

§9.2 一元回归分析

本节讨论回归函数是一元线性函数或可线性化函数的情况.

一.一元线性回归模型

若回归函数是线性函数

$$\mu(x_1, x_2, \dots, x_k) = b_0 + b_1 x_1 + b_2 x_2 + \dots + b_k x_k$$

其中 b_0, b_1, \dots, b_k 是未知常数, 称为**线性回归问题**.

若 Y 关于 X 的回归函数为

$$\mu(x) = E(Y | X = x) = a + bx$$

有一元线性回归模型:

$$Y = a + bx + \varepsilon, \quad \varepsilon \sim N(0, \sigma^2)$$

其中 a 、 b 、 σ^2 为未知参数, 且

a — 回归常数(又称截距)

b — 回归系数(又称斜率)

ε — 随机误差（随机扰动项）

若随机误差 $\varepsilon \sim N(0, \sigma^2)$, 称为**一元线性正态回归模型**.

ε_i 是第 i 次观察时的随机误差, 有

- 1) $E(\varepsilon_i) = 0, \quad D(\varepsilon_i) = \sigma^2, \quad i = 1, 2, \dots, n;$
- 2) $\varepsilon_1, \dots, \varepsilon_n$ 相互独立.

回归假定

二.一元线性回归模型的参数估计

用观察值 (x_i, y_j) a 、 b 、 σ^2 进行估计.

对所假定的回归模型进行检验 (回归系数 b)

对自变量 X 的一组值 x_1, x_2, \dots, x_n 做 n 次独立试验, 得独立观察值 y_1, y_2, \dots, y_n .

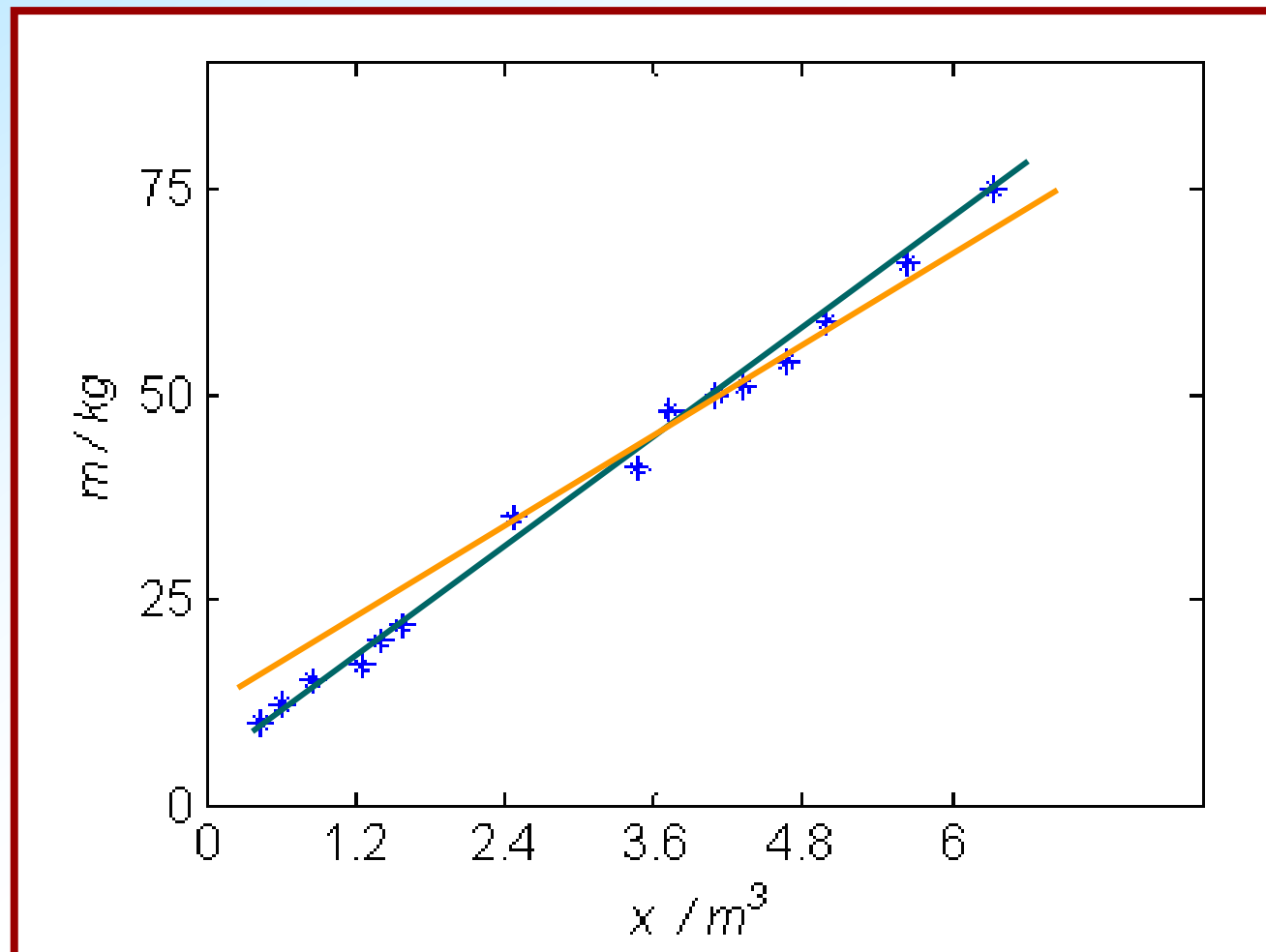
问题 如何依据观察值

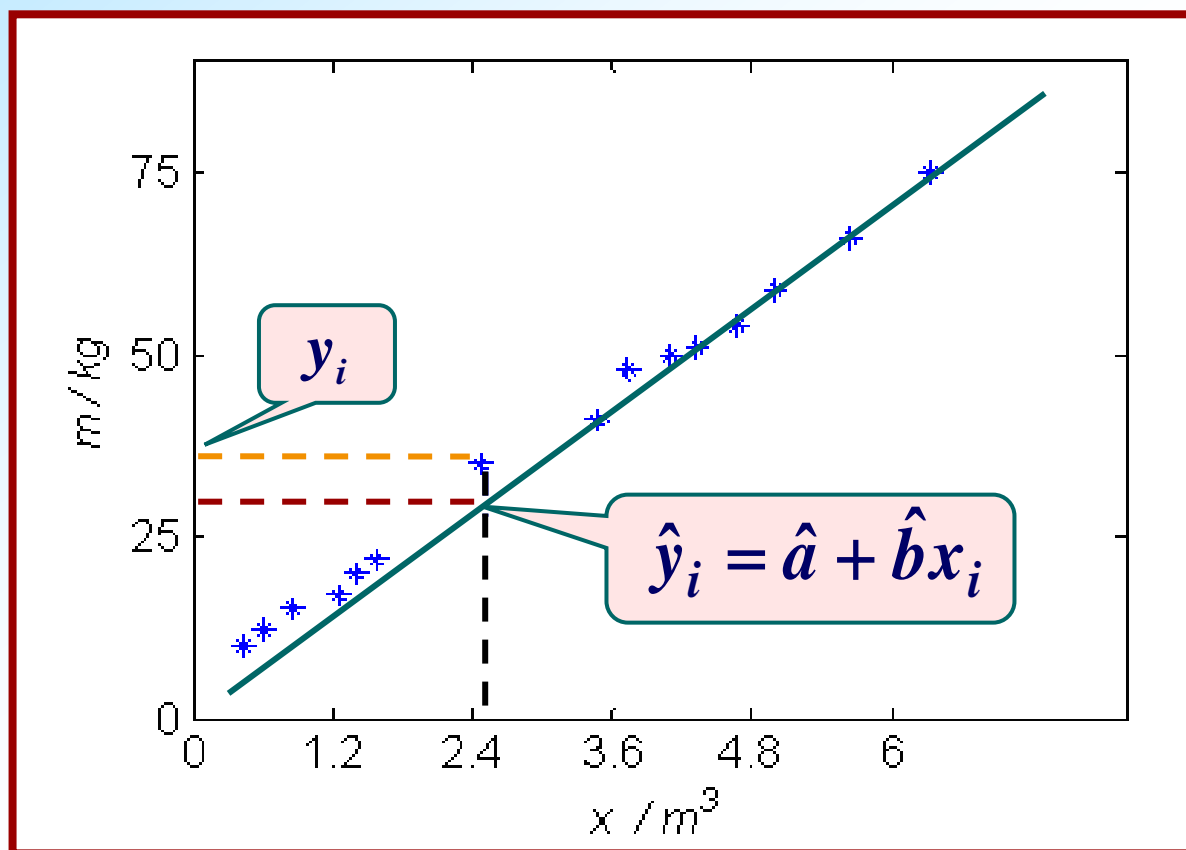
$$(x_i, y_i), i=1, 2, \dots, n.$$

求 a 、 b 的估计值 \hat{a}, \hat{b} ?

记 y_i 的估计值为 $\hat{y}_i = \hat{a} + \hat{b}x_i$

称为**回归值**.





对所有的 i ，应使偏差 $y_i - \hat{y}_i$ 都尽可能小，
有三种思路：

1) 使误差总和 $\sum (y_i - \hat{y}_i)$ 最小;

缺点: 可能正负误差抵消

2) 使误差绝对值之和 $\sum |y_i - \hat{y}_i|$ 最小;

缺点: 数学处理困难

3) 使误差的平方和 $\sum (y_i - \hat{y}_i)^2$ 最小.

结论 应选 a 、 b 的估计使离差(误差)平方和:

$$Q = \sum_{i=1}^n (y_i - \hat{y}_i)^2 = \sum_{i=1}^n (y_i - \hat{a} - \hat{b}x_i)^2$$

达最小.

$$\text{令 } Q(a, b) = \sum_{i=1}^n (y_i - a - bx_i)^2$$

应用最小二乘法 分别对 a, b 求一阶偏导,
并建立方程组:

$$\begin{cases} \sum_{i=1}^n (y_i - a - bx_i) = 0 \\ \sum_{i=1}^n (y_i - a - bx_i)x_i = 0 \end{cases}$$

或

$$\begin{cases} na + (\sum_{i=1}^n x_i)b = \sum_{i=1}^n y_i & (1) \\ (\sum_{i=1}^n x_i)a + (\sum_{i=1}^n x_i^2)b = \sum_{i=1}^n x_i y_i & (2) \end{cases}$$

称为**正规方程组**，由克莱姆法则，解得：

$$\begin{cases} \hat{b} = \frac{l_{xy}}{l_{xx}} \\ \hat{a} = \bar{y} - \hat{b}\bar{x} \end{cases}$$

其中

$$\bar{x} = \frac{1}{n} \sum_{i=1}^n x_i, \quad \bar{y} = \frac{1}{n} \sum_{i=1}^n y_i$$

$$l_{xy} = \sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y}) = \sum_{i=1}^n x_i y_i - n\bar{x} \bar{y}$$

$$l_{xx} = \sum_{i=1}^n (x_i - \bar{x})^2 = \sum_{i=1}^n x_i^2 - n\bar{x}^2$$

由回归假定

$$E(\varepsilon_i)=0, \quad D(\varepsilon_i)=\sigma^2, \quad i=1,2, \dots, n;$$

有 $\sigma^2=D(\varepsilon)=E(\varepsilon^2)$, 故

$$\frac{1}{n} \sum_{i=1}^n \varepsilon_i^2$$

是 σ^2 的矩估计量.

可证明 σ^2 的无偏估计量为

$$\hat{\sigma}^2 = \frac{1}{n-2} \sum_{i=1}^n (y_i - \hat{y}_i)^2$$

代入 $\hat{y}_i = \hat{a} + \hat{b}x_i$, 得

$$\hat{\sigma}^2 = \frac{1}{n-2} \sum_{i=1}^n (l_{yy} - \hat{b}^2 l_{xx})$$

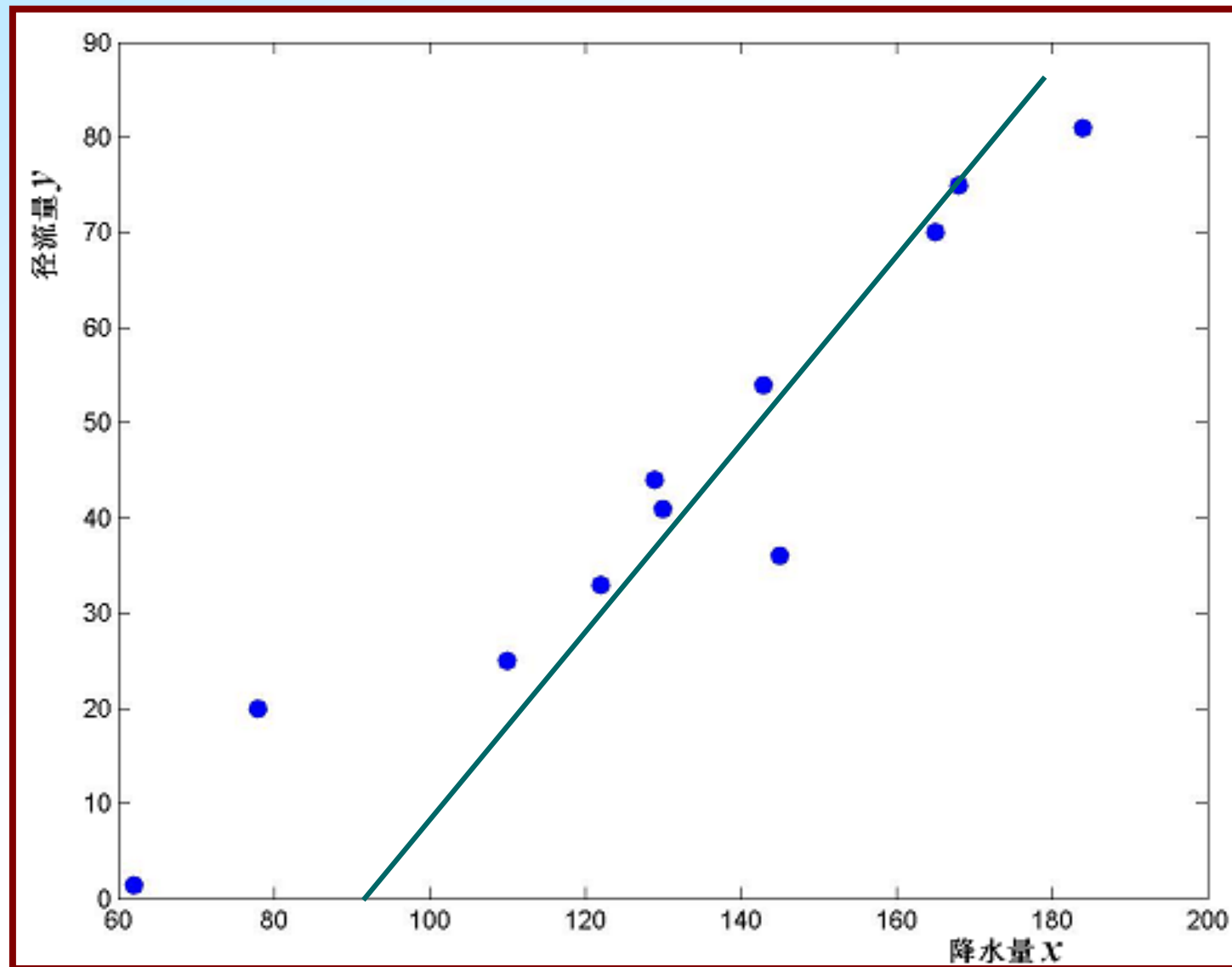
其中

$$l_{yy} = \sum_{i=1}^n (y_i - \bar{y})^2 = \sum_{i=1}^n y_i^2 - n\bar{y}^2$$

例9.2.1 流经某地区的降雨量 X 和该地河流的径流量 Y 的观察值如下表,

降雨量 x_i :	110	184	145	122	165	143	78
径流量 y_i :	25	81	36	33	70	54	20
	129	62	130	168	1436	(Σ)	
	44	1.41	41	75	480.4	(Σ)	

求 Y 关于 X 的(经验)线性回归方程,试估计降雨量为200时径流量为多少?



解 $n=11, \bar{x} = 130.5, \bar{y} = 43.7$

$$l_{xx} = \sum_{i=1}^{11} (x_i - \bar{x})^2 = 13768.7$$

$$\begin{aligned} l_{xy} &= \sum_{i=1}^n x_i y_i - n\bar{x}\bar{y} \\ &= 71424.8 - 62731.35 = 8693.45 \end{aligned}$$

$$\begin{cases} \hat{b} = \frac{l_{xy}}{l_{xx}} = \frac{8693.45}{13768.7} = 0.63 \\ \hat{a} = \bar{y} - \hat{b}\bar{x} = 43.7 - 0.63 \times 130.5 = -38.5 \end{cases}$$

所求经验回归方程为

$$\hat{y} = \hat{a} + \hat{b}x = -38.5 + 0.693x$$

降雨量为200时的径流量值为

$$\hat{y}(200) = 0.693 \times 200 - 38.5 = 100.1$$

$$\text{又 } l_{yy} = \sum_{i=1}^n (y_i - 43.7)^2 = 6050.59$$

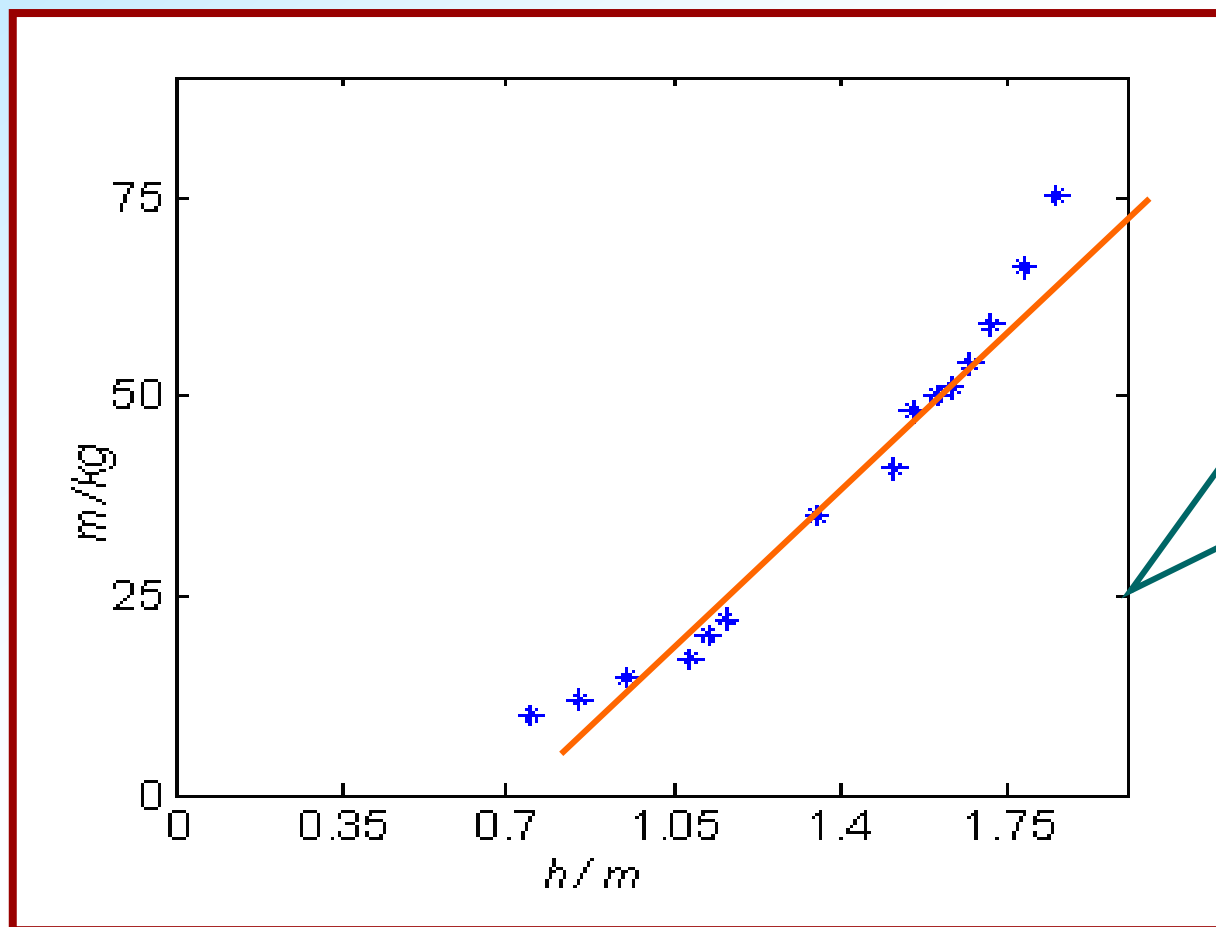
随机误差的方差 σ^2 的估计为

$$\hat{\sigma}^2 = \frac{1}{n-2} \sum_{i=1}^n (l_{yy} - \hat{b}^2 l_{xx})$$
$$= (6050.59 - 0.63^2 \times 13768.7) / 11 = 53.25$$

问题

随机变量 Y 与 X 间是否存在线性相关关系？
是否能由数据散布图完全确定回归函数？

续 例9.1.1 身高体重关系



身高 h 和
体重 m 无
明显的
线性相
关关系.

形式地估计回归系数和回归常数,并建立经验线性回归方程无实际意义.

2. 一元线性回归的假设检验 (相关系数法)

相关系数法是基于试验数据检验随机变量间线性相关关系是否显著的一种方法.

相关系数

$$\rho_{XY} = \frac{E\{[X - E(X)][Y - E(Y)]\}}{\sqrt{D(X)}\sqrt{D(Y)}}$$

是表征随机变量 Y 与 X 的线性相关程度的数字特征.

样本相关系数:

$$\begin{aligned}\hat{\rho}_{XY} = R &= \frac{\frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})}{\sqrt{\frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})^2} \sqrt{\frac{1}{n} \sum_{i=1}^n (y_i - \bar{y})^2}} \\ &= \frac{l_{xy}}{\sqrt{l_{xx}} \sqrt{l_{yy}}}\end{aligned}$$

作为 ρ_{XY} 的估计值.

有 $|R| \leq 1$ ，统计量 R 描述了 X 与 Y 间的线性相关关系的密切程度.

1) 当 $R = 0 \longrightarrow l_{xy} = 0$

$\longrightarrow \hat{b} = \frac{l_{xy}}{l_{xx}} = 0$

\longrightarrow 经验回归方程形为 $\hat{y} = \hat{a}$

说明自变量 X 的变化不会引起因变量 Y 的变化(不相关).

$$2) \quad |R| = 1 \longrightarrow \frac{l_{xy}}{\sqrt{l_{xx}} \sqrt{l_{yy}}} = \pm 1$$

$$\longrightarrow l_{xy} = \pm \sqrt{l_{xx}} \sqrt{l_{yy}}$$

$$\longrightarrow \hat{b} = \frac{l_{xy}}{l_{xx}} = \frac{\pm \sqrt{l_{xx}} \sqrt{l_{yy}}}{l_{xx}} = \pm \sqrt{\frac{l_{yy}}{l_{xx}}} \neq 0,$$

可认为X与Y间存在线性相关关系.

结论 1) $|R|$ 越接近于1, X 与 Y 间的线性相关关系越显著;

2) $|R|$ 越靠近于0, X 与 Y 间的线性相关关系越不显著.

根据附表6 相关系数临界值表, 有

判别准则 给定显著性水平 α (0.05, 0.01)

当 $|R| > R_{\alpha}(n-2)$,

认为 X 与 Y 之间的线性相关关系显著

当 $|R| \leq R_{\alpha}(n-2)$,

认为 X 与 Y 之间的线性相关关系不显著

例9.2.2(续前例)利用相关系数显著性检验法, 检验降雨量 X 和径流量 Y 的线性相关关系是否显著.

解 X 与 Y 的样本相关系数为

$$R = \frac{l_{XY}}{\sqrt{l_{XX}} \sqrt{l_{YY}}}$$

$$= \frac{8693.45}{\sqrt{13768.7} \sqrt{6050.58}} = 0.952$$

查表得

$$R_{\alpha}(n-2) = R_{0.01}(9) = 0.735 < 0.952 = R$$

可认为 X 与 Y 的线性相关关系显著.

例

3. 非线性回归问题的线性化处理

§9.4

例、为制定在服装标准, 调查了一组女青年的身高 X 与裤长 Y 的数据, 经计算得 $(x_i, y_i), i = 1, 2, \dots, 30$

$$\sum_{i=1}^{30} x_i = 4797, \quad \sum_{i=1}^{30} y_i = 3068, \quad \sum_{i=1}^{30} x_i^2 = 767949,$$

$$\sum_{i=1}^{30} y_i^2 = 314112, \quad \sum_{i=1}^{30} x_i y_i = 491124$$

试求： 1. 裤长 Y 对身高 X 的经验回归直线;
2. 用相关系数检验法, 在显著性水平
 $\alpha=0.01$ 下检验回归方程的显著性.

解：1. 由已知

$$\bar{x} = \frac{1}{30} \sum_{i=1}^{30} x_i = 159.9, \quad \bar{y} = \frac{1}{30} \sum_{i=1}^{30} y_i = 102.3,$$

$$l_{xx} = \sum_{i=1}^{30} x_i^2 - \frac{1}{30} \left[\sum_{i=1}^{30} x_i \right]^2 = 908.7,$$

$$l_{xy} = \sum_{i=1}^{30} x_i y_i - \frac{1}{30} \left(\sum_{i=1}^{30} x_i \right) \left(\sum_{i=1}^{30} y_i \right) = 550.8$$

$$l_{yy} = \sum_{i=1}^{30} y_i^2 - \frac{1}{30} \left[\sum_{i=1}^{30} y_i \right]^2 = 357.87,$$

$$\hat{b} = \frac{l_{xy}}{l_{xx}} = \frac{550.8}{908.7} = 0.61,$$

$$\hat{a} = \bar{y} - \hat{b}\bar{x} = 102.3 - 0.61 \times 159.9 = 5.4$$

Y 关于 X 的经验回归方程为

$$\hat{y} = \hat{a} + \hat{b}x = 5.4 + 0.61x$$

2. 样本相关系数 $R = \frac{l_{xy}}{\sqrt{l_{xx}} \sqrt{l_{yy}}} = 0.966$

查表得 $R_{0.01}(28)=0.463$, 因为 $|R| > R_{0.01}(28)$,

可认为裤长 Y 与身高 X 之间的线性相关关系显著。#