

“回归分析”与“时间序列分析”

中国地震局地球物理研究所图书馆

顾功叙 赠

专家赠书/刊纪念

统计预报短训班

一九七三年四月

回 归 分 析

目 录

§ 1、回归模型	(1· 1— 3)
§ 2、回归系数的最小二乘估计	(1· 3—10)
§ 3、回归问题的方差分析与统计检验	
(1)回归效果的分析——方差分析	(1· 10—13)
(2) β_i 与 Y 的估计精度——区间估计	(1· 13—15)
(3)各变量的重要性—— β_i 的显著性检验	(1· 15—18)
§ 4、逐步回归的基本思想——正交筛选原则	
(1)逐步回归的基本思想	(1· 18—19)
(2)两个衡量变量重要性的指标	(1· 19—21)
(3)应用正交变换原则的逐步回归 ——正交筛选法	(1· 21—29)
§ 5、双重检验的逐步回归	
(1)双重检验问题的提出	(1· 29—30)
(2)双重检验的逐步回归	(1· 30—34)
(3)双重检验的逐步回归的计算步骤	(1· 34—36)
§ 6、逐步回归的一些问题	(1· 36—38)
附录: 计算框图与数值例子	(1· 39—46)
参考文献	(1· 47—48)

回归分析

§1. 回归模型

已知量 y 随着 m 个自变量 x_1, x_2, \dots, x_m 的变化而变化, 要有如下线性关系:

$$y = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \dots + \beta_m x_m + \varepsilon \quad (1.1)$$

这一式子通常称为 y 关于 x_1, x_2, \dots, x_m 的“回归方程”, 式中的 $\beta_0, \beta_1, \beta_2, \dots, \beta_m$ 为待定参数, 也称“回归系数”, ε 为随机变量, 也称“残差”, 这是因为 ε 可以看作关于 y 中无法用 x_1, x_2, \dots, x_m 表示的各种因素或随机因素组成的残差。

在实际问题中, 我们常假数量 x_1, x_2, \dots, x_m, y 有 N 组观测数据:

$$(x_{t1}, x_{t2}, \dots, x_{tm}, y_t) \quad (t=1, 2, \dots, N)$$

它们满足回归方程, 即满足下面的 N 个式子:

$$y_t = \beta_0 + \beta_1 x_{t1} + \beta_2 x_{t2} + \dots + \beta_m x_{tm} + \varepsilon_t \quad (1.2)$$
$$(t=1, 2, \dots, N)$$

这时随机变量 ε_t ($t=1, 2, \dots, N$) 表示各次观测的残差, 通常它们满足如下假设:

① 观测残差没有系统性的倚倚现象。即 ε_t 的数学期望全为零

$$E(\varepsilon_t) = 0 \quad (t=1, 2, \dots, N) \quad (1.3)$$

② 这些观测残差相互独立的, 并且具有相同的方差。因此 ε_t 的协方差可表为

$$V(\varepsilon_t, \varepsilon_\tau) = \delta_{t\tau} \sigma^2 \quad (t, \tau=1, 2, \dots, N) \quad (1.4)$$

这里的 δ_{tc} 是 Kronecker 符号, 当 $t=c$ 时, 它等于 1, $t \neq c$ 时等于 0。 σ^2 是公共方差。当 $t=c$ 时, $V(E_t, E_t)$ 就是 E_t 的方差, 可简记为 $V(E_t)$ 。

③ 观测误差服从正态分布。有了这个假设, 就可以顺利地应用一些统计检验方法。

如果用符号 $N(\mu, \sigma^2)$ 表示数学期望为 μ , 方差为 σ^2 的正态分布, 则上述三个假设可简单地概括为“误差 E_t 相互独立地服从正态分布 $N(0, \sigma^2)$ ”。

在回归分析中, 需要解决的主要问题是:

① 根据给出的 N 组观测数据, 确定各回归系数 β_i 的估值 b_i ($i=0, 1, 2, \dots, m$)。同时对各 b_i 作统计检验, 以便指出这些估值的可靠性。

② 对自变量 (X_1, X_2, \dots, X_m) 的任一观测值, 求被变量 y 的取值, 并且指出预报的精度。这就需要给出误差 E 的未知方差 σ^2 的估值。

回归分析在一般数据处理, 经验公式, 曲线拟合, 以及各类预报问题中都有广泛的应用。但是随着应用领域不同, 回归分析的重点也常常有差别: 例如在一般经验公式与曲线拟合中, 人们往往只注意求出各 b_i 值; 在各类预报问题中, 人们显然更关心回归方程在预报中的精度 (主要是 σ^2 的大小)。此外为建立较好的回归方程, 人们往往需要对所有可供使用的自变量作统计检验 (即对 b_i 作显著性检验), 以判别哪些是重要的, 哪些是不重要的, 从而为变量的筛选提供依据。

本文将从应用的角度简要地讨论上述各类问题, 重点放在预报问题中各预报因子的筛选上, 特别将介绍在筛选变量中已有广泛应用的“正交筛选法”和“逐步回归法”。在内容安排上, 为突出重点, 我们将把部分内容以补充“说明”形式给出在各节之后。

本节补充说明:

说明 1.1: 在上述回归模型中, 自变量 X_i 被假定为确定性变量, 而观测误差 ε_t 是相互独立地服从正态分布 $N(0, \sigma^2)$, 因此 y_t 也是相互独立地服从正态分布的随机变量, 其数学期望和方差分别是:

$$E(y_t) = \beta_0 + \beta_1 X_{t1} + \beta_2 X_{t2} + \dots + \beta_m X_{tm} \quad (1.5)$$

$$V(y_t) = \sigma^2 \quad (1.6)$$

(1.5) 或

$$E(y) = \beta_0 + \beta_1 X_1 + \beta_2 X_2 + \dots + \beta_m X_m \quad (1.7)$$

可以代替 (1.2) 或 (1.1), 它们都称为回归方程。从这一形式的回归方程开始讨论, 就可以不涉及误差 ε 的问题, 很多统计书籍都这样做。

说明 1.2: 上述回归模型, 严格地说, 应称为“多元线性正态回归模型”。所谓“多元”, 是指 y 依赖于不止一个的自变量 ($m > 1$); 所谓“线性”, 是指回归方程是关于自变量 X_i 的线性函数; 所谓“正态”, 是指 ε 或 y 是正态分布的随机变量。今后我们可以看到, 线性的含义可以推广为只是关于 β_i 的线性函数, 而不一定是 X_i 的线性函数。

说明 1.3: 在一般回归分析中, 我们总要假设观测次数比回归系数的个数大 (即 $N > m + 1$), 而且任一自变量不能用其它自变量线性表出。在正交筛选和逐步回归中, 后一假设不是必需的。

§2. 回归系数的最小二乘估计

假设我们有某一种方法可以得到 b_i 的估值, 则 y 的观测值可表为

$$y_t = b_0 + b_1 X_{t1} + b_2 X_{t2} + \dots + b_m X_{tm} + \varepsilon_t \quad (2.1)$$

上式也称回归方程，或加“经验”二字，用以与(1.1)的“理论”形式相区别，这里 e_t 是误差 E_t 的估值，常称为“残差”或“剩余”。令 y_t^* 为 y_t 的估值，即

$$y_t^* = b_0 + b_1 \chi_{t1} + b_2 \chi_{t2} + \dots + b_m \chi_{tm} \quad (2.2)$$

$$e_t = y_t - y_t^* \quad (2.3)$$

b_i 的最优估值可以用最小二乘法得到，即 b_i 的确定将使残差平方和 Q 达到最小

$$Q = \sum_t e_t^2 = \sum_t (y_t - y_t^*)^2 \quad (\text{致})$$

$$= \sum_t (y_t - (b_0 + b_1 \chi_{t1} + b_2 \chi_{t2} + \dots + b_m \chi_{tm}))^2 \quad (2.4)$$

由微分分析的极值原理，知各 b_i 满足如下方程组

$$\begin{cases} \frac{\partial Q}{\partial b_0} = 0 \\ \frac{\partial Q}{\partial b_1} = 0 \\ \vdots \\ \frac{\partial Q}{\partial b_m} = 0 \end{cases}$$

即 Q 对各 b_i 的确微商全为零。

$$\text{由 } \frac{\partial Q}{\partial b_0} = 0 \quad \text{得}$$

$$\sum_t (y_t - (b_0 + b_1 \chi_{t1} + b_2 \chi_{t2} + \dots + b_m \chi_{tm}))(-1) = 0 \quad (2.6)$$

$$\text{或 } \sum_t y_t - (N b_0 + b_1 \sum_t \chi_{t1} + b_2 \sum_t \chi_{t2} + \dots + b_m \sum_t \chi_{tm}) = 0$$

(致): 求和号 \sum 下常添有下标 t 或 i ，其中 t (或 i) 表示从 1 到 N 的和，即 $\sum_{t=1}^N$ ； i (或 j) 表示从 1 到 m 的和，即 $\sum_{i=1}^m$

把观测值的平均值为

$$\begin{cases} \bar{x}_i = \frac{1}{N} \sum_x x_{ti} & (i=1, 2, \dots, m) \\ \bar{y} = \frac{1}{N} \sum_x y_x \end{cases} \quad (2.7)$$

(2.6)可表为

$$b_0 = \bar{y} - (b_1 \bar{x}_1 + b_2 \bar{x}_2 + \dots + b_m \bar{x}_m) \quad (2.8)$$

以 (2.8) 代入 (2.5), 且令

$$\begin{cases} x'_{ti} = x_{ti} - \bar{x}_i & (i=1, 2, \dots, m, \\ y'_x = y_x - \bar{y} & (x=1, 2, \dots, N) \end{cases} \quad (2.9)$$

将 $Q = \sum_x \{y'_x - (b_1 x'_{t1} + b_2 x'_{t2} + \dots + b_m x'_{tm})\}^2 \quad (2.10)$

这时 (2.5) 可化为

$$\begin{cases} \frac{\partial Q}{\partial b_1} = \frac{1}{2} \sum_x \{y'_x - (b_1 x'_{t1} + b_2 x'_{t2} + \dots + b_m x'_{tm})\} (-x'_{t1}) = 0 \\ \frac{\partial Q}{\partial b_2} = \frac{1}{2} \sum_x \{y'_x - (b_1 x'_{t1} + b_2 x'_{t2} + \dots + b_m x'_{tm})\} (-x'_{t2}) = 0 \\ \dots \\ \frac{\partial Q}{\partial b_m} = \frac{1}{2} \sum_x \{y'_x - (b_1 x'_{t1} + b_2 x'_{t2} + \dots + b_m x'_{tm})\} (-x'_{tm}) = 0 \end{cases} \quad (2.11)$$

经整理并引入如下记号

$$\begin{cases} S_{ij} = S_{ji} = \sum_x x'_{ti} x'_{tj} = \sum_x (x_{ti} - \bar{x}_i)(x_{tj} - \bar{x}_j) \\ S_{iy} = \sum_x x'_{ti} y'_x = \sum_x (x_{ti} - \bar{x}_i)(y_x - \bar{y}) \end{cases} \quad (2.12)$$

(i, j=1, 2, \dots, m)

(2.11)可化为

$$\begin{cases} S_{11} b_1 + S_{12} b_2 + \dots + S_{1m} b_m = S_{1y} \\ S_{21} b_1 + S_{22} b_2 + \dots + S_{2m} b_m = S_{2y} \\ \dots \\ S_{m1} b_1 + S_{m2} b_2 + \dots + S_{mm} b_m = S_{my} \end{cases} \quad (2.13)$$

上式常称为“正规方程”，当观测数据变化的情况下， S_{ij} 与 S_{ji} 可由 (2.12) 一一算出，因此正规方程是以 b_1, b_2, \dots, b_m 为未知数的 m 个联立方程式，例如可用消去法求解（详见 §5）。解出 b_1, b_2, \dots, b_m 后，由 (2.8) 便可将 b_0 。

本节补充说明：

[说明 2.1]： 在上述回归模型中， b_0 与 b_1, b_2, \dots, b_m 出现的形式不一致，给处理带来麻烦。如果引入一个形式变量 X_0 ，其观测值恒为 1，即 $X_{t0} \equiv 1$ ，则回归方程可写作

$$Y_t = b_0 X_{t0} + b_1 X_{t1} + \dots + b_m X_{tm} + e_t \quad (2.14)$$

类似 (2.11) ~ (2.13) 的推导，可得

$$\frac{\partial Q}{\partial b_i} = \frac{1}{2} \sum_x [Y_t - (b_0 X_{t0} + b_1 X_{t1} + \dots + b_m X_{tm})] (-X_{ti}) = 0 \quad (i=0, 1, 2, \dots, m) \quad (2.15)$$

$$\begin{cases} a_{ij} = a_{ji} = \sum_x X_{ti} X_{tj} \\ a_{iy} = \sum_x X_{ti} Y_t \end{cases} \quad (i, j = 0, 1, 2, \dots, m) \quad (2.16)$$

可得正规方程

$$\begin{cases} a_{00} b_0 + a_{01} b_1 + \dots + a_{0m} b_m = a_{0y} \\ a_{10} b_0 + a_{11} b_1 + \dots + a_{1m} b_m = a_{1y} \\ \dots \dots \dots \dots \dots \dots \\ a_{m0} b_0 + a_{m1} b_1 + \dots + a_{mm} b_m = a_{my} \end{cases} \quad (2.17)$$

上式是以 b_0, b_1, \dots, b_m 为未知数的 $m+1$ 个联立方程式，解之便得全部 b_i 。虽然，这样的处理方法要比本节正文中更为简便，不少统计著作就是这样做。

[说明 2.2]： 当 $b_0 = \bar{y} - (b_1 \bar{X}_1 + b_2 \bar{X}_2 + \dots + b_m \bar{X}_m)$ (见 (2.8)) 可知当各自变量的观测数据的平均值 \bar{X}_i, \bar{y} 全为零，则 $b_0 = 0$ ，这时可以认为回归方程中根本没有 b_0 这一项，即

$$Y_t = b_1 X_{t1} + b_2 X_{t2} + \dots + b_m X_{tm} + e_t \quad (2.18)$$

在实际问题中 \bar{X}_i, \bar{y} 全为零是不可能的，但是可以通过变量变换

$$X'_{ti} = X_{ti} - \bar{X}_i \quad y'_t = y_t - \bar{y}$$

使新的变量 X'_i, y' 的平均值全为零，于是仍然可以得到不含 b_0 的回归方程

$$y'_t = b'_1 X'_{t1} + b'_2 X'_{t2} + \dots + b'_m X'_{tm} + e_t \quad (2.19)$$

这就是正文中的处理方法，从正文的推导来看，(用(2.8)~(2.10))这里的 b'_1, b'_2, \dots, b'_m 分别等于原回归模型(有 b_0 的)

$$y_t = b_0 + b_1 X_{t1} + b_2 X_{t2} + \dots + b_m X_{tm} + e_t$$

中的 b_1, b_2, \dots, b_m 。

[说明2.3]: 本节正文中把 b_0 与 b_1, b_2, \dots, b_m 分开来处理的方法，虽不及[说明2.1]的简便，但可以附带得到一些重要信息，因此也有很多统计著作这样处理。事实上，各变量之间的协方差(这时暂把各 X_{ti} 看作是随机变量)可用 $\frac{1}{N-1} S_{ij}$ 估计，而相关系数 Y_{ij} 可表为

$$Y_{ij} = \frac{S_{ij}}{\sqrt{S_{ii}} \sqrt{S_{jj}}} = \frac{\sum_t (X_{ti} - \bar{X}_i)(X_{tj} - \bar{X}_j)}{\sqrt{\sum_t (X_{ti} - \bar{X}_i)^2} \sqrt{\sum_t (X_{tj} - \bar{X}_j)^2}} \quad (2.20)$$

在实际计算中，故感到 Y_{ij} 是冗繁的，各 Y_{ij} 的误差要比 S_{ij} 小，如果用 Y_{ij} 代替 S_{ij} 来求解正规方程的话，捨入误差一般要小一些，在电子计算机上计算时更是这样。要注意的是，这样得到的解 b'_i 与原来的解 b_i 不同，其关系式是

$$b'_i = \frac{\sqrt{S_{ii}}}{\sqrt{S_{yy}}} b_i \quad (i=1, 2, \dots, m) \quad (2.21)$$

其中

$$S_{yy} = \sum_t (y'_t)^2 = \sum_t (y_t - \bar{y})^2 \quad (2.22)$$

这一事实可作如下结论：

(1) 在正规方程 (2.13) 中, 第 i 个方程的各方项同除以 $\sqrt{S_{ii}}\sqrt{S_{yy}}$ 其解不变, 这时有

$$\frac{S_{i1}}{\sqrt{S_{ii}}\sqrt{S_{yy}}} b_1 + \frac{S_{i2}}{\sqrt{S_{ii}}\sqrt{S_{yy}}} b_2 + \dots + \frac{S_{im}}{\sqrt{S_{ii}}\sqrt{S_{yy}}} b_m = \frac{S_{iy}}{\sqrt{S_{ii}}\sqrt{S_{yy}}} \quad (i=1, 2, \dots, m)$$

(2) 上述方程中的第 i 个方程的分子分母同乘以 $\sqrt{S_{ii}}$, 其解不变, 这时有

$$\frac{S_{i1}}{\sqrt{S_{ii}}\sqrt{S_{11}}} \left(\frac{\sqrt{S_{11}}}{\sqrt{S_{yy}}} b_1 \right) + \frac{S_{i2}}{\sqrt{S_{ii}}\sqrt{S_{22}}} \left(\frac{\sqrt{S_{22}}}{\sqrt{S_{yy}}} b_2 \right) + \dots$$

$$\dots + \frac{S_{im}}{\sqrt{S_{ii}}\sqrt{S_{mm}}} \left(\frac{\sqrt{S_{mm}}}{\sqrt{S_{yy}}} b_m \right) = \frac{S_{iy}}{\sqrt{S_{ii}}\sqrt{S_{yy}}} \quad (i=1, 2, \dots, m)$$

由 (2.20) 和 (2.21) 得

$$Y_{i1} b_1 + Y_{i2} b_2 + \dots + Y_{im} b_m = Y_{iy} \quad (2.23)$$

$$(i=1, 2, \dots, m)$$

或

$$\begin{cases} Y_{11} b_1 + Y_{12} b_2 + \dots + Y_{1m} b_m = Y_{1y} \\ Y_{21} b_1 + Y_{22} b_2 + \dots + Y_{2m} b_m = Y_{2y} \\ \dots \\ Y_{m1} b_1 + Y_{m2} b_2 + \dots + Y_{mm} b_m = Y_{my} \end{cases} \quad (2.24)$$

(说明 2.4): 关于 \bar{X}_i 与 S_{ij} 的计算问题。

① \bar{X}_i 的计算 (包括 Y_i):

$$\bar{X}_i = \frac{1}{N} \sum_x X_{xi} = \frac{1}{N} \sum_x (X_{xi} - C_i) + C_i \quad (2.25)$$

其中 C_i 是与 x 无关的常量。因此当某 x 的观测值的前几位几乎相同时, 可先减去其公共部分 C_i 再求和, 最后才把 C_i 加回去。这种计算法对于手算非常合适, 可以节省计算时间。对电子计算机来说, 这一算法虽然不能起到节省计算时间的作用, 却可以提高计算精度。这作为第一步, 先按下式算出 C_i 。

$$C_i = \frac{1}{N} \sum_x X_{xi} \quad (2.26)$$

C_i 已是 X_i 的均值, 但可能有较大的捨入误差, 为此作第二步新称

$$\bar{X}_i = \frac{1}{N} \sum_{t=1}^n (X_{ti} - C_i) + C_i \quad (2.27)$$

其中右端均求和项集中了均值 C_i 的新称过程中的主要捨入误差, 因此这样的“二步称法”可以得到更精确的 \bar{X}_i , 特别当观测数据位数较多, 而机口的字长较短时更是如此。[文献12] [文献13]

② S_{ij} 的新称 (包括 S_{iy} 与 S_{44}):

$$\begin{aligned} S_{ij} &= \sum_{t=1}^n (X_{ti} - \bar{X}_i)(X_{tj} - \bar{X}_j) \\ &= \sum_{t=1}^n X_{ti} X_{tj} - N \bar{X}_i \bar{X}_j \end{aligned} \quad (2.28)$$

新称时用第二等式可以节省时间, 但改用电子计算机新称时, 第一等式的精度要比第二等式高。[文献12]、[文献13]

[说明2.5]: 引入矩阵符号

$$S = \begin{pmatrix} S_{11} & S_{12} & \cdots & S_{1m} \\ S_{21} & S_{22} & \cdots & S_{2m} \\ \vdots & \vdots & \ddots & \vdots \\ S_{m1} & S_{m2} & \cdots & S_{mm} \end{pmatrix} \quad b = \begin{pmatrix} b_1 \\ b_2 \\ \vdots \\ b_m \end{pmatrix} \quad g = \begin{pmatrix} S_{1y} \\ S_{2y} \\ \vdots \\ S_{my} \end{pmatrix} \quad (2.29)$$

则正规方程 (2.13) 可表为

$$S b = g \quad (2.30)$$

在 [说明1.3] 的假设下, S 矩阵满秩 (其行列式的值不等于零), b 有唯一解, 这时 S 矩阵的逆矩阵 S^{-1} 存在, 记 S^{-1} 为 C , 则

$$C = \begin{pmatrix} C_{11} & C_{12} & \cdots & C_{1m} \\ C_{21} & C_{22} & \cdots & C_{2m} \\ \vdots & \vdots & \ddots & \vdots \\ C_{m1} & C_{m2} & \cdots & C_{mm} \end{pmatrix} \quad (2.31)$$

b 可表为

$$b = S^{-1}g \quad \text{或} \quad b = Cg$$

即
$$b_i = \sum_j C_{ij} g_j \quad (2.32)$$

[说明 2.6]: 对于经验公式与曲线拟合一类问题, 人们只关心 b_i 的估计, 因此其后的内容 (特别是 §3) 也许是多余的。从估计规范看, 确定 b_i 的过程与解线性矛盾方程组的最小二乘法无异。

§3. 回归问题的方差分析与统计检验

(1) 回归效果的分析——方差分析

y 的总平方和 S_{yy} 可表为

$$\begin{aligned} S_{yy} &= \sum_t (y_t - \bar{y})^2 = \sum_t [(y_t - y_t^*) + (y_t^* - \bar{y})]^2 \\ &= \sum_t (y_t - y_t^*)^2 + \sum_t (y_t^* - \bar{y})^2 \quad (\text{注1}) \end{aligned} \quad (3.1)$$

上式中右端第一部分是前面已定义过的残差平方和 Q , 第二部分称为回归平方和, 记为 U [注2]。

[注1]: 这要因为由 (2.8), (2.9) 有

$$\begin{aligned} & \sum_t (y_t - y_t^*) (y_t^* - \bar{y}) \\ &= \sum_t [y_t - (b_1 x_{t1} + \dots + b_m x_{tm})] (b_1 x_{t1} + \dots + b_m x_{tm}) \\ &= b_1 \sum_t [y_t - (b_1 x_{t1} + \dots + b_m x_{tm})] (x_{t1}) + \dots \\ & \dots + b_m \sum_t [y_t - (b_1 x_{t1} + \dots + b_m x_{tm})] (x_{tm}) = 0 \end{aligned}$$

最后一等号是因为式中的每一个求和项都为零（见(2.11)）。

〔註2〕：配估计值 y_t^* 之均值

$$\bar{y}^* = \frac{1}{N} \sum_t y_t^*$$

由(2.8)易知 $\bar{y}^* = \bar{y}$ ，故

$$U = \sum_t (y_t^* - \bar{y})^2 = \sum_t (y_t^* - \bar{y}^*)^2$$

是回归方程的误差平方和。

即总的误差平方和 S_{yy} 可分解为残差平方和与回归平方和两部分，可记为

$$S_{yy} = Q + U \quad (3.2)$$

Q 的名称可用前面的定义(2.4)或(2.10)进行，当 S_{yy} 已算出的情况下，用下列名称则更方便

$$Q = S_{yy} - \sum_i b_i S_{iy} \quad (3.3)$$

这是因为

$$\begin{aligned} Q &= \sum_t (y_t - y_t^*)^2 = \sum_t (y_t - y_t^*) [(y_t - \bar{y}) - (y_t^* - \bar{y})] \\ &= \sum_t (y_t - y_t^*) (y_t - \bar{y}) - \sum_t (y_t - y_t^*) (y_t^* - \bar{y}) \end{aligned}$$

上式第二部分由上页〔註〕知其为零，故

$$\begin{aligned} Q &= \sum_t (y_t - y_t^*) (y_t - \bar{y}) = \sum_t [(y_t - \bar{y}) - (y_t^* - \bar{y})] (y_t - \bar{y}) \\ &= \sum_t (y_t - \bar{y})^2 - \sum (b_1 X'_{t1} + \dots + b_m X'_{tm}) y'_t \\ &= S_{yy} - (b_1 S_{1y} + b_2 S_{2y} + \dots + b_m S_{my}) \end{aligned}$$

对比(3.2)与(3.3)，易知回归平方和也可记为

$$U = \sum_i b_i S_{iy} \quad (3.4)$$

对给定的观测值而言, S_{yy} 是不变的, 因此 $Q+U$ 是常量, Q 大则 U 小, 反之 Q 小则 U 大, 因此 Q 与 U 都可以用于衡量回归的效果, 但是 Q 或 U 与 S_{yy} 一样, 都没有量纲的, 即与 y 的单位有关, 大小不好判别。为此, 可采用如下定义的无量纲指标:

$$R^2 = \frac{U}{S_{yy}} \quad \text{或} \quad R = \sqrt{1 - \frac{Q}{S_{yy}}} \quad (3.5)$$

回归平方和 U 实质上是回归方程中的全部自变量的“方差贡献”, 因此 R^2 就是这种贡献在总和中所占的比例。 R 称为“复(或全)相关系数”, 它是全部自变量与 y 的相关, 为强调这一点, 有些书籍把 R 记作 $r_{y \cdot 12 \dots m}$ 。另知 $0 \leq R \leq 1$ 。显然, 复相关系数越接近 1, 则回归效果越好。

尽管我们常用 R 作为总的回归效果的一个重要指标, 但是我们要清醒地看到, R 是与回归方程中变量个数 m 以及观测次数 N 有关的。当 m 较大而 N 不那么大时, 常有较大的 R (即较小的 Q)。特别当 $m+1=N$ 时, 即使这 m 个变量与 y 风马牛不相干, 亦必然有 $R=1$ (即 $Q=0$)。这就是说, 在实际问题中, 要注意 m 与 N 的比例适当。一般人认为 N 至少是 m 的 5 (或 10) 倍以上。当然, 这样的要求, 并不是称 R 所特有的, 在整个回归分析中都有此需要。

下面给出另一个衡量回归效果的指标:

$$F = \frac{U/m}{Q/(N-m-1)} \quad (3.6)$$

这是两个方差比, 它放大了 m , N 的作用, 实际比 R 更为合理。重要的是在 $\beta_i \equiv 0$ ($i=1, 2, \dots, m$) 的假设下, 这个统计量服从 F 分布, 自由度分别为 m 与 $N-m-1$ 。因此可用它来检验这 m 个变量的回归效果。例如在给定的显著水平 α 后, 可从 F 分布表中查出相应的临界值 F_α , 当统计的 $F > F_\alpha$, 则各 β_i 全为零的假设不成立, 即回归效果显著。

上述内容，可归纳成一张“方差分析表”

方差分析表

来源	平方和	自由度	方差	F 检验
回归	$U = \sum_x (y_x^* - \bar{y})^2 = \sum_0 b_i S_{iy}$	m	U/m	$F = \frac{U/m}{Q/(N-m-1)}$
残差	$Q = \sum_x (y_x - y_x^*)^2$	$N-m-1$	$Q/(N-m-1)$	
总和	$S_{yy} = \sum_x (y_x - \bar{y})^2$	$N-1$		

(2) β_i 与 y 的估计精度——区间估计

在 §2 中，已给出 β_i 的估计值 b_i ，要给出 y_t 的估计值 y_t^* 。但是这种估计只给出一个值，并没有给出估计值的精度的知识，因此称为点估计。这里将给出 β_i 与 y 的区间估计，由估计区间的长短，便可知道它们的精度。

首先我们指出 b_i 是从正态分布的随机变量，且有

均值 $E(b_i) = \beta_i$ (3.7)

方差 $V(b_i) = C_{ii} \sigma^2$ (3.8)

协方差 $V(b_i, b_j) = C_{ij} \sigma^2$ ($i, j = 1, 2, \dots, m$) (3.9)

这里 $C = (C_{ij})$ 是正规方程的系数矩阵 $S = (S_{ij})$ 的逆矩阵^[注]。 σ^2 是 y_t (或 ϵ_t) 的公共方差。一般 σ^2 未知，我们可以给出它的无偏估计 S_y^2 ，可知：

$$S_y^2 = \frac{Q}{N-m-1} \quad (3.10)$$

有关 b_i 和 S_y 的一些结论的证明，这里不再给出，有兴趣的读者，可从 [文献 2] 的前几页中看到。

下面我们将使用这些结论给出 β_i 与 y 的区间估计。

[注]：注意此处有关 b_i 的结论，虽然没有把 b_0 包括进去（人们往往不太注意 b_0 的性质），但只需作小的修改，即把 (C_{ij}) 理解为 [规范 2.1] 中的 $m+1$ 阶矩阵 (a_{ij}) 的逆，则有关结论对 b_0 也是正确的。

因而 b_i 服从正态分布 $N(\beta_i, C_{ii}\sigma^2)$, 因此

$$z_i = \frac{b_i - \beta_i}{\sqrt{C_{ii}}\sigma}$$

是标准正态量 $N(0, 1)$ 。如果 σ 已知, 则对给定的显著水平 α , 可给出 β_i 的区间估计

$$\left| \frac{b_i - \beta_i}{\sqrt{C_{ii}}\sigma} \right| < N_{\alpha/2} \quad (3.12)$$

或

$$b_i - N_{\alpha/2}\sqrt{C_{ii}}\sigma < \beta_i < b_i + N_{\alpha/2}\sqrt{C_{ii}}\sigma$$

式中 $N_{\alpha/2}$ 为正态分布表上相应于 $\alpha/2$ 的临界值〔注〕。对于典型值 $\alpha=0.05$, $N_{0.025}=1.96$ 。在实际问题中, 一般 σ 是未知的, 这时可用它的估计值 S_y 替代, 但按统计理论,

$$t_i = \frac{b_i - \beta_i}{\sqrt{C_{ii}}S_y} \quad (3.13)$$

已不再是正态分布, 而是 t 分布, 自由度为 $N-m-1$, 因此 β_i 的区间估计为

$$\left| \frac{b_i - \beta_i}{\sqrt{C_{ii}}S_y} \right| < t_{\alpha/2} \quad (3.14)$$

式中 $t_{\alpha/2}$ 为 t 分布表上相应于 $\alpha/2$ 的临界值〔注〕。

但是当 t 分布的自由度较大时, t 分布很接近正态分布。因此在大子样问题中, 可直接把 (3.12) 式中的 σ 代之以 S_y 作 β_i 的估计。

由〔说明 1.1〕知 y 服从正态分布

$$N(\beta_0 + \beta_1 X_1 + \dots + \beta_m X_m, \sigma^2)$$

类似上述 β_i 的区间估计的讨论, 可知: 与自变量中的任意一观测

〔注〕: 由于 $|b_i - \beta_i|$ 包括 $b_i - \beta_i > 0$ 和 $b_i - \beta_i < 0$ 两种情况, 根据分布的对称性, 因此临界值相等于 $\alpha/2$ 。

测值 (X_1, X_2, \dots, X_m) 相对变动的 y 的区间估计为

$$\left| \frac{y - (\beta_0 + \beta_1 X_1 + \dots + \beta_m X_m)}{\sigma} \right| < N_{\alpha/2} \quad (3.15)$$

(显著水平为 α)

使用上式时要求 σ 与各 β_i 都已知, 这在实际问题中是几乎不可能的。为了给出区间估计, 我们也可以按统计理论构造新的统计量, 但统计量相当复杂。在实际统计中, 只要样本足够大, 人们往往用 b_i 代替 β_i , 用 S_y 代替 σ , 得如下近似式

$$\left| \frac{y - (b_0 + b_1 X_1 + \dots + b_m X_m)}{S_y} \right| < N_{\alpha/2} \quad (3.16)$$

或

$$y^* - N_{\alpha/2} S_y < y < y^* + N_{\alpha/2} S_y$$

上述分析表明, β_i 与 y 的估计精度直接依赖于 S_y , 因此 S_y 也是常用于衡量总的回归效果的一个重要指标。由于

$$S_y^2 = \frac{Q}{N - m - 1}$$

要使 S_y 减小 (即提高估计精度), 则必须降低 Q 或 m 。众所周知, 残差平方和 Q 是随着回归方程中变量个数 m 的增加而减少的, 因此要求同时降低 Q 和 m 是矛盾的。解决矛盾的方法虽然是对给定的自变量逐一鉴别, 去粗取精, 选用重要的变量, 这正是下面的讨论内容。

(3) 各变量的重要性—— β_i 的显著性检验。

类似本节 (1) 的讨论, 考察回归方程中的任一变量 (例如 X_i) 的方差贡献。为明确起见, 令 $Q^{(m)}$ 表示 m 个变量组成的回归方程的残差平方和, $Q^{(m-1)}$ 表示 m 个变量中去掉任一自变量 (例如 X_i) 后所得的残差平方和。令

$$V_i = Q^{(m-1)} - Q^{(m)} \quad (3.17)$$

显然 V_i 就是 X_i 在这 m 个变量的回归方程中的方差贡献。仿