

变量之间的关系一般可分为**确定性**与**非确定性**两种:

- (1) 确定性关系是指变量之间的关系可以用函数关系来表达;
- (2) 而非确定性的关系即所谓的相关关系. 回归分析是研究相关关系的一种数学工具.

## 一、回归的含义

### 1. 自变量与因变量的定义:

**【补例1】**人的体重Y与身高X之间存在着相关关系.

由图4-2可作如下假设:

$$y = \beta_0 + \beta_1 x + \varepsilon$$

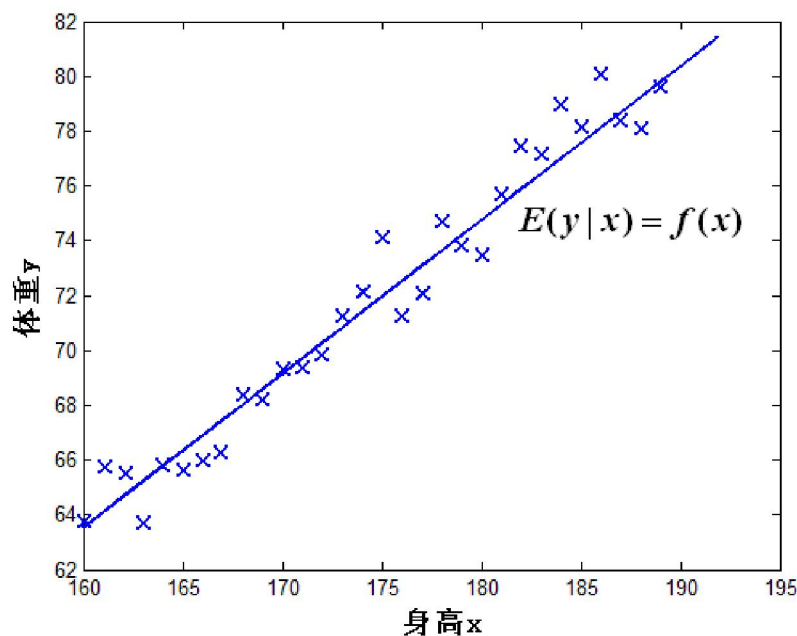


图4-2 三十个男子的身高和体重

设随机变量 $y$ 与 $x$ 之间存在着某种相关关系. 这里:

①  $x$ 是可控或可精确观察的变量, 故把 $x$ 看成普通变量, 称 $x$ 为自变量(预报变量、回归变量);

②  $y$ 称为因变量(响应变量).

若对每一确定 $x$ 值,  $E(y|x) = f(x)$ 存在, 则称  $f(x)$  为 $y$ 关于 $x$ 的回归函数.

**【补例2】**人的体重 $Y$ 与身高 $X$ 、性别 $S$ 、地区 $C$ 之间存在着相关关系.

因此还可作如下假设:

$$y = f(x, s, c) + \varepsilon$$

## 2. 回归分析的任务与目的:

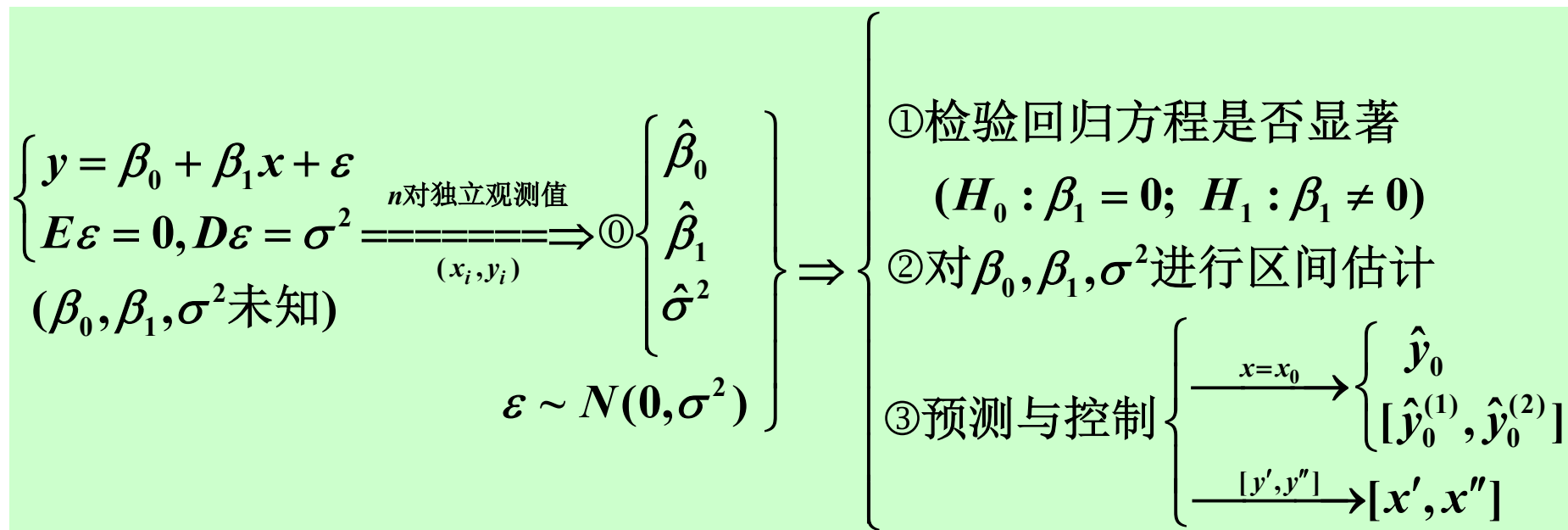
任务: 研究自变量之变动对因变量之变动的影晌程度;

目的: 根据自变量的变化来估计或预测因变量的变化情况.

## 3. 回归分析的内容:

- (1) 确定回归模型;
- (2) 根据样本估计未知参数;
- (3) 检验回归方程与各自变量的显著性;
- (4) 利用自变量的值来估计和预测因变量.

## 二、本节内容与思路



### 三、一元线性回归模型

下面的讨论中, 自变量 $x$ 为非随机变量, 而因变量 $y$ 为随机变量.

#### 1. 总体模型:

$$\begin{cases} y = \beta_0 + \beta_1 x + \varepsilon \\ E\varepsilon = 0, D\varepsilon = \sigma^2 < \infty (\sigma^2 \text{未知}) \end{cases} \quad (4.3)$$

回归系数——固定的未知参数 $\beta_0, \beta_1$ .

$$E(y|x) = \beta_0 + \beta_1 x \quad \text{——} y \text{对} x \text{的回归函数}$$

$$\hat{y} = \hat{\beta}_0 + \hat{\beta}_1 x \quad \text{——} y \text{对} x \text{的回归方程}$$

$$(\text{拟合方程或预报方程}) \quad (4.4)$$

回归直线—— $\hat{y} = \hat{\beta}_0 + \hat{\beta}_1 x$ 所表示的直线 (拟合直线);

回归值——对固定的 $x$ , 相应 $y$ 的估计值 $\hat{y} = \hat{\beta}_0 + \hat{\beta}_1 x$  (拟合值或预报值).

## 2. 样本模型:

假设有  $n$  组独立观测值  $(x_i, y_i) (i = 1, 2, \dots, n)$ , 则由(4.3)有

$$\begin{cases} y_i = \beta_0 + \beta_1 x_i + \varepsilon_i, & i = 1, 2, \dots, n \\ E\varepsilon_i = 0, D\varepsilon_i = \sigma^2, & \text{且 } \varepsilon_1, \varepsilon_2, \dots, \varepsilon_n \text{ 相互独立} \end{cases} \quad (4.5)$$

## 四、最小二乘估计及统计性质

### 1. $\beta_0, \beta_1$ 的最小二乘估计:

#### (1) 最小二乘法的原理:

$$\text{误差平方和: } Q(\beta_0, \beta_1) = \sum_{i=1}^n \varepsilon_i^2 = \sum_{i=1}^n (y_i - \beta_0 - \beta_1 x_i)^2$$

$$Q(\hat{\beta}_0, \hat{\beta}_1) = \min_{\beta_0, \beta_1} Q(\beta_0, \beta_1)$$

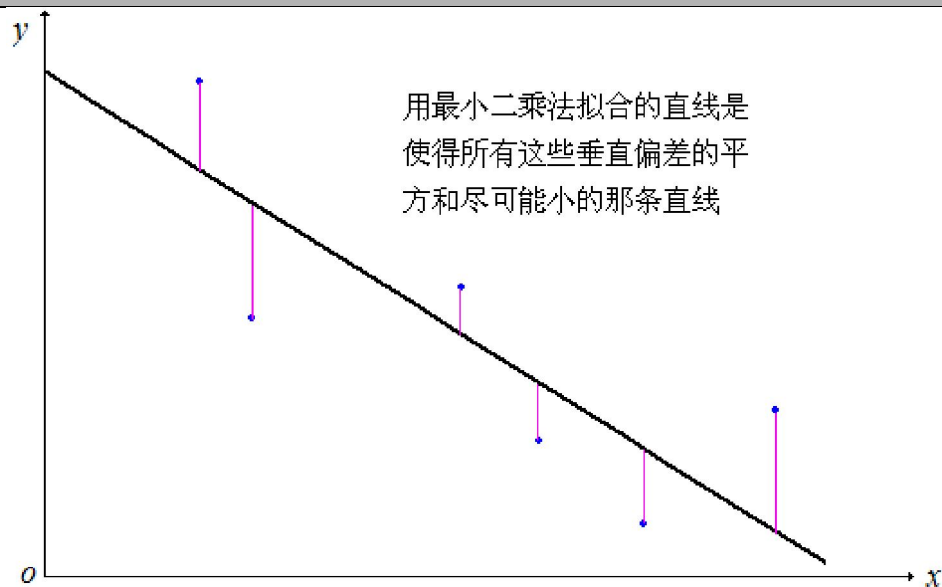


图4-5 最小二乘法拟合的直线

(2)  $\beta_0, \beta_1$  的最小二乘估计 (Least Squares Estimation, 简称LS估计):

$$\hat{\beta}_1 = \frac{L_{xy}}{L_{xx}}, \quad \hat{\beta}_0 = \bar{y} - \hat{\beta}_1 \bar{x} \quad (4.11)(4.12)$$

$$\bar{x} = \frac{1}{n} \sum_{i=1}^n x_i, \quad \bar{y} = \frac{1}{n} \sum_{i=1}^n y_i, \quad L_{xx} = \sum_{i=1}^n (x_i - \bar{x})^2, \quad L_{xy} = \sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y}).$$

### (3) $y$ 对 $x$ 的回归方程:

$$\hat{y} = \hat{\beta}_0 + \hat{\beta}_1 x \quad \text{—— (4.13)}$$

【★例4.1(P<sub>138</sub>)】为研究温度对某个化学过程生产量的影响, 收集数据如下. 求 $y$ 对 $x$ 的回归方程, 并作拟合曲线图与观测数据的散点图.

$x$	-5	-4	-3	-2	-1	0	1	2	3	4	5
$y$	1	5	4	7	10	8	9	13	14	13	18

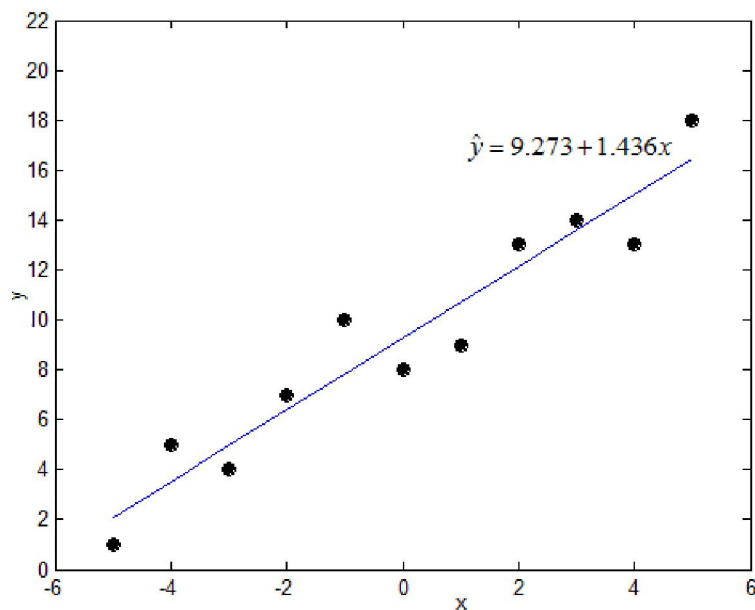


图4-4 数据散点图和拟合直线



## 2. LS估计 $\hat{\beta}_0$ 和 $\hat{\beta}_1$ 的性质:

(1)  $\hat{\beta}_0$ 和 $\hat{\beta}_1$ 是 $y_1, y_2, \dots, y_n$ 的线性函数.

(2) 定理4.1  $E(\hat{\beta}_0) = \beta_0, \quad E(\hat{\beta}_1) = \beta_1,$

$$D(\hat{\beta}_0) = \left( \frac{1}{n} + \frac{\bar{x}^2}{L_{xx}} \right) \sigma^2, \quad D(\hat{\beta}_1) = \frac{1}{L_{xx}} \sigma^2.$$

说明1: 安排试验时应注意:

- ①  $x_1, \dots, x_n$ 可正可负时, 尽可能使 $\bar{x} = 0$ , 此时 $D(\hat{\beta}_0)$ 最小;
- ②  $x_1, \dots, x_n$ 越分散越好, 即使 $L_{xx}$ 越大越好;
- ③ 试验次数 $n$ 不能太小.

### 3. $\sigma^2$ 的无偏估计:

残差: 
$$e_i = y_i - \hat{y}_i$$

残差平方和: 
$$Q_e = Q(\hat{\beta}_0, \hat{\beta}_1) = \sum_{i=1}^n (y_i - \hat{y}_i)^2$$

**定理4.2** 
$$EQ_e = (n-2)\sigma^2 \Rightarrow E\left(\frac{Q_e}{n-2}\right) = \sigma^2.$$

**结论1:**  $\hat{\sigma}_e^2 = \frac{Q_e}{n-2}$  (剩余方差、残差的方差) 是  $\sigma^2$  的无偏估计.

## 五、回归方程的显著性检验和回归系数的置信区间

**定理4.3:** 若假定  $\varepsilon \sim N(0, \sigma^2)$ , 则模型(4.5)可写成

$$\begin{cases} y_i = \beta_0 + \beta_1 x_i + \varepsilon_i, & i = 1, 2, \dots, n \\ \varepsilon_i \sim N(0, \sigma^2), & \text{且 } \varepsilon_1, \varepsilon_2, \dots, \varepsilon_n \text{ 相互独立} \end{cases} \quad (4.17)$$

$$\left. \begin{aligned} \hat{\beta}_0 &\sim N\left(\beta_0, \left(\frac{1}{n} + \frac{\bar{x}^2}{L_{xx}}\right)\sigma^2\right) \\ \hat{\beta}_1 &\sim N\left(\beta_1, \frac{1}{L_{xx}}\sigma^2\right) \\ \frac{Q_e}{\sigma^2} &\sim \chi^2(n-2) \\ \bar{y}, \hat{\beta}_1, Q_e &\text{相互独立, } \hat{\beta}_0 \text{与 } Q_e \text{独立} \end{aligned} \right\} \Rightarrow \begin{cases} \frac{(\hat{\beta}_0 - \beta_0) / \sqrt{\frac{1}{n} + \frac{\bar{x}^2}{L_{xx}}}}{\sqrt{Q_e / (n-2)}} \sim t(n-2) \\ \frac{(\hat{\beta}_1 - \beta_1) / \sqrt{1/L_{xx}}}{\sqrt{Q_e / (n-2)}} \sim t(n-2) \\ \frac{Q_e}{\sigma^2} \sim \chi^2(n-2) \end{cases}$$

## (一) 回归方程的显著性检验

检验 $y$ 与 $x$ 之间有无线性关系  $y = \beta_0 + \beta_1 x + \varepsilon$ , 即检验假设:

$$H_0: \beta_1 = 0; \quad H_1: \beta_1 \neq 0 \quad (4.18)$$

### 1. 总离差平方和的分解式:

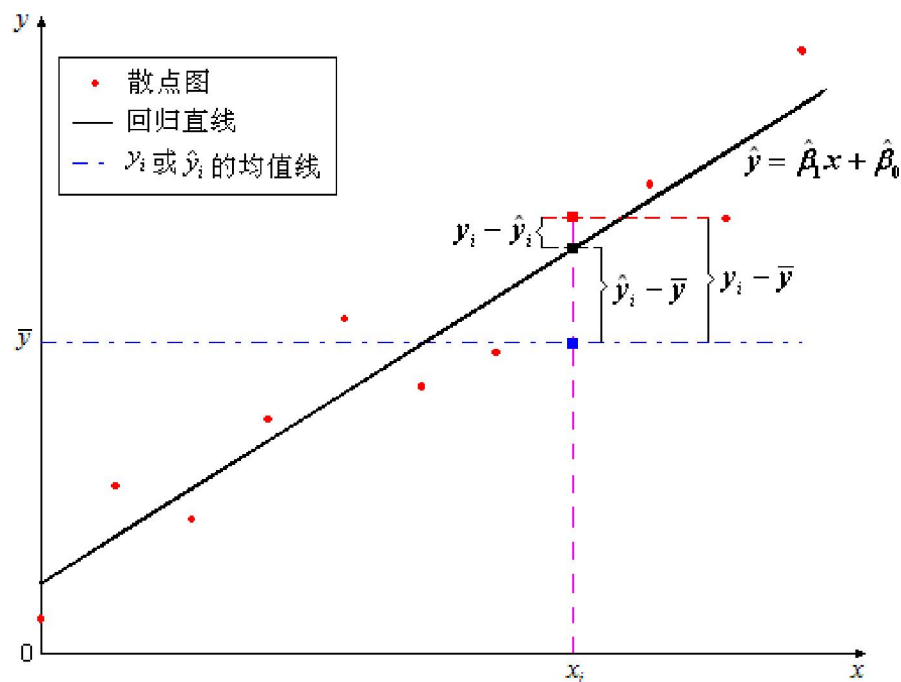


图4-7 总离差平方和分解式的几何意义

$$L_{yy} = \sum_{i=1}^n (y_i - \bar{y})^2 \text{——总离差平方和}$$

$$Q_e = \sum_{i=1}^n (y_i - \hat{y}_i)^2 \text{——残差平方和}$$

$$U = \sum_{i=1}^n (\hat{y}_i - \bar{y})^2 = \hat{\beta}_1^2 L_{xx} = L_{xy}^2 / L_{xx}$$

——回归平方和

则

$$L_{yy} = Q_e + U$$

## 2. 分解式中各元素的统计意义:

$$EQ_e = (n-2)\sigma^2, \quad EU = \sigma^2 + \beta_1^2 L_{xx}$$

$Q_e$ ——反映了误差引起数据  $y_1, \dots, y_n$  的波动程度大小;

$U$ ——除反映了误差的作用外, 还反映了回归因子  $x$  对  $y$  的线性影响.

## 3. 回归方程的显著性检验:

### (1) F检验法

① 定理4.4  $Q_e$  与  $U$  独立, 且  $U/\sigma^2 \overset{H_0 \text{为真}}{\sim} \chi^2(1)$ .

② 检验统计量: 
$$F = \frac{U}{Q_e/(n-2)} = \frac{\hat{\beta}_1^2 L_{xx}}{Q_e/(n-2)} \overset{H_0 \text{为真}}{\sim} F(1, n-2).$$

③  $H_0: \beta_1 = 0$  的拒绝域:  $F > F_{1-\alpha}(1, n-2)$ .

## (2) t检验法:

① 检验统计量: 
$$T = \frac{\sqrt{L_{xx}} \hat{\beta}_1}{\sqrt{Q_e / (n-2)}} \overset{H_0 \text{为真}}{\sim} t(n-2).$$

②  $H_0: \beta_1 = 0$  的拒绝域:  $|T| > t_{1-\alpha/2}(n-2).$

【★例4.2(P<sub>153</sub>)】为研究温度对某个化学过程的生产量的影响, 收集到如下数据:

x	-5	-4	-3	-2	-1	0	1	2	3	4	5
y	1	5	4	7	10	8	9	13	14	13	18

并用最小二乘法得到拟合直线为:  $\hat{y} = 1.436x + 9.273$ . 现要求在正态分布下分别用t检验法和F检验法, 检验回归方程效果是否显著 ( $\alpha = 0.05$ ).

### (3) r检验法

$$\textcircled{1} \quad \text{由 } r^2 = \frac{U}{L_{yy}} = \frac{L_{xy}^2}{L_{xx}L_{yy}} \Rightarrow r = \frac{L_{xy}}{\sqrt{L_{xx}L_{yy}}} = \frac{\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})}{\sqrt{\sum_{i=1}^n (x_i - \bar{x})^2 \sum_{i=1}^n (y_i - \bar{y})^2}}$$

——样本相关系数 (4.23)

$$\textcircled{2} \quad H_0: \beta_1 = 0 \text{ 的拒绝域: } |r| > r_{1-\alpha}.$$

其中  $r_{1-\alpha}$  通过查“相关系数检验临界值表”得到.

**注意1:** r检验法临界值可由F检验法推得.

**说明2:** ★若拒绝 $H_0$ , 则认为 $y$ 与 $x$ 存在线性关系, 所求回归方程有意义;

★否则此回归方程无意义. 此时, 可能有如下几种情况:

- $x$ 对 $y$ 没有显著影响, 此时应丢掉 $x$ ;
- $x$ 对 $y$ 有显著影响, 但该影响不是线性的, 改用非线性回归;
- 除 $x$ 外, 还有其它不可忽略的自变量对 $y$ 有显著影响, 从而削弱了 $x$ 对 $y$ 的影响, 此时改用多元线性回归.



## (二) 回归系数的置信区间( 记 $\hat{\sigma}_e = \sqrt{Q_e/(n-2)}$ )

1.  $\beta_0$  的置信水平为  $1-\alpha$  的置信区间:

$$\frac{(\hat{\beta}_0 - \beta_0) / \sqrt{\frac{1}{n} + \frac{\bar{x}^2}{L_{xx}}}}{\sqrt{Q_e/(n-2)}} \sim t(n-2) \Rightarrow \left[ \hat{\beta}_0 \pm t_{1-\alpha/2}(n-2) \cdot \hat{\sigma}_e \cdot \sqrt{\frac{1}{n} + \frac{\bar{x}^2}{L_{xx}}} \right].$$

2.  $\beta_1$  的置信水平为  $1-\alpha$  的置信区间:

$$\frac{(\hat{\beta}_1 - \beta_1) / \sqrt{1/L_{xx}}}{\sqrt{Q_e/(n-2)}} \sim t(n-2) \Rightarrow \left[ \hat{\beta}_1 \pm t_{1-\alpha/2}(n-2) \cdot \hat{\sigma}_e \cdot \sqrt{\frac{1}{L_{xx}}} \right].$$

3.  $\sigma^2$  的置信水平为  $1-\alpha$  的置信区间:

$$\frac{Q_e}{\sigma^2} \sim \chi^2(n-2) \Rightarrow \left[ \frac{Q_e}{\chi_{1-\alpha/2}^2(n-2)}, \frac{Q_e}{\chi_{\alpha/2}^2(n-2)} \right].$$

【★例4.3(P<sub>156</sub>)】 为研究温度对某个化学过程的生产量的影响, 收集到如下数据:

x	-5	-4	-3	-2	-1	0	1	2	3	4	5
y	1	5	4	7	10	8	9	13	14	13	18

求回归系数 $\beta_0, \beta_1$ 的置信区间( $\alpha = 0.05$ ).

## 六、预测与控制

设y与x满足模型 $\begin{cases} y = \beta_0 + \beta_1 x + \varepsilon \\ \varepsilon \sim N(0, \sigma^2) \end{cases}$  (4.17). 并已得回归方程 $\hat{y} = \hat{\beta}_1 x + \hat{\beta}_0$ .

1. 预测: 对固定的x值预测相应的y值.

令 $x_0$ 为x的一固定值, 且 $\begin{cases} y_0 = \beta_0 + \beta_1 x_0 + \varepsilon_0 \\ \varepsilon_0 \sim N(0, \sigma^2) \end{cases}$ . 设 $y_0, y_1, \dots, y_n$ 相互独立,

求:

### (1) $y_0$ 的预测值:

$\hat{y}_0 = \hat{\beta}_0 + \hat{\beta}_1 x_0$ , 且  $\hat{y}_0$  是  $Ey_0$  的无偏估计.

### (2) $y_0$ 的置信水平为 $1-\alpha$ 的预测区间:

$$T = \frac{y_0 - \hat{y}_0}{\hat{\sigma}_e \sqrt{1 + \frac{1}{n} + \frac{(x_0 - \bar{x})^2}{L_{xx}}}} \sim t(n-2) \Rightarrow [\hat{y}_0 \pm \delta(x_0)] \quad (4.31)$$

其中,  $\delta(x_0) = \hat{\sigma}_e \cdot t_{1-\alpha/2}(n-2) \cdot \sqrt{1 + \frac{1}{n} + \frac{(x_0 - \bar{x})^2}{L_{xx}}}$ .

### 说明3:

- $\hat{\sigma}_e$  越小, 预测区间越窄, 预测就越精确;
- $x_0$  越靠近  $\bar{x}$ , 预测精度也就越高.

### (3) $y = \beta_0 + \beta_1 x + \varepsilon$ 的预测区间:

$$[\hat{y} \pm \delta(x)] \quad (4.32)$$

特别, 当  $n$  很大而  $|x - \bar{x}|$  很小时,  $y$  的  $1-\alpha$  的预测区间近似为:

$$[\hat{y} \pm \hat{\sigma}_e u_{1-\alpha/2}] \quad (4.33)$$

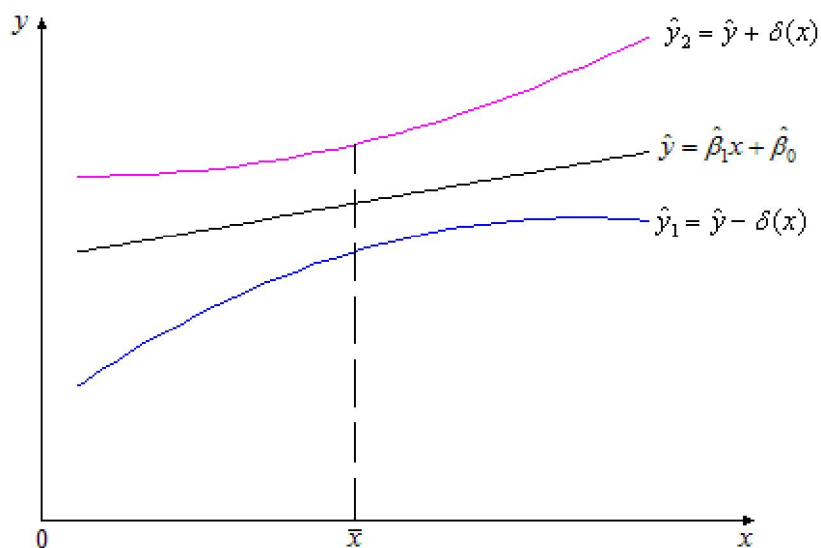


图4-8  $y$  的预测带(4.32)

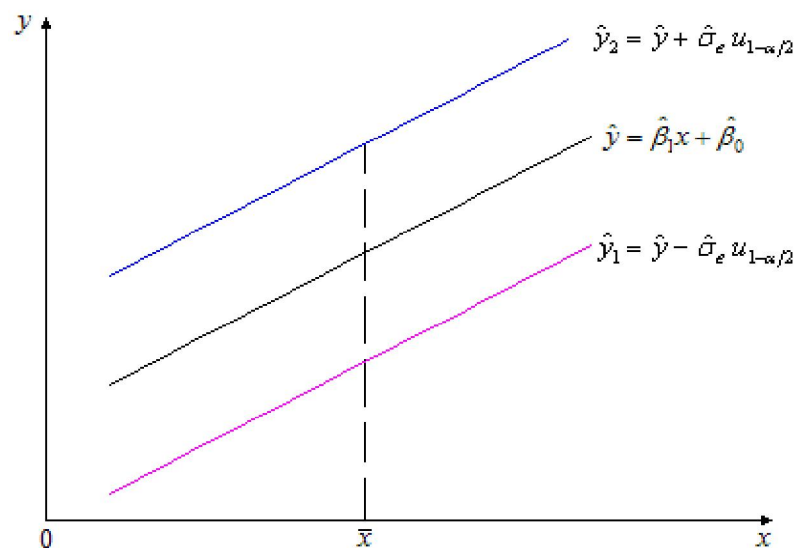


图4-9  $y$  的近似预测带(4.33)

【★例4.4(P<sub>159</sub>)】 为研究温度对某个化学过程的生产量的影响, 收集到如下数据:

x	-5	-4	-3	-2	-1	0	1	2	3	4	5
y	1	5	4	7	10	8	9	13	14	13	18

取  $x_0 = 3$ , 求  $y_0$  的预测值与置信水平为  $1 - \alpha = 0.95$  的预测区间.

2. 控制: 控制  $x$  的值以便把  $y$  的值控制在指定的范围内.

若要  $y = \beta_0 + \beta_1 x + \varepsilon$  的值以  $1 - \alpha$  的概率落在指定区间  $(y', y'')$  之内, 求自变量  $x$  的控制范围  $(x', x'')$ .

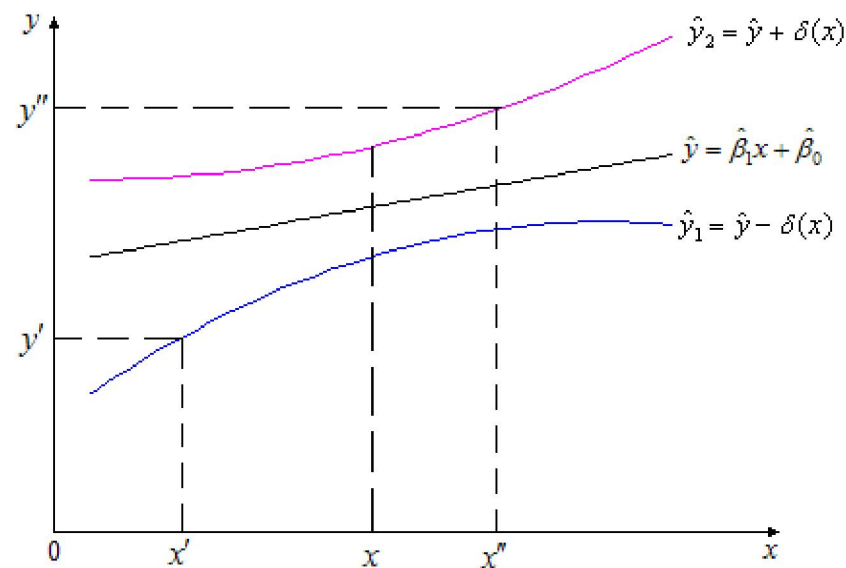


图4-10  $x$  的控制范围

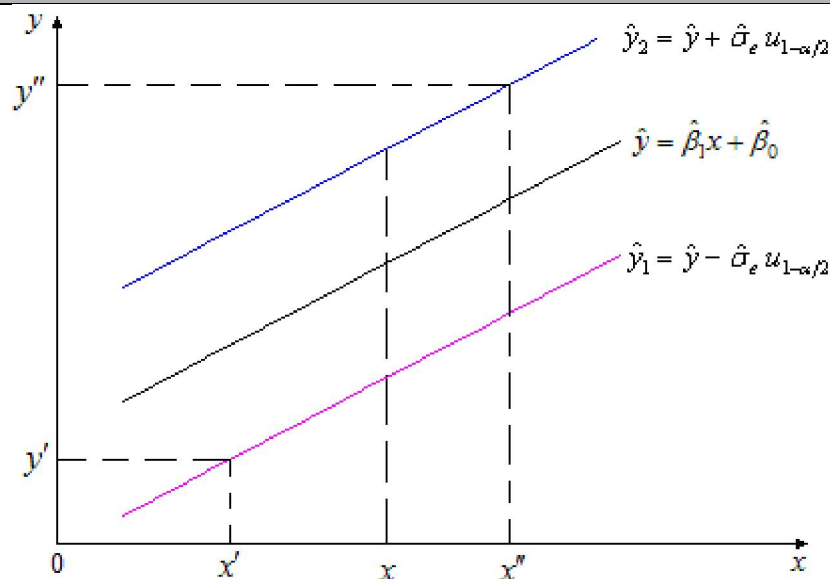


图4-11  $x$ 的近似控制范围 ( $n$ 很大且 $|x - \bar{x}|$ 很小)

对图4-11, 由 
$$\begin{cases} y' = \hat{\beta}_1 x + \hat{\beta}_0 - \hat{\sigma}_e u_{1-\alpha/2} \\ y'' = \hat{\beta}_1 x + \hat{\beta}_0 + \hat{\sigma}_e u_{1-\alpha/2} \end{cases} \Rightarrow \begin{cases} x' = \frac{1}{\hat{\beta}_1} (y' - \hat{\beta}_0 + \hat{\sigma}_e u_{1-\alpha/2}) \\ x'' = \frac{1}{\hat{\beta}_1} (y'' - \hat{\beta}_0 - \hat{\sigma}_e u_{1-\alpha/2}) \end{cases}.$$

当  $\hat{\beta}_1 > 0$  时,  $x$  控制为  $(x', x'')$ ; 当  $\hat{\beta}_1 < 0$  时,  $x$  控制为  $(x'', x')$ .

**注意2:** 为了实现上述控制, 必须使区间  $(y', y'')$  的长度大于  $2\hat{\sigma}_e u_{1-\alpha/2}$ .