

# 回归分析及其试验设计

华东师范大学出版社

# 回归分析及其试验设计

茆诗松 丁元 编著  
周纪芗 吕乃刚

华东师范大学出版社

## 回归分析及其试验设计

茆诗松 丁元 编著  
周纪芗 吕乃刚

---

华东师范大学出版社出版

(上海市中山北路 3663 号)

---

新华书店上海发行所发行 华东师大印刷厂印刷

开本 850×1168 1/32 印张 12 273 千字  
1981 年 10 月第 2 版 1986 年 6 月第 1 次印刷  
印数 12,001—18,000

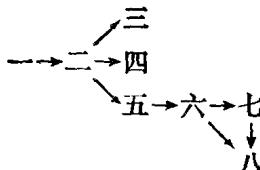
书号：13135·003 定价：1.20 元  
2.50 元

## 再 版 前 言

回归分析是研究随机现象中变量之间关系的一种数理统计方法，它在工农业生产和科学实验中有着广泛的应用。近几年来，广大科技人员运用回归分析方法，在进行数据处理、寻求经验公式、制定某些产品的新标准、探索新工艺与新配方、研究气象与地震的预报、及提取自动控制中的数学模型等方面，都取得了可喜的成绩，积累了新的经验。为了在工农业生产和科学实验中进一步推广回归分析方法，并在实践基础上进行理论研究，使之更好地为工农业生产服务，我们编写了本书。

全书共分八章。前二章通过生产中的实际问题，较详细地介绍了回归分析中的参数估计、统计检验和预报控制等问题。第三章阐述了逐步回归分析方法。第四章是讲多项式回归分析方法。这些方法在挑选因子或简化计算方面都是行之有效的。最后四章是向读者介绍回归的试验设计(简称回归设计)。这方面内容是近代发展起来的，在国内已有应用，并获得了一定的成果。可以预计，在我国工农业生产中它将会得到更多的应用与发展。为了适应教学的需要，我们在每章后面都配置了思考题与练习题。为了减少读者阅读本书理论部分的困难，我们编写了矩阵代数知识作为本书的附录。

这八章的内部联系是：



读者可根据自己的需要选择阅读的顺序。

书中采用的例子大部分是我们近几年来所接触到的课题，在

解决这些课题的过程中得到了有关单位的大力支持和协助，但由于我们接触实际还不够，也有部分例子是选自有关单位编写的书刊，在此一起表示深切的谢意。

在本书编著过程中得到魏宗舒教授、曹锡华教授以及原上海师大概率统计教研室其他同志的大力支持和帮助，在此表示深切感谢。

本书由上海教育出版社初版以来受到广大读者的关心与支持。很多读者给我们提出了宝贵意见和指出了一些错误。还有些读者通过书信和来访与编者一起讨论他们在实际工作中遇到的问题。在此我们表示由衷的感谢。

这次由华东师范大学出版社再版，我们除了对已发现的错误作了修正外，还增加了§ 6.7几个定理的证明一节和第八章 混料试验设计。这无论对于内容的完善，还是对于实践都是必要的。

本书再版，由于我们的水平有限，实践经验不足，不免仍会有不少错误。请读者批评指正。

编 者

1981年10月

# 目 录

|   |           |
|---|-----------|
| <b>第一章 一元线性回归 .....</b>                 | <b>1</b>  |
| § 1.1 什么是回归分析 .....                     | 1         |
| § 1.2 一元线性回归的数学模型 .....                 | 2         |
| § 1.3 参数 $\beta_0, \beta$ 的最小二乘估计 ..... | 3         |
| § 1.4 回归方程的显著性检验 .....                  | 7         |
| § 1.5 重复试验情况 .....                      | 12        |
| § 1.6 利用回归方程进行预报和控制 .....               | 18        |
| § 1.7 可化为线性回归的例子 .....                  | 27        |
| § 1.8 回归直线的简便求法 .....                   | 30        |
| 思考题与练习题 .....                           | 32        |
| <b>第二章 多元线性回归 .....</b>                 | <b>37</b> |
| § 2.1 多元线性回归的数学模型 .....                 | 37        |
| § 2.2 参数 $\beta$ 的最小二乘估计 .....          | 38        |
| § 2.3 线性回归数学模型的其他形式 .....               | 44        |
| § 2.4 回归方程的显著性检验 .....                  | 51        |
| § 2.5 回归系数的显著性检验 .....                  | 55        |
| § 2.6 利用回归方程进行预报和控制 .....               | 62        |
| § 2.7 多元线性回归的计算程序 .....                 | 67        |
| 思考题与练习题 .....                           | 71        |
| <b>第三章 逐步回归分析 .....</b>                 | <b>74</b> |
| § 3.1 “最优”回归方程的选择 .....                 | 74        |
| § 3.2 逐步回归分析的数学模型 .....                 | 78        |
| § 3.3 线性代数的有关知识 .....                   | 81        |
| § 3.4 逐步回归中的基本公式 .....                  | 88        |
| § 3.5 逐步回归的具体步骤 .....                   | 93        |
| § 3.6 逐步回归计算的框图与程序 .....                | 101       |
| 思考题与练习题 .....                           | 108       |

|                           |     |
|---------------------------|-----|
| <b>第四章 多项式回归与正交多项式</b>    | 109 |
| § 4.1 多项式回归               | 109 |
| § 4.2 正交多项式的应用            | 110 |
| § 4.3 多元正交多项式回归的例子        | 122 |
| 思考题与练习题                   | 131 |
| 附录：正交多项式表( $N=2\sim 30$ ) | 132 |
| <b>第五章 回归的正交设计</b>        | 141 |
| § 5.1 什么是回归设计             | 141 |
| § 5.2 一次回归的正交设计           | 142 |
| § 5.3 单纯形计划               | 151 |
| § 5.4 交互效应与部分实施           | 153 |
| § 5.5 一次回归正交设计的应用         | 157 |
| § 5.6 二次回归的正交设计           | 162 |
| § 5.7 二次回归正交设计的统计分析       | 174 |
| § 5.8 二次回归正交设计的计算程序       | 186 |
| 思考题与练习题                   | 189 |
| <b>第六章 回归的旋转设计</b>        | 191 |
| § 6.1 旋转性条件               | 191 |
| § 6.2 二次旋转设计              | 195 |
| § 6.3 二次旋转组合设计中 $m_e$ 的选择 | 201 |
| § 6.4 二次旋转设计的统计分析         | 211 |
| § 6.5 时间漂移与正交区组           | 219 |
| § 6.6 三次旋转设计              | 230 |
| § 6.7 几个定理的证明             | 241 |
| 思考题与练习题                   | 253 |
| <b>第七章 回归的 D-最优设计</b>     | 255 |
| § 7.1 回归模型与计划概念的拓广        | 255 |
| § 7.2 密集椭球体与 D-最优设计       | 259 |
| § 7.3 等价定理及其应用            | 263 |
| § 7.4 构造 D-最优计划的数值方法      | 284 |
| § 7.5 等价定理的证明             | 293 |
| § 7.6 饱和 D-最优设计           | 297 |

|                                |     |
|--------------------------------|-----|
| <b>第八章 混料试验设计</b>              | 303 |
| § 8.1 混料试验设计的概念                | 303 |
| § 8.2 单形格子设计                   | 307 |
| § 8.3 单形重心设计                   | 314 |
| § 8.4 极端顶点设计                   | 319 |
| § 8.5 最优设计与渐近最优设计              | 326 |
| § 8.6 混料试验的正交设计与旋转设计           | 336 |
| <b>附录 矩阵代数</b>                 | 351 |
| <b>附表 一、正态分布表</b>              | 356 |
| 二、 $t$ 分布的双侧分位数( $t_\alpha$ )表 | 360 |
| 三、 $F$ 检验的临界值( $F_\alpha$ )表   | 362 |

# 第一章 一元线性回归

## § 1.1 什么是回归分析

在生产斗争和科学实验中，经常遇到一些同处于一个统一体中的变量。在这个统一体中，这些变量是相互联系、相互制约的，也就是说，它们之间客观上存在着一定的关系。为了深入了解事物的本质，往往需要找出描述这些变量之间依存关系的数学表达式。在微积分中，我们研究了完全确定的函数关系（如图 1.1）。然而，在许多实际问题中，不是由于变量之间的关系比较复杂，使我们无法得到精确的数学表达式；就是由于生产或试验过程中不可避免地存在着误差的影响，而使它们之间的关系具有某种不确定性（如图 1.2）。

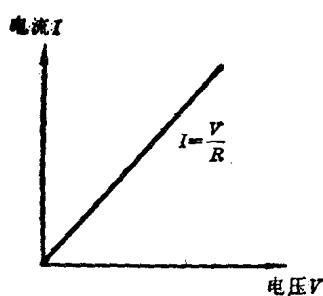


图 1.1

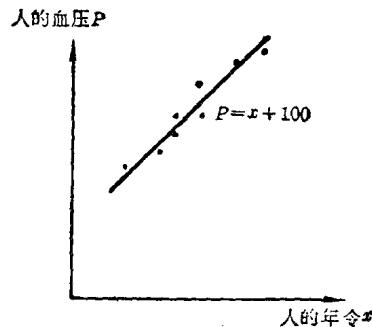


图 1.2

因此，需要我们用统计方法，在大量的试验和观察中，寻找隐藏在上述随机性后面的统计规律性。这类统计规律称为回归关系，有关回归关系的计算方法和理论通称为回归分析，它是数理统计的一个重要分支，在生产和科研中有着广泛的应用。譬如求经验公式，找出产量或质量指标与生产条件的关系，确定最佳生产条件，预报气象与病虫害，制定自动控制中的数学模型等等，都要用到回归分析的工具。

**回归分析的主要内容是：**

- (1) 从一组数据出发，确定这些变量间的定量关系式；
- (2) 对这些关系式的可信程度进行统计检验；
- (3) 从影响着某一个量的许多变量中，判断哪些变量的影响是显著的，哪些是不显著的；
- (4) 利用所求得的关系式对生产过程进行预报和控制；
- (5) 近代又出现，根据回归的分析方法，特别是进行预报和控制所提出的要求，选择试验点，对试验进行某种设计；
- (6) 寻求点数较少，且具有较好统计性质的回归设计方法。

关于这些内容的详细情况，将在以后各章逐步介绍。

回归分析所研究的数学模型主要是线性回归模型和多项式回归模型。后者可以化为前者，但后者本身也有一些特殊的方法。本章先讨论一元线性回归模型。

## § 1.2 一元线性回归的数学模型

一元回归处理的是两个变量之间的关系，即两个变量  $x$  和  $y$  间若存在一定的关系，则通过试验，分析所得数据，找出两者之间关系的经验公式。假如两个变量的关系是线性的，那就是一元线性回归分析所研究的对象。

先看一个例子。

[例 1] 上海市某生产队是三麦高产单位，在摸索高产经验的过程中，贫下中农经过多次试验，总结出一种根据小麦基本苗数推算成熟期有效穗数的方法。1973 年他们在五块田上进行了对比试验，在同样的肥料和管理水平下，取得如下数据：

| 试 验 号<br>$\alpha$ | 播 种 量<br>(11 月 18 日) | 基 本 苗 数 $x_\alpha$<br>(12 月 19 日) | 有 效 穗 数 $y_\alpha$<br>(5 月 5 日) |
|-------------------|----------------------|-----------------------------------|---------------------------------|
| 1                 | 25(斤/亩)              | 15(万/亩)                           | 39.4(万/亩)                       |
| 2                 | 30                   | 25.8                              | 42.9                            |
| 3                 | 35                   | 30                                | 41.0                            |
| 4                 | 40                   | 36.6                              | 43.1                            |
| 5                 | 45                   | 44.4                              | 49.2                            |

为了研究这些数据中所蕴藏的规律性，他们把第  $\alpha$  块田的基本苗数  $x_\alpha$  作为横坐标，有效穗数  $y_\alpha$  作为纵坐标，描出各点(图 1.3). 可以发现，这些点大致都落在一条直线附近。这就是说，变量  $x$  和  $y$  之间的关系可以基本上看作是线性关系；这些点与直线的偏离，是由试验过程中其他一些随机因素的影响而引起的。因此，表 1.1 中的数据可以假设具有如下的结构式：

$$y_\alpha = \beta_0 + \beta x_\alpha + \varepsilon_\alpha, \\ \alpha = 1, 2, \dots, N. \quad (1.1)$$

其中  $\varepsilon_1, \varepsilon_2, \dots, \varepsilon_N$  分别表示其他随机因素对有效穗数  $y_\alpha$  影响的总和，一

般假设它们是一组相互独立，且服从同一正态分布  $N(0, \sigma^2)$  的随机变量(本书中，对  $\varepsilon_\alpha, \alpha = 1, 2, \dots, N$ ，都作这样的假定，以后一般不再另行说明)。变量  $x$  可以是随机变量，也可以是一般变量，我们只讨论它是一般变量的情况，即它是可以精确测量或严格控制的变量。在上述这些条件下，变量  $y$  是服从正态分布  $N(\beta_0 + \beta x_\alpha, \sigma^2)$  的随机变量。 $(1.1)$  式就是一元线性回归的数学模型。在例 1 中  $N = 5$ 。

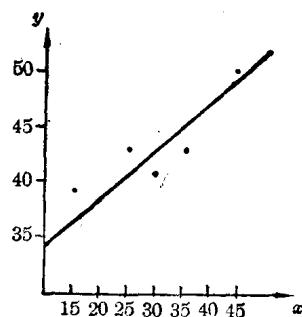


图 1.3

### §1.3 参数 $\beta_0, \beta$ 的最小二乘估计

我们采用最小二乘法来估计 $(1.1)$ 式中的参数  $\beta_0, \beta$ 。

设  $b_0$  和  $b$  分别是参数  $\beta_0$  和  $\beta$  的最小二乘估计，于是得到一元线性回归的回归方程

$$\hat{y} = b_0 + bx. \quad (1.2)$$

$b_0, b$  又叫做回归方程的回归系数。对于每一个  $x_\alpha$ ，由方程 $(1.2)$  可以确定一个回归值  $\hat{y}_\alpha = b_0 + bx_\alpha$ 。这个回归值  $\hat{y}_\alpha$  与实际观察值  $y_\alpha$  之差  $y_\alpha - \hat{y}_\alpha = y_\alpha - b_0 - bx_\alpha$ ，刻划了  $y_\alpha$  与回归直线  $\hat{y} = b_0 + bx$  的偏离程度。一个很自然的想法是，对于所有的  $x_\alpha$ ，若  $\hat{y}_\alpha$  与  $y_\alpha$  的偏离越小，则就认为直线和所有的试验点拟合得越好。显然，全部观

察值  $y_\alpha$  与回归值  $\hat{y}_\alpha$  的偏离平方和

$$Q(b_0, b) = \sum_{\alpha=1}^N (y_\alpha - \hat{y}_\alpha)^2 = \sum_{\alpha=1}^N (y_\alpha - b_0 - bx_\alpha)^2 \quad (1.3)$$

刻画了全部观察值与回归直线的偏离程度。所谓最小二乘法，就是使得

$$Q(b_0, b) = \text{最小}$$

的一种确定  $b_0$  和  $b$  的方法。因此，用最小二乘法配出的直线  $\hat{y} = b_0 + bx$  是这样一条直线，它和点  $(x_\alpha, y_\alpha)$ ,  $\alpha = 1, 2, \dots, N$  的偏离是一切直线中最小的。由于  $Q(b_0, b)$  是  $b_0$  和  $b$  的二次函数，又是非负的，所以它的最小值总是存在的。根据微积分学中的极值原理，要求的估计值  $b_0$  和  $b$  是下列方程组的解：

$$\begin{cases} \frac{\partial Q}{\partial b_0} = -2 \sum_{\alpha} (y_\alpha - b_0 - bx_\alpha) = 0, \\ \frac{\partial Q}{\partial b} = -2 \sum_{\alpha} (y_\alpha - b_0 - bx_\alpha) x_\alpha = 0. \end{cases} \quad (1.4)$$

其中和式  $\sum_{\alpha}$  表示对  $\alpha$  从 1 到  $N$  求和，今后不再说明。方程组 (1.4) 称为正规方程组，它还可以写成如下的形式：

$$\begin{cases} \sum_{\alpha} (y_\alpha - \hat{y}_\alpha) = 0, \\ \sum_{\alpha} (y_\alpha - \hat{y}_\alpha) x_\alpha = 0. \end{cases} \quad (1.5)$$

解正规方程组 (1.4)，得

$$\begin{cases} b_0 = \bar{y} - b\bar{x}, \\ b = \frac{\sum_{\alpha} x_\alpha y_\alpha - \frac{1}{N} (\sum_{\alpha} x_\alpha) (\sum_{\alpha} y_\alpha)}{\sum_{\alpha} x_\alpha^2 - \frac{1}{N} (\sum_{\alpha} x_\alpha)^2}. \end{cases} \quad (1.6)$$

其中  $\bar{x} = \frac{1}{N} \sum_{\alpha} x_\alpha$ ,  $\bar{y} = \frac{1}{N} \sum_{\alpha} y_\alpha$ . 假如把  $b_0 = \bar{y} - b\bar{x}$  代入 (1.2)，可得回归方程的另一形式

$$\hat{y} - \bar{y} = b(x - \bar{x}).$$

由此可见，回归直线 (1.2) 是通过点  $(\bar{x}, \bar{y})$  的，明确这一点，对回归直线的作图是有帮助的。

由公式(1.6)求回归方程的具体计算,通常是列表进行的.

例1的计算在表1.1和表1.2中进行.表1.2中

$$l_{xx} = \sum_{\alpha} (x_{\alpha} - \bar{x})^2 = \sum_{\alpha} x_{\alpha}^2 - \frac{1}{N} (\sum_{\alpha} x_{\alpha})^2,$$

$$l_{yy} = \sum_{\alpha} (y_{\alpha} - \bar{y})^2 = \sum_{\alpha} y_{\alpha}^2 - \frac{1}{N} (\sum_{\alpha} y_{\alpha})^2,$$

$$l_{xy} = \sum_{\alpha} (x_{\alpha} - \bar{x})(y_{\alpha} - \bar{y}) = \sum_{\alpha} x_{\alpha} y_{\alpha} - \frac{1}{N} (\sum_{\alpha} x_{\alpha})(\sum_{\alpha} y_{\alpha}).$$

其中  $l_{yy}$  的值在计算回归方程时并不需要,但在进一步分析中要经常用到,因此也顺便计算出来.

表1.1 回归直线方程计算表(I)

| 编 号      | $x$   | $y$   | $x^2$   | $y^2$   | $xy$    |
|----------|-------|-------|---------|---------|---------|
| 1        | 15.0  | 39.4  | 225.00  | 1552.36 | 591.00  |
| 2        | 25.8  | 42.9  | 665.64  | 1840.41 | 1106.82 |
| 3        | 30.0  | 41.0  | 900.00  | 1681.00 | 1230.00 |
| 4        | 36.6  | 43.1  | 1339.56 | 1857.61 | 1577.46 |
| 5        | 44.4  | 49.2  | 1971.36 | 2420.64 | 2184.48 |
| $\Sigma$ | 151.8 | 215.6 | 5101.56 | 9352.02 | 6689.76 |

表1.2 回归直线方程计算表(II)

|                                    |  |  |
|------------------------------------|--|--|
| $\sum x = 151.8$                   | $\sum y = 215.6$   | $N = 5$                                  |
| $\bar{x} = 30.36$                  | $\bar{y} = 43.12$  |  |
| $\sum x^2 = 5101.56$               | $\sum y^2 = 9352.02$   | $\sum xy = 6689.76$                      |
| $\frac{1}{N} (\sum x)^2 = 4608.64$ | $\frac{1}{N} (\sum y)^2 = 9296.67$                             | $\frac{1}{N} (\sum x)(\sum y) = 6545.61$ |
| $l_{xx} = 492.92$                  | $l_{yy} = 55.35$   | $l_{xy} = 144.15$                        |
|                                    | $b = \frac{l_{xy}}{l_{xx}} = \frac{144.15}{492.92} = 0.29$     |  |
|                                    | $b_0 = \bar{y} - b\bar{x} = 43.12 - 0.29 \times 30.36 = 34.32$ |  |
|                                    | $\hat{y} = 34.32 + 0.29x$                                      | (1.7)                                    |

下面,我们进一步研究最小二乘估计  $b_0$ ,  $b$  的统计性质.

因为  $y_{\alpha}$ ,  $\alpha = 1, 2, \dots, N$  是  $N$  个相互独立的随机变量,且

$$E(y_{\alpha}) = E(\beta_0 + \beta x_{\alpha} + \varepsilon_{\alpha}) = \beta_0 + \beta x_{\alpha},$$

所以它们算术平均数的平均值

$$E(\bar{y}) = E\left(\frac{1}{N} \sum_{\alpha} y_{\alpha}\right) = \frac{1}{N} \sum_{\alpha} E y_{\alpha} = \beta_0 + \beta \bar{x}.$$

最小二乘估计  $b_0, b$  是诸  $y_{\alpha}$  的线性函数，因而它们也是正态随机变量。

由

$$E(y_{\alpha}) = \beta_0 + \beta x_{\alpha},$$

$$E(\bar{y}) = \beta_0 + \beta \bar{x},$$

可得  $b, b_0$  的平均值：

$$\begin{aligned} E(b) &= \frac{\sum_{\alpha} (x_{\alpha} - \bar{x}) E(y_{\alpha} - \bar{y})}{\sum_{\alpha} (x_{\alpha} - \bar{x})^2} \\ &= \frac{\sum_{\alpha} (x_{\alpha} - \bar{x}) [(\beta_0 + \beta x_{\alpha}) - (\beta_0 + \beta \bar{x})]}{\sum_{\alpha} (x_{\alpha} - \bar{x})^2} \\ &= \beta, \end{aligned} \tag{1.8}$$

$$E(b_0) = E(\bar{y} - b \bar{x}) = (\beta_0 + \beta \bar{x}) - \beta \bar{x} = \beta_0. \tag{1.9}$$

所以  $b_0, b$  分别是  $\beta_0, \beta$  的无偏估计。这是最小二乘估计的一个重要性质。

由此可以推出

$$E(\hat{y}) = E(b_0 + bx) = \beta_0 + \beta x = E(y).$$

这表明  $\hat{y}$  是  $E(y)$  的无偏估计，即回归值  $\hat{y}$  可看作是某一点实际观察值  $y$  的平均值。

现在来计算  $b_0, b$  的方差。注意到

$$\begin{aligned} b &= \frac{\sum_{\alpha} (x_{\alpha} - \bar{x})(y_{\alpha} - \bar{y})}{\sum_{\alpha} (x_{\alpha} - \bar{x})^2} = \frac{\sum_{\alpha} (x_{\alpha} - \bar{x}) y_{\alpha} - \bar{y} \sum_{\alpha} (x_{\alpha} - \bar{x})}{\sum_{\alpha} (x_{\alpha} - \bar{x})^2} \\ &= \frac{\sum_{\alpha} (x_{\alpha} - \bar{x}) y_{\alpha}}{\sum_{\alpha} (x_{\alpha} - \bar{x})^2} - \sum_{\alpha} \frac{(x_{\alpha} - \bar{x})}{\sum_{\alpha} (x_{\alpha} - \bar{x})^2} y_{\alpha}, \end{aligned}$$

由

$$D(y_{\alpha}) = D(\beta_0 + \beta x_{\alpha} + \varepsilon_{\alpha}) = \sigma^2,$$

立即得

$$\begin{aligned} D(b) &= \sum_{\alpha} \left[ -\frac{x_{\alpha} - \bar{x}}{\sum_{\alpha} (x_{\alpha} - \bar{x})^2} \right]^2 D(y_{\alpha}) = \frac{\sum_{\alpha} (x_{\alpha} - \bar{x})^2}{[\sum_{\alpha} (x_{\alpha} - \bar{x})^2]^2} \sigma^2 \\ &= \frac{\sigma^2}{\sum_{\alpha} (x_{\alpha} - \bar{x})^2}. \end{aligned} \quad (1.10)$$

大家知道，方差的大小表示随机变量取值波动的大小。式(1.10)表明，回归系数  $b$  的波动大小不仅与误差的方差  $\sigma^2$  有关，而且还取决于观测数据中自变量  $x$  波动的程度。如果  $x$  值波动较大（即比较分散），则  $b$  的波动就较小，也就是估计比较精确。反之，若原始数据  $x$  是在一个较小的范围内取得的，则  $\beta$  的估计就不会精确。这些对安排试验有一定的指导意义。

类似地可以求得估计量  $b_0$  的方差。因为

$$b_0 = \bar{y} - b\bar{x} = \sum_{\alpha} \left[ \frac{1}{N} - \frac{\bar{x}(x_{\alpha} - \bar{x})}{\sum_{\alpha} (x_{\alpha} - \bar{x})^2} \right] y_{\alpha},$$

故有

$$D(b_0) = \sigma^2 \left[ \frac{1}{N} + \frac{\bar{x}^2}{\sum_{\alpha} (x_{\alpha} - \bar{x})^2} \right]. \quad (1.11)$$

由此可知，回归系数  $b_0$  的方差不仅与  $\sigma$  和  $x$  的波动大小有关，而且还同观察数据的个数  $N$  有关。数据越多，且  $x$  值越分散，估计量  $b_0$  就越精确。

## § 1.4 回归方程的显著性检验

回归方程(1.2)是求出来了，但它是否基本上符合变量  $y$  与  $x$  之间的客观规律呢？用它来根据自变量  $x$  的值预报因变量  $y$  的值，效果如何？这就需要对变量  $y$  与  $x$  之间是否是线性关系，进行统计检验。

我们知道，观察值  $y_1, y_2, \dots, y_N$  之间的差异，是由两个方面的原因引起的：(1)自变量  $x$  取值的不同；(2)其他因素（包括试验误差）的影响。为了检验这两个方面的影响哪一个是主要的，首先

就必须把它们所引起的差异，从  $y$  总的差异中分解出来。

$N$  个观察值之间的差异，可用观察值  $y_a$  与其算术平均值  $\bar{y}$  的偏差平方和来表示，称为总的偏差平方和，记作

$$S_{\text{总}} = \sum_a (y_a - \bar{y})^2 = l_{yy}. \quad (1.12)$$

因为

$$\begin{aligned} S_{\text{总}} &= \sum_a (y_a - \bar{y})^2 = \sum_a [(y_a - \hat{y}_a) + (\hat{y}_a - \bar{y})]^2 \\ &= \sum_a (\hat{y}_a - \bar{y})^2 + \sum_a (y_a - \hat{y}_a)^2 + 2 \sum_a (y_a - \hat{y}_a)(\hat{y}_a - \bar{y}), \end{aligned}$$

由(1.5)可知，交叉项

$$\begin{aligned} \sum_a (y_a - \hat{y}_a)(\hat{y}_a - \bar{y}) &= \sum_a (y_a - \hat{y}_a)(b_0 + b x_a - \bar{y}) \\ &= (b_0 - \bar{y}) \sum_a (y_a - \hat{y}_a) + b \sum_a (y_a - \hat{y}_a) x_a \\ &= 0, \end{aligned}$$

于是我们获得了总的偏差平方和的分解公式

$$\sum_a (y_a - \bar{y})^2 = \sum_a (\hat{y}_a - \bar{y})^2 + \sum_a (y_a - \hat{y}_a)^2, \quad (1.13)$$

或者写成

$$S_{\text{总}} = S_{\text{回}} + S_{\text{剩}}. \quad (1.14)$$

其中

$$S_{\text{回}} = \sum_a (\hat{y}_a - \bar{y})^2 \quad (1.15)$$

称为回归平方和，它是由自变量  $x$  的变化而引起的，它的大小（在与误差相比的意义下）反映了自变量  $x$  的重要程度；

$$S_{\text{剩}} = \sum_a (y_a - \hat{y}_a)^2 \quad (1.16)$$

称为剩余平方和，它是由试验误差以及其他未加控制的因素引起的，它的大小反映了试验误差及其他因素对试验结果的影响。这样一来，通过平方和分解公式(1.13)，就把对  $N$  个观察值的两种影响从数量上基本区分开来了。

现在我们回到统计检验问题上来。如果变量  $y$  与  $x$  之间无线性关系，那么模型(1.1)中一次项系数  $\beta=0$ ；反之， $\beta \neq 0$ 。所以，要检验两个变量之间是否有线性关系，归根结蒂就是要检验  $\beta$  是

否为零。而这一点可以通过比较  $S_{\text{总}}$  与  $S_{\text{固}}$  来实现。为此我们先引进下列分解定理，然后建立检验“ $\beta=0$ ”这个假设的统计量。

**分解定理<sup>[4](\*)</sup>** 若  $Q=Q_1+Q_2+\cdots+Q_k$ ，其中  $Q$  为  $\chi^2(f)$  变量， $Q_i$  为正态随机变量的平方和，它的自由度为  $f_i$ ，则  $Q_i$  间相互独立，且服从  $\chi^2(f_i)$  分布的充要条件是

$$f=f_1+f_2+\cdots+f_k. \quad (1.17)$$

如今的平方和分解公式(1.13)相当于  $Q=Q_1+Q_2$ ，假如我们能证明  $S_{\text{总}}/\sigma^2$  为  $\chi^2(f)$  变量， $S_{\text{固}}/\sigma^2$ ， $S_{\text{固}}/\sigma^2$  分别为  $\chi^2(f_1)$ ， $\chi^2(f_2)$  变量，并且  $f=f_1+f_2$ ，那末运用分解定理就可证得  $S_{\text{总}}$  与  $S_{\text{固}}$  相互独立。

首先，我们证明

$$Q=\frac{1}{\sigma^2} S_{\text{总}}=\frac{1}{\sigma^2} \sum_{\alpha} (y_{\alpha}-\bar{y})^2 \sim \chi^2(N-1)^{(**)}.$$

事实上

$$\begin{aligned} Q &= \sum_{\alpha} \left( \frac{y_{\alpha}-\bar{y}}{\sigma} \right)^2 = \sum_{\alpha} \left( \frac{y_{\alpha}-\beta_0-\bar{y}+\beta_0}{\sigma} \right)^2 \\ &= \sum_{\alpha} (y_{\alpha}^*-\bar{y}^*)^2, \end{aligned}$$

其中  $y_{\alpha}^*=\frac{y_{\alpha}-\beta_0}{\sigma}$ ，在假设“ $\beta=0$ ”成立下， $y_{\alpha}^*$  是  $y_{\alpha}$  的标准化变量，即  $y_{\alpha}^* \sim N(0, 1)$ 。我们对随机变量  $y_1^*, y_2^*, \dots, y_N^*$  进行如下正交线性变换：

$$\begin{pmatrix} z_1 \\ z_2 \\ z_3 \\ \vdots \\ z_N \end{pmatrix} = \begin{pmatrix} \frac{1}{\sqrt{N}} & \frac{1}{\sqrt{N}} & \frac{1}{\sqrt{N}} & \frac{1}{\sqrt{N}} & \cdots & \frac{1}{\sqrt{N}} \\ \frac{1}{\sqrt{1 \times 2}} & \frac{-1}{\sqrt{1 \times 2}} & 0 & 0 & \cdots & 0 \\ \frac{1}{\sqrt{2 \times 3}} & \frac{1}{\sqrt{2 \times 3}} & \frac{-2}{\sqrt{2 \times 3}} & 0 & \cdots & 0 \\ \cdots & \cdots & \cdots & \cdots & \cdots & \cdots \\ \frac{1}{\sqrt{N(N-1)}} & \frac{1}{\sqrt{N(N-1)}} & \frac{1}{\sqrt{N(N-1)}} & \frac{1}{\sqrt{N(N-1)}} & \cdots & \frac{-(N-1)}{\sqrt{N(N-1)}} \end{pmatrix} \begin{pmatrix} y_1^* \\ y_2^* \\ y_3^* \\ \vdots \\ y_N^* \end{pmatrix}$$

(\*) 上角码<sup>[4]</sup>，意为参阅书末参考文献[4]，下同。

(\*\*)  $Q \sim \chi^2(N-1)$  表示随机变量  $Q$  服从自由度为  $N-1$  的  $\chi^2$  分布，以下类同。