1i} \beta_1 + x_{2i} \beta_2 + \varepsilon_i \quad i = 1, 2, \dots, n \quad (18)$$

其中 $x_{1i} \sim N(1, 1)$, $x_{2i} \sim U(0, 1)$, (β_1, β_2) 的真实值为 $(3, 2)$ 。为了考察模型误差的分布对估计和检验结果影响, 考虑如下两种情况(1) $\varepsilon_i \sim N(0, 0.5^2)$; (2) $\varepsilon_i \sim U(-\sqrt{3}/2, \sqrt{3}/2)$, 表1中分别用 N 和 U 表示。缺失的机制采用安佰玲等(2013) [10]中的设置, 当 $|x_{1i} - 1| + |x_{2i} - 0.5| \leq 1$ 时

$\Delta = p(\delta = 1 | x_1 = x_{1i}, x_2 = x_{2i}) = 0.8 + 0.2(|x_{1i} - 1| + |x_{2i} - 0.5|)$; 其余情况等于 0.9。

随机约束设定为 $\beta_1 + \beta_2 = 5 + \eta$, $E\eta = 0$, $\text{Var}(\eta) = 0.49$, 即 $A = (1, 1)$, $b = 5$, $\Omega = 0.49$ 。

针对模型(18), 基于上面的各种缺失机制和误差分布, 我们取样本量 n 分别为 50、100 和 150, 在每一种设定下分别求取 (β_1, β_2) 的基于完整数据分析方法的估计(表1中用 C 表示, 下同)和随机约束估计

(C-SR), 基于单点插补方法的估计(I)和随机约束估计(I-SR)。每种情况重复计算500次, 以这些估计量的均方误差(EMSE)来衡量其表现,

$$\text{EMSE}(\beta^*) = \frac{1}{500} \sum_{k=1}^{500} \sum_{j=1}^2 (\beta_{kj}^* - \beta_j)^2$$

其中 β_{kj}^* 是参数 β_j 的第 k 次重复时的估计值, 模拟结果见表 1。

Table 1. EMSEs of the estimators

表 1. 不同估计量的 MSE

β	$n = 50$		$n = 100$		$n = 150$	
	N	U	N	U	N	U
C	0.0305	0.0342	0.0146	0.0169	0.0112	0.0101
C-RS	0.0253	0.0281	0.0133	0.0152	0.0105	0.0095
I	0.0321	0.0330	0.0158	0.0155	0.0100	0.0105
I-RS	0.0266	0.0274	0.0144	0.0142	0.0094	0.0099

从模拟结果可以看出: (1) 这四类估计量的均方误差都随着样本量的增加而减小, 此外, 估计值对于误差分布的改变几乎没有变化。(2) 同样设置下, 单点插补方法的均方误差一般小于完整估计方法。(3) 同样设置下, 考虑了约束条件的随机约束估计的表现优于没有考虑约束条件时的估计。

5. 定理的证明

在给出定理的证明之前, 我们先给出下面条件。

条件 1: $\Sigma = E(\delta_1 X_1 X_1^T)$ 为正定矩阵。

条件 2: $E(\varepsilon | X) = 0$, $E(|\varepsilon|^3 | X) < \infty$ 。

引理 1. 如果前面的假设条件成立, $\hat{\beta}_c$ 和 $\hat{\beta}_l$ 都是渐进正态的, 二者都满足

$$\sqrt{n}(\hat{\beta} - \beta) \xrightarrow{d} N(0, \sigma^2 \Omega^{-1}),$$

其中 $\hat{\beta}$ 为 $\hat{\beta}_c$ 或 $\hat{\beta}_l$ 。

证明: 该引理即为杨徐佳等(2011) [9]中的定理 1。

引理 2. 如果前面的假设条件成立, 有

$$\frac{1}{n} \bar{X}^T \bar{X} \xrightarrow{p} \Sigma, \quad \frac{1}{n} X^T X \xrightarrow{p} EX_1 X_1^T$$

定理 1 的证明: 定义 $S_{SR} = \bar{X}^T \bar{X} + A^T \Omega^{-1} A$, $\Delta = \text{diag}(\delta_1, \delta_2, \dots, \delta_n)$, 则由 $\hat{\beta}_c^{RS}$ 的定义可得

$$\begin{aligned} \hat{\beta}_c^{RS} &= S_{SR}^{-1} (\bar{X}^T \bar{Y} + A^T \Omega^{-1} b) = S_{SR}^{-1} (X^T \Delta Y + A^T \Omega^{-1} b) \\ &= S_{SR}^{-1} [X^T \Delta X \beta + X^T \Delta \varepsilon + A^T \Omega^{-1} (A \beta + \eta)] \\ &= S_{SR}^{-1} [(X^T \Delta X + A^T \Omega^{-1} A) \beta + X^T \Delta \varepsilon + A^T \Omega^{-1} \eta] \\ &= \beta + S_{SR}^{-1} X^T \Delta \varepsilon + S_{SR}^{-1} A^T \Omega^{-1} \eta \end{aligned}$$

从而有

$$\sqrt{n}(\hat{\beta}_C^{RS} - \beta) = \left(\frac{S_{SR}}{n}\right)^{-1} \frac{1}{\sqrt{n}} X^T \Delta \varepsilon + \left(\frac{S_{SR}}{n}\right)^{-1} \frac{1}{\sqrt{n}} A^T \Omega^{-1} \eta$$

由引理 2, 结合 $\frac{1}{n} A^T \Omega^{-1} A = o_p(1)$, 可得:

$$\frac{S_{SR}}{n} = \frac{1}{n} (\bar{X}^T \bar{X} + A^T \Omega^{-1} A) \xrightarrow{p} \Sigma, \quad \frac{1}{\sqrt{n}} X^T \Delta \varepsilon \xrightarrow{D} N(0, \sigma^2 \Sigma)$$

另一方面由 $E A^T \Omega^{-1} \eta = 0$ 和 $Cov(A^T \Omega^{-1} \eta) = A^T \Omega^{-1} A$, 可得

$$\frac{1}{\sqrt{n}} A^T \Omega^{-1} \eta = o_p(1)$$

结合上面的结论, 由 Slutsky 定理, 可得

$$\sqrt{n}(\hat{\beta}_C^{RS} - \hat{\beta}) \xrightarrow{D} N(0, \sigma^2 \Sigma^{-1})$$

定理 2 的证明: 由 $Y_i^* = \delta_i Y_i + (1 - \delta_i) X_i^T \hat{\beta}_C$, 可得

$$\begin{aligned} Y^* &= \Delta Y + (I_n - \Delta) X \hat{\beta}_C \\ &= \Delta Y + (I_n - \Delta) X (X^T \Delta X)^{-1} X^T \Delta Y \\ &= \Delta Y + X (X^T \Delta X)^{-1} X^T \Delta Y - \Delta X (X^T \Delta X)^{-1} X^T \Delta Y \end{aligned}$$

从而有

$$\begin{aligned} \beta_I^{RS} &= (X^T X + A^T \Omega^{-1} A)^{-1} (X^T Y^* + A^T \Omega^{-1} b) \\ &= (X^T X + A^T \Omega^{-1} A)^{-1} \left[X^T (\Delta Y + X (X^T \Delta X)^{-1} X^T \Delta Y - \Delta X (X^T \Delta X)^{-1} X^T \Delta Y) + A^T \Omega^{-1} b \right] \\ &= (X^T X + A^T \Omega^{-1} A)^{-1} \left[X^T X (X^T \Delta X)^{-1} X^T \Delta Y + A^T \Omega^{-1} b \right] \\ &= (X^T X + A^T \Omega^{-1} A)^{-1} \left[X^T X (X^T \Delta X)^{-1} X^T \Delta (X \beta + \varepsilon) + A^T \Omega^{-1} (A \beta + \eta) \right] \\ &= \beta + (X^T X + A^T \Omega^{-1} A)^{-1} \left[X^T X (X^T \Delta X)^{-1} X^T \Delta \varepsilon + A^T \Omega^{-1} \eta \right] \end{aligned}$$

由引理 2 以及 $\frac{1}{n} A^T \Omega^{-1} A \xrightarrow{p} 0$ 可得

$$\begin{aligned} \sqrt{n}(\hat{\beta}_{PM} - \beta) &= \left(\frac{X^T X + A^T \Omega^{-1} A}{n}\right)^{-1} X^T X (X^T \Delta X)^{-1} \frac{1}{\sqrt{n}} X^T \Delta \varepsilon + \left(\frac{X^T X + A^T \Omega^{-1} A}{n}\right)^{-1} \frac{1}{\sqrt{n}} A^T \Omega^{-1} \eta \\ &= (X^T \Delta X)^{-1} \frac{1}{\sqrt{n}} X^T \Delta \varepsilon + o_p(1) \end{aligned}$$

从而可得

$$\sqrt{n}(\hat{\beta}_I^{RS} - \hat{\beta}) \xrightarrow{D} N(0, \sigma^2 \Sigma^{-1}).$$

基金项目

中国劳动关系学院教育教学改革立项项目(JG1406); 2020 年度安徽高等学校自然科学基金项目(KJ2020A1200)。

参考文献

- [1] Toutenburg, H. (1982) *Prior Information in Linear Models*. Wiley, New York.
- [2] Rao, C.R. and Toutenburg, H. (1999) *Linear Models: Least Squares and Alternatives*. 2nd Edition, Springer, Berlin.
- [3] Little, R.J.A. and Rubin, D.B. (1987) *Statistical Analysis with Missing Data*. John Wiley & Sons, Inc., New York.
- [4] Cheng, P.E. (1990) Nonparametric Estimation of Mean Functionals with Data Missing at Random. *Journal of the American Statistical Association*, **89**, 81-87. <https://doi.org/10.1080/01621459.1994.10476448>
- [5] Chu, C.K. and Cheng, P.E. (1995) Nonparametric Regression Estimation with Missing Data. *Journal of Statistical Planning and Inference*, **48**, 85-99. [https://doi.org/10.1016/0378-3758\(94\)00151-K](https://doi.org/10.1016/0378-3758(94)00151-K)
- [6] Wang, Q.H. and Rao, J.N.K. (2002) Empirical Likelihood-Based Inference under Imputation for Missing Response Data. *The Annals of Statistics*, **30**, 896-924. <https://doi.org/10.1214/aos/1028674845>
- [7] Wang, Q.H. and Rao, J.N.K. (2004) Empirical Likelihood for Linear Regression Models under Imputation for Missing Responses. *Canadian Journal of Statistics*, **29**, 597-608. <https://doi.org/10.2307/3316009>
- [8] Qin, Y.S., Li, L. and Lei, Q.Z. (2009) Empirical Likelihood for Linear Regression Models with Missing Responses. *Statistics and Probability Letters*, **79**, 1391-1396. <https://doi.org/10.1016/j.spl.2009.03.002>
- [9] 杨徐佳, 于倩倩, 王森. 因变量缺失下线性回归模型的估计与检验[J]. 淮北煤炭师范学院学报(自然科学版), 2011, 32(1): 24-28.
- [10] 安佰玲. 线性回归模型在因变量缺失下的约束估计[J]. 统计与决策, 2013(11): 19-21.
- [11] Durbin, J.A. (1953) A Note on Regression When There Is Extraneous Information about One of the Coefficients. *Journal of the American Statistical Association*, **48**, 799-808. <https://doi.org/10.1080/01621459.1953.10501201>
- [12] Theil, H. and Goldberger, A.S. (1961) On Pure and Mixed Statistical Estimation in Economics. *International Economic Review*, **2**, 65-78. <https://doi.org/10.2307/2525589>
- [13] Theil, H. (1963) On the Use of Incomplete Prior Information in Regression Analysis. *Journal of the American Statistical Association*, **58**, 401-414. <https://doi.org/10.1080/01621459.1963.10500854>