

§ 4.3 随机解释变量问题

教材129页 4.3 内生解释变量问题

一、随机解释变量问题

对于模型

$$Y_i = \beta_0 + \beta_1 X_{1i} + \beta_2 X_{2i} + \dots + \beta_k X_{ki} + \mu_i$$

基本假设:解释变量 X_1, X_2, \dots, X_k 是确定性变量。

如果存在一个或多个随机变量作为解释变量,则称原模型出现**随机解释变量问题**。

假设 X_2 为随机解释变量。对于随机解释变量问题,分三种不同情况:

1. 随机解释变量与随机误差项独立
(Independence)

$$Cov(X_2, \mu) = E(X_2 \mu) = E(X_2)E$$

2. 随机解释变量与随机误差项同期无关
(contemporaneously uncorrelated), 但异期相关。

$$Cov(X_{2i}, \mu_i) = E(X_{2i} \mu_i)$$

$$Cov(X_{2i}, \mu_{i-s}) = E(X_{2i} \mu_{i-s}) \quad s \neq 0$$

3. 随机解释变量与随机误差项同期相关
(contemporaneously correlated)。

$$Cov(X_{2i}, \mu_i) = E(X_{2i} \mu_i)$$

二、实际经济问题中的随机解释变量问题

在实际经济问题中，经济变量往往都具有随机性。

但是在单方程计量经济学模型中，凡是外生变量都被认为是确定性的。

随机解释变量问题：

- ✓ 模型设定错误
- ✓ 测量误差
- ✓ 联立性

模型设定错误

- ✓ 模型设定错误是导致内生性最常见的原因
- ✓ 模型设定错误往往表现为相关变量的缺失
- ✓ 缺失变量成为错误设定模型误差项的一部分
- ✓ 当缺失变量和模型中其他变量相关时，就会导致这些变量的内生性。

工资模型

$$\log(\text{wage}) = \beta_0 + \beta_1 \text{exper} + \beta_2 (\text{exper})^2 + \beta_3 \text{edu} + \beta_4 \text{abl} + \mu$$

其中：**wage**：工资 **exper**：工作年限

edu：受教育年限 **abl**：工作能力

由于工作能力不可观测，变量**abl**不得不丢掉，建立模型

$$\log(\text{wage}) = \beta_0 + \beta_1 \text{exper} + \beta_2 (\text{exper})^2 + \beta_3 \text{edu} + \mu'$$

其中： $\mu' = \mu + \text{abl}$

一般**abl**和**edu**存在正相关性（受教育年限越长，工作能力越强），所以模型存在着内生性

测量误差

- ✓ 非系统性的因变量测量误差不会引起解释变量的内生性
- ✓ 解释变量测量误差会带来内生性

例1 考虑以下模型

$$Y_t^* = \alpha_0 + \alpha_1 X_t + \mu_t$$

由于 Y_t^* 不可直接观测，所以使用一个可观测变量 Y_t

$$Y_t = Y_t^* + \varepsilon_t,$$

ε_t 表示 Y_t^* 中的测量误差。建立 Y_t 关于 X_t 的回归模型

$$\hat{Y}_t = \hat{\beta}_0 + \hat{\beta}_1 X_t$$

假设 $E(\mu_t) = E(\varepsilon_t) = 0$, $Cov(X_t, \mu_t) = 0$,

$Cov(X_t, \varepsilon_t) = 0$ 即 Y_t^* 的测量误差与 X_t 不相关,

$Cov(\mu_t, \varepsilon_t) = 0$ 即方程误差与测量误差不相关。

讨论 $\hat{\beta}_1$ 的性质。

例2 考虑以下模型

$$Y_t = \beta_0 + \beta_1 X_t^* + \mu_t$$

假设解释变量 X_t^* 的实测值 X_t 与之有偏误： $X_t = X_t^* + \varepsilon_t$,

其中 ε_t 是具有零均值，无序列相关，且与 X_t^* 和 μ_t 不相关的随机变量。

建立 Y_t 关于 X_t 的回归模型，试讨论模型的内生性。

当用滞后被解释变量作为模型的解释变量时，也会产生随机解释变量问题

例如：

(1) 耐用品存量调整模型：

耐用品的存量 Q_t 由前一个时期的存量 Q_{t-1} 和当期收入 I_t 共同决定：

$$Q_t = \beta_0 + \beta_1 I_t + \beta_2 Q_{t-1} + \mu_t \quad t=1, \dots, T$$

这是一个滞后被解释变量作为解释变量的模型。

如果模型不存在随机误差项的序列相关性，那么随机解释变量 Q_{t-1} 只与 μ_{t-1} 相关，与 μ_t 不相关，属于上述的第2种情况。

(2) 合理预期的消费函数模型

合理预期理论认为消费 C_t 是由对收入的预期 Y_t^e 所决定的：

$$C_t = \beta_0 + \beta_1 Y_t^e + \mu_t$$

预期收入 Y_t^e 与实际收入 Y 间存如下关系的假设

$$Y_t^e = (1 - \lambda)Y_t + \lambda Y_{t-1}^e$$

容易推出

$$\begin{aligned} C_t &= \beta_0 + \beta_1(1 - \lambda)Y_t + \beta_1\lambda Y_{t-1}^e + \mu_t \\ &= \beta_0 + \beta_1(1 - \lambda)Y_t + \lambda(C_{t-1} - \beta_0 - \mu_{t-1}) + \mu_t \\ &= \beta_0(1 - \lambda) + \beta_1(1 - \lambda)Y_t + \lambda C_{t-1} + \mu_t - \lambda\mu_{t-1} \end{aligned}$$

C_{t-1} 是一随机解释变量，且与 $(\mu_t - \lambda\mu_{t-1})$ 高度相关。
属于上述第3种情况。

三、随机解释变量的后果

计量经济学模型一旦出现随机解释变量，且与随机扰动项相关的话，如果仍采用OLS法估计模型参数，不同性质的随机解释变量会产生不同的后果。

渐进无偏性和一致性的

- ▶ 如果一个随机变量的精确抽样分布很难得到，那近这里
- ▶ 当模型满足OLS的假定条件时，其参数的OLS估计偏性和有效性
- ▶ 有时OLS估计量并不具有这种特征，但随样本容量有了这种特征
- ▶ 随着随样本容量的增加，随机变量的分布称为渐应的统计特性称为渐进特性
- ▶ 即渐进无偏性和渐进有

► 设 $\hat{\beta}^n$ 是参数 β 的估计量，其中 n 为

► 设其数学期望值为 $E(\hat{\beta}^n)$ ，方差为 $\text{Var}(\hat{\beta}^n) = E[\hat{\beta}^n - E(\hat{\beta}^n)]^2$

► 所谓渐进分布是指，当样本容量 $n \rightarrow \infty$ 时，上量序列分别收敛到一定分布。均值、方差有以下

$$\lim_{n \rightarrow \infty} E(\hat{\beta}^n) = E(\hat{\beta})$$

$$\lim_{n \rightarrow \infty} \text{Var}(\hat{\beta}^n) = E[\hat{\beta} - E(\hat{\beta})]^2$$

► $E(\hat{\beta})$ 和 $E[\hat{\beta} - E(\hat{\beta})]^2$ 分别为 $\hat{\beta}^n$ 的渐进期望值

(1) 渐近无偏

▶ 如果 $\lim_{n \rightarrow \infty} E(\hat{\beta}^n) = \beta$ 称 $\hat{\beta}^n$ 为 β 的渐近无偏估计。
量 n 充分大时 $\hat{\beta}^n$ 的均值趋向于总体参数 β 。

▶ 如果小样本估计量是有偏的，但其估计量具有

▶ 我们就可以增加样本，以优化估

(2) 一致

- ▶ 所谓一致性估计是指对于任意给定的任意，
满足

$$\lim_{n \rightarrow \infty} P\{|\hat{\beta}^n - \beta| < \varepsilon\} = 1$$

- ▶ 即当样本容量 n 充分大时 $\hat{\beta}^n$ 值趋于总体真实，
近1，记为

$$P \lim_{n \rightarrow \infty} \hat{\beta}^n = \beta$$

简记为

一致估计的充分条件为：

$$\lim_{n \rightarrow \infty} E(\hat{\beta}^n) = \beta \quad \text{和} \quad \lim_{n \rightarrow \infty} \text{Var}(\hat{\beta}^n) = 0$$

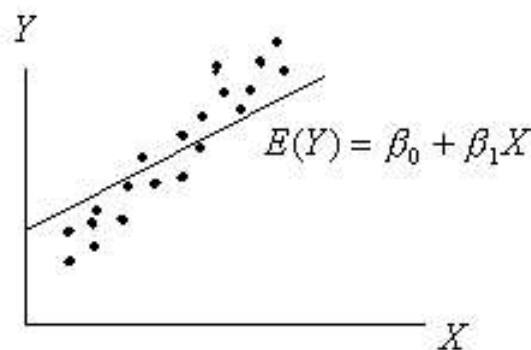
► $P \lim$ 表示概率极限。概率极限有以下运算法则

$$P \lim(C_1 X_1 + C_2 X_2) = C_1 P \lim(X_1) + C_2 P \lim(X_2)$$

$$P \lim(X_1 X_2) = P \lim(X_1) P \lim(X_2)$$

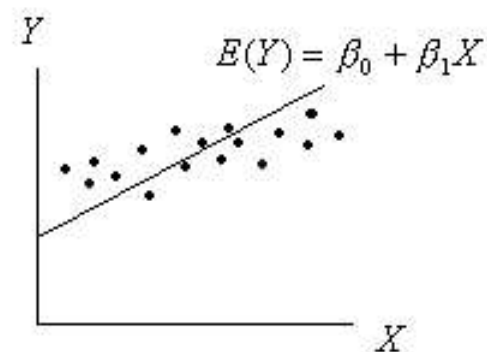
$$P \lim \left(\frac{X_1}{X_2} \right) = \frac{P \lim X_1}{P \lim X_2}$$

- 随机解释变量与随机误差项相关图



(a) 正相关

拟合的样本回归线可能低估截距项，而高估斜率项。



(b) 负相关

拟合的样本回归线高估截距项，而低估斜率项。

对一元线性回归模型：

$$Y_t = \beta_0 + \beta_1 X_t + \mu_t$$

OLS估计量为

$$\hat{\beta}_1 = \frac{\sum x_t y_t}{\sum x_t^2} = \beta_1 + \frac{\sum x_t \mu_t}{\sum x_t^2}$$

随机解释变量 X 与随机项 μ 的关系不同，参数OLS估计量的统计性质也会不同。

1、如果 X 与 μ 相互独立，得到的参数估计量仍然是无偏、一致估计量。

证明

2、如果 \mathbf{X} 与 μ 同期不相关，异期相关，得到的参数估计量有偏、但却是一致的。

$$E(\hat{\beta}_1) = \beta_1 + E\left(\frac{\sum x_t \mu_t}{\sum x_t^2}\right) = \beta_1 + \sum E(k_t \mu_t)$$

k_t 的分母中包含不同期的 \mathbf{X} ；由异期相关性知： k_t 与 μ_t 相关，因此，

$$E(\hat{\beta}_1) \neq \beta_1$$

但是

$$P \lim_{n \rightarrow \infty} \left(\beta_1 + \frac{\sum x_t \mu_t}{\sum x_t^2} \right) = \beta_1 + \frac{P \lim(\frac{1}{n} \sum x_t \mu_t)}{P \lim(\frac{1}{n} \sum x_t^2)}$$

3、如果 \mathbf{X} 与 μ 同期相关，得到的参数估计量有偏、且非一致。

$$E(\hat{\beta}_1) = \beta_1 + E\left(\sum \frac{x_t}{\sum x_t^2} \mu_t\right) = \beta_1 + \sum E(k_t \mu_t)$$

k_t 的分母中包含不同期的 \mathbf{X} ，因此

$$E(\hat{\beta}_1) \neq \beta_1$$

并且

$$P \lim_{n \rightarrow \infty} \left(\beta_1 + \frac{\sum x_t \mu_t}{\sum x_t^2} \right) = \beta_1 + \frac{P \lim(\frac{1}{n} \sum x_t \mu_t)}{P \lim(\frac{1}{n} \sum x_t^2)}$$

对一元线性回归模型：

$$Y_t = \beta_0 + \beta_1 X_t + \mu_t$$

1、如果 X 与 μ 相互独立，得到的参数估计量仍然是无偏、一致估计量。

2、如果 X 与 μ 同期不相关，异期相关，得到的参数估计量有偏、但却是一致的。

3、如果 X 与 μ 同期相关，得到的参数估计量有偏、且非一致。

注意：

在多元线性回归模型中，

内生解释变量回归系数的OLS估计不是一致估计，

并且会引起与之具有相关性的外生解释变量回归系数OLS估计的不一致性。

四、工具变量法

模型中出现随机解释变量且与随机误差项相关时，OLS估计量是有偏的。

如果随机解释变量与随机误差项异期相关，则可以通过增大样本容量的办法来得到一致的估计量；

但如果是同期相关，即使增大样本容量也无济于事。这时，最常用的估计方法是**工具变量法**（**Instrument variables**）。

1、工具变量的选取

工具变量：在模型估计过程中被作为工具使用，以替代模型中与随机误差项相关的随机解释变量。

选择为工具变量的变量必须满足以下条件：

- (1) 与所替代的随机解释变量高度相关；
- (2) 与随机误差项不相关；
- (3) 与模型中其它解释变量不相关，以避免出现多重共线性。

2、工具变量的应用

以一元回归模型的离差形式为例说明如下：

$$y_i = \beta_1 x_i + \mu_i$$

用OLS估计模型，相当于用 x_i 去乘模型两边、对 i 求和、再略去 $\sum x_i \mu_i$ 项后得到正规方程：

$$\sum x_i y_i = \beta_1 \sum x_i^2$$

解得

$$\hat{\beta}_1 = \frac{\sum x_i y_i}{\sum x_i^2} \quad (*)$$

由于 $\text{Cov}(X_i, \mu_i) = 0$ ，意味着大样本下

$$(\sum x_i \mu_i) / n \rightarrow 0$$

表明大样本下OLS估计量

$$\hat{\beta}_1 = \frac{\sum x_i y_i}{\sum x_i^2}$$

具有一致性。

然而，如果 x_i 与 μ_i 相关，即使在大样本下，也不存在 $(\sum x_i \mu_i) / n \rightarrow 0$ ，则OLS估计量

$$\hat{\beta}_1 = \frac{\sum x_i y_i}{\sum x_i^2}$$

不具有 consistency。

如果选择 Z 为 X 的**工具变量**，那么在上述估计过程可改为：

$$\sum z_i y_i = \beta_1 \sum z_i x_i + \sum z_i \mu_i$$

利用 $E(z_i \mu_i) = 0$ ，在大样本下可得到：

$$\tilde{\beta}_1 = \frac{\sum z_i y_i}{\sum z_i x_i}$$

关于 β_0 的估计，仍用 $\tilde{\beta}_0 = \bar{Y} - \tilde{\beta}_1 \bar{X}$ 完成。

这种求模型参数估计量的方法称为**工具变量法** (**instrumental variable method**)，相应的估计量称为**工具变量法估计量** (**instrumental variable (IV) estimator**)。

2、工具变量的应用

- 对于多元线性模型

$$Y_i = \beta_0 + \beta_1 X_{1i} + \beta_2 X_{2i} + \cdots + \beta_k X_{ki} + \mu_i$$

$i=1,2,\dots,n$

- 用普通最小二乘法估计模型，最后归结为求解一个关于参数估计量的正规方程组：

$$\left\{ \begin{array}{l} \Sigma(\hat{\beta}_0 + \hat{\beta}_1 X_{1i} + \hat{\beta}_2 X_{2i} + \cdots + \hat{\beta}_k X_{ki}) = \Sigma Y_i \\ \Sigma(\hat{\beta}_0 + \hat{\beta}_1 X_{1i} + \hat{\beta}_2 X_{2i} + \cdots + \hat{\beta}_k X_{ki}) X_{1i} = \Sigma Y_i X_{1i} \\ \Sigma(\hat{\beta}_0 + \hat{\beta}_1 X_{1i} + \hat{\beta}_2 X_{2i} + \cdots + \hat{\beta}_k X_{ki}) X_{2i} = \Sigma Y_i X_{2i} \\ \vdots \\ \Sigma(\hat{\beta}_0 + \hat{\beta}_1 X_{1i} + \hat{\beta}_2 X_{2i} + \cdots + \hat{\beta}_k X_{ki}) X_{ki} = \Sigma Y_i X_{ki} \end{array} \right.$$

- 该方程组也可以看作为矩方法的结果。用每个解释变量分别乘以模型的两边，并对所有样本点求和：

$$\left\{ \begin{array}{l} \Sigma y_i = \Sigma(\beta_0 + \beta_1 x_{1i} + \beta_2 x_{2i} + \cdots + \beta_k x_{ki} + \mu_i) \\ \Sigma y_i x_{1i} = \Sigma(\beta_0 + \beta_1 x_{1i} + \beta_2 x_{2i} + \cdots + \beta_k x_{ki} + \mu_i) x_{1i} \\ \Sigma y_i x_{2i} = \Sigma(\beta_0 + \beta_1 x_{1i} + \beta_2 x_{2i} + \cdots + \beta_k x_{ki} + \mu_i) x_{2i} \\ \vdots \\ \Sigma y_i x_{ki} = \Sigma(\beta_0 + \beta_1 x_{1i} + \beta_2 x_{2i} + \cdots + \beta_k x_{ki} + \mu_i) x_{ki} \end{array} \right.$$

- 然后再对方程的两边求期望：

$$\left\{ \begin{array}{l} E(\Sigma y_i) = E(\Sigma(\beta_0 + \beta_1 x_{1i} + \beta_2 x_{2i} + \cdots + \beta_k x_{ki} + \mu_i)) \\ E(\Sigma y_i x_{1i}) = E(\Sigma(\beta_0 + \beta_1 x_{1i} + \beta_2 x_{2i} + \cdots + \beta_k x_{ki} + \mu_i) x_{1i}) \\ E(\Sigma y_i x_{2i}) = E(\Sigma(\beta_0 + \beta_1 x_{1i} + \beta_2 x_{2i} + \cdots + \beta_k x_{ki} + \mu_i) x_{2i}) \\ \vdots \\ E(\Sigma y_i x_{ki}) = E(\Sigma(\beta_0 + \beta_1 x_{1i} + \beta_2 x_{2i} + \cdots + \beta_k x_{ki} + \mu_i) x_{ki}) \end{array} \right.$$

- 利用下列条件得到的:

$$E(\mu_i) = 0$$

$$E(\mu_i x_{ji}) = 0, j = 1, 2, \dots, k$$

$$E(\beta_j) = \hat{\beta}_j, j = 0, 1, 2, \dots, k$$

- 如果 \mathbf{x}_2 为随机变量，且与随机误差项相关；选择 \mathbf{z} 作为它的工具变量。在应该用 \mathbf{x}_2 乘方程两边时，不用 \mathbf{x}_2 ，而用 \mathbf{z} 。

$$\left\{ \begin{array}{l} \Sigma y_i = \Sigma(\beta_0 + \beta_1 x_{1i} + \beta_2 x_{2i} + \cdots + \beta_k x_{ki} + \mu_i) \\ \Sigma y_i x_{1i} = \Sigma(\beta_0 + \beta_1 x_{1i} + \beta_2 x_{2i} + \cdots + \beta_k x_{ki} + \mu_i) x_{1i} \\ \Sigma y_i \mathbf{z}_i = \Sigma(\beta_0 + \beta_1 x_{1i} + \beta_2 x_{2i} + \cdots + \beta_k x_{ki} + \mu_i) \mathbf{z}_i \\ \vdots \\ \Sigma y_i x_{ki} = \Sigma(\beta_0 + \beta_1 x_{1i} + \beta_2 x_{2i} + \cdots + \beta_k x_{ki} + \mu_i) x_{ki} \end{array} \right.$$

- 得到采用工具变量法的正规方程组：

$$\left\{ \begin{array}{l} \Sigma y_i = \Sigma (\hat{\beta}_0 + \hat{\beta}_1 x_{1i} + \hat{\beta}_2 x_{2i} + \cdots + \hat{\beta}_k x_{ki}) \\ \Sigma y_i x_{1i} = \Sigma (\hat{\beta}_0 + \hat{\beta}_1 x_{1i} + \hat{\beta}_2 x_{2i} + \cdots + \hat{\beta}_k x_{ki}) x_{1i} \\ \Sigma y_i z_i = \Sigma (\hat{\beta}_0 + \hat{\beta}_1 x_{1i} + \hat{\beta}_2 x_{2i} + \cdots + \hat{\beta}_k x_{ki}) z_i \\ \vdots \\ \Sigma y_i x_{ki} = \Sigma (\hat{\beta}_0 + \hat{\beta}_1 x_{1i} + \hat{\beta}_2 x_{2i} + \cdots + \hat{\beta}_k x_{ki}) x_{ki} \end{array} \right.$$

- 求解该方程组即可得到关于原模型参数的工具变量法估计量。

- 对于矩阵形式:

$$Y = XB + N$$

采用工具变量法（假设 x_2 与随机项相关，用工具变量 z 替代）得到的正规方程组为：

$$Z'Y = Z'X\hat{B}$$

对于矩阵形式：

$$\mathbf{Y} = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\mu}$$

采用工具变量法（假设 \mathbf{X}_2 与随机项相关，用工具变量 \mathbf{Z} 替代）得到的正规方程组为：

$$\mathbf{Z}'\mathbf{Y} = \mathbf{Z}'\mathbf{X}\boldsymbol{\beta}$$

参数估计量为：

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

其中

$$\mathbf{Z}' = \begin{bmatrix} 1 & 1 & \cdots & 1 \\ X_{11} & X_{12} & \cdots & X_{1n} \\ Z_1 & Z_2 & \cdots & Z_n \\ \vdots & & & \\ X_{k1} & X_{k2} & \cdots & X_{kn} \end{bmatrix}$$

称为工具变量矩阵

参数估计量为：

$$\hat{B} = (Z'X)^{-1} Z'Y$$

其中

$$Z' = \begin{bmatrix} 1 & 1 & \cdots & 1 \\ x_{11} & x_{12} & \cdots & x_{1n} \\ z_1 & z_2 & \cdots & z_n \\ \vdots & & & \\ x_{k1} & x_{k2} & \cdots & x_{kn} \end{bmatrix}$$

通常，对于没有选择另外的变量作为工具变量的解释变量，可以认为用自身作为工具变量。于是被称为**工具变量矩阵**。

3、工具变量法估计量是一致估计量

一元回归中，工具变量法估计量为

$$\tilde{\beta}_1 = \frac{\sum z_i(\beta_1 x_i + \mu_i)}{\sum z_i x_i} = \beta_1 + \frac{\sum z_i \mu_i}{\sum z_i x_i}$$

两边取概率极限得：

$$P\lim(\tilde{\beta}_1) = \beta_1 + \frac{P\lim \frac{1}{n} \sum z_i \mu_i}{P\lim \frac{1}{n} \sum z_i x_i}$$

如果工具变量 Z 选取恰当，即有

$$P\lim \frac{1}{n} \sum z_i \mu_i = \text{cov}(Z_i, \mu_i) = 0 \quad P\lim \frac{1}{n} \sum z_i x_i = \text{cov}(Z_i, X_i) \neq 0$$

因此：

$$P\lim(\tilde{\beta}_1) = \beta_1$$

多元回归情况下，工具变量法估计量仍是一致估计量

注意：

1、在小样本下，工具变量法估计量仍是有偏的。

$$E\left(\frac{1}{\sum z_i x_i} \sum z_i \mu_i\right) \neq E\left(\frac{1}{\sum z_i x_i}\right) E\left(\sum z_i \mu_i\right) = 0$$

2、工具变量并没有替代模型中的解释变量，只是在估计过程中作为“工具”被使用。

上述工具变量法估计过程可等价地分解成下面的两步最小二乘法（TSLS）回归：

第一步，用OLS法进行X关于工具变量Z的回归：

$$\hat{X}_i = \hat{\alpha}_0 + \hat{\alpha}_1 Z_i$$

第二步，以第一步得到的 \hat{X} 为解释变量，进行如下OLS回归：

$$\hat{Y}_i = \tilde{\beta}_0 + \tilde{\beta}_1 \hat{X}_i$$

容易验证仍有：

$$\tilde{\beta}_1 = \frac{\sum z_i y_i}{\sum z_i x_i}$$

因此，工具变量法仍是Y对X的回归，而不是对Z的回归。

3、如果模型中有两个以上的随机解释变量与随机误差项相关，就必须找到两个以上的工具变量。但是，一旦工具变量选定，它们在估计过程被使用的次序不影响估计结果。

4、如果1个随机解释变量可以找到多个互相独立的工具变量

可以用TSLS法找出最优工具变量

两阶段最小二乘法：TSLS

一个内生自变量

$$Y = \beta_0 + \beta_1 X_1 + \beta_2 X_2 + \beta_3 X_3 + u$$

X_1 为内生变量， X_2 和 X_3 为外生变量， Z_1 、 Z_2 为 X_1 的工具变量。

两阶段最小二乘步骤：

第一阶段（first stage）：以内生变量为因变量，所有外生变量为自变量做回归

$$X_1 = \alpha_0 + \alpha_1 Z_1 + \alpha_2 Z_2 + \alpha_3 X_2 + \alpha_4 X_3 + v$$

得拟合值

$$\hat{X}_1 = \hat{\alpha}_0 + \hat{\alpha}_1 Z_1 + \hat{\alpha}_2 Z_2 + \hat{\alpha}_3 X_2 + \hat{\alpha}_4 X_3$$

第二阶段（second stage）：将 \hat{x}_1 作为 x_1 的工具变量，对模型 $Y = \beta_0 + \beta_1 \hat{X}_1 + \beta_2 X_2 + \beta_3 X_3 + u$ 实施工具变量估计

两阶段最小二乘法：TSLS

多个内生自变量

$$Y = \beta_0 + \beta_1 X_1 + \beta_2 X_2 + \beta_3 X_3 + u$$

X_1 、 X_2 为内生变量, X_3 为外生变量, Z_1 和 Z_2 为 X_1 的工具变量, W 为 X_2 的工具变量

第一阶段（first stage）：

分别以内生变量 X_1 和 X_2 为因变量，以所有外生变量 Z_1 、 Z_2 、 W 和 X_3 为自变量进行回归，即

$$X_1 = \alpha_0 + \alpha_1 Z_1 + \alpha_2 Z_2 + \alpha_3 W + \alpha_4 X_3 + v$$

$$X_2 = \gamma_0 + \gamma_1 Z_1 + \gamma_2 Z_2 + \gamma_3 W + \gamma_4 X_3 + \varepsilon$$

两阶段最小二乘法：TSLS

第二阶段：

以 \hat{X}_1 和 \hat{X}_2 代替 X_1 和 X_2 对原模型进行OLS估计，即对模型

$$Y = \beta_0 + \beta_1 \hat{X}_1 + \hat{\beta}_2 X_2 + \beta_3 X_3 + u$$

进行OLS估计，得出回归系数的一致估计。

5、如果1个随机解释变量可以找到多个互相独立的工具变量，人们希望充分利用这些工具变量的信息，就形成了广义矩方法（**Generalized Method of Moments, GMM**）。

- 在**GMM**中，矩条件大于待估参数的数量，于是如何求解成为它的核心问题。
- **GMM**是近20年计量经济学理论方法发展的重要方向之一。
- **IV**是**GMM**的一个特例。

5、OLS可以看作工具变量法的一种特殊情况。

6、要找到与随机扰动项不相关而又与随机解释变量相关的工具变量并不是一件很容易的事

五、解释变量的内生性检验 ——豪斯曼检验

- ◆ 自变量若内生，OLS估计会不一致；
- ◆ 自变量若外生，盲目用工具变量会降低有效性
- ◆ 故需要检验自变量是否内生。

$$Y = \beta_0 + \beta_1 X_1 + \beta_2 X_2 + \beta_3 X_3 + u$$

X_1 和 X_2 为外生变量，对 X_3 的内生性检验，设 Z 为 X_3 的工具变量，将 X_3 对 X_1 、 X_2 和 Z 回归

$$X_3 = \gamma_0 + \gamma_1 X_1 + \gamma_2 X_2 + \gamma_3 Z + v$$

若有内生性，则是 u 和 v 之间有关系：

$$u = \rho v + \varepsilon$$

$$H_0 : \rho = 0; \quad H_1 : \rho \neq 0$$

将上述 u 和 v 的关系代入原模型

$$Y = \beta_0 + \beta_1 X_1 + \beta_2 X_2 + \beta_3 X_3 + \rho v + \varepsilon$$

v 不可观测，用 \hat{v} 代替（从模型 $X_3 = \gamma_0 + \gamma_1 X_1 + \gamma_2 X_2 + \gamma_3 Z + v$ 中估计而来），即最终估计模型：

$$Y = \beta_0 + \beta_1 X_1 + \beta_2 X_2 + \beta_3 X_3 + \rho \hat{v} + \varepsilon$$

并检验 ρ 是否为0。

六、案例——中国居民人均消费函数

例4.4.1 在例2.5.1的中国居民人均消费函数的估计中，采用OLS估计了下面的模型：

$$CONSP = \beta_0 + \beta_1 GDPP + \mu$$

由于：居民人均消费支出（CONSP）与人均国内生产总值（GDPP）相互影响，因此，

容易判断**GDPP**与 **μ** 同期相关（往往是正相关），OLS估计量有偏并且是非一致的（低估截距项而高估计斜率项）。

OLS估计结果:

$$\widehat{CONSP}=201.11+0.3862GDPP$$

$$(13.51) \quad (53.47)$$

$$R^2=0.9927 \quad F=2859.23 \quad DW=0.5503 \quad SSR=23240.7$$

如果用 $GDPP_{t-1}$ 为工具变量, 可得如下工具变量法估计结果:

$$\widehat{CONSP}=212.45+0.3817GDPP$$

$$(14.84) \quad (56.04)$$

$$R^2=0.9937 \quad F=3140.58 \quad DW=0.6691 \quad SSR=18366.5$$

例1 . 某国的政府税收 T (单位: 百万美元)、国内生产总值 GDP (单位: 车数量 Z (单位: 百万辆) 的观测数据如下表所示:

序号	T	GDP	Z
1	3	4	5
2	2	1	2
3	5	7	6
4	6	8	7
7	7	8	6
8	9	11	7
9	8	10	7

要求: 试以汽车数量 Z 作为国内生产总值 GDP 的工具变量,

例2? 对于模型

$$Y_t = \beta_0 + \beta_1 X_{1t} + \beta_2 X_{2t} + \beta_3 Y_{t-1} + \mu_t \quad \text{中,}$$

假设 Y_{t-1} 与 μ_t 相关。为了消除该相关性，采用工具变量法：

先求 Y_t 关于 X_{1t} 与 X_{2t} 回归，得到 \hat{Y}_t ，再作如下回归：

$$Y_t = \beta_0 + \beta_1 X_{1t} + \beta_2 X_{2t} + \beta_3 \hat{Y}_{t-1} + \mu_t$$

试问：这一方法能否消除原模型中 Y_{t-1} 与 μ_t 的相关性？为什么？

- 例3 书149页第4题（缺了后两小问）

例3对于模型

$$Y_t = \beta_0 + \beta_1 X_t^* + \mu_t$$

假设解释变量 X_t^* 的实测值 X_t 与之有偏误： $X_t = X_t^* + e_t$,

其中 e_t 是具有零均值，无序列相关，且与 X_t^* 和 μ_t 不相关的随机变量。
试问：

（1）能否将 $X_t^* = X_t - e_t$ 代入原模型，使之变换成 $Y_t = \beta_0 + \beta_1 X_t + v_t$ 后进行估计？其中 v_t 为变换后模型的随机干扰项。

（2）进一步假设 μ_t 与 e_t 之间，以及它们与 X_t^* 之间无异期相关，那么 $E(X_{t-1}v_t) = 0$ 成立吗？ X_t 与 X_{t-1} 相关吗？

（3）由（2）的结论，你能寻找什么样的工具变量对变换后的模型进行估计？

例5? 在二元回归模型 $Y = \beta_0 + \beta_1 X_1 + \beta_2 X_2 + \varepsilon$ 中，
 X_1 为内生变量， X_2 为外生变量。变量 Z_1 ， Z_2 为 X_1 的工具变量。
TSLS估计法的第一步是将 X_1 对 Z_1 ， Z_2 和 X_2 进行回归，
并以回归的拟合值 \hat{X}_1 作为 X_1 的工具变量。
为什么第一步回归中要包含 X_2 ? 不包含 X_2 会有什么影响?

作业

- 书 149页 5

§ 4.4 模型设定偏误问题

- 一、模型设定偏误的类型
- 二、模型设定偏误的后果
- 三、模型设定偏误的检验

一、模型设定偏误的类型

- 模型设定偏误主要有两大类：
 - (1) 关于解释变量选取的偏误，主要包括漏选相关变量和多选无关变量，
 - (2) 关于模型函数形式选取的偏误。

1、相关变量的遗漏 (omitting relevant variables)

- 例如，如果“正确”的模型为

$$Y = \beta_0 + \beta_1 X_1 + \beta_2 X_2 + \mu$$

而我们将模型设定为

$$Y = \alpha_0 + \alpha_1 X_1 + v$$

即设定模型时漏掉了一个相关的解释变量。

这类错误称为遗漏相关变量。

2、无关变量的误选 (including irrevelant variables)

- 例如，如果

$$Y = \beta_0 + \beta_1 X_1 + \beta_2 X_2 + \mu$$

仍为“真”，但我们将模型设定为

$$Y = \alpha_0 + \alpha_1 X_1 + \alpha_2 X_2 + \alpha_3 X_3 + \mu$$

即设定模型时，多选了一个无关解释变量。

3、错误的函数形式 (wrong functional form)

- 例如，如果“真实”的回归函数为

$$Y = AX_1^{\beta_1} X_2^{\beta_2} e^{\mu}$$

但却将模型设定为

$$Y = \beta_0 + \beta_1 X_1 + \beta_2 X_2 + v$$

二、模型设定偏误的后果

- 当模型设定出现偏误时，模型估计结果也会与“实际”有偏差。这种偏差的性质与程度与模型设定偏误的类型密切相关。

1、遗漏相关变量偏误

采用遗漏相关变量的模型进行估计而带来的偏误称为遗漏相关变量偏误（omitting relevant variable bias）。

设正确的模型为

$$Y = \beta_0 + \beta_1 X_1 + \beta_2 X_2 + \mu$$

却对

$$Y = \alpha_0 + \alpha_1 X_1 + v$$

进行回归，得

$$\hat{\alpha}_1 = \frac{\sum x_{1i} y_i}{\sum x_{1i}^2}$$

将正确模型 $Y = \beta_0 + \beta_1 X_1 + \beta_2 X_2 + \mu$ 的离差形式

$$y_i = \beta_1 x_{1i} + \beta_2 x_{2i} + \mu_i - \bar{\mu}$$

代入 $\hat{\alpha}_1 = \frac{\sum x_{1i} y_i}{\sum x_{1i}^2}$ 得

$$\begin{aligned}\hat{\alpha}_1 &= \frac{\sum x_{1i} y_i}{\sum x_{1i}^2} = \frac{\sum x_{1i} (\beta_1 x_{1i} + \beta_2 x_{2i} + \mu_i - \bar{\mu})}{\sum x_{1i}^2} \\ &= \beta_1 + \beta_2 \frac{\sum x_{1i} x_{2i}}{\sum x_{1i}^2} + \frac{\sum x_{1i} (\mu_i - \bar{\mu})}{\sum x_{1i}^2}\end{aligned}$$

$$\begin{aligned}\hat{\alpha}_1 &= \frac{\sum x_{1i} y_i}{\sum x_{1i}^2} = \frac{\sum x_{1i} (\beta_1 x_{1i} + \beta_2 x_{2i} + \mu_i - \bar{\mu})}{\sum x_{1i}^2} \\ &= \beta_1 + \beta_2 \frac{\sum x_{1i} x_{2i}}{\sum x_{1i}^2} + \frac{\sum x_{1i} (\mu_i - \bar{\mu})}{\sum x_{1i}^2}\end{aligned}$$

(1) 如果漏掉的 X_2 与 X_1 相关，则上式中的第二项在小样本下求期望与大样本下求概率极限都不会为零，从而使得OLS估计量在小样本下有偏，在大样本下非一致。

(2) 如果 X_2 与 X_1 不相关，则 α_1 的估计满足无偏性与一致性；但这时 α_0 的估计却是有偏的。

(3) 随机扰动项 μ 的方差估计 $\hat{\sigma}^2$ 也是有偏的。

(4) $\hat{\alpha}_1$ 的方差是真实估计量 $\hat{\beta}_1$ 的方差的有偏估计。

由 $Y = \alpha_0 + \alpha_1 X_1 + v$ 得

$$Var(\hat{\alpha}_1) = \frac{\sigma^2}{\sum x_{1i}^2}$$

由 $Y = \beta_0 + \beta_1 X_1 + \beta_2 X_2 + \mu$ 得

$$Var(\hat{\beta}_1) = \sigma^2 \frac{\sum x_{2i}^2}{\sum x_{1i}^2 \sum x_{2i}^2 - (\sum x_{1i} x_{2i})^2} = \frac{\sigma^2}{\sum x_{1i}^2 (1 - r_{x_1 x_2}^2)}$$

如果 X_2 与 X_1 相关，显然有 $Var(\hat{\alpha}_1) \neq Var(\hat{\beta}_1)$

如果 X_2 与 X_1 不相关，也有 $Var(\hat{\alpha}_1) \neq Var(\hat{\beta}_1)$

(5) 通常的置信区间和假设检验程序对于所估计参数的统计显著性容易导出错误性的结论。

(6) 基于不正确模型做出的预测及预测区间都是不可靠的。

2、包含无关变量偏误

采用包含无关解释变量的模型进行估计带来的偏误，称为包含无关变量偏误（including irrelevant variable bias）。

设
$$Y = \alpha_0 + \alpha_1 X_1 + v \quad (*)$$

为正确模型，但却估计了

$$Y = \beta_0 + \beta_1 X_1 + \beta_2 X_2 + \mu \quad (**)$$

如果 $\beta_2 = 0$ ，则 (**) 与 (*) 相同，因此，可将 (**) 式视为以 $\beta_2 = 0$ 为约束的 (*) 式的特殊形式。

由于所有的经典假设都满足，因此对

$$Y = \beta_0 + \beta_1 X_1 + \beta_2 X_2 + \mu \quad (**)$$

式进行OLS估计，可得到无偏且一致的估计量。

注意：由于 $\beta_2 = 0$ ，因此， $E(\hat{\beta}_2) = 0$ 。

但是，OLS估计量却不具有最小方差性。

$Y = \alpha_0 + \alpha_1 X_1 + v$ 中 X_1 的方差：

$$Var(\hat{\alpha}_1) = \frac{\sigma^2}{\sum x_{1i}^2}$$

$Y = \beta_0 + \beta_1 X_1 + \beta_2 X_2 + \mu$ 中 X_1 的方差：

$$Var(\hat{\beta}_1) = \frac{\sigma^2}{\sum x_{1i}^2 (1 - r_{x_1 x_2}^2)}$$

当 X_1 与 X_2 完全线性无关时： $Var(\hat{\alpha}_1) = Var(\hat{\beta}_1)$

否则： $Var(\hat{\beta}_1) > Var(\hat{\alpha}_1)$

总结:

- ✓ 在多选无关解释变量的情形下，OLS估计是**无偏**且**一致**的估计量
- ✓ 随机干扰项的方差也能被正确估计
- ✓ 通常的置信区间和假设检验程序仍然有效
- ✓ 但是OLS估计往往是无效的

3、错误函数形式的偏误

当选取了错误函数形式并对其进行估计时，带来的偏误称**错误函数形式偏误**（wrong functional form bias）。

容易判断，这种**偏误是全方位的**。

例如，如果“真实”的回归函数为

$$Y = AX_1^{\beta_1} X_2^{\beta_2} e^{\mu}$$

却估计线性式

$$Y = \beta_0 + \beta_1 X_1 + \beta_2 X_2 + v$$

显然，两者的参数具有完全不同的经济含义，且估计结果一般也是不相同的。

三、模型设定偏误的检验

1、检验是否含有无关变量

可用t 检验与F检验完成。

检验的基本思想:如果模型中误选了无关变量,则其系数的真值应为零。因此,只须对无关变量系数的显著性进行检验。

t检验: 检验某1个变量是否应包括在模型中;

F检验: 检验若干个变量是否应同时包括在模型中

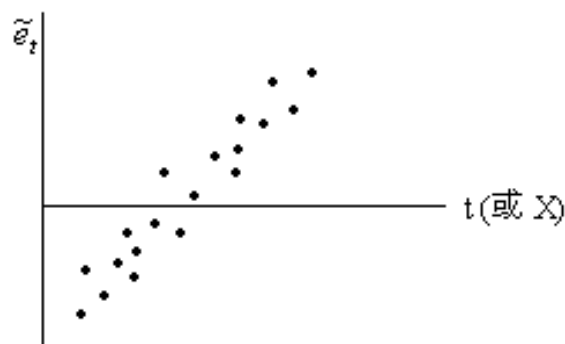
2、检验是否有相关变量的遗漏或函数形式设定偏误

(1) 残差图示法

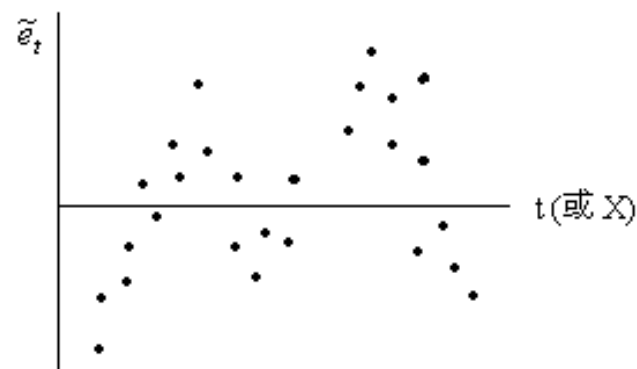
对所设定的模型进行OLS回归，得到估计的残差序列 \tilde{e}_t ；

做出 \tilde{e}_t 与时间 t 或某解释变量 X 的散点图，考察 \tilde{e}_t 是否有规律地在变动，以判断是否遗漏了重要的解释变量或选取了错误的函数形式。

- 残差序列变化图

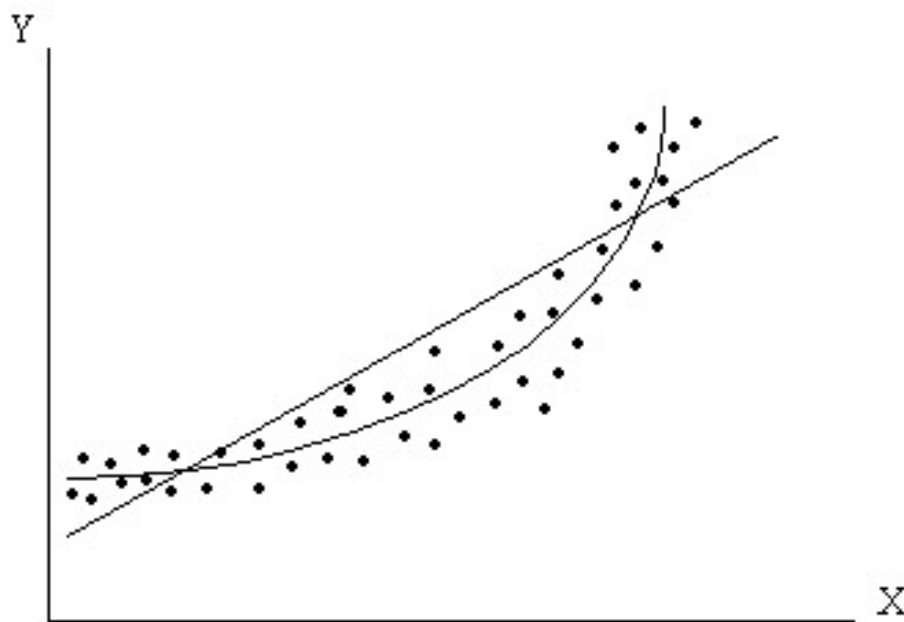


(a) 趋势变化：
模型设定时可能遗漏了一随着时间的推移而持续上升的变量



(b) 循环变化：
模型设定时可能遗漏了一随着时间的推移而呈现循环变化的变量

- 模型函数形式设定偏误时残差序列呈现正负交替变化



图示：一元回归模型中，真实模型呈幂函数形式，但却选取了线性函数进行回归。

(2) 一般性设定偏误检验

但更准确更常用的判定方法是拉姆齐(Ramsey)于1969年提出的所谓**RESET 检验** (regression error specification test)。

基本思想:

如果事先知道遗漏了哪个变量, 只需将此变量引入模型, 估计并检验其参数是否显著不为零即可;

问题是不知道遗漏了哪个变量, 需寻找一个替代变量 Z , 来进行上述检验。

RESET检验中, 采用所设定模型中被解释变量 Y 的估计值 \hat{Y} 的若干次幂来充当该“替代”变量。

例如，先估计 $Y = \alpha_0 + \alpha_1 X_1 + v$ 得

$$\hat{Y} = \hat{\alpha}_0 + \hat{\alpha}_1 X_1$$

再用通过残差项 \tilde{e}_t 与估计的 \hat{Y} 的图形判断引入 \hat{Y} 的若干次幂充当“替代”变量。

如 \tilde{e}_t 与 Y 的图形呈现曲线形变化时，回归模型可选为：

$$Y = \beta_0 + \beta_1 X_1 + \gamma_1 \hat{Y}^2 + \gamma_2 \hat{Y}^3 + \mu$$

再根据第三章第五节介绍的增加解释变量的F检验来判断是否增加这些“替代”变量。

若仅增加一个“替代”变量，也可通过t检验来判断。

RESET检验也可用来检验函数形式设定偏误的问题。

例如，在一元回归中，假设真实的函数形式是非线性的，用泰勒定理将其近似地表示为多项式

$$Y = \beta_0 + \beta_1 X_1 + \beta_2 X_1^2 + \beta_3 X_1^3 + \cdots + \mu \quad (*)$$

因此，如果设定了线性模型，就意味着遗漏了相关变量 X_1^2 、 X_1^3 ，等等。

因此，在一元回归中，可通过检验(*)式中的各高次幂参数的显著性来判断是否将非线性模型误设成了线性模型。

对多元回归，非线性函数可能是关于若干个或全部解释变量的非线性，这时可按遗漏变量的程序进行检验。

例如，估计 $Y = \beta_0 + \beta_1 X_1 + \beta_2 X_2 + \mu$

但却怀疑真实的函数形式是非线性的。

这时，只需以估计出的 \hat{Y} 的若干次幂为“替代”变量，进行类似于如下模型的估计

$$Y = \beta_0 + \beta_1 X_1 + \beta_2 X_2 + \gamma_1 \hat{Y}^2 + \gamma_2 \hat{Y}^3 + \mu$$

再判断各“替代”变量的参数是否显著地不为零即可。

案例：中国商品进口模型

我们主要研究中国商品进口与国内生产总值的关系。（下表）。

表 4.2.1 1978~2001 年中国商品进口与国内生产总值

	国内生产总值 GDP (亿元)	商品进口 M (亿美元)		国内生产总值 GDP (亿元)	商品进口 M (亿美元)
1978	3624.1	108.9	1990	18547.9	533.5
1979	4038.2	156.7	1991	21617.8	637.9
1980	4517.8	200.2	1992	26638.1	805.9
1981	4862.4	220.2	1993	34634.4	1039.6
1982	5294.7	192.9	1994	46759.4	1156.1
1983	5934.5	213.9	1995	58478.1	1320.8
1984	7171.0	274.1	1996	67884.6	1388.3
1985	8964.4	422.5	1997	74462.6	1423.7
1986	10202.2	429.1	1998	78345.2	1402.4
1987	11962.5	432.1	1999	82067.46	1657
1988	14928.3	552.7	2000	89442.2	2250.9
1989	16909.2	591.4	2001	95933.3	2436.1

资料来源：《中国统计年鉴》（1995、2000、2002）。

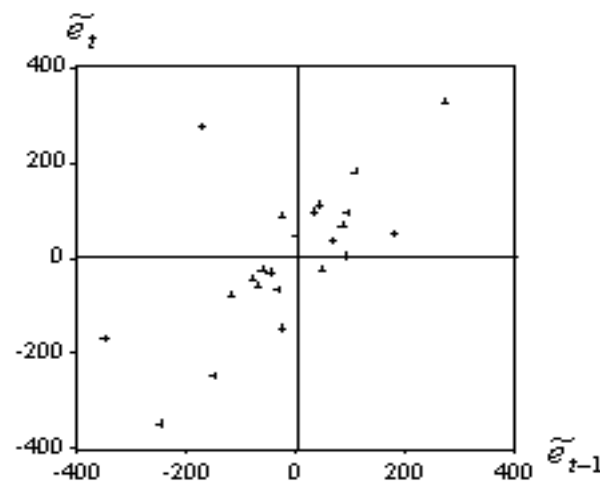
1. 通过OLS法建立如下中国商品进口方程:

$$\hat{M}_t = 152.91 + 0.02GDP_t$$

(2.32) (20.12)

$$R^2=0.948 \quad \bar{R}^2=0.946 \quad SE=154.9 \quad DW=0.628$$

2. 进行序列相关性检验。



上例估计了中国商品进口M与GDP的关系，并发现具有强烈的一阶自相关性。

然而，由于仅用GDP来解释商品进口的变化，明显地遗漏了诸如商品进口价格、汇率等其他影响因素。因此，序列相关性的主要原因可能就是建模时遗漏了重要的相关变量造成的。

下面进行RESET检验。

用原回归模型估计出商品进口序列

$$\hat{M}_t = 152.91 + 0.020GDP_t$$

$$R^2=0.9484$$

在原回归模型中加入 \hat{M}_t^2 、 \hat{M}_t^3 后重新进行估计，得：

$$\tilde{M}_t = -3.860 + 0.072GDP - 0.0028\hat{M}_t^2 + 8.59E-07\hat{M}_t^3$$

(-0.085) (8.274) (-6.457) (6.692)

$$R^2=0.9842$$

$$F = \frac{(R_U^2 - R_R^2)/q}{(1 - R_U^2)/(n - (k + q + 1))} = \frac{(0.984 - 0.948)/2}{(1 - 0.984)/(24 - 4)} = 22.5$$

在 $\alpha=5\%$ 下，查得临界值 $F_{0.05}(2, 20)=3.49$

判断：拒绝原模型与引入新变量的模型可决系数无显著差异的假设，表明原模型确实存在遗漏相关变量的设定偏误。

(3) 线性模型与双对数线性模型的选择

无法通过判定系数的大小来辅助决策，因为在两类模型中被解释变量是不同的。

为了在两类模型中比较，可用Box-Cox变换：

第一步，计算Y的样本几何均值。

$$\tilde{Y} = (Y_1 Y_2 \cdots Y_n)^{1/n} = \exp\left(\frac{1}{n} \sum \ln Y_i\right)$$

第二步，用得到的样本几何均值去除原被解释变量Y，得到被解释变量的新序列Y*。

$$Y_i^* = Y_i / \tilde{Y}$$

第三步，用 Y^* 替代 Y ，分别估计双对数线性模型与线性模型。并通过比较它们的残差平方和是否有显著差异来进行判断。

Zarembka（1968）提出的检验统计量为：

$$\frac{1}{2}n \ln\left(\frac{RSS_2}{RSS_1}\right)$$

其中， RSS_1 与 RSS_2 分别为对应的较小的残差平方和与较大的残差平方和， n 为样本容量。

可以证明：该统计量在两个回归的残差平方和无差异的假设下服从自由度为1 的 χ^2 分布。

因此，拒绝原假设时，就应选择 RSS_1 的模型。

例5.3.2 在 § 4.3 中国商品进口的例中，
采用线性模型： **$R^2=0.948$** ；
采用双对数线性模型： **$R^2=0.973$** ，
但不能就此简单地判断双对数线性模型优于线性模型。下面进行Box-Cox变换。

计算原商品进口样本的几何平均值为：

$$\tilde{M} = \exp\left(\frac{1}{n} \sum \ln(M_t)\right) = 583.12$$

计算出新的商品进口序列：

$$M_t^* = M_t / \tilde{M}$$

以 M_t^* 替代 M_t ，分别进行双对数线性模型与线性模型的回归，得：

$$\ln(\hat{M}_t^*) = -1.3565 + 0.7836 \ln GDP_t \quad RSS_1 = 0.5044$$

$$\hat{M}_t^* = 0.2622 + 0.000035 GDP_t \quad RSS_2 = 1.5536$$

于是， $\frac{1}{n} \ln\left(\frac{RSS_2}{RSS_1}\right) = \frac{1}{n} \times 24 \times 1.1249$:

在 $\alpha=5\%$ 下，查得临界值 $\chi^2_{0.05}(1)=3.841$

判断：拒绝原假设，表明双对数线性模型确实“优于”线性模型。

例1 (空)在回归模型中由于自变量的_____

小二乘法估计引起的后果有:(1)估计量是有偏的;(2)方差变大。由于自变量的_____ ,进行最小二乘法估

计时,估计量是有偏的,且方差变大。

例 2、一个估计某行业 ECO 薪水的回归模型如下

$$\ln(\text{salary}) = \beta_0 + \beta_1 \ln(\text{sales}) + \beta_2 \ln(\text{mktval}) + \beta_3 \text{profmarg} \\ + \beta_4 \text{ceoten} + \beta_5 \text{comten} + \mu$$

其中，salary 为年薪 sales 为公司的销售收入，mktval 为公司的市值，profmarg 为利润占销售额的百分比，ceoten 为其就任当前公司 CEO 的年数，comten 为其在该公司的年数。一个有 177 个样本数据集的估计得到 $R^2=0.353$ 。若添加 ceoten^2 和 comten^2 后， $R^2=0.375$ 。问：此模型中是否有函数设定的偏误？

例3 18 假设真实模型是

$$Y_i = \beta_1 + \beta_2 X_i + \beta_3 X_i^2 + \beta_4 X_i^3 + u_i$$

但你估计了

$$Y_i = \alpha_1 + \alpha_2 X_i + v_i$$

如果你利用 Y 在 $X = -3, -2, -1, 0, 1, 2, 3$ 外的观测并估计了 “ α

例 4、假设真实模型是：

$$Y_i = \beta_1 X_i + u_i$$

但没有拟和这个过原点回归，却例行拟和了通常带有截距项的模型：

$$Y_i = \alpha_0 + \alpha_1 X_i + v_i。$$

评述这一设定误差的后果。

例 5、假设真实模型是：

$$Y_i = \alpha_0 + \alpha_1 X_i + v_i$$

但拟和的模型为：

$$Y_i = \beta_1 X_i + u_i。$$

评述这一设定误差的后果。

例 6、假设真实模型是：

$$Y_i = \beta_1 + \beta_2 X_{2i} + u_i \quad (1)$$

而我们增加了一个“无关”变量 X_3 到模型中去（“无关”指变量 X_3 的真实系数 β_3 为 0），并估计了

$$Y_i = \beta_1 + \beta_2 X_{2i} + \beta_3 X_{3i} + v_i \quad (2)$$

- a. 模型（2）的可决系数和修正的可决系数会不会比模型（1）的大？
- b. 从模型（2）得到的 β_1, β_2 的估计值是无偏的吗？
- c. 无关变量 X_3 的引入对 $\hat{\beta}_1$ 和 $\hat{\beta}_2$ 的方差有影响吗？

例、设家庭每月消费支出为被解释变量 Y ，每月可支配收入为解释变量 X 。现有某社区 10 个家庭每月消费与收入的有关数据如下（单位：元）：

$$\bar{X} = 6800, \sum (X_i - \bar{X})^2 = 5.28 \times 10^7, \bar{Y} = 4604,$$

$$\sum (Y_i - \bar{Y})^2 = 19220640, \sum (X_i - \bar{X})(Y_i - \bar{Y}) = 30416000$$

- (1) 用最小二乘法估计样本回归方程 $\hat{Y}_i = \hat{\beta}_1 + \hat{\beta}_2 X_i$ 。
- (2) 求 $\hat{\beta}_2$ 的标准差。
- (3) 求可决系数 R^2 。
- (4) $\alpha = 0.05$ ，检验假设 $\beta_2 = 0$ 。
- (5) 假定家庭月可支配收入为 20000 元，预测家庭月消费支出，并给出置信度为 95% 的置信区间。