

第十讲 最小二乘法

2023.11.24

$g_n(\theta) = 0$ 是一个合理的估计方程的必要条件是 $Eg_n(\theta) = 0$

多变量线性回归模型

多变量回归模型

多变量/多重线性回归模型(multiple linear regression model)总体模型:

$$y = \beta_0 + x_1\beta_1 + \dots + x_{p-1}\beta_{p-1} + \varepsilon = \beta_0 + \mathbf{x}^\top \mathbf{b} + \varepsilon, \quad \varepsilon \sim (0, \sigma^2), \varepsilon \text{与} \mathbf{x} \text{独立.}$$

假设样本 $(y_i, x_{i1}, \dots, x_{i,p-1})$, $i = 1, 2, \dots, n$ 独立, 来自于上述总体模型, 即

$$y_i = \beta_0 + x_{i1}\beta_1 + \dots + x_{i,p-1}\beta_{p-1} + \varepsilon_i = \beta_0 + \mathbf{x}_i^\top \mathbf{b} + \varepsilon_i, \quad (1)$$

ε_i iid $\sim (0, \sigma^2)$, ε_i 与 \mathbf{x}_i 独立.

记所有响应变量 $\mathbf{y} = (y_1, \dots, y_n)^\top$, 所有误差 $\boldsymbol{\varepsilon} = (\varepsilon_1, \dots, \varepsilon_n)^\top$,

所有回归系数 $\boldsymbol{\beta} = (\beta_0, \beta_1, \dots, \beta_{p-1})^\top = \begin{pmatrix} \beta_0 \\ \mathbf{b} \end{pmatrix}$ 。所有自变量(包含1)

组成的矩阵称为设计阵(design matrix):

$$\mathbf{X}_{n \times p} = \begin{pmatrix} 1 & x_{11} & \cdots & x_{1,p-1} \\ 1 & x_{21} & \cdots & x_{2,p-1} \\ \vdots & \vdots & \ddots & \vdots \\ 1 & x_{n1} & \cdots & x_{n,p-1} \end{pmatrix} = \begin{pmatrix} 1 & \mathbf{x}_1^\top \\ 1 & \mathbf{x}_2^\top \\ \vdots & \vdots \\ 1 & \mathbf{x}_n^\top \end{pmatrix} = \begin{pmatrix} \tilde{\mathbf{x}}_1^\top \\ \tilde{\mathbf{x}}_2^\top \\ \vdots \\ \tilde{\mathbf{x}}_n^\top \end{pmatrix}$$

矩阵-向量形式

模型 (1) 以矩阵-向量表达为

$$\mathbf{y}_{n \times 1} = X_{n \times p} \boldsymbol{\beta}_{p \times 1} + \boldsymbol{\varepsilon}_{n \times 1}, \quad \boldsymbol{\varepsilon} \sim (\mathbf{0}, \sigma^2 I_n), \quad \boldsymbol{\varepsilon} \text{ 与 } X \text{ 独立} \quad (2)$$

注：模型也可写作：

$$\mathbf{y} = \mathbf{u} + \boldsymbol{\varepsilon}, \quad \mathbf{u} \in C(X)$$

记 X 的各列为 $X = (\mathbf{1}, \mathbf{x}_{(1)}, \dots, \mathbf{x}_{(p-1)})$, 模型为

$$\mathbf{y} = \mathbf{u} + \boldsymbol{\varepsilon} = \mathbf{1}\beta_0 + \mathbf{x}_{(1)}\beta_1 + \dots + \mathbf{x}_{(p-1)}\beta_{p-1} + \boldsymbol{\varepsilon}$$

例1. 假设 (y_i, x_i) , $i = 1, 2, \dots, n$ 独立, 满足简单线性回归模型:

$$y_i = a + bx_i + \varepsilon_i,$$

记 $\boldsymbol{\beta} = \begin{pmatrix} a \\ b \end{pmatrix}$, $\mathbf{y} = (y_1, \dots, y_n)^\top$, 设计阵 $X = \begin{pmatrix} 1 & x_1 \\ \vdots & \vdots \\ 1 & x_n \end{pmatrix} = (\mathbf{1}, \mathbf{x})$

模型为 $\mathbf{y} = X\boldsymbol{\beta} + \boldsymbol{\varepsilon} = \mathbf{1}a + \mathbf{x}b + \boldsymbol{\varepsilon}$

最小二乘法 (LS: Least Squares)

最小化误差平方和

为了估计参数 $\boldsymbol{\beta}$ ，最小二乘法最小化误差平方和：

$$\min_{\boldsymbol{\beta} \in R^p} \sum \varepsilon_i^2 = \min_{\boldsymbol{\beta} \in R^p} \sum (y_i - \beta_0 - \mathbf{x}_i^T \mathbf{b})^2 = \min_{\boldsymbol{\beta} \in R^p} \|\boldsymbol{\varepsilon}\|^2 = \min_{\boldsymbol{\beta} \in R^p} \|\mathbf{y} - X\boldsymbol{\beta}\|^2$$

使得平方和达到最小值的 $\hat{\boldsymbol{\beta}}$ 称为最小二乘(LS)估计。

定理1. 假设线性回归模型满足Gauss - Markov(GM)假设：

$$\mathbf{y}_{n \times 1} = X_{n \times p} \boldsymbol{\beta}_{p \times 1} + \boldsymbol{\varepsilon}_{n \times 1}, \quad \boldsymbol{\varepsilon} \sim (\mathbf{0}, \sigma^2 I_n), \quad \boldsymbol{\varepsilon} \text{ 与 } X \text{ 独立}$$

最小二乘估计为 $\hat{\boldsymbol{\beta}} = (X^T X)^{-1} X^T \mathbf{y}$ ，满足正则方程（估计方程）

$$X^T \boldsymbol{\varepsilon} = X^T (\mathbf{y} - X\boldsymbol{\beta}) = \mathbf{0}$$

若 X 列满秩，则LS估计 $\hat{\boldsymbol{\beta}} = (X^T X)^{-1} X^T \mathbf{y}$ 唯一，且 $\hat{\boldsymbol{\beta}}$ 是 $\boldsymbol{\beta}$ 的无偏估计。

无偏性是LS估计最重要的性质，无偏性说明在GM假设成立的条件下，基于LS方法能得到参数正确的估计。LS估计的另一个重要性质是在所有线性无偏估计中其方差最小（GM定理，后面）

证法1:
投影

证：由第9讲定理2(投影的最小二乘性质, P23),

$$\min_{\beta \in R^p} \|\mathbf{y} - X\beta\|^2 = \min_{\mathbf{u} \in C(X)} \|\mathbf{y} - \mathbf{u}\|^2$$

最小值在 \mathbf{u} 等于 \mathbf{y} 在 $C(X)$ 上的正交投影时达到最小, 投影

$$\hat{\mathbf{y}} = P_X \mathbf{y} = X \underbrace{(X^T X)^{-1} X^T \mathbf{y}} = X \underline{\hat{\beta}}$$

$$\Rightarrow LS估计 \hat{\beta} = (X^T X)^{-1} X^T \mathbf{y}$$

$\mathbf{y} - \hat{\mathbf{y}} \perp C(X)$, 特别地与 X 的每一列正交, 满足正则方程:

$$X^T (\mathbf{y} - \hat{\mathbf{y}}) = X^T (\mathbf{y} - X\hat{\beta}) = \mathbf{0}.$$

若 X 列满秩, 即若 $X^T X$ 可逆, 则 $\hat{\beta} = (X^T X)^{-1} X^T \mathbf{y}$ 唯一.

因为 ε 与 X 独立, 所以

$$\begin{aligned} E(\hat{\beta} | X) &= (X^T X)^{-1} X^T E(\mathbf{y} | X) = (X^T X)^{-1} X^T E(X\beta + \varepsilon | X) \\ &= \beta + (X^T X)^{-1} X^T E(\varepsilon | X) = \beta + (X^T X)^{-1} X^T E(\varepsilon) = \beta \end{aligned}$$

$$\Rightarrow E(\hat{\beta}) = E(E(\hat{\beta} | X)) = \beta.$$

引理. (矩阵向量求导)

(1) 若 $\mathbf{a}, \mathbf{x} \in R^n$, 则 $\frac{\partial(\mathbf{a}^\top \mathbf{x})}{\partial \mathbf{x}} = \mathbf{a}$;

(2) 若 A 为 $n \times n$ 对称矩阵, $\mathbf{x} \in R^n$, 则 $\frac{\partial(\mathbf{x}^\top A \mathbf{x})}{\partial \mathbf{x}} = 2A\mathbf{x}$.

证法2:
求导

目标函数 $Q(\boldsymbol{\beta}) = \|\mathbf{y} - X\boldsymbol{\beta}\|^2 = \mathbf{y}^\top \mathbf{y} - 2\mathbf{y}^\top X\boldsymbol{\beta} + \boldsymbol{\beta}^\top X^\top X\boldsymbol{\beta}$,

$$\frac{\partial Q}{\partial \boldsymbol{\beta}} = -2(\mathbf{y}^\top X)^\top + 2X^\top X\boldsymbol{\beta} = -2X^\top(\mathbf{y} - X\boldsymbol{\beta}) = -2X^\top \boldsymbol{\varepsilon}$$

令之为0, 得正则方程

$$X^\top \boldsymbol{\varepsilon} = X^\top(\mathbf{y} - X\boldsymbol{\beta}) = 0 \Leftrightarrow X^\top X\boldsymbol{\beta} = X^\top \mathbf{y}$$

$$\Rightarrow \text{LS估计 } \hat{\boldsymbol{\beta}} = (X^\top X)^{-1} X^\top \mathbf{y}.$$

注: 也可以展开平方和求导 (避免矩阵向量求导)

记 X 的各行为 $X = (\tilde{\mathbf{x}}_1, \dots, \tilde{\mathbf{x}}_n)^\top$,

$$Q(\boldsymbol{\beta}) = \|\mathbf{y} - X\boldsymbol{\beta}\|^2 = \sum_{i=1}^n (y_i - \tilde{\mathbf{x}}_i^\top \boldsymbol{\beta})^2, \quad \text{令 } \frac{\partial Q}{\partial \boldsymbol{\beta}} = -2 \sum_{i=1}^n \tilde{\mathbf{x}}_i (y_i - \tilde{\mathbf{x}}_i^\top \boldsymbol{\beta}) = 0$$

$$\text{即 } \sum_{i=1}^n \tilde{\mathbf{x}}_i \tilde{\mathbf{x}}_i^\top \boldsymbol{\beta} = \sum_{i=1}^n \tilde{\mathbf{x}}_i y_i \Rightarrow \hat{\boldsymbol{\beta}} = \left(\sum_{i=1}^n \tilde{\mathbf{x}}_i \tilde{\mathbf{x}}_i^\top \right)^{-1} \left(\sum_{i=1}^n \tilde{\mathbf{x}}_i y_i \right) = (X^\top X)^{-1} X^\top \mathbf{y}$$

正则方程

LS估计是正则方程

$$X^T \boldsymbol{\varepsilon} = X^T (\mathbf{y} - X\boldsymbol{\beta}) = \mathbf{0}$$

的解，该方程与模型假设 $X \perp \boldsymbol{\varepsilon}$ 是相容的，这是因为

$$X \perp \boldsymbol{\varepsilon} \Rightarrow E(X^T \boldsymbol{\varepsilon}) = \mathbf{0} \Rightarrow \text{方程 } X^T \boldsymbol{\varepsilon} = \mathbf{0} \text{ 合理}$$

一般地，假设基于 n 个独立样本的估计方程为

$$g_n(\boldsymbol{\theta}) = 0$$

该方程的解 $\hat{\boldsymbol{\theta}}$ 为参数 $\boldsymbol{\theta}$ 的估计（比如似然方程 $g_n(\boldsymbol{\theta}) = \frac{\partial \log L(\boldsymbol{\theta})}{\partial \boldsymbol{\theta}} = 0$ ）。
 $\hat{\boldsymbol{\theta}}$ 是无偏估计或近似无偏估计的一个必要条件是

$$E g_n(\boldsymbol{\theta}) = 0$$

直观上， $g_n(\hat{\boldsymbol{\theta}}) = 0$ 且 $E(\hat{\boldsymbol{\theta}}) = \boldsymbol{\theta} \Rightarrow g_n(\boldsymbol{\theta}) \approx 0$ ，平均来看 $E g_n(\boldsymbol{\theta}) = 0$ 。

细节：假设 $g_n(\boldsymbol{\theta}) = 0$ 的解为 $\hat{\boldsymbol{\theta}}$ ，Taylor 展开

$$\begin{aligned} 0 &= g_n(\hat{\boldsymbol{\theta}}) \approx g_n(\boldsymbol{\theta}) + g_n'(\boldsymbol{\theta})(\hat{\boldsymbol{\theta}} - \boldsymbol{\theta}) \\ &\approx g_n(\boldsymbol{\theta}) + G_n(\boldsymbol{\theta})(\hat{\boldsymbol{\theta}} - \boldsymbol{\theta}) \end{aligned}$$

其中 $g_n'(\boldsymbol{\theta})$ 一般是 n 项独立和，由大数定律 $g_n'(\boldsymbol{\theta}) \approx E[g_n'(\boldsymbol{\theta})] \triangleq G_n(\boldsymbol{\theta})$

$\Rightarrow \hat{\boldsymbol{\theta}} - \boldsymbol{\theta} \approx G_n(\boldsymbol{\theta})^{-1} g_n(\boldsymbol{\theta})$ ， $E(\hat{\boldsymbol{\theta}} - \boldsymbol{\theta}) \approx G_n(\boldsymbol{\theta})^{-1} E(g_n(\boldsymbol{\theta})) = 0$ 。

LS估计-矩方法观点 (moment method)

模型假设: $y_i = \beta_0 + \mathbf{x}_i^\top \boldsymbol{\beta} + \varepsilon_i = \tilde{\mathbf{x}}_i^\top \boldsymbol{\beta} + \varepsilon_i$, ε_i 与 $\tilde{\mathbf{x}}_i$ 独立, $\varepsilon_i \sim (0, \sigma^2)$ 。
 \Leftrightarrow 以矩阵-向量表达: $\mathbf{y} = X\boldsymbol{\beta} + \boldsymbol{\varepsilon}$, $\boldsymbol{\varepsilon} \sim (0, \sigma^2 I_n)$, $\boldsymbol{\varepsilon}$ 与 X 独立。

矩方法

因为 $E(\tilde{\mathbf{x}}_i \varepsilon_i) = 0$, 矩方法令样本矩 $\sum_{i=1}^n \tilde{\mathbf{x}}_i \varepsilon_i / n$ 等于 0:

$$\text{令 } \sum_{i=1}^n \tilde{\mathbf{x}}_i \varepsilon_i = X^\top \boldsymbol{\varepsilon} = X^\top (\mathbf{y} - X\boldsymbol{\beta}) = 0,$$

即LS的正则方程 $X^\top \mathbf{y} = X^\top X\boldsymbol{\beta} = 0 \Rightarrow \hat{\boldsymbol{\beta}} = (X^\top X)^{-1} X^\top \mathbf{y}$

LS方法或矩方法也可简单地描述为:

模型 $\mathbf{y} = X\boldsymbol{\beta} + \boldsymbol{\varepsilon}$ 两边同时左乘 X^\top

$$X^\top \mathbf{y} = X^\top X\boldsymbol{\beta} + X^\top \boldsymbol{\varepsilon}$$

并舍弃 $X^\top \boldsymbol{\varepsilon}$ 一项(令 $X^\top \boldsymbol{\varepsilon} = \mathbf{0}$), 得正则方程 $X^\top \mathbf{y} = X^\top X\boldsymbol{\beta}$ 。

除了同时左乘 X^T ，是否还有其它选择？

显然方程两边同时左乘 $p \times n$ 矩阵 $Z^T = f(X^T)$ ：

$$Z^T \mathbf{y} = Z^T X \boldsymbol{\beta} + Z^T \boldsymbol{\varepsilon}$$

并令 $Z^T \boldsymbol{\varepsilon} = 0$ ，即 $Z^T \mathbf{y} = Z^T X \boldsymbol{\beta}$ ，也是一个合理的方程，这是因为 $\boldsymbol{\varepsilon}$ 与 $Z^T = f(X^T)$ 独立 $\Rightarrow E(Z^T \boldsymbol{\varepsilon}) = 0$ 。

事实上，任何外生的 Z （即 Z 与 $\boldsymbol{\varepsilon}$ 独立）都满足 $E(Z^T \boldsymbol{\varepsilon}) = 0$ ，导致合理的方程：

$$Z^T \boldsymbol{\varepsilon} = 0 \Leftrightarrow Z^T \mathbf{y} = Z^T X \boldsymbol{\beta}$$

进一步，若 $C(Z^T X) = C(Z^T)$ ，则方程有解；

若 $Z^T X$ 可逆，则有唯一解

$$\tilde{\boldsymbol{\beta}} = (Z^T X)^{-1} Z^T \mathbf{y}$$

它是无偏的。

$$\begin{aligned} E(\tilde{\boldsymbol{\beta}} | Z, X) &= (Z^T X)^{-1} Z^T E(X \boldsymbol{\beta} + \boldsymbol{\varepsilon} | Z, X) \\ &= \boldsymbol{\beta} + (Z^T X)^{-1} Z^T E(\boldsymbol{\varepsilon} | Z, X) = \boldsymbol{\beta}. \end{aligned}$$

有趣的是，即使线性模型的Gauss-Markov假设不成立（一般如此!），上述矩方法仍能产生合理的估计方程，并得到具有优良性质的估计。

工具变量法

假设 $\mathbf{y} = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\varepsilon}$ 中 \mathbf{X} 与 $\boldsymbol{\varepsilon}$ 不独立，此时正则方程 $\mathbf{X}^\top \boldsymbol{\varepsilon} = \mathbf{0}$ 不合理。

假设存在外生的 \mathbf{Z} （称为工具变量）满足

$$\mathbf{Z} \perp \boldsymbol{\varepsilon}, \text{ 且 } \mathbf{Z} \text{ 与 } \mathbf{X} \text{ 有关}$$

则 $E(\mathbf{Z}^\top \boldsymbol{\varepsilon}) = \mathbf{0}$ 成立，那么下述估计方程是合理的

$$\mathbf{Z}^\top \boldsymbol{\varepsilon} = \mathbf{Z}^\top (\mathbf{y} - \mathbf{X}\boldsymbol{\beta}) = 0$$

由此解得 $\tilde{\boldsymbol{\beta}} = (\mathbf{Z}^\top \mathbf{X})^{-1} \mathbf{Z}^\top \mathbf{y}$ ，称为工具变量最小二乘估计(参见附录)。



Ill. Niklas Elmehed © Nobel Prize Outreach.
David Card
Prize share: 1/2



Ill. Niklas Elmehed © Nobel Prize Outreach.
Joshua D. Angrist
Prize share: 1/4



Ill. Niklas Elmehed © Nobel Prize Outreach.
Guido W. Imbens
Prize share: 1/4

2021诺贝尔经济奖:

J. Angrist, G. Imbens, D. Card.

对于特定的经济社会问题，发现了一些工具变量。

LS估计 - 解方程的观点

我们也可以将线性模型参数估计问题看作是解近似方程的问题，这与矩方法类似，但可以借鉴解方程的方法处理传统矩方法或LS方法难以处理的问题，比如 $p > n$ 的情形。

超定方程
 $n > p$

$\mathbf{y} = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\varepsilon}$ 可看作是超定方程(over-determined system):

$$\mathbf{y}_{n \times 1} = \mathbf{X}_{n \times p} \boldsymbol{\beta}_{p \times 1}, \quad n \text{ 个方程} > p \text{ 个未知量}$$

只有当 $\mathbf{y} \in C(\mathbf{X})$ 时才有解(一般不可能)。

- 我们可求最优近似解使误差 $\|\mathbf{y} - \mathbf{X}\boldsymbol{\beta}\|^2$ 最小，这是LS方法。
- 也可以设法减少方程个数，即通过整合方程(比如合并某些方程、删除冗余方程,称为Pre-conditioning)的方式减少方程个数，这等价于方程两边同时做某个 \mathbf{Z}^T -线性变换($\mathbf{Z}^T: k \times n$)，得

$$\mathbf{Z}^T \mathbf{y} = \mathbf{Z}^T \mathbf{X} \boldsymbol{\beta},$$

由前述讨论知 \mathbf{Z} 的选取不能太随意： \mathbf{Z} 应与 \mathbf{X} 有关，且 \mathbf{Z} 外生(随机)。

欠定方程
 $n < p$

$p > n$ 时, 模型 $\mathbf{y}_{n \times 1} = X_{n \times p} \boldsymbol{\beta}_{p \times 1} + \boldsymbol{\varepsilon}_{n \times 1}$ 中求解参数 $\boldsymbol{\beta}$ 可看作是解欠定方程(underdetermined system)

$$\mathbf{y}_{n \times 1} = X_{n \times p} \boldsymbol{\beta}_{p \times 1}, \quad n < p,$$

通常有无穷多解, 为了求出有意义的解, 通常对解施加某些限制, 比如

- 不定方程: 数论中限制线性方程的解为正整数或有理数;
- 主成分回归: 极小化欧氏模长 $\|\boldsymbol{\beta}\|_2$, $\tilde{\boldsymbol{\beta}} = X^+ \mathbf{y}$;
- 压缩感知或lasso: 极小化 $\|\boldsymbol{\beta}\|_0$ (非0的个数)或其放松 $\|\boldsymbol{\beta}\|_1$ 。

误差方差的LS估计

拟合值
残差

- 投影 $\hat{\mathbf{y}} = P_X \mathbf{y} = X\hat{\boldsymbol{\beta}} = X(X^\top X)^{-1} X^\top \mathbf{y}$, 称为拟合值向量。
- $\mathbf{e} = \mathbf{y}^\perp = \mathbf{y} - \hat{\mathbf{y}} = (I_n - P_X)\mathbf{y}$ 称为残差向量。

$\hat{\boldsymbol{\beta}}$ 满足正则方程: $X^\top (\mathbf{y} - X\hat{\boldsymbol{\beta}}) = X^\top \mathbf{e} = 0$, 即 $\mathbf{e} \perp C(X)$,
特别地 $\mathbf{e} \perp \mathbf{1}$, 即 $\bar{e} = \mathbf{1}^\top \mathbf{e} / n = 0$ (样本均值为0)

误差方
差估计

残差平方和: $RSS = \|\mathbf{e}\|^2 = \|\mathbf{y} - X\hat{\boldsymbol{\beta}}\|^2 = \mathbf{y}^\top (I_n - P_X)\mathbf{y}$

σ^2 的LS估计定义为: $\hat{\sigma}^2 = \frac{RSS}{n-p} = \frac{1}{n-p} \|\mathbf{e}\|^2$ 。

LS估计的统计性质

模型假设 $\varepsilon_1, \dots, \varepsilon_n \text{ iid } \sim (0, \sigma^2)$ 反映了我们对误差的认知:

$$\varepsilon_1, \dots, \varepsilon_n \text{ 大致相同, } \varepsilon_i \approx 0 \text{ 且 } |\varepsilon_i| \approx \sigma$$

即使没有这些概率模型假设, 我们也能操作LS。有了这些假设, 我们可评价LS方法, 比如, LS估计 $\hat{\boldsymbol{\beta}}$ 在平均意义下是否等于或接近真正的 $\boldsymbol{\beta}$? 其精度/方差如何? 逼近误差RSS大概多大?

第6讲引理1(2)

$$\begin{aligned} \text{RSS} &= s_{yy} - s_{xy}^2/s_{xx} \\ &= s_{\varepsilon\varepsilon} - s_{x\varepsilon}^2/s_{xx} \end{aligned}$$

从投影的观点是显然的

例如, $\mathbf{y} = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\varepsilon}$ 中,

$$\text{RSS} = \|\mathbf{e}\|^2 = \mathbf{y}^\top (\mathbf{I}_n - \mathbf{P}_X) \mathbf{y} = (\mathbf{X}\boldsymbol{\beta} + \boldsymbol{\varepsilon})^\top (\mathbf{I}_n - \mathbf{P}_X) (\mathbf{X}\boldsymbol{\beta} + \boldsymbol{\varepsilon}) = \boldsymbol{\varepsilon}^\top (\mathbf{I}_n - \mathbf{P}_X) \boldsymbol{\varepsilon},$$

平均来看 (当X给定时):

$$E(\text{RSS}) = E\boldsymbol{\varepsilon}^\top (\mathbf{I}_n - \mathbf{P}_X) \boldsymbol{\varepsilon} = \text{tr}((\mathbf{I}_n - \mathbf{P}_X) \text{var}(\boldsymbol{\varepsilon})) = (n - r)\sigma^2$$

其中 $r = \text{tr}(\mathbf{P}_X) = \text{rank}(X)$.

引理 (下页)

引理1: 若 $\mathbf{x} \sim (\boldsymbol{\mu}, \Sigma)$, 则 $E(\mathbf{x}^T A \mathbf{x}) = \boldsymbol{\mu}^T A \boldsymbol{\mu} + \text{tr}(A \Sigma)$.

证明: $E(\mathbf{x}^T A \mathbf{x}) = E \text{tr}(\mathbf{x}^T A \mathbf{x}) = E \text{tr}(A \mathbf{x} \mathbf{x}^T) = \text{tr} E(A \mathbf{x} \mathbf{x}^T) = \text{tr} A E(\mathbf{x} \mathbf{x}^T)$,

因为 $\Sigma = E(\mathbf{x} - \boldsymbol{\mu})(\mathbf{x} - \boldsymbol{\mu})^T = E(\mathbf{x} \mathbf{x}^T) - \boldsymbol{\mu} \boldsymbol{\mu}^T$

所以 $E(\mathbf{x}^T A \mathbf{x}) = \text{tr} A (\Sigma + \boldsymbol{\mu} \boldsymbol{\mu}^T) = \text{tr}(A \Sigma) + \text{tr} A \boldsymbol{\mu} \boldsymbol{\mu}^T = \text{tr}(A \Sigma) + \boldsymbol{\mu}^T A \boldsymbol{\mu}$

无偏性

定理1. 假设线性模型 $\mathbf{y} = X_{n \times p} \boldsymbol{\beta} + \boldsymbol{\varepsilon}$, $E \boldsymbol{\varepsilon} = 0$, $\text{var}(\boldsymbol{\varepsilon}) = \sigma^2 I_n$, 其中 $\boldsymbol{\varepsilon}$ 与 X 独立。假设 X 是列满秩的 ($n \geq p$, 且 $X^T X$ 可逆), 则

(1) LS估计的无偏性: $E(\hat{\boldsymbol{\beta}}) = \boldsymbol{\beta}$.

(2) LS估计的方差: $\text{var}(\hat{\boldsymbol{\beta}} | X) = \sigma^2 (X^T X)^{-1}$.

(3) 误差方差估计的无偏性: $E(\hat{\sigma}^2) = \sigma^2$.

证明: 因为 $\boldsymbol{\varepsilon}$ 与 X 独立, 所以

- $E(\mathbf{y} | X) = E(X \boldsymbol{\beta} + \boldsymbol{\varepsilon} | X) = X \boldsymbol{\beta} + E(\boldsymbol{\varepsilon} | X) = X \boldsymbol{\beta} + E(\boldsymbol{\varepsilon}) = X \boldsymbol{\beta}$
- $\text{var}(\mathbf{y} | X) = \text{var}(X \boldsymbol{\beta} + \boldsymbol{\varepsilon} | X) = \text{var}(\boldsymbol{\varepsilon} | X) = \text{var}(\boldsymbol{\varepsilon}) = \sigma^2 I_n$

$$(1) E(\hat{\boldsymbol{\beta}} | X) = E\left(\left(X^T X\right)^{-1} X^T \mathbf{y} | X\right) = \left(X^T X\right)^{-1} X^T [E(\mathbf{y} | X)] = \boldsymbol{\beta} \Rightarrow E(\hat{\boldsymbol{\beta}}) = \boldsymbol{\beta}.$$

$$(2) \text{var}(\hat{\boldsymbol{\beta}} | X) = \text{var}\left(\left(X^T X\right)^{-1} X^T \mathbf{y} | X\right) \\ = \left(X^T X\right)^{-1} X^T [\text{var}(\mathbf{y} | X)] X \left(X^T X\right)^{-1} = \sigma^2 \left(X^T X\right)^{-1}$$

$$(3) \text{前一页我们已经证明了 } E(RSS | X) = (n - p)\sigma^2,$$

$$\Rightarrow E(RSS) = (n - p)\sigma^2$$

$$\Rightarrow E(\hat{\sigma}^2) = E(RSS / (n - p)) = \sigma^2.$$

$\hat{\boldsymbol{\beta}}$ 的方差估计

在方差公式 $\text{var}(\hat{\boldsymbol{\beta}} | X) = \sigma^2 (X^T X)^{-1}$ 中代入 (Plug-in) σ^2 的估计, 即得到 $\hat{\boldsymbol{\beta}}$ 方差的估计:

$$\widehat{\text{var}(\hat{\boldsymbol{\beta}} | X)} = \hat{\sigma}^2 (X^T X)^{-1}$$

例1(续). 假设 (y_i, x_i) , $i = 1, 2, \dots, n$ 独立, 满足简单线性回归模型:

$$y_i = a + bx_i + \varepsilon_i,$$

第6讲命题1、2中我们已经求得LS估计及其方差, 这里我们以矩阵向量形式再次计算如下:

$$\text{记 } \boldsymbol{\beta} = \begin{pmatrix} a \\ b \end{pmatrix}, \mathbf{y} = (y_1, \dots, y_n)^\top, \text{ 设计阵 } X = \begin{pmatrix} 1 & x_1 \\ \vdots & \vdots \\ 1 & x_n \end{pmatrix} = (\mathbf{1}, \mathbf{x}), \mathbf{x} \neq \mathbf{1}c,$$

模型为 $\mathbf{y} = X\boldsymbol{\beta} + \boldsymbol{\varepsilon} = \mathbf{1}a + \mathbf{x}b + \boldsymbol{\varepsilon}$, LS估计

$$\begin{aligned} \hat{\boldsymbol{\beta}} &= \begin{pmatrix} \hat{a} \\ \hat{b} \end{pmatrix} = (X^\top X)^{-1} X^\top \mathbf{y} = \begin{pmatrix} n & n\bar{x} \\ n\bar{x} & \sum x_i^2 \end{pmatrix}^{-1} \begin{pmatrix} n\bar{y} \\ \sum x_i y_i \end{pmatrix} \\ &= \frac{1}{n(\sum x_i^2 - n\bar{x}^2)} \begin{pmatrix} \sum x_i^2 & -n\bar{x} \\ -n\bar{x} & n \end{pmatrix} \begin{pmatrix} n\bar{y} \\ \sum x_i y_i \end{pmatrix} = \begin{pmatrix} \bar{y} - \bar{x}s_{xy}/s_{xx} \\ s_{xy}/s_{xx} \end{pmatrix} \\ \text{var}(\hat{\boldsymbol{\beta}} | X) &= \sigma^2 (X^\top X)^{-1} = \frac{\sigma^2}{ns_{xx}} \begin{pmatrix} \sum x_i^2 & -n\bar{x} \\ -n\bar{x} & n \end{pmatrix} \\ &= \sigma^2 \begin{pmatrix} 1/n + \bar{x}^2/s_{xx} & -\bar{x}/s_{xx} \\ -\bar{x}/s_{xx} & 1/s_{xx} \end{pmatrix} \end{aligned}$$

定义：对任何两个对称 $n \times n$ 矩阵 A, B ，若 $A - B \geq 0$ (非负定)，则称在Loewner偏序意义下 A 不小于 B ，记作 $A \geq B$ 。

性质：若 $A \geq B$ ，则对任何 $k \times n$ 矩阵 C ， $CAC^T \geq CBC^T$

LS估计的最优性

定理2(Gauss - Markov定理). 线性模型 $\mathbf{y} = X\boldsymbol{\beta} + \boldsymbol{\varepsilon}$ 中, 假设 $X_{n \times p}$ 列满秩, 则 $\boldsymbol{\beta}$ 的最小二乘估计 $\hat{\boldsymbol{\beta}} = (X^T X)^{-1} X^T \mathbf{y}$ 是最优无偏线性估计 (BLUE, best linear unbiased estimate,), 即对任何 $\boldsymbol{\beta}$ 的线性无偏估计 $\tilde{\boldsymbol{\beta}} = C\mathbf{y}$

$$\text{var}(\tilde{\boldsymbol{\beta}}) \geq \text{var}(\hat{\boldsymbol{\beta}}),$$

其中 C 是仅与 X 有关的 $p \times n$ 常数矩阵。

证明：由 $\tilde{\boldsymbol{\beta}} = C\mathbf{y}$ 的无偏性(注意 C 仅与 X 有关)

$$\boldsymbol{\beta} = E(\tilde{\boldsymbol{\beta}} | X) = E(C\mathbf{y} | X) = CE(\mathbf{y} | X) = CX\boldsymbol{\beta},$$

上式对任何 $\boldsymbol{\beta}$ 成立, 故 $CX = I_p$ 。

因为 $\text{var}(\tilde{\boldsymbol{\beta}}) = E\text{var}(\tilde{\boldsymbol{\beta}} | X) + \text{var}(E(\tilde{\boldsymbol{\beta}} | X)) = E\text{var}(\tilde{\boldsymbol{\beta}} | X)$,
 $\text{var}(\hat{\boldsymbol{\beta}}) = E\text{var}(\hat{\boldsymbol{\beta}} | X)$, 我们只需要证明:
 $\text{var}(\tilde{\boldsymbol{\beta}} | X) = C \text{var}(\mathbf{y} | X) C^T = \sigma^2 C C^T \geq \text{var}(\hat{\boldsymbol{\beta}} | X) = \sigma^2 (X^T X)^{-1}$ 。
 即 $C C^T \geq (X^T X)^{-1}$ 。

因为 $CX = I_p$, $P_X \leq I_n$, 所以

$$(X^T X)^{-1} = CX(X^T X)^{-1}X^T C^T = CP_X C^T \leq CC^T.$$

所以 $\text{var}(\tilde{\boldsymbol{\beta}} | X) \geq \text{var}(\hat{\boldsymbol{\beta}} | X)$, $\text{var}(\tilde{\boldsymbol{\beta}}) \geq \text{var}(\hat{\boldsymbol{\beta}})$ 。

推论：假设 $A_{k \times p}$ 是仅与 X 有关的常数矩阵，则 $A\hat{\boldsymbol{\beta}}$ 是 $A\boldsymbol{\beta}$ 的 BLUE。
 特别地， $\hat{\beta}_k$ 是 β_k 的 BLUE， $\hat{\beta}_k - \hat{\beta}_j$ 是 $\beta_k - \beta_j$ 的 BLUE，等等。

不确定性
原理

对于简单线性模型 $y_i = a + bx_i + \varepsilon_i$ 的任一无偏估计 \tilde{b} ，由GM定理

$$\text{var}(\tilde{b} | x) \geq \text{var}(\hat{b} | x) = \sigma^2 / s_{xx}$$

记样本方差 $s_x^2 = s_{xx} / (n-1)$ ，上式等价于

$$\text{var}(\tilde{b} | x) \times s_x^2 \geq \sigma^2 / (n-1),$$

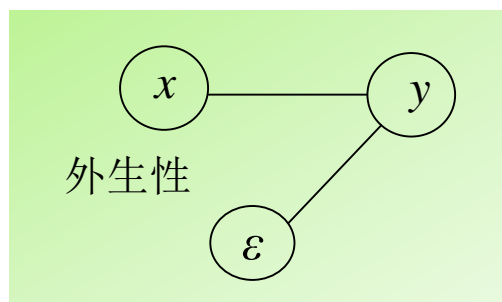
这说明左端两个方差不可能同时很小，即 x 和 \tilde{b} 不可能同时精确测量 (Uncertainty principle, 不确定性原理, 测不准原理).

附录: 工具变量法

工具变量法 (Instrumental Variable method) 试图在线性模型中误差与自变量不独立的情况下, 求解回归系数的无偏估计, 发现因果。

外生变量

如果线性模型 $y = a + bx + \varepsilon$ 中 x 不是研究对象本身固有而是外界随机赋予的, 称 x 是外生的 (exogenous), 此时 x 与 ε 独立。外生性是推断 x 和 y 之间因果关系的关键条件。



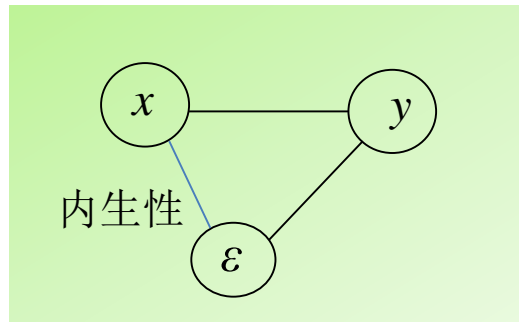
例2. 第1讲霍乱案例中, 比较两个自来水公司客户的问题可表述成:

$$y = a + bx + \varepsilon$$

其中 y 为霍乱状态, x 为饮用水来源 (是否污染), Snow 医生论证了 x 是外生变量 (随机取值, 自然试验), 因而 b 的LS估计是无偏的。

内生变量

观察研究中，模型 $y = a + bx + \varepsilon$ 中 x 是研究对象本身固有的，一般与 ε 不独立，称 x 是内生的 (endogenous)。此时如何推断 x 和 y 之间因果关系？



当 x 是内生变量时，可以尝试在回归模型中控制所有与 x, y 都有关系的干扰因素（非常困难）。计量经济学家发明了工具变量法 (Wright, 1928)，其关键在于发现恰当的自然试验 (natural experiment)，当然，这也是非常困难的。

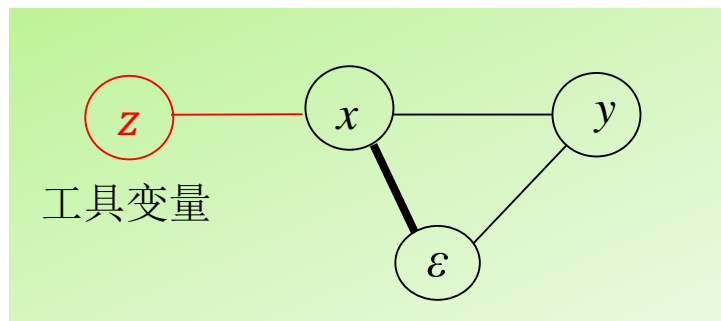
工具变量

线性模型 $y = a + bx + \varepsilon$, $\varepsilon \sim (0, \sigma^2)$, x 与 ε 不独立（内生）。

假设存在一个自然试验(natural experiment)产生的外生变量 z , 满足

z 与 ε 独立（外生）， z 与 x 相关，

z 称为工具变量（instrumental variable, IV）。



工具变量 z 满足条件：

- (1) z 与 ε 独立；
- (2) z 与 x 相关；

x 与 ε 之间有联系，但 z 与 ε 独立。

工具变量法在研究 y 与 x 的关系时，以 $\hat{x} = P_z x$ 替代 x , 换言之，利用 z 将 x 中与 ε 有关的成分消除/清洗掉。

工具变量 最小二乘

假设模型 $y_i = a + bx_i + \varepsilon_i$, $\varepsilon_i \sim (0, \sigma^2)$, ε_i 和 x_i 相关, $i = 1, \dots, n$ 。

假设 z_i 是工具变量, 满足条件: (1) z_i 与 ε_i 独立; (2) z_i 与 x_i 相关;

记 $\mathbf{y} = (y_i)^\top$, $\mathbf{x} = (x_i)^\top$, $\mathbf{z} = (z_i)^\top$, $\boldsymbol{\varepsilon} = (\varepsilon_i)^\top$, $\mathbf{Z} = (\mathbf{1}, \mathbf{z})$, $\mathbf{X} = (\mathbf{1}, \mathbf{x})$,

模型: $\mathbf{y} = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\varepsilon} = \mathbf{1}a + \mathbf{x}b + \boldsymbol{\varepsilon}$ 。

$\boldsymbol{\varepsilon}$ 与 \mathbf{Z} 独立 $\Rightarrow E(\mathbf{Z}^\top \boldsymbol{\varepsilon}) = 0$ 。

令 $\mathbf{Z}^\top \boldsymbol{\varepsilon} = \begin{pmatrix} \mathbf{1}^\top (\mathbf{y} - \mathbf{1}a - \mathbf{x}b) \\ \mathbf{z}^\top (\mathbf{y} - \mathbf{1}a - \mathbf{x}b) \end{pmatrix} = 0$, 得工具变量最小二乘估计 (IVLS):

$$\tilde{b}_{IVLS} = \frac{\sum (z_i - \bar{z})(y_i - \bar{y})}{\sum (z_i - \bar{z})(x_i - \bar{x})} = \frac{s_{yz}}{s_{xz}}, \quad \tilde{a}_{IVLS} = \bar{y} - \tilde{b}\bar{x}.$$

等价地, 也可从总体模型求出参数的表达, 再代入样本:

总体模型 $y = a + bx + \varepsilon$, 因为 z 与 ε 独立, 所以

$$0 = \text{cov}(\varepsilon, z) = \text{cov}(y - a - bx, z) = \text{cov}(y, z) - b \text{cov}(x, z)$$

所以, $b = \text{cov}(y, z) / \text{cov}(x, z)$, 代入样本方差: $\tilde{b} = s_{yz} / s_{xz}$ 。

性质： $\hat{b}_{IV} = \frac{s_{yz}}{s_{xz}}$ 渐近无偏，即 $E(\hat{b}_{IV}) \rightarrow b, n \rightarrow \infty$

两步估计1

注意到IVLS有下述表示，

$$\hat{b}_{IV} = \frac{s_{yz}}{s_{xz}} = \frac{s_{yz} / s_{zz}}{s_{xz} / s_{zz}},$$

所以，IVLS估计可由下述两步LS方法得到：

(1) $lm(x \sim z): x = c + dz + \varepsilon^{(1)} \Rightarrow LS$ 估计 $\hat{d} = s_{xz} / s_{zz}$

(2) $lm(y \sim z): y = e + fz + \varepsilon^{(2)} \Rightarrow LS$ 估计 $\hat{f} = s_{yz} / s_{zz}$

$$\Rightarrow \hat{b}_{IV} = \hat{f} / \hat{d}.$$

两步估计2

注意到IVLS估计也可表示如下

$$\hat{b}_{IV} = \frac{s_{yz}}{s_{xz}} = \frac{s_{yz}}{s_{zz}(s_{xz}/s_{zz})} = \frac{s_{yz}}{s_{zz}\hat{d}} = \frac{\sum(\hat{d} z_i - \hat{d} \bar{z})(y_i - \bar{y})}{\sum(\hat{d} z_i - \hat{d} \bar{z})^2} = \frac{s_{y\hat{x}}}{s_{\hat{x}\hat{x}}}$$

其中 $\hat{d} = s_{xz} / s_{zz}$, $\hat{x}_i = \hat{d} z_i$

所以, IVLS估计看作是下述两步LS估计:

$$(1) \text{ } lm(x \sim z): x = c + dz + \varepsilon^{(1)} \Rightarrow \text{LS估计 } \hat{d} = s_{xz} / s_{zz},$$

$$\text{令 } \hat{x}_i = \hat{d} z_i$$

$$(2) \text{ } lm(y \sim \hat{x}): y = e + f \hat{x} + \varepsilon^{(2)} \Rightarrow \text{LS估计 } \hat{f} = s_{y\hat{x}} / s_{\hat{x}\hat{x}} \text{ 这就是 } \hat{b}_{IV}$$

该观点反映了工具变量法的本质:

- (1) 回归 $y \sim x$ 不正确, 因为 x 是内生的 (x 与 ε 相关);
- (2) 以 \hat{x} 替代 x (\hat{x} 与 ε 独立);
- (3) $y \sim \hat{x}$

例3. 2021诺贝尔经济奖获得者Angrist研究了教育程度（ x ，受教育时长，月）是否与收入（ y ）存在因果关系(Angrist and Krueger, 1991)。考虑线性模型

$$y = a + bx + \varepsilon,$$

与收入有关的其它因素（比如能力，家庭因素，努力程度等）无法全部测量或精确测量到，我们将它们放到 ε 中。显然，这些因素与自变量 x 有关，因此基于上述简单模型无法正确地（无偏）估计 x 的效应。

Angrist and Krueger 注意到小学入学时所有儿童年龄最大有1年的差距（当年12月达到6岁即可入学）

出生季节	1 (10-12月)	2 (1-3月)	3 (4-6月)	4 (7-9月)
入学年龄	5 $\frac{3}{4}$ -6	6 -6 $\frac{1}{4}$	6 $\frac{1}{4}$ - 6 $\frac{1}{2}$	6 $\frac{1}{2}$ - 6 $\frac{3}{4}$

而义务教育法规定青少年在16岁生日之前必须在校学习，所以**对于16岁生日当天离开学校的那些学生而言，其受教育时长是由其生日 z 决定的**，所以生日与 x 有关，但与 ε 无关（天然试验），所以 z 可作为工具变量

$$y = \text{earnings}$$

$$x = \text{years of education}$$

$$z = \text{出生季节}$$

例4. 研究服兵役与否 x (0或1) 对生活质量 y 的影响, 建立模型

$$y = a + bx + \varepsilon$$

其中 ε 代表与 y 有关的个人和环境因素。志愿征兵制下入伍是个人的选择, 因此 x 是内生的, 基于LS方法得到的 b 的估计是有偏的/不正确的。

Angrist (1990) 发现了一个天然试验: 1970's 越战期间, 美国推行了基于“抽签”的强制征兵制度 (draft lottery)。每个19-26岁符合兵役条件的男性被随机分配了一个号码RSN (1-365)。征兵前国家公布一个数字 T , 小于 T 的人将被强制列入入伍候选, 再进行体检等其它程序。1970, 71, 72年分别 $T=195, 125, 95$ 。变量 $z = 1_{(RSN < T)}$ 是工具变量, 它与兵役状态 x 相关, 但与研究对象独立。

参考文献:

Angrist, Joshua D. 1990. "Lifetime Earnings and the Vietnam Era Draft Lottery: Evidence from Social Security Administrative Records." *American Economic Review*. June, 80:3, pp. 313–36.

Angrist, Joshua D. and Alan B. Krueger. 1991. "Does Compulsory School Attendance Affect Schooling and Earnings?" *Quarterly Journal of Economics*. November, 106:4, pp. 979–1014.

Angrist, J. D., and Krueger, A. B. 2001. "Instrumental Variables and the Search for Identification: From Supply and Demand to Natural Experiments," *The Journal of Economic Perspectives* (15:4), pp. 69-85.

其它工具变量的例子(Angrist and Krueger 2001)

Natural and Randomized Experiments

<i>Outcome Variable</i>	<i>Endogenous Variable</i>	<i>Source of Instrumental Variable(s)</i>	<i>Reference</i>
<i>1. Natural Experiments</i>			
Labor supply	Disability insurance replacement rates	Region and time variation in benefit rules	Gruber (2000)
Labor supply	Fertility	Sibling-Sex composition	Angrist and Evans (1998)
Education, Labor supply	Out-of-wedlock fertility	Occurrence of twin births	Bronars and Grogger (1994)
Wages	Unemployment insurance tax rate	State laws	Anderson and Meyer (2000)
Earnings	Years of schooling	Region and time variation in school construction	Duflo (2001)
Earnings	Years of schooling	Proximity to college	Card (1995)
Earnings	Years of schooling	Quarter of birth	Angrist and Krueger (1991)
Earnings	Veteran status	Cohort dummies	Imbens and van der Klaauw (1995)
Earnings	Veteran status	Draft lottery number	Angrist (1990)
Achievement test scores	Class size	Discontinuities in class size due to maximum class-size rule	Angrist and Lavy (1999)
College enrollment	Financial aid	Discontinuities in financial aid formula	van der Klaauw (1996)
Health	Heart attack surgery	Proximity to cardiac care centers	McClellan, McNeil and Newhouse (1994)
Crime	Police	Electoral cycles	Levitt (1997)
Employment and Earnings	Length of prison sentence	Randomly assigned federal judges	Kling (1999)
Birth weight	Maternal smoking	State cigarette taxes	Evans and Ringel (1999)

常见的内生性原因

(0) 观察研究基本都有内生性。

(1) 丢失变量(Omitted variable,回归方程中没有控制相关变量)

假设正确模型为: $y = a + bx + cz + \delta$, $\delta \perp x$ (\perp 代表独立, 下同)

其中 z 与 x 相关。如果我们没有测量, 或者没有在上述

模型中控制 z , 那么工作模型 $y = a + bx + \varepsilon$ 中 $\varepsilon = cz + \delta$ 与 x 不独立。

(2) 因果颠倒(reverse causation,响应变量是自变量的原因)

假设正确模型为 $y = a + bx + \varepsilon$, $\varepsilon \perp x$

但实际操作中工作模型取为: $x = c + dy + \delta$,

从正确模型我们知道 $x = -a/b - y/b - \varepsilon/b$, 所以工作模型中 $\delta = -\varepsilon/b$ 与 y 有关。

(3) 自变量带误差模型 (Error in Variable, EV模型)

正确模型: $y = a + bx_0 + \varepsilon_0$, $x_0 \perp \varepsilon_0$. 假设对于 x_0 的测量有误差, 即我们

只能测量到: $x = x_0 + \delta$, 其中 $x_0 \perp \delta$. 则

$$y = a + b(x - \delta) + \varepsilon_0 = a + bx + (\varepsilon_0 - b\delta) \triangleq a + bx + \varepsilon$$

显然 $x = x_0 + \delta$ 与 $\varepsilon = \varepsilon_0 - b\delta$ 都含 δ , 它们相关。