

金融计量学

(Financial Econometrics)

主讲教师: 任飞

商学院金融系

办公地点: 商学院1601

email: olrenfei@163.com

第二章 最小二乘法和线性回归模型

- 第一节 最小二乘法的基本属性
- 一、有关回归的基本介绍
- 二、参数的最小二乘估计
- 三、最小二乘估计量的性质和分布
- 第二节 一元线性回归模型的统计检验
- 一、拟合优度检验
- 二、假设检验
- 第三节 多变量线性回归模型的统计检验
- 一、多变量模型的简单介绍
- 二、拟合优度检验
- 三、假设检验

第四节 预测

- 一、预测的概念和类型
- 二、预测的评价标准
- 第五节 模型选择
- 一、"好"模型具有的特性
- 二、用于预测的模型的选择

第二章 最小二乘法和线性回归模型

- 第一节 最小二乘法的基本属性
- 一、有关回归的基本介绍
- 二、参数的最小二乘估计
- 三、最小二乘估计量的性质和分布

一、有关回归的基本介绍

金融、经济变量之间的关系,大体上可以分为两种:

- (1) 函数关系: $Y = f(X_1, X_2,, X_p)$, 其中Y的值是由 X_i (i = 1, 2 p)所唯一确定的。
- (2) 相关关系: $Y = f(X_1, X_2,, X_p)$, 其中Y的值不能由 X_i (i = 1, 2 p)精确地唯一确定。

有关回归的基本介绍

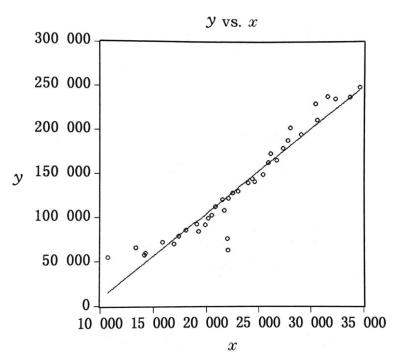


图2-1 货币供应量(y)和GDP(x)散点图

对于变量间的相关关系,可以根据 大量的统计资料,找出它们在数量 变化方面的规律(即"平均"的规 律),这种统计规律所揭示的关系 就是回归关系(regressive relationship) ,所表示的数学方程 就是回归方程(regression equation) 或回归模型(regression model)。

有关回归的基本介绍

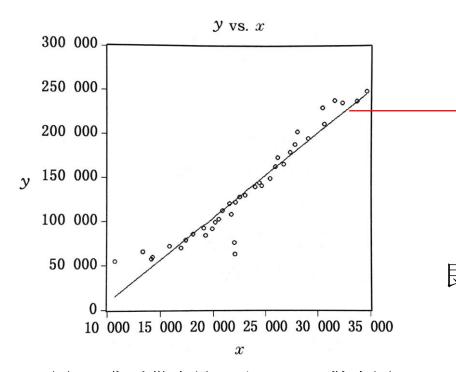


图2-1 货币供应量(y)和GDP(x)散点图

$$y = \alpha + \beta x$$
 (2.1)
与实际值存在误差

$$y = \alpha + \beta x + u \qquad (2.2)$$

即:
$$y_t = \alpha + \beta x_t + u_t$$
 (2.3)

其中t (= 1,2,3,,T) 表示观测数。

一、有关回归的基本介绍

$$y_t = \alpha + \beta x_t + u_t \quad (2.3)$$

- 其中 y_t 被称作 因变量(dependent variable) 被解释变量(explained variable) 结果变量(effect variable)
- x_t 被称作 自变量(independent variable) 解释变量(explanatory variable) 原因变量(causal variable)
- α 、 β 为参数(parameters)/回归系数(regression coefficients);
- u_t 通常被称为随机误差项(stochastic error term)/随机扰动项 (random disturbance term),简称误差项。

一、有关回归的基本介绍

$$y_t = \alpha + \beta x_t + u_t \quad (2.3)$$

为什么将 u_t 包含在模型中?

- (1) 有些变量是观测不到的或者是无法度量的,又或者影响因变量 y_t 的因素太多;
- (2) $\text{在}y_t$ 的度量过程中会发生偏误,这些偏误在模型中是表示不出来的;
- (3) 外界随机因素对 y_t 的影响也很难模型化,比如:恐怖事件、自然灾害、设备故障等。

(一) 方法介绍

- 本章所介绍的是普通最小二乘法(ordinary least squares, 简记OLS);
- 最小二乘法的基本原则是: 最优拟合直线应该使各点到直线的距离的和最小, 也可表述为距离的平方和最小。
- 假定根据这一原理得到的 α 、 β 估计值为 $\hat{\alpha}$ 、 $\hat{\beta}$,则直线可表示为 $\hat{y}_t = \hat{\alpha} + \hat{\beta}x_t$

(一) 方法介绍

- 直线上的 y_t 值,记为 \hat{y}_t ,称为拟合值(fitted value),实际值与拟合值的差,记为 \hat{u}_t ,称为残差(residual),可以看作是随机误差项 u_t 的估计值。
- 根据OLS的基本原则,使直线与各散点的距离的平方和最小,实际上是使残差平方和(residual sum of squares, 简记RSS) $\sum_{t=1}^{T} \hat{u}_t^2$ 最小,即最小化:

$$RSS = \sum_{t=1}^{T} (y_t - \hat{y}_t)^2 = \sum_{t=1}^{T} (y_t - \hat{\alpha} - \hat{\beta}x_t)^2$$
 (2.4)

(一) 方法介绍

• 根据最小化的一阶条件,将式**2.4**分别对 $\hat{\alpha}$ 、 $\hat{\beta}$ 求偏导,并令其为零,即可求得结果如下:

$$\hat{\beta} = \frac{\sum x_t y_t - T\bar{x}\bar{y}}{\sum x_t^2 - T\bar{x}^2}$$
 (2.5)

$$\hat{\alpha} = \bar{y} - \hat{\beta}\bar{x} \tag{2.6}$$

OLS系数估计量的推导过程(附录A2-1): 残差平方和:

$$RSS = \sum_{t=1}^{T} (y_t - \hat{y}_t)^2 = \sum_{t=1}^{T} (y_t - \hat{\alpha} - \hat{\beta}x_t)^2$$

对 $\hat{\alpha}$ 、 $\hat{\beta}$ 取值,使得残差平方和最小,因此分别求RSS对 $\hat{\alpha}$ 、 $\hat{\beta}$ 的偏导,并令其为零,

$$\frac{\partial RSS}{\partial \hat{\alpha}} = -2\sum_{t} (y_t - \hat{\alpha} - \hat{\beta}x_t) = 0$$
 (A2.1)

$$\frac{\partial RSS}{\partial \hat{\beta}} = -2\sum_{t} x_{t} (y_{t} - \hat{\alpha} - \hat{\beta}x_{t}) = 0 \quad (A2.2)$$

OLS系数估计量的推导过程:

将式(A2.1)展开可得:

$$\sum y_t - T\hat{\alpha} - \hat{\beta} \sum x_t = 0$$
 (A2.3)

而
$$\sum y_t = T\bar{y}$$
、 $\sum x_t = T\bar{x}$,所以
$$T\bar{y} - T\hat{\alpha} - T\hat{\beta}\bar{x} = 0$$
 (A2.4)

即:

$$\bar{y} - \hat{\alpha} - \hat{\beta}\bar{x} = 0 \tag{A2.5}$$

所以,

$$\hat{\alpha} = \bar{y} - \hat{\beta}\bar{x} \tag{A2.6}$$

OLS系数估计量的推导过程:

由式 (A2.2) 可得:

$$\sum_{t} x_t (y_t - \hat{\alpha} - \hat{\beta} x_t) = 0$$
 (A2.7)

将式(A2.6)代入式(A2.7)得:

$$\sum_{t} x_t \left(y_t - \bar{y} + \hat{\beta}\bar{x} - \hat{\beta}x_t \right) = 0 \tag{A2.8}$$

OLS系数估计量的推导过程:

展开可得:

$$\sum_{t} x_{t} y_{t} - \bar{y} \sum_{t} x_{t} + \hat{\beta} \bar{x} \sum_{t} x_{t} - \hat{\beta} \sum_{t} x_{t}^{2} = 0 \quad (A2.9)$$

$$\sum_{t} x_{t} y_{t} - T\bar{x}\bar{y} + \hat{\beta}T\bar{x}^{2} - \hat{\beta} \sum_{t} x_{t}^{2} = 0$$
 (A2.10)

因此,

$$\hat{\alpha} = \bar{y} - \hat{\beta}\bar{x} \qquad \qquad \hat{\beta} = \frac{\sum x_t y_t - T\bar{x}\bar{y}}{\sum x_t^2 - T\bar{x}^2} \quad (A2.11)$$

(二)一些基本概念

1. 总体(the population)和样本(the sample)

总体是指待研究变量的所有数据集合,可以是有限的,也可以是无限的;

样本是总体的一个子集。

(二)一些基本概念

- 2. 总体回归方程(the population regression function,简记PRF),样本回归方程(the sample regression function,简记SRF)。
- 总体回归方程 (PRF) 表示变量之间的真实关系,其中的 α 、 β 值是真实值,方程为:

$$y_t = \alpha + \beta x_t + u_t \tag{2.7}$$

样本回归方程(SRF)是根据所选样本估算的变量之间的关系函数,方程为:

$$\hat{y}_t = \hat{\alpha} + \hat{\beta}x_t \tag{2.8}$$

(二)一些基本概念

- 2. 总体回归方程(the population regression function,简记PRF),样本回归方程(the sample regression function,简记SRF)。
- 根据样本回归方程,总体回归方程(2.7)可以写为:

$$y_{t} = \hat{\alpha} + \hat{\beta}x_{t} + \hat{u}_{t}$$
 (2.9)
模型拟合值 \hat{y}_{t} 残差项

(二)一些基本概念

- 3. 线性关系
- 对线性的解释一: y = x的线性函数,比如, $y = \alpha + \beta x$ 。
- 对线性的解释二: y是参数的一个线性函数,它可以不是变量x的线性函数。比如, $y = \alpha + \beta x^2$ 就是一个线性回归模型, 但 $y = \alpha + \sqrt{\beta} x$ 则不是。
- 在本课程中,线性回归一词总是对指参数β为线性的一种 回归(即参数只以一次方出现),对解释变量x则可以是 或不是线性的。

(二)一些基本概念

3. 线性关系

有些模型看起来不是线性回归,但经过一些基本代数变换可以转换成线性回归模型。例如,

(1)
$$y_t = Ax_t^{\beta} e^{u_t} \longrightarrow ln(y_t) = ln(A) + \beta ln(x_t) + u_t$$

$$\Rightarrow Y_t = ln(y_t), \quad \alpha = ln(A), \quad X_t = ln(x_t)$$

$$Y_t = \alpha + \beta X_t + u_t$$

(2)
$$y_t = \alpha + \frac{\beta}{z_t} + u_t$$

$$\Rightarrow x_t = \frac{1}{z_t},$$

$$y_t = \alpha + \beta x_t + u_t$$

(二)一些基本概念

4. 估计量 (estimator) 和估计值 (estimate)

估计量: 计算系数的方程

估计值: 估计出来的系数的数值

- (一)经典线性回归模型(CLRM)的基本假设
 - (1) $E(u_t) = 0$, 即误差项具有零均值;
 - (2) $var(u_t) = \sigma^2 < \infty$,即误差项具有常数方差,且对于所有x值是有限的;
 - (3) $cov(u_i, u_i) = 0$,即误差项之间在统计意义上相互独立;
 - (4) $cov(u_t, x_t) = 0$,即误差项与变量x无关;
 - (5) $u_t \sim N(0, \sigma^2)$, 即误差项服从正态分布。

随机误差项满足假设条件(1)-(5)的线性回归模型即为CLRM(the classic linear regression model)。

(二)最小二乘估计量的性质

如果满足假设(1)一(4),由最小二乘法得到的 $\hat{\alpha}$ 、 $\hat{\beta}$ 具有一些特性,它们是最优线性无偏估计量(Best Linear Unbiased Estimators,简记BLUE)

线性 (linear): $\hat{\alpha}$ 、 $\hat{\beta}$ 与随机变量y之间是线性函数关系;

无偏 (unbiased): 平均而言,实际得到的 $\hat{\alpha}$ 、 $\hat{\beta}$ 值与其真实值是一致的;

最优(best): 在所有线性无偏估计量里,OLS估计量 $\hat{\beta}$ 具有最小方差。

(二)最小二乘估计量的性质

BLUE的另一种表述方法: OLS估计量具有一致性、无偏性和有效性。

1. 一致性(consistency): 随着样本容量的增大,该估计量会逐渐接近参数的真实值。

假定 $X \sim N(u_x, \sigma^2)$,从该正态总体中抽取一容量为n的随机样本。考虑X的两个估计量:

$$ar{X} = \sum rac{X_i}{n}$$
 $X^* = \sum rac{X_i}{n+1}$ $E(ar{X}) = u_x$ $E(X^*) = rac{n}{n+1}u_x \neq u_x$

 X^* 是X的有偏估计量,但随着样本容量(n)的增加, $E(X^*)$ 将近似于真实的 u_x ,即 X^* 具有一致性。

- (二)最小二乘估计量的性质
- 2. 无偏性(unbiasedness): 估计量的均值等于总体回归参数真值。

$$\text{iff: } \hat{\beta} = \frac{\sum x_t y_t - T\bar{x}\bar{y}}{\sum x_t^2 - T\bar{x}^2} = \frac{\sum (x_t - \bar{x})(y_t - \bar{y})}{\sum (x_t - \bar{x})^2} = \frac{\sum (x_t - \bar{x})y_t}{\sum (x_t - \bar{x})^2} = \frac{\sum (x_t - \bar{x})(\alpha + \beta x_t + u_t)}{\sum (x_t - \bar{x})^2}$$

可以证明
$$\frac{\sum (x_t - \bar{x})}{\sum (x_t - \bar{x})^2} = 0$$
, $\frac{\sum (x_t - \bar{x})x_t}{\sum (x_t - \bar{x})^2} = 1$

故
$$\hat{\beta} = \beta + \frac{\sum (x_t - \bar{x})u_t}{\sum (x_t - \bar{x})^2}$$
,

$$E(\hat{\beta}) = E\left(\beta + \frac{\sum (x_t - \bar{x})u_t}{\sum (x_t - \bar{x})^2}\right) = \beta + \frac{\sum (x_t - \bar{x})E(u_t)}{\sum (x_t - \bar{x})^2} = \beta$$

同理可得,
$$E(\hat{\alpha}) = E(\bar{y} - \hat{\beta}\bar{x}) = E(\alpha + (\beta - \hat{\beta})\bar{x} + \bar{u}) = \alpha$$

- (二) 最小二乘估计量的性质
- 3. 有效性(efficiency): 在所有线性无偏估计量里,OLS估计量 $\hat{\beta}$ 具有最小方差。

$$var(\hat{\beta}) = var(\beta + \frac{\sum (x_t - \bar{x})u_t}{\sum (x_t - \bar{x})^2}) = \sum \frac{(x_t - \bar{x})^2 var(u_t)}{\sum (x_t - \bar{x})^2} = \frac{\sigma^2}{\sum (x_t - \bar{x})^2}$$

$$\exists \exists k_t = \frac{x_t - \bar{x}}{\sum (x_t - \bar{x})^2}$$

$$var(\hat{\alpha}) = var(\bar{y} - \hat{\beta}\bar{x}) = var(\frac{1}{n}\sum y_{t} - \sum k_{t}y_{t}\bar{x}) = var(\sum y_{t}(\frac{1}{n} - \bar{x}k_{t})) = \sigma^{2}\sum(\frac{1}{n} - k_{t}\bar{x})^{2} = \sigma^{2}\sum(\frac{1}{n^{2}} - 2\frac{1}{n}k_{t}\bar{x} + k_{t}^{2}\bar{x}^{2}) = \sigma^{2}(\frac{1}{n} - 2\frac{1}{n}\bar{x}\sum k_{t} + \bar{x}^{2}\sum k_{t}^{2}) = \sigma^{2}(\frac{1}{n} + \bar{x}^{2}\frac{1}{\sum(x_{t} - \bar{x})^{2}}) = \sigma^{2}(\frac{\sum(x_{t} - \bar{x})^{2} + n\bar{x}^{2}}{n\sum(x_{t} - \bar{x})^{2}}) = \sigma^{2}(\frac{\sum x_{t}^{2}}{n\sum(x_{t} - \bar{x})^{2}})$$

-

三、最小二乘估计量的性质和分布

- (二)最小二乘估计量的性质
- 3. 有效性 (efficiency): 在所有线性无偏估计量里,OLS估计量 $\hat{\beta}$ 具有最小方差。

假设 $\hat{\beta}$ *是其它估计方法得到的关于 β 的线性无偏估计:

$$\hat{\beta}^* = \sum c_t y_t$$

设 $c_t = k_t + d_t$, $d_t > 0$ 为不全为零的常数,则很容易证明 $var(\hat{\beta} *) \geq var(\hat{\beta})$

同理,可以证明α的最小二乘估计â具有最小方差。

- (三) OLS估计量的方差、标准差和其概率分布
- 1. OLS估计量的标准差: 衡量估计量精确程度

$$SE(\hat{\alpha}) = s \sqrt{\frac{\sum x_t^2}{T \sum (x_t - \bar{x})^2}} = s \sqrt{\frac{\sum x_t^2}{T[(\sum x_t^2) - T\bar{x}^2]}}$$
 (2.21)

$$SE(\hat{\beta}) = s \sqrt{\frac{1}{T\sum(x_t - \bar{x})^2}} = s \sqrt{\frac{1}{\sum x_t^2 - T\bar{x}^2}}$$
 (2.22)

其中, s是残差的估计标准差。

估计量的标准差由解释变量x的真实观察值、样本容量T、 残差估计标准差s决定。

s是残差的估计标准差,通常被称为回归标准差,若用 σ^2 表示干扰项的真实方差

随机变量 u_t 的方差

$$var(u_t) = E[(u_t) - E(u_t)]^2$$
 (2.23)

根据 $E(u_t) = 0$,所以,

$$var(u_t) = E[u_t^2]$$
 (2.24)

$$E[u_t^2] = \frac{1}{T} \sum u_t^2 \tag{2.25}$$

即:

$$s^2 = \frac{1}{T} \sum u_t^2 \tag{2.26}$$

事实上,我们不可能得到 $\sum u_t^2$ 的值,因此只能以样本值 $\sum \hat{u}_t^2$ 来代替,即:

$$s^2 = \frac{1}{T} \sum \hat{u}_t^2 \tag{2.27}$$

但是这一估计量是有偏的, σ^2 的无偏估计量应该是:

$$s^2 = \frac{\sum \hat{u}_t^2}{T - 2} \tag{2.28}$$

$$s = \sqrt{\frac{\sum \hat{u}_t^2}{T - 2}} \tag{2.29}$$

- (三) OLS估计量的方差、标准差和其概率分布
 - 1. 0LS估计量的标准差具有如下性质:
 - (1) 样本容量T越大,参数估计值的标准差越小;
 - (2) $SE(\hat{\alpha})$ 和 $SE(\hat{\beta})$ 都取决于 s^2 ;
 - (3) 参数估计值的方差与 $\sum (x_t \bar{x})^2$ 成反比;
 - (4) $\sum x_t^2$ 项只影响截距的标准差,不影响斜率的标准差。

(三) OLS估计量的方差、标准差和其概率分布

2. OLS估计量的概率分布:

给定假定条件(5),即 $u_t \sim N(0, \sigma^2)$,则 y_t 也服从正态分布,系数估计量也是服从正态分布的。

$$\hat{\alpha} \sim N(\alpha, var(\alpha))$$
 $\hat{\beta} \sim N(\beta, var(\beta))$

 $\hat{\alpha}$ 、 $\hat{\beta}$ 的标准正态分布为:

$$\frac{\widehat{\alpha} - \alpha}{\sqrt{var(\alpha)}} \sim N(0,1) \qquad \frac{\widehat{\beta} - \beta}{\sqrt{var(\beta)}} \sim N(0,1)$$

(三) OLS估计量的方差、标准差和其概率分布

2. OLS估计量的概率分布:

用样本的标准差去替代总体标准差会产生不确定性, $\frac{\hat{\alpha} - \alpha}{SE(\hat{\alpha})}$ 、 $\frac{\hat{\beta} - \beta}{SE(\hat{\beta})}$ 将服从自由度为T-2的t分布:

$$\frac{\hat{\alpha} - \alpha}{SE(\hat{\alpha})} \sim t_{T-2} \tag{2.32}$$

$$\frac{\hat{\beta} - \beta}{SE(\hat{\beta})} \sim t_{T-2} \tag{2.33}$$

第二章 最小二乘法和线性回归模型

第二节 一元线性回归模型的统计检验

- 一、拟合优度检验
- 二、假设检验

一、拟合优度(goodness of fit statistics)检验

拟合优度可用R²表示:

$$R^2 = \frac{ESS}{TSS} = 1 - \frac{RSS}{TSS} \in [0,1],$$

 R^2 越大,说明回归线拟合程度越好。

 $TSS = \sum (y_t - \bar{y})^2$: 总平方和(the total sum of squares),反映了 y相对于其均值的波动性,这一平方和可以分成两部分:

$$\sum (y_t - \bar{y})^2 = \sum (\hat{y}_t - \bar{y})^2 + \sum \hat{u}_t^2 \qquad (2.36)$$

 $ESS = \sum (\hat{y}_t - \bar{y})^2$: 回归平方和(the explained sum of squares),即被模型解释的部分;

 $RSS = \sum \hat{u}_t^2$: 残差平方和(the residual sum of squares),即不能被模型解释的部分。



一、拟合优度(goodness of fit statistics)检验

 R^2 作为拟合优度的一个衡量标准存在的问题:

增加了一个解释变量以后, R²只会增大而不会减小。

解决方法: 使用调整过的 R^2 代替未调整过的 R^2 。

$$\overline{R^2} = 1 - \left[\frac{T-1}{T-K}(1-R^2)\right] \tag{2.40}$$

其中, T为样本容量, K为自变量个数。

 $\overline{R^2}$ 可以用来决定是否应该将某个解释变量包含在模型中。如果 $\overline{R^2}$ 增大,则包含该变量;反之则不包含。

假设检验的目的:确定根据样本得到的系数估计值是否与理论值相符或处在理论区间内。

基本思想: 概率性质的反证法。

H₀: 零假设 (null hypothesis) 或原假设

H₁: 备择假设 (alternative hypothesis)

为了检验原假设 H_0 是否正确,先假定这个假设是正确的。如果由此推出一个不合理的结果,则表明"假设 H_0 正确"是错误的,即原假设错误,因此要拒绝原假设。

概率性质的反证法的根据是**小概率事件原理**。该原理认为"小概率事件在一次实验中几乎是不可能发生的"。在原假设 H_0 下构造一个事件(即检验统计量),这个事件在"原假设 H_0 是正确的"的条件下是一个小概率事件,如果该事件发生了,说明"原假设 H_0 是正确的"是错误的,因为不应该出现的小概率事件出现了,应该拒绝原假设 H_0 。

例. 对于某一回归结果 β 的真实值是否等于0.5,提出如下双侧检验:

 $H_0: \beta = 0.5$

 $H_1: \beta \neq 0.5$

假设检验方法:

- (1) 置信区间检验法(confidence interval approach)
- (2) 显著性检验法(test of significance approach) t检验:对单个变量系数的显著性检验; F检验,对多个变量系数的联合显著性检验。

(一) *t*检验

t检验是显著性检验法(test of significance approach)中最常用的一种方法,是对单个变量系数的显著性检验。

基本原理:如果参数的假设值与估计值差别很大,就会导致小概率事件的发生,从而导致我们拒绝参数的假设值。

(一) *t*检验

t检验的主要步骤:

- (1) 用OLS方法回归方程 $y_t = \alpha + \beta x_t + u_t$,得到 β 的估计值 $\hat{\beta}$ 及其标准差SE($\hat{\beta}$)。
- (2) H_0 : $\beta = \beta^*$, H_1 : $\beta \neq \beta^*$ (双侧检验),则建立的统计量 $t_{sta} = \frac{\hat{\beta} \beta^*}{SE(\hat{\beta})}$ 服从自由度为T 2的t分布。
- (3)选择一个显著性水平(通常是5%),在t分布中确定拒绝区域和非拒绝区域,如图2-2。
- (4) 当检验统计值 t_{sta} 的绝对值大于临界值时,它就落在拒绝区域,因此拒绝 H_0 ,而接受 H_1 。反之则相反。

(一) *t*检验

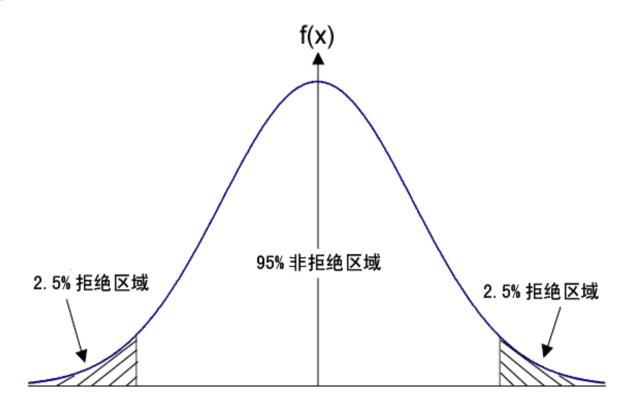


图2-2 双侧检验拒绝区域和非拒绝区域分布

(二) 置信区间法

置信区间检验的主要步骤(所建立的零假设同t检验):

- (1) 用OLS方法回归方程 $y_t = \alpha + \beta x_t + u_t$,得到 β 的估计值 $\hat{\beta}$ 及其标准差SE($\hat{\beta}$)。
- (2) 选择一个显著性水平(通常是5%),查t分布表,获得自由度为T-2的临界值 t_{crit} 。
 - (3) 所建立的置信区间为:

$$(\hat{\beta} - t_{crit}SE(\hat{\beta}), \hat{\beta} + t_{crit}SE(\hat{\beta}))$$
 (2.41)

(4) 如果零假设值 β *落在置信区间外,拒绝 H_0 : $\beta = \beta^*$; 反之,则不能拒绝。

(三) t检验与置信区间检验的关系

在显著性检验法下,当 t_{crit} 的绝对值小于临界值时,即:

$$-t_{crit} \le \frac{\widehat{\beta} - \beta^*}{SE(\widehat{\beta})} \le +t_{crit}$$
 (2.42)

时,我们不能拒绝原假设。

对式(2.42)变形,我们可以得到:

$$\hat{\beta} - t_{crit} SE(\hat{\beta}) \le \beta^* \le \hat{\beta} + t_{crit} SE(\hat{\beta})$$
 (2.43)

可以看到,式(2.43)恰好是置信区间法的置信区间式(2.41), 因此,实际上t检验与置信区间法提供的结果是完全一样的。

(四)第一类错误和第二类错误

第一类错误:如果有一个零假设在5%的显著性水平下被拒绝了,则这个拒绝就有可能是不正确的;

第二类错误:得到95%的一个置信区间,落在这个区间的零假设都不能拒绝,当接受一个零假设的时候也可能犯错误,因为回归系数的真实值可能是该区间内的另外一个值。

在选择显著性水平时人们面临抉择:降低犯第一类错误的概率就会增加犯第二类错误的概率。

在计量经济学中,我们通常选择相当小的显著性水平和犯第一类错误的概率。

(五) *p*值

*p*值: 计量经济结果对应的精确的显著性水平,一个0.07的*p* 值说明有关系数在0.07水平统计显著。

p值度量的是犯第一类错误的概率,即拒绝正确的零假设的概率。p值越小,拒绝零假设时就越放心。

第二章 最小二乘法和线性回归模型

第三节 多变量线性回归模型的统计检验

- 一、多变量模型的简单介绍
- 二、拟合优度检验
- 三、假设检验

考察下面这个方程:

$$y_t = \beta_1 + \beta_2 x_{2t} + \beta_3 x_{3t} + \dots + \beta_k x_{kt} + \mu_t \quad t = 1, 2, 3, \dots, T \quad (2.44)$$

对y产生影响的解释变量共有k-1(x_{2t} , x_{3t} ,..., x_{kt})个,系数(β_1 , β_2 ,..., β_k)分别衡量了解释变量对因变量y的边际影响的程度。

方程(2.44)的矩阵形式为:

$$Y = X\beta + \mu \quad (2.46)$$

$$Y = \begin{pmatrix} y_1 \\ y_2 \\ \dots \\ y_T \end{pmatrix} \quad X = \begin{pmatrix} 1 & x_{21} & \dots & x_{k1} \\ 1 & x_{22} & \dots & x_{k2} \\ \dots & \dots & \dots & \dots \\ 1 & x_{2T} & \dots & x_{kT} \end{pmatrix} \qquad \beta = \begin{pmatrix} \beta_1 \\ \beta_2 \\ \dots \\ \beta_k \end{pmatrix} \qquad \mu = \begin{pmatrix} \mu_1 \\ \mu_2 \\ \dots \\ \mu_T \end{pmatrix}$$

多变量OLS回归系数估计量的推导过程 残差项向量:

$$\hat{\boldsymbol{\mu}} = \begin{bmatrix} \hat{\boldsymbol{\mu}}_1 \\ \hat{\boldsymbol{\mu}}_2 \\ \dots \\ \hat{\boldsymbol{\mu}}_T \end{bmatrix} \tag{2.47}$$

$$RSS = \hat{\mu}'\hat{\mu} = \begin{bmatrix} \hat{\mu}_1 & \hat{\mu}_2 & \dots & \hat{\mu}_T \end{bmatrix} \begin{bmatrix} \hat{\mu}_1 \\ \hat{\mu}_2 \\ \dots \\ \hat{\mu}_T \end{bmatrix} = \sum \hat{\mu}_t^2 \qquad (2.48)$$

多变量OLS回归系数估计量的推导过程

RSS也可以用待估计向量 $\hat{\beta}$ 表示:

$$RSS = \hat{\mu}'\hat{\mu} = (Y - X\hat{\beta})'(Y - X\hat{\beta}) = Y'Y - \hat{\beta}'X'Y - Y'X\hat{\beta} + \hat{\beta}'X'X\hat{\beta}$$
$$\hat{\beta}'X'Y = Y'X\hat{\beta}$$

所以,

$$RSS = Y'Y - 2\hat{\beta}'X'Y + \hat{\beta}'X'X\hat{\beta}'$$

求RSS对 $\hat{\beta}$ 的导数,并令其为零,可得:

$$\frac{\partial RSS}{\partial \hat{\beta}} = -2X'Y + 2X'X\hat{\beta} = 0$$

多变量OLS回归系数估计量的推导过程 所以,

$$X'Y = X'X\hat{\beta}$$

可得:

$$\hat{\beta} = \begin{bmatrix} \hat{\beta}_1 \\ \hat{\beta}_2 \\ \dots \\ \hat{\beta}_T \end{bmatrix} = (X'X)^{-1}X'Y \qquad (2.49)$$

同样我们可以得到多变量模型残差的样本方差:

$$s^2 = \frac{\hat{\mu}'\hat{\mu}}{T - k} \tag{2.50}$$

以及参数的协方差矩阵:

$$var(\hat{\beta}) = s^2 (X'X)^{-1}$$
 (2.51)

二、拟合优度检验

与双变量模型类似, R²定义如下:

$$R^2 = \frac{ESS}{TSS} \tag{2.53}$$

可以证明:

$$ESS = \beta_2 \sum y_t x_{2t} + \beta_3 \sum y_t x_{3t} + \dots + \beta_k \sum y_t x_{kt}$$
 (2.54)

因此,

$$R^{2} = \frac{\beta_{2} \sum y_{t} x_{2t} + \beta_{3} \sum y_{t} x_{3t} + \dots + \beta_{k} \sum y_{t} x_{kt}}{\sum y_{t}^{2}}$$
(2.55)

 R^2 的值在0~1之间, R^2 越接近于1,说明估计的回归直线拟合得越好。

(一) *t*检验

在多元回归模型中, t统计量:

$$t_1 = \frac{\hat{\beta}_1 - \beta_1^*}{\text{SE}(\hat{\beta}_1)}$$
$$t_2 = \frac{\hat{\beta}_2 - \beta_2^*}{\text{SE}(\hat{\beta}_2)}$$

$$t_k = \frac{\hat{\beta}_k - \beta_k^*}{\text{SE}(\hat{\beta}_k)}$$

均服从自由度为(n-k)的t分布。后续检验过程与双变量线性回归模型的检验过程一样。

(二) F检验

F检验是用来检验有关部分回归系数的联合检验。 考虑如下多元回归模型:

 $y = \beta_1 + \beta_2 x_2 + \beta_3 x_3 + \dots + \beta_k x_k + u$ (2.57) 该模型称为**无约束回归模型**(unrestricted regression),因 为关于回归系数没有任何限制。

(二) F检验

假设我们想检验其中q个回归系数是否同时为零,为此改写公式为:

$$y = \beta_1 + \beta_2 x_2 + \dots + \beta_{k-q} x_{k-q} + \beta_{k-q+1} x_{k-q+1} + \dots + \beta_k x_k + u$$
 (2.58)

建立假设 H_0 : $\beta_{k-q+1} = \beta_k = 0$

则修正的模型将变为**有约束回归模型**(restricted regression):

$$y = \beta_1 + \beta_2 x_2 + \dots + \beta_{k-q} x_{k-q} + u \tag{2.59}$$



(二) F检验

检验的统计量为:

$$\frac{(RSS_R - RSS_{UR})/q}{RSS_{UR}/(N-K)} \tag{2.60}$$

如果零假设为真,则上式中的统计量将服从分子自由度为q、分母自由度为(N-K)的F分布。

联合F检验所检验的是一组变量是否显著。

(二) F检验

F检验与 R^2 有密切的联系。

由
$$R^2 = 1 - \frac{RSS}{TSS}$$
,则有:

$$R_{UR}^2 = 1 - \frac{RSS_{UR}}{TSS_{UR}}, \quad R_R^2 = 1 - \frac{RSS_R}{TSS_R}$$
 (2.61)

两个统计量具有相同的因变量,因此 $TSS_{UR} = TSS_R$ 检验的统计量可以写成:

$$F_{q,N-k} = \frac{(R_{UR}^2 - R_R^2)/q}{(1 - R_{UR}^2)/(N - k)}$$
(2.62)



随着人工智能、大数据等新兴信息技术在金融领域的广泛应用,金融科技应运而生并取得快速发展。金融科技公司利用其在信息技术上的比较优势,对传统金融机构进行"赋能",缓解传统金融机构面临的信息不对称难题,使金融更好的服务实体经济。在数量上缓解企业的融资约束,在质量上提高部门之间的信贷配置效率,进而提高企业全要素生产率。

研究假说:

地区金融科技发展水平越高,企业全要素生产率就越高。



样本选取:2011年-2018年沪深A股上市公司,最终得到15761个公司-年度观测值

回归模型:

$$TFP_{i,t} = \beta_0 + \beta_1 Fintech_{m,t-1} + \beta_2 SIZE_t + \varepsilon_{i,t}$$

 $TFP_{i,t}$: 企业i在t年的全要素生产率。

 $Fintech_{m,t-1}$: 地区m在第t-1年的金融科技发展水平,用地区金融科技公司数量测度。

 $SIZE_t$: 第t年企业规模,年末总资产取自然对数。

注:此例对原始模型进行了简化,只保留了关键解释变量与部分控制变量。

Eview操作步骤:导入数据,首先,选择以Eviews打开数据文件ex2.1.xls,出现如下对话框。

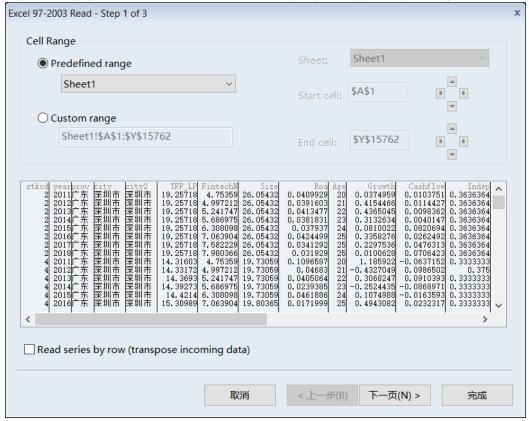


图2-3 数据导入窗口

Eview操作步骤: 导入数据

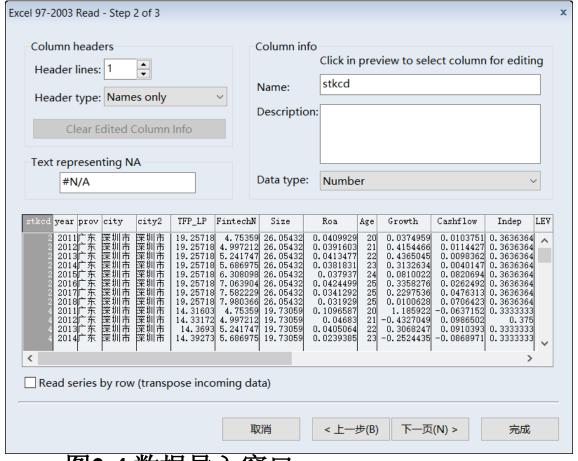


图2-4 数据导入窗口

Eview操作步骤:导入数据

Excel 97-2003 Read - Step	o 3 of 3	х
Import method Create new workfili	Basic structure Dated Panel Structure of the Data to be Imported Frequency: Annual	V
Import options Rename Series Frequency Conversion	Panel identifier series Cross section	
	Date series: year	
	Cancel <back next=""></back>	Finish

图2-5 数据导入窗口

数据导入后的工作文件

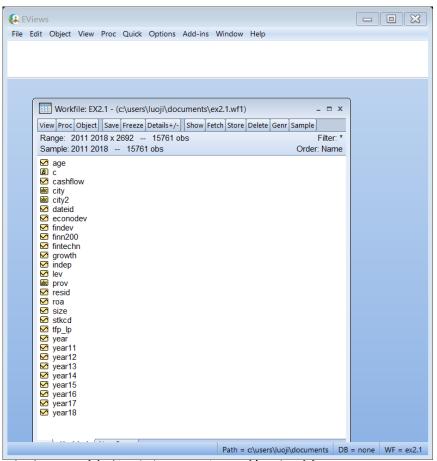


图2-6数据导入后工作文件

Eview操作步骤: Quick→Estimate Equation

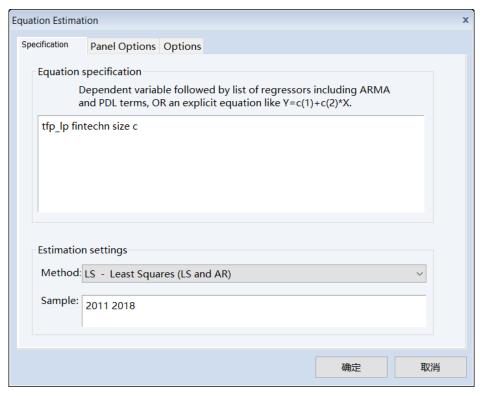


图2-7 方程设定窗口

Eview操作步骤: 在 "Panel Option"菜单栏,选择 "Cross-section"为 "Fixed"

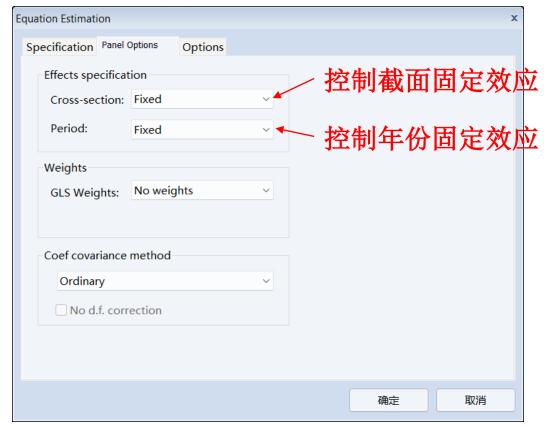


图2-8 方程设定窗口

回归结果:

$$TFP_{i,t} = 5.9129 + 0.0469Fintech_{m,t-1} + 0.4711SIZE_t$$
 $(0.1457) (0.0043)$ (0.0070) $t = (40.573) (10.849)$ (67.558) $p = (0.0000) (0.0000)$ (0.0000) (0.0000) $\bar{R}^2 = 0.9317$, $\bar{R}^2 = 0.9176$, $F_{2700.13060} = 66.0233$

自由度df = 15761 - (3 + 2691 + 7) = 13060 $F_{2700,13060}$ 对应的p值为0.00,因此拒绝所有回归系数全为0的零假设,回归方程整体显著。

总样本数 自变量数 截面虚拟变量数 日期虚拟变量数



回归结果:

$$TFP_{i,t} = 5.9129 + 0.0469Fintech_{m,t-1} + 0.4711SIZE_t$$
 $(0.1457) \ (0.0043) \ t = (40.573) \ (10.849) \ p = (0.0000) \ (0.0000)$ (0.0000) (0.0000) (0.0000) $\bar{R}^2 = 0.9317$, $\bar{R}^2 = 0.9176$, $F_{2700,13060} = 66.0233$

Fintech的估计系数在1%水平上显著为正,表面金融科技发展对企业全要素生产率存在显著正向影响,地区金融科技发展越好,当地企业的全要素生产率就越高,假说1得证。 SIZE的估计系数在1%水平上显著为正,表面企业规模越大,全要素生产率越高,因为大企业具有较高的规模效应。

第二章 最小二乘法和线性回归模型

第四节 预测

- 一、预测的概念和类型
- 二、预测的评价标准

一、预测的概念和类型

(一) 预测的概念

根据金融经济变量的过去和现在的发展规律,借助计量模型,对其未来的发展趋势和状况进行描述和分析,形成科学的假设和判断。

一、预测的概念和类型

(二)预测原理

条件期望(conditional expectations):在所有已知的t期的信息的条件下,Y的t+1期的条件期望值记作 $E(Y_{t+1}|I_t)$ 在t期对Y的下一期的所有预测值中,Y的条件期望 $E(Y_{t+1}|I_t)$ 是最优的(即具有最小方差)。

因此,在t期,对因变量Y的下一期(即t+1期)的预测值:

$$f_{t,1} = E(Y_{t+1} \mid I_t) \tag{2.65}$$

- (三)预测的类型
 - 1. 无条件预测和有条件预测
 - **无条件预测**:在预测模型中所有解释变量的值都是已知的这一条件下所进行的预测。

$$Y_t = \alpha + \beta_1 X_{1,t-2} + \beta_2 X_{2,t-3} + u_t$$

有条件预测:预测模型中某些解释变量的值时未知的。

$$Y_t = \alpha + \beta_1 X_{1,t} + \beta_2 X_{2,t} + u_t$$

(三)预测的类型

2. 样本内预测和样本外预测

样本内预测:用全部观测值来估计模型,然后用估计得到的模型对其中的一部分观测值进行预测。

例: 以某上市公司股票1997年1月-**2004年12月**的月收益率数据估计模型,利用模型对2004年6月-2004年12月的数据进行预测。

》 **样本外预测**:将全部观测值分为两部分,一部分用来估计模型,然后用估计得到的模型对另一部分数据进行预测。

例:以某上市公司股票1997年1月-**2004年5月**的月收益率数据估计模型,利用模型对2004年6月-2004年12月的数据进行预测。

- (三)预测的类型
 - 3. 事前预测和事后模拟
 - **事前预测**:指在不知道因变量真实值的情况下对其的预测。
 - **事后模拟**:指已经获得要预测的值的实际值,进行预测是为了评价模型的好坏。

- (三)预测的类型
 - 4. 一步向前预测和多步向前预测
 - **一步向前预测**:指仅对下一期的变量值进行预测。
 - 多步向前预测: 指对下一期值、更下一期值进行预测。如在t期对t+1期、t+2期、…、t+r期的值进行预测。如果 $r \leq 5$,则为短期预测(short-term forecast);如果r > 5,则为长期预测(long-term forecast)。

在将模型的预测结果应用于实践前,可以通过**事后模拟**对**预测结果的精确性**进行判断。

有不同的统计量可以量化预测变量与它对应的数据序列的接近程度:

- (一) 平均预测误差平方和平均预测误差绝对值
- (二) Theil 不相等系数

(一) 平均预测误差平方和平均预测误差绝对值

平均预测误差平方和(mean squared error,MSE)定义为:

$$MSE = \frac{1}{T} \sum_{t=1}^{T} (y_t^s - y_t^{\alpha})^2$$
 (2.66)

其中, y_t^s 表示 y_t 的预测值, y_t^{α} 表示 y_t 的实际值,T表示时段数。 **平均预测误差绝对值**(mean absolute error,MAE)定义为:

$$MAE = \frac{1}{T} \sum_{t=1}^{T} |y_t^s - y_t^{\alpha}|$$
 (2.67)

MSE和MAE度量的是误差的绝对大小,误差越大,说明模型的预测效果越不理想。

(二) Theil 不相等系数

$$U = \frac{\sqrt{\frac{1}{T} \sum_{t=1}^{T} (y_t^s - y_t^{\alpha})^2}}{\sqrt{\frac{1}{T} \sum_{t=1}^{T} (y_t^s)^2} + \sqrt{\frac{1}{T} \sum_{t=1}^{T} (y_t^{\alpha})^2}}$$
(2.68)

需要注意,U的分子就是MSE的平方根,而分母使得U总在0和1之间。

如果U = 0,则对所有的t,都有 $y_t^s = y_t^{\alpha}$,且完全拟合;

如果U=1,则模型的预测能力最差。

因此,Theil不等系数度量的是误差的相对大小。

(二) Theil 不相等系数

MSE可以分解成如下形式:

$$\frac{1}{T} \sum_{t=1}^{T} (y_t^s - y_t^{\alpha})^2 = (\bar{y}^s - \bar{y}^{\alpha})^2 + (\sigma_s - \sigma_{\alpha})^2 + 2(1 - \rho)\sigma_s\sigma_{\alpha} \quad (2.69)$$

其中, \bar{y}^s 、 \bar{y}^α 、 σ_s 、 σ_α 分别是序列 y_t^s 和 y_t^α 的平均值和标准差, ρ 是它们的相关系数。

$$\rho = (1/\sigma_s \sigma_\alpha T) \sum (y_t^s - \bar{y}^s)(y_t^\alpha - \bar{y}^\alpha)$$

定义不相等比例如下:

U的偏误比例,表示系统误差。

$$U^{M} = \frac{(\bar{y}^{s} - \bar{y}^{\alpha})^{2}}{(1/T)\sum(y_{t}^{s} - y_{t}^{\alpha})^{2}}$$
(2.70)

U的方差比例,表示模型中的变量重复其实际变化程度的能力。

$$U^{S} = \frac{(\sigma_{S} - \sigma_{\alpha})^{2}}{(1/T)\sum(y_{t}^{S} - y_{t}^{\alpha})^{2}}$$
(2.71)

U的协方差比例,度量的是非系统误差。

$$U^{C} = \frac{2(1-\rho)\sigma_{S}\sigma_{\alpha}}{(1/T)\sum(y_{t}^{S} - y_{t}^{\alpha})^{2}}$$
(2.72)

$$U^M + U^S + U^C = 1$$

第二章 最小二乘法和线性回归模型

第五节 模型选择

- 一、"好"模型具有的特性
- 二、用于预测的模型的选择

一、"好"模型具有的特性

- 一般而言,一个"好"模型应具有以下特征:
- 节省性 (parsimony): 应在相对精确反映现实的基础上尽可能地简单。
- 可识别性 (identifiability): 估计的参数要有唯一确定值。
- 高拟合性(goodness of fit): 单方程中调整的 \bar{R}^2 ,应尽可能地高。
- **理论一致性**(theoretical consistency):估计值符号与经济理论相符。
- **预测能力**(predictive power): 对未来有较强的预测能力。

二、用于预测的模型的选择

在判断样本内模型的拟合度时, R²是一个有效的标准。

因为*R*²将随着模型解释变量的增多而不断增加,按照此标准将不会得到最佳的预测模型。

因此必须对由于解释变量增多而造成自由度丢失施加一个惩罚项,其中的一个标准就是:

$$\bar{R}^2 = 1 - \left[\frac{T-1}{T-K} (1-R^2) \right]$$

随着解释变量的个数K的增大, \bar{R}^2 将减小。

二、用于预测的模型的选择

对自由度丢失惩罚更为严格的标准:

Akaike的信息准则(Akaike information criterion,AIC):

$$AIC = \ln(\hat{\sigma}^2) + \frac{2K}{T} \tag{2.73}$$

Schwarz的信息准则(Schwarz information criterion,SC):

$$SC = \ln(\hat{\sigma}^2) + \frac{K}{T}(\ln T) \tag{2.74}$$

其中, $\hat{\sigma}^2$ 是方程随机误差项方差的估计值,K是解释变量的个数,T是样本容量。

无论从AIC标准还是SC标准,度量值越低,模型的预测会更好。

本章小结

本章内容在计量经济学中是最基础也是最重要的部分。在这一章中,我们首先介绍了最小二乘法及其估计量的性质和分布。在此基础上我们对一元线性回归模型的统计检验进行了详细讨论,接着将模型扩展,讨论了多元线性回归模型。在用模型进行预测时,主要有两种情况:即有条件预测和无条件预测。最后一小节我们简单介绍了模型的选择。

本章小结

- 最小二乘法的基本原理和计算方法。
- > 经典线性回归模型的基本假定。
- **BLUE**统计量的性质。
- 检验和置信区间检验的原理及步骤。
- > 多变量模型的回归系数的检验。
- > 预测的类型及评判预测的标准。
- > 好模型具有的特征。