

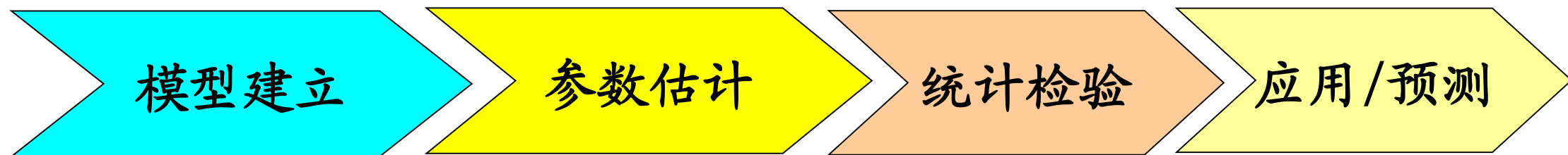
# 计量经济学基础

## 第三章 多元线性回归模型

# 主要内容

- 多元线性回归模型概述
- 多元线性回归模型的参数估计
- 多元线性回归模型的统计检验
- 常用计量模型结构参数的经济含义

# 本章内容的逻辑体系



总体回归模型/方程  
样本回归模型/方程  
模型的假设条件

参数的最小二乘估计  
随机误差项的方差估计量  
最小二乘估计量的特性

总离差平方和分解公式;  
多元样本可决系数;  
三个平方和的计算公式;  
修正的可决系数  
方程的显著性检验;  
解释变量的显著性检验;  
 $t$ 检验、区间估计

点预测: 内插预测、外推预测  
区间预测: 个值的区间预测、均值的区间预测

## § 3.1 模型的建立及假设条件

一、多元线性回归模型的基本概念

二、多元线性回归模型的基本假定

# 一、多元线性回归模型的基本概念

## ➤ 多元线性总体回归模型:

假设被解释变量 $Y$ 是多个解释变量 $X_2, X_3, \dots, X_k$ 和随机误差项 $u$ 的线性函数:

$$Y_i = \beta_1 + \beta_2 X_{2i} + \beta_3 X_{3i} + \dots + \beta_k X_{ki} + u_i$$

则称上式为多元总体线性回归模型。其中: $k$ 为待估参数个数

设  $(X_{2i}, X_{3i}, \dots, X_{ki}; Y_i)$ ,  $i = 1, 2, \dots, n$  是来自总体的观测值, 将其代入总体回归模型:

$$\begin{cases} Y_1 = \beta_1 \times 1 + \beta_2 X_{21} + \beta_3 X_{31} + \dots + \beta_k X_{k1} + u_1 \\ Y_2 = \beta_1 \times 1 + \beta_2 X_{22} + \beta_3 X_{32} + \dots + \beta_k X_{k2} + u_2 \\ \vdots \\ Y_n = \beta_1 \times 1 + \beta_2 X_{2n} + \beta_3 X_{3n} + \dots + \beta_k X_{kn} + u_n \end{cases}$$
  

$$\begin{bmatrix} Y_1 \\ Y_2 \\ \vdots \\ Y_n \end{bmatrix}_{n \times 1} = \begin{bmatrix} 1 & X_{21} & X_{31} & \dots & X_{k1} \\ 1 & X_{22} & X_{32} & \dots & X_{k2} \\ \vdots & \vdots & \vdots & & \vdots \\ 1 & X_{2n} & X_{3n} & \dots & X_{kn} \end{bmatrix}_{n \times k} \begin{bmatrix} \beta_1 \\ \beta_2 \\ \vdots \\ \beta_k \end{bmatrix}_{k \times 1} + \begin{bmatrix} u_1 \\ u_2 \\ \vdots \\ u_n \end{bmatrix}_{n \times 1}$$

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

则多元线性回归模型的矩阵表达式为:  $\mathbf{Y} = \mathbf{X} \boldsymbol{\beta} + \mathbf{U}$

描述被解释变量  $Y$  的条件均值与解释变量  $X_2, X_3, \dots, X_k$  之间线性关系的方程：

$$E(Y | X) = \beta_1 + \beta_2 X_2 + \beta_3 X_3 + \dots + \beta_k X_k$$

称为**多元总体回归方程**。其矩阵表达式为： $E(Y | X) = X\beta$

**注：**只要满足  $E(U | X) = 0$ ，则对总体回归模型  $Y = X\beta + U$  两边取条件数学期望即可得到总体回归方程：

$$Y = X\beta + U \Rightarrow E(Y | X) = E(X\beta + U | X) = X\beta + E(U | X) = X\beta$$

总体回归方程的含义是：给定各自变量  $X$  的取值时因变量  $Y$  的平均值

总体回归方程中的  $\beta_j$  称为**偏回归系数**，表示在其他解释变量保持不变的情况下， $X_j$  的单位变化所引起的  $Y$  的条件均值  $E(Y | X_j)$  的变化（与偏导数含义等同）。

若只有样本信息  $(X_{2i}, X_{3i}, \dots, X_{ki}; Y_i)$ ，则表述样本信息集中自变量与因变量间数量关系的表达式  $Y_i = \hat{\beta}_1 + \hat{\beta}_2 X_{2i} + \hat{\beta}_3 X_{3i} + \dots + \hat{\beta}_k X_{ki} + e_i$  称为**多元样本回归模型**。

其矩阵表述式为： $Y = X\hat{\beta} + e$

利用样本观测值  $(X_{2i}, X_{3i}, \dots, X_{ki}; Y_i)$  , 可对总体回归方程中的参数  $\beta_1, \beta_2, \dots, \beta_k$  进行估计。设其估计值为  $\hat{\beta}_1, \hat{\beta}_2, \dots, \hat{\beta}_k$  , 则由此得到的估计方程  $\hat{Y}_i = \hat{\beta}_1 + \hat{\beta}_2 X_{2i} + \hat{\beta}_3 X_{3i} + \dots + \hat{\beta}_k X_{ki}$ 。该方程即为**样本回归方程**或**经验回归方程**。其矩阵表述式为:  $\hat{Y} = X\hat{\beta}$

### 概 念

### 简单线性回归

总体回归方程:  $E(Y | X_i) = \beta_1 + \beta_2 X_i$

样本回归方程:  $\hat{Y}_i = \hat{\beta}_1 + \hat{\beta}_2 X_i$

总体回归模型:  $Y_i = \beta_1 + \beta_2 X_i + u_i$

样本回归模型:  $Y_i = \hat{\beta}_1 + \hat{\beta}_2 X_i + e_i$

### 多元线性回归

$E(Y | X) = X\beta$

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

$Y = X\beta + U$

$Y = X\hat{\beta} + e$



## 二、多元线性回归模型的基本假定

假定1: 零均值假定: 在给定 $X$ 的条件下,  $u_i$ 的条件期望为零, 即:

$$E(u_i | X_2, X_3, \dots, X_k) = 0 \Leftrightarrow E(U | X) = 0$$

假定2: 同方差假定: 在给定 $X$ 的条件下,  $u_i$ 的条件方差为某个常数, 即:

$$\text{Var}(u_i | X_2, X_3, \dots, X_k) = \sigma^2$$

假定3: 无自相关假定: 随机扰动项 $u_i$ 的逐次值互不相关, 即:

$$\text{Cov}(u_i, u_j | X_2, X_3, \dots, X_k) = 0 \quad i \neq j$$

定当随机扰动项满足同方差和无自相关假定时, 其协方差矩阵为:

$$\text{Var}(U | X) = E(UU' | X) = \begin{pmatrix} \sigma^2 & \dots & 0 \\ \vdots & \ddots & \vdots \\ 0 & \dots & \sigma^2 \end{pmatrix} = \sigma^2 I$$

反之, 只要有一个假定被违背, 则随机扰动项的协方差矩阵就不能简化为 $\sigma^2 I$

**假定4: 解释变量 $X$ 是非随机的**, 或者虽然 $X$  是随机的但与扰动项 $u$  不相关, 即:

$$Cov(X_{ij}, u_i | X_2, X_3, \dots, X_k) = 0 \quad j = 1, 2, \dots, k$$

**假定5: 无多重共线性假定:** 各解释变量之间不存在线性关系, 或各个解释变量观测值之间线性无关。

注: 无多重共线性假定是指系数矩阵 $X$  中的 $k$  个列向量  $(I \ X_2 \ X_3 \ \dots \ X_k)$  线性无关, 即系数矩阵 $X$ 为列满秩矩阵。

若 $Rank(X) = k \Rightarrow Rank(X'X)_{k \times k} = k \Rightarrow X'X$ 为满秩方阵 $\Rightarrow (X'X)^{-1}$ 存在。反之, 若无多重共线性假定被违背, 即 $Rank(X) < k$ , 则 $(X'X)^{-1}$ 不存在。

**最低样本容量要求:** 计量分析中, 样本容量不得少于待估参数个数, 即 $n \geq k$ , 否则 $(X'X)^{-1}$ 不存在。

**假定6: 正态性假定:** 随机扰动项服从均值为零, 方差为 $\sigma^2$  的正态分布

$$U|X \sim N(0, \sigma^2 I_n)$$

## § 3.2 最小二乘法

- 1、参数的最小二乘估计
- 2、最小二乘估计量的性质
- 3、随机误差项方差 $\sigma^2$ 的估计量

# 1、参数的最小二乘估计

- 随机抽取 $n$ 组观测值 $(X_{2i}, X_{3i}, \dots, X_{ki}; Y_i)$ ,  $i=1, 2, \dots, n$ , 若样本回归方程的参数估计值已经得到, 即:

$$\hat{Y}_i = \hat{\beta}_1 + \hat{\beta}_2 X_{2i} + \hat{\beta}_3 X_{3i} + \dots + \hat{\beta}_k X_{ki}$$

- 则样本观察值与样本回归值之间的残差平方和为:

$$\begin{aligned} Q(\hat{\beta}_1, \hat{\beta}_2, \dots, \hat{\beta}_k) &= \sum e_i^2 = \sum (Y_i - \hat{Y}_i)^2 \\ &= \sum (Y_i - \hat{\beta}_1 - \hat{\beta}_2 X_{2i} - \hat{\beta}_3 X_{3i} - \dots - \hat{\beta}_k X_{ki})^2 \end{aligned}$$

# 1、参数的最小二乘估计

- 最小二乘准则：使残差平方和最小。即求解下列无约束极值问题：

$$\begin{aligned} & \underset{\hat{\beta}_1, \hat{\beta}_2, \dots, \hat{\beta}_k}{\text{Min}} Q(\hat{\beta}_1, \hat{\beta}_2, \dots, \hat{\beta}_k) \\ &= \underset{\hat{\beta}_1, \hat{\beta}_2, \dots, \hat{\beta}_k}{\text{Min}} \sum \left( Y_i - \hat{\beta}_1 - \hat{\beta}_2 X_{2i} - \hat{\beta}_3 X_{3i} - \dots - \hat{\beta}_k X_{ki} \right)^2 \end{aligned}$$

- 由多元函数极值的必要条件知： $\hat{\beta}_1, \hat{\beta}_2, \dots, \hat{\beta}_k$  应满足：

$$\frac{\partial Q}{\partial \hat{\beta}_i} = 0, \quad i = 1, 2, \dots, k$$

由极值问题  $\underset{\hat{\beta}_1, \hat{\beta}_2, \dots, \hat{\beta}_k}{Min} Q = \sum (Y_i - \hat{\beta}_1 - \hat{\beta}_2 X_{2i} - \hat{\beta}_3 X_{3i} - \dots - \hat{\beta}_k X_{ki})^2$  的一阶条件即可得到

参数估计值应满足的正规方程组：

$$\begin{cases} \sum [Y_i - (\hat{\beta}_1 + \hat{\beta}_2 X_{2i} + \hat{\beta}_3 X_{3i} + \dots + \hat{\beta}_k X_{ki})] = 0 \Rightarrow \sum e_i = 0 \\ \sum X_{2i} [Y_i - (\hat{\beta}_1 + \hat{\beta}_2 X_{2i} + \hat{\beta}_3 X_{3i} + \dots + \hat{\beta}_k X_{ki})] = 0 \Rightarrow \sum X_{2i} e_i = 0 \\ \vdots \\ \sum X_{ki} [Y_i - (\hat{\beta}_1 + \hat{\beta}_2 X_{2i} + \hat{\beta}_3 X_{3i} + \dots + \hat{\beta}_k X_{ki})] = 0 \Rightarrow \sum X_{ki} e_i = 0 \end{cases}$$

上述 $k$ 个中间结论可用矩阵重新表述为：

$$\begin{bmatrix} 1 & 1 & \dots & 1 \\ X_{21} & X_{22} & \dots & X_{2n} \\ \vdots & \vdots & \vdots & \vdots \\ X_{k1} & X_{k2} & \dots & X_{kn} \end{bmatrix} \begin{bmatrix} e_1 \\ e_2 \\ \vdots \\ e_n \end{bmatrix} = \begin{bmatrix} 0 \\ 0 \\ \vdots \\ 0 \end{bmatrix}, \text{ 即 } X'e = 0$$

为利用条件  $X'e = 0$ , 在样本回归模型  $Y = X\hat{\beta} + e$  两边同时左乘  $X'$  可得:

$$X'Y = X'X\hat{\beta} + X'e \Rightarrow X'Y = X'X\hat{\beta}$$

若模型满足无多重共线性假定, 则  $X'X$  为满秩方阵, 其逆矩阵  $(X'X)^{-1}$  存在, 故在等式  $X'Y = X'X\hat{\beta}$  两边同时左乘  $(X'X)^{-1}$  可得:  $(X'X)^{-1} X'Y = (X'X)^{-1} X'X\hat{\beta} = \hat{\beta}$ , 故多元线性回归模型的最小二乘估计量为:  $\hat{\beta} = (X'X)^{-1} X'Y$

注: 在推导出多元线性回归模型  $OLS$  估计量  $\hat{\beta} = (X'X)^{-1} X'Y$  的过程中, 除解释变量为确定性变量和无多重共线性假定外并没有用到其它4个假定。即其它4个假定是否满足并不影响  $OLS$  估计量的表达式。

## 2、最小二乘估计量的性质

1、线性性：结构参数 $\beta$ 的最小二乘估计量 $\hat{\beta}=(X'X)^{-1}X'Y$ 总可以表示为被解释变量的线性函数（即 $\hat{\beta}$ 是 $n$ 个独立正态随机变量的某一个线性组合，因此 $\hat{\beta}$ 必然服从正态分布）

线性性是不证自明的。因为 $\hat{\beta}=(X'X)^{-1}X'Y \Leftrightarrow$

$$\begin{pmatrix} \hat{\beta}_1 \\ \hat{\beta}_2 \\ \vdots \\ \hat{\beta}_k \end{pmatrix} = \begin{pmatrix} a_{11} & a_{12} & \cdots & a_{1n} \\ a_{21} & a_{22} & \cdots & a_{2n} \\ \vdots & \vdots & \ddots & \vdots \\ a_{k1} & a_{k2} & \cdots & a_{kn} \end{pmatrix} \begin{pmatrix} Y_1 \\ Y_2 \\ \vdots \\ Y_n \end{pmatrix} = \begin{pmatrix} \sum_{i=1}^n a_{1i} Y_i \\ \sum_{i=1}^n a_{2i} Y_i \\ \vdots \\ \sum_{i=1}^n a_{ki} Y_i \end{pmatrix}$$

2、无偏性：最小二乘估计量 $\hat{\beta}$ 的数学期望等于参数的真实值 $\beta$ ，即 $E(\hat{\beta}) = \beta$

证明：  $\hat{\beta} = (X'X)^{-1}X'Y = (X'X)^{-1}X'(X\beta + U) = (X'X)^{-1}X'X\beta + (X'X)^{-1}X'U = \beta + (X'X)^{-1}X'U$   
 $\Rightarrow E\hat{\beta} = E[\beta + (X'X)^{-1}X'U] = \beta + (X'X)^{-1}X'E(U) = \beta$



**3、有效性（最小方差性）：**在结构参数 $\beta$ 的所有线性无偏估计量中，由最小二乘法所得到的 $OLS$ 估计量 $\hat{\beta}=(X'X)^{-1}X'Y$ 的方差是最小的。

***Gauss-Markov*定理：**在满足经典假定的情况下， $OLS$ 估计量 $\hat{\beta}$ 是 $\beta$ 的**最优线性无偏估计量**(*Best Linear Unbiased Estimator*, 简称***BLUE***估计量)

$$\begin{aligned} \text{OLS估计量}\hat{\beta}\text{的协方差矩阵: } Var(\hat{\beta}) &= E \left\{ \begin{pmatrix} \hat{\beta}_1 - \beta_1 \\ \hat{\beta}_2 - \beta_2 \\ \vdots \\ \hat{\beta}_k - \beta_k \end{pmatrix} \begin{pmatrix} \hat{\beta}_1 - \beta_1 & \hat{\beta}_2 - \beta_2 & \cdots & \hat{\beta}_k - \beta_k \end{pmatrix} \right\} \\ &= E \left[ (\hat{\beta} - \beta)(\hat{\beta} - \beta)' \right]. \end{aligned}$$

$$\text{因 } \hat{\beta} = (X'X)^{-1} X'Y = (X'X)^{-1} X'(X\beta + U) = \beta + (X'X)^{-1} X'U$$

$$\Rightarrow \hat{\beta} - \beta = (X'X)^{-1} X'U \Rightarrow (\hat{\beta} - \beta)' = U'X(X'X)^{-1}$$

$$\Rightarrow \text{Var}(\hat{\beta}) = E\left[(X'X)^{-1} X'UU'X(X'X)^{-1}\right] = (X'X)^{-1} X'E(UU')X(X'X)^{-1}$$

式中 $E(UU')$ 为随机扰动项的协方差矩阵, 若同方差和无自相关假定满足, 则 $E(UU') = \sigma^2 I_n$

$$\Rightarrow \text{Var}(\hat{\beta}) = (X'X)^{-1} X'\sigma^2 I X(X'X)^{-1} = \sigma^2 (X'X)^{-1} X'X(X'X)^{-1} = \sigma^2 (X'X)^{-1}$$

故由OLS估计量的性质, 可知 $\hat{\beta} \sim N\left(\beta, \sigma^2 (X'X)^{-1}\right)$

### 3、随机误差项方差 $\sigma^2$ 的估计量

尽管我们知道 $\hat{\beta} \sim N(\beta, \sigma^2 (X'X)^{-1})$ ，但 $\sigma^2$ 未知，因此这是一个方差未知的正态分布，仍然需要利用样本方差来估计 $\sigma^2$ 。

可以证明， $\sigma^2$ 的无偏估计量为 $\hat{\sigma}^2 = \frac{e'e}{n-k} = \frac{\sum_{i=1}^n e_i^2}{n-k}$

式中 $n$ 为样本容量， $k$ 为模型中待估参数个数， $n-k$ 为残差平方和RSS的自由度。

因为在计算RSS时要受到 $\begin{cases} \sum e_i = 0 \\ \sum X_{2i} e_i = 0 \\ \vdots \\ \sum X_{ki} e_i = 0 \end{cases}$   $k$ 个约束，其自由度相应减少 $k$ 个。

$\hat{\sigma} = \sqrt{\frac{\sum_{i=1}^n e_i^2}{n-k}}$  称为回归标准误，它衡量了样本点偏离样本回归线的平均偏离程度。

# 回归系数的方差及其估计值

由于  $\hat{\beta} \sim N(\beta, \sigma^2 (X'X)^{-1})$ , 令  $(C_{ij}) = (X'X)^{-1}$ , 则可得  $\hat{\beta}_j (j=1, 2, \dots, k)$  的方差为  $Var(\hat{\beta}_j) = \sigma^2 (X'X)^{-1}_{jj} = \sigma^2 C_{jj}$ 。其中  $C_{jj}$  为  $(X'X)^{-1}$  主对角线上的元素。即若  $\sigma^2$  已知:

$$\hat{\beta}_j \sim N(\beta_j, \sigma^2 C_{jj}) \Rightarrow \frac{\hat{\beta}_j - \beta_j}{\sigma \sqrt{C_{jj}}} \sim N(0, 1)$$

若  $\sigma^2$  未知, 可用其无偏估计量  $\hat{\sigma}^2$  代替, 由此可得  $\hat{\beta}_j$  的方差估计量为:

$$S_{\hat{\beta}_j}^2 = \hat{\sigma}^2 C_{jj} = \frac{\sum e_i^2}{n-k} C_{jj} (j=1, 2, \dots, k)$$

故当  $\sigma^2$  未知, 并用样本方差  $\hat{\sigma}^2$  代替总体方差  $\sigma^2$  后:

$$\hat{\beta}_j \sim N(\beta_j, S_{\hat{\beta}_j}^2) \Rightarrow \frac{\hat{\beta}_j - \beta_j}{S_{\hat{\beta}_j}} = \frac{\hat{\beta}_j - \beta_j}{\hat{\sigma} \sqrt{C_{jj}}} \sim t(n-k)$$

## § 3.4 多元线性回归模型的统计检验

- 一、参数的显著性检验 ( $t$ 检验)
- 二、参数的置信区间
- 三、拟合优度检验
- 四、方程的显著性检验( $F$ 检验)

# 一、参数的显著性检验 (t检验)

## ➤1、t统计量

从前面的分析知道, 参数 $\beta_j$  的  $OLS$ 估计量 $\hat{\beta}_j$  服从正态分布:  $\hat{\beta}_j \sim N(\beta_j, S_{\hat{\beta}_j}^2)$

其中 $S_{\hat{\beta}_j}^2 = \hat{\sigma}^2 C_{jj'}$ ,  $\hat{\sigma} = \sqrt{\frac{\sum_{i=1}^n e_i^2}{n-k}}$ ,  $C_{jj'}$ 为 $(X'X)^{-1}$  主对角线上的元素。

由此可构造如下  $t$  统计量:  $t_j = \frac{\hat{\beta}_j - \beta_j}{S_{\hat{\beta}_j}} \sim t(n-k)$

特别地, 若 $\beta_j$  的真实值为0, 则  $t$  统计量为简化为:  $t_j = \frac{\hat{\beta}_j}{S_{\hat{\beta}_j}} \sim t(n-k)$

## 参数的显著性检验的具体步骤:

1: 提出假设

$$H_0: \beta_j = 0 (j = 1, 2, \dots, k)$$

$$H_1: \beta_j \neq 0$$

2: 在 $H_0$ 成立的前提下构造检验统计量:

$$t = \frac{\hat{\beta}_j - \beta_j}{S_{\hat{\beta}_j}} = \frac{\hat{\beta}_j}{S_{\hat{\beta}_j}} \sim t(n-k)$$

3: 给定显著性水平 $\alpha$ , 查表获得临界值 $t_{\alpha/2}(n-k)$ , 确定拒绝域  $|t| > t_{\alpha/2}(n-k)$

4: 代入样本信息, 计算出检验统计量值 $t = \hat{\beta}_j / S_{\hat{\beta}_j}$ , 并与临界值比较:

(1)若 $|t| > t_{\alpha/2}(n-k)$ , 落入拒绝域, 拒绝 $H_0$ , 即 $\beta_j$ 显著异于零

(2)若 $|t| < t_{\alpha/2}(n-k)$ , 未落入拒绝域, 不能拒绝 $H_0$ , 即 $\beta_j$ 不显著

## 二、回归系数的置信区间

- 回归系数的置信区间主要用来考察：在一次抽样中所估计出的参数值离其真实值有多“近”。

■ 在变量的显著性检验中已经知道：

$$t_j = \frac{\hat{\beta}_j - \beta_j}{S_{\hat{\beta}_j}} \sim t(n-k)$$

容易推出：在 $(1-\alpha)$ 的置信水平下 $\beta_j$ 的置信区间是

$$\left( \hat{\beta}_j - t_{\alpha/2}(n-k) \times S_{\hat{\beta}_j}, \hat{\beta}_j + t_{\alpha/2}(n-k) \times S_{\hat{\beta}_j} \right)$$



### 三、拟合优度检验

- 1、可决系数与调整的可决系数

总离差平方和的分解:  $TSS = ESS + RSS$

$$TSS = \sum y_i^2 = \sum (Y_i - \bar{Y})^2 \quad - - \text{总离差平方和}$$

$$ESS = \sum \hat{y}_i^2 = \sum (\hat{Y}_i - \bar{Y})^2 \quad - - \text{回归平方和}$$

$$RSS = \sum e_i^2 = \sum (Y_i - \hat{Y}_i)^2 \quad - - \text{残差平方和}$$

总离差平方和的分解的离差形式为:

$$\sum y_i^2 = \sum \hat{y}_i^2 + \sum e_i^2$$

### 三、拟合优度检验

在多元线性回归模型中，总离差平方和的分解公式仍然成立，即：

$$\text{总离差平方和}(TSS) = \text{回归平方和}(ESS) + \text{残差平方和}(RSS)$$

衡量了因变量（围绕  
均值取值）的波动性

衡量了在因变量的波  
动性中，能够由自变  
量的波动（ $x_i^2$ ）进行  
解释的部分

衡量了在因变量的波  
动性中，不能由自变  
量的波动，只能由其  
它随机影响因素的波  
动进行解释的部分

由此看来，似乎仍然可以利用可决系数  $R^2 = \frac{ESS}{TSS} = 1 - \frac{RSS}{TSS}$  来度量拟合优度检验的结果

**问题：**在应用中发现，如果往模型中增加一个解释变量， $R^2$ 往往增大。这就给人一个错觉：要使得模型拟合得好，只要增加解释变量即可。

**注：**若新增解释变量与被解释变量的协方差为零（即新增解释变量与被解释变量无线性相关性），则新增解释变量不改变模型的 $R^2$ 值。否则 $R^2$ 值将随模型中解释变量个数的增加而严格递增。

为惩罚通过往模型中增加(无关)解释变量而提高拟合优度的行为，必须对 $R^2$ 进行修正。

在样本容量 $n$ 既定的情况下，增加解释变量(即 $k$ 增加)必定使得模型自由度减少，所以直观的修正思路是：将残差平方和与总离差平方和分别除以各自的自由度，以剔除变量

个数对拟合优度的影响：
$$\bar{R}^2 = 1 - \frac{RSS/(n-k)}{TSS/(n-1)}$$

在上式中，若样本既定，则 $TSS/(n-1)$ 值既定。往模型中增加1个解释变量通常会使 $RSS$ 的值下降，但同时 $RSS$ 的自由度也在下降，因此并不能保证 $RSS/(n-k)$ 的值一定下降，从而 $\bar{R}^2$ 的值上升。

由于 $\bar{R}^2$ 是通过 $R^2$ 进行自由度修正之后得到的，因此二者必然具有内在联系，可以证明：

$$\bar{R}^2 = 1 - (1 - R^2) \frac{n-1}{n-k}$$

证明：因 $R^2 = 1 - \frac{RSS}{TSS} \Rightarrow \frac{RSS}{TSS} = 1 - R^2$ 。将其代入 $\bar{R}^2$ 的定义式可得：

$$\bar{R}^2 = 1 - \frac{RSS/(n-k)}{TSS/(n-1)} = 1 - \frac{RSS}{TSS} \frac{n-1}{n-k} = 1 - (1 - R^2) \frac{n-1}{n-k}$$

**补充结论：**当往模型中新增解释变量对应的 $t$  统计量值大于1 时， $\bar{R}^2$  值会增加，否则减少。

## 方差分析表

变差来源	平方和	自由度	均方差
回归平方和	$ESS = \sum (\hat{Y}_i - \bar{Y})^2$	$k-1$	$ESS / (k-1)$
残差平方和	$RSS = \sum (Y_i - \hat{Y}_i)^2$	$n-k$	$RSS / (n-k)$
总离差平方和	$TSS = \sum (Y_i - \bar{Y})^2$	$n-1$	$TSS / (n-1)$

例：对于一个四元线性回归模型，  
已知可决系数  $R^2 = 0.85$

变差来源	平方和	自由度
回归平方和 (ESS)		
残差平方和 (RSS)	120	35
总变差 (TSS)		

由四元线性回归模型  $\Rightarrow k = 5$

由RSS的自由度为35  $\Rightarrow n - k = 35 \Rightarrow n = 40$

$\Rightarrow$  TSS的自由度为  $n - 1 = 39$

$\Rightarrow$  ESS的自由度为  $k - 1 = 4$

$$\text{由 } RSS=120, R^2=1-\frac{RSS}{TSS}=0.85 \Rightarrow 1-\frac{120}{TSS}=0.85 \Rightarrow TSS=800 \Rightarrow ESS=680$$

$$\text{由 } TSS=800, ESS=680, RSS=120 \Rightarrow \begin{cases} \bar{R}^2 = 1 - \frac{RSS / n - k}{TSS / n - 1} = 1 - \frac{120 / 35}{800 / 39} = 0.833 \\ F = \frac{ESS / k - 1}{RSS / n - k} = \frac{680 / 4}{120 / 35} = 49.583 \end{cases}$$

$$\text{也可由 } \bar{R}^2 = 1 - (1 - R^2) \frac{n - 1}{n - k} = 1 - (1 - 0.85) \frac{39}{35} = 0.833, \quad F = \frac{R^2}{1 - R^2} \frac{n - k}{k - 1} = \frac{0.85}{1 - 0.85} \frac{35}{4} = 49.583$$

### 三、回归方程的显著性检验(F检验)

对多元线性回归模型： $Y_i = \beta_1 + \beta_2 X_{2i} + \beta_3 X_{3i} + \cdots + \beta_k X_{ki} + u_i \ (i = 1, 2, \cdots, n)$

进行方程显著性检验的目的在于判断模型中**被解释变量与解释变量之间的线性关系在总体上是否显著成立**，等价于**检验模型中的所有自变量联合起来是否会对因变量产生显著的线性影响**。

从技术上看，方程的显著性检验就是检验模型中的结构参数 $\beta_j \ (j = 2, 3, \cdots, k)$ 是否全部同时等于0。若是，表明自变量联合起来不会对因变量产生显著的线性影响，也意味着自变量与因变量之间的线性关系在总体上不成立。

注意：方程显著性检验中的“线性”同样是指**参数线性**。例如：如果如下变量非线性模型 $Y = \beta_1 + \beta_2 X + \beta_3 X^2 + \beta_4 \ln X + u$ 通过了方程的显著性检验，则表明模型中的3个自变量 $X$ 、 $X^2$ 和 $\ln X$ 与因变量 $Y$ 之间存在显著的线性关系。

# 方程显著性检验(F检验)的具体步骤

$$Y_i = \beta_1 + \beta_2 X_{2i} + \beta_3 X_{3i} + \cdots + \beta_k X_{ki} + u_i$$

## 一、提出假设

$H_0$ :  $\beta_2 = \beta_3 = \cdots = \beta_k = 0$  (方程/模型的线性性在总体上不成立)

$H_1$ : 至少有一个  $\beta_j \neq 0$  ( $j = 2, 3, \cdots, k$ )

二、在  $H_0$  成立的前提下构造检验统计量  $F = \frac{ESS / (k-1)}{RSS / (n-k)} \sim F(k-1, n-k)$

三、给定显著性水平  $\alpha$ , 查表求得临界值  $F_\alpha(k-1, n-k)$ , 确定拒绝域  $F > F_\alpha(k-1, n-k)$

四、代入样本信息计算出检验统计量值  $F = \frac{ESS / (k-1)}{RSS / (n-k)}$ , 并与临界值  $F_\alpha(k-1, n-k)$  比较:

(1) 若  $F > F_\alpha(k-1, n-k)$ , 落入拒绝域, 拒绝  $H_0$ , 即方程/模型的线性性在总体上显著成立

(2) 若  $F < F_\alpha(k-1, n-k)$ , 未落入拒绝域, 接受  $H_0$ , 即方程/模型的线性性在总体上不成立



# 拟合优度检验与方程的显著性检验 (F检验)

由  $F$  统计量的定义式可以看出,  $F$  检验与拟合优度检验 (可决系数) 有密切联系, 二者都建立在因变量总离差平方和分解的基础上。

事实上,  $F$  统计量可由可决系数  $R^2$  推算: 
$$F = \frac{R^2 / (k-1)}{(1-R^2) / (n-k)} = \frac{R^2}{1-R^2} \frac{n-k}{k-1}$$

证明: 因  $R^2 = \frac{ESS}{TSS} = 1 - \frac{RSS}{TSS} \Rightarrow \begin{cases} ESS = R^2 \cdot TSS \\ RSS = (1-R^2) \cdot TSS \end{cases}$ , 将其代入  $F$  统计量的定义式可得:

$$F = \frac{ESS / (k-1)}{RSS / (n-k)} = \frac{R^2 \cdot TSS / (k-1)}{(1-R^2) \cdot TSS / (n-k)} = \frac{R^2 / (k-1)}{(1-R^2) / (n-k)} = \frac{R^2}{1-R^2} \frac{n-k}{k-1}$$

由上式可以看出,  $F$  是  $R^2$  的单调增函数, 即  $R^2$  越大,  $F$  统计量值越大。由此可知:

**方程的显著性检验和拟合优度检验具有内在一致性**

**内在逻辑:** 多元线性回归模型的  $R^2$  值越大, 对样本点的拟合程度越高, 说明模型的线性设定越合理, 也就越容易通过方程的显著性检验。

## 参数显著性检验与方程显著性检验之间的关系

1、检验目的不同：参数显著性检验用于检验结构参数的真实值是否显著异于零，方程显著性检验用于检验方程 / 模型的线性性在总体上是否显著成立。

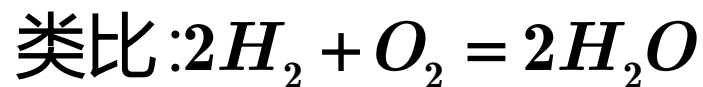
另一等价的表述则没有这么明显的区别：参数的显著性检验是在假定其它自变量保持不变时，检验某一自变量单独变动时是否会对因变量产生显著影响。而方程的显著性检验则是检验模型中所有自变量联合起来是否会对因变量产生显著的线性影响。（此时所有自变量均可自由变动）

2、对简单线性回归模型  $Y = \beta_1 + \beta_2 X + u$  而言，参数的显著性检验与方程的显著性检验等价。

若  $\beta_2$  通过了参数的显著性检验，即  $\beta_2$  显著异于零，则意味着自变量  $X$  与因变量  $Y$  之间存在显著的线性关系。反之，若简单线性回归模型  $Y = \beta_1 + \beta_2 X + u$  通过了方程的显著性检验，即  $X$  与  $Y$  之间存在显著的线性关系，则  $\beta_2$  必然显著异于零。

对多元线性回归模型而言：

- 1) 即使模型通过了方程的显著性检验，也并不意味着每一个结构参数都能通过参数的显著性检验
- 2) 只要有一个结构参数通过了参数的显著性检验，必然意味着模型能够通过方程的显著性检验
- 3) 即使模型中每一个结构参数都通不过参数的显著性检验，模型也依然有可能通过方程的显著性检验。



# 补充：四种常用计量模型结构参数的经济意义

1、线性模型：  $Y = \beta_1 + \beta_2 X + u$   $\Rightarrow$   $E(Y | X) = \bar{Y} = \beta_1 + \beta_2 X$

对应的总体  
回归方程为

$$\beta_2 = \frac{d\bar{Y}}{dX} \stackrel{\text{离散化}}{=} \frac{\Delta \bar{Y}}{\Delta X} \text{ (自变量的绝对变化所带来的因变量平均值的绝对变化)}$$

$$\Rightarrow \Delta \bar{Y} = \beta_2 \Delta X \stackrel{\text{令 } \Delta X=1}{\Rightarrow} \Delta \bar{Y} = \beta_2 \Rightarrow \beta_2 \text{ 的含义:}$$

自变量  $X$  每增加1个单位，因变量  $Y$  平均将增加  $\beta_2$  个单位（边际量）。

2、双对数线性模型： $Ln Y = \beta_1 + \beta_2 Ln X + u$   $\Rightarrow$   $E(Ln Y | X) = \beta_1 + \beta_2 Ln X$  对应的总体  
回归方程为

$$\text{由于 } E(Ln Y | X) = \frac{\sum Ln Y_i}{n} = Ln \left[ \left( \prod Y_i \right)^{1/n} \right] = Ln \bar{Y} \text{ (}\bar{Y}\text{为几何平均数),}$$

故双对数线性模型的总体回归方程可简写为： $Ln \bar{Y} = \beta_1 + \beta_2 Ln X$

$$\text{则 } \beta_2 = \frac{dLn \bar{Y}}{dLn X} = \frac{d\bar{Y} / \bar{Y}}{dX / X} \xrightarrow{\text{离散化}} \beta_2 = \frac{\Delta \bar{Y} / \bar{Y}}{\Delta X / X}$$

$\beta_2$ 即为自变量的相对变化所带来的因变量的平均相对变化

$$\text{由 } \beta_2 = \frac{\Delta \bar{Y} / \bar{Y}}{\Delta X / X} \Rightarrow \frac{\Delta \bar{Y}}{\bar{Y}} = \beta_2 \frac{\Delta X}{X} \xrightarrow{\text{令 } \Delta X / X = 1\%} \frac{\Delta \bar{Y}}{\bar{Y}} = \beta_2 \% \Rightarrow \beta_2 \text{ 的含义:}$$

自变量  $X$  每增长1%，因变量  $Y$  平均将增长  $\beta_2 \%$  (弹性)。

3、对数-线性模型（增长率模型）： $Ln Y = \beta_1 + \beta_2 X + u$  对应的总体  
回归方程为  $\Rightarrow Ln \bar{Y} = \beta_1 + \beta_2 X$

$$\beta_2 = \frac{dLn \bar{Y}}{dX} = \frac{d\bar{Y} / \bar{Y}}{dX} \xrightarrow{\text{离散化}} \beta_2 = \frac{\Delta \bar{Y} / \bar{Y}}{\Delta X}$$

$\beta_2$ 即为自变量的绝对变化所带来的因变量的平均相对变化

由上式可知： $\frac{\Delta \bar{Y}}{\bar{Y}} = \beta_2 \Delta X \xrightarrow{\text{令}\Delta X=1} \frac{\Delta \bar{Y}}{\bar{Y}} = \beta_2 = \beta_2 \times 100\% = 100\beta_2\%$ ，故 $\beta_2$ 的含义是：

自变量 $X$ 每增加1个单位，因变量 $Y$ 平均将增长 $100\beta_2\%$ (增长率)。

4、线性-对数： $Y = \beta_1 + \beta_2 \ln X + u$  对应的总体  
回归方程为  $\Rightarrow \bar{Y} = \beta_1 + \beta_2 \ln X$

$$\beta_2 = \frac{d\bar{Y}}{d\ln X} = \frac{d\bar{Y}}{dX/X} \xrightarrow{\text{离散化}} \beta_2 = \frac{\Delta \bar{Y}}{\Delta X/X}$$

$\beta_2$ 即为自变量的相对变化所带来的因变量的平均绝对变化

由上式可知： $\Delta \bar{Y} = \beta_2 \frac{\Delta X}{X} \xrightarrow{\text{令 } \Delta X/X = 1\%} \Delta \bar{Y} = \beta_2 \times 1\% = \frac{\beta_2}{100}$ ，故 $\beta_2$ 的含义是：

自变量  $X$  每增长1%，因变量  $Y$  平均将增加  $\frac{\beta_2}{100}$  个单位。

例：下表根据不同的模型回归了消费支出 $Y$ (元)与收入 $X$ (元)之间的关系：

模型	截距	斜率	参数含义
双对数线性模型 $\ln Y = \beta_1 + \beta_2 \ln X + u$	0.673	0.901	收入每增长1%，消费支出平均将增长0.901%
	$Se=(0.333)$	(0.033)	
对数－线性模型 $\ln Y = \beta_1 + \beta_2 X + u$	9.139	0.00253	收入每增加1元，消费支出平均将增长0.253%
	$Se=(0.040)$	(1.3E-06)	
线性－对数模型 $Y = \beta_1 + \beta_2 \ln X + u$	- 198826.8	21475.06	收入每增长1%，消费支出平均将增加214.75元
	$Se=(9657.49)$	(947.602)	
线性模型 $Y = \beta_1 + \beta_2 X + u$	2372.630	0.623	收入每增加1元，消费支出平均将增加0.623元
	$Se=(546.027)$	(0.018)	