

应用数学丛书

多元统计分析
及计算程序



林少宫 袁蒲佳 申鼎煊 编著

华中工学院出版社

多元统计分析及计算程序

林少宫 袁蒲佳 申鼎煊 编著

华中工学院出版社

多元统计分析及计算程序

林少宫 袁蒲佳 申鼎煊 编著

责任编辑 龙纯曼

华中工学院出版社出版发行

（武昌喻家台）

新华书店湖北发行所经销

华中工学院出版社沔阳印刷厂印刷

开本：787×1092 1/32 印张 9.525 字数：197*000

1987年1月第1版 1987年1月第1次印刷

印数：1—1,500

ISBN 7-5609-0023-2/0·3

统一书号：13255—044 定价：1.65 元

内 容 简 介

本书较系统地介绍了多元统计分析和它的应用，并配有用标准 BASIC 语言编写的计算程序。内容包括有多元回归分析、方差分析、判别分析、聚类分析、主成分分析、因子分析、典型相关分析、对应分析、非线性映射及时间序列分析等。各章都有程序使用说明，框图和实例。

本书可供地质、气象、经济、医学、生物、农业、体育等方面的科技工作者和研究人员及学校师生参考。

作者，提供一套便于在微型电算机上进行多元分析的计算程序。诚然，在早先出版过的一些多元统计计算的书籍中，这些算法程序并不完全缺乏，但大多是用ALGOL语言编写的，特别是针对国产DJS—21机来实现的。在微型电算机日趋普及的今天，这些程序虽然仍有可借鉴的意义，但因微型机上均不配备ALGOL语言，故而不能直接使用。为此，作者尝试用BASIC语言重新编写这些算法。BASIC语言最为普及，容易学会，又具有一定的会话功能，可以说任何微型机上都可以使用它。现在，为了使这些程序更具有普遍适用性，作者有意把所编程序限制在BASIC语言的最基本语法内，而舍弃了一些较为特殊的部分如扩展BASIC的内容。这样看起来程序的某些方面如输出格式方面显得不够完美，但它能在已拥有的任何微型机上运行而不必去更改程序，从而对用户提供了极大便利。这也正是作者的至诚愿望。

本书取名为多元统计分析及计算程序，而不取名为多元统计分析的计算程序，是因为作者还抱着另一种愿望，即希望读者通过本书的阅读，能对多元统计分析的原理、模型、问题和方法有一个基本的了解。也就是说，撇开算法程序不谈，对于具备单元统计分析和矩阵代数基础知识的读者，本书也可作为应用多元统计分析的一种手册式的入门读物。

作者在编写本书的过程中得到了华中工学院数学系陆传务教授的热情支持和具体帮助，计算机系刘健教授在百忙中为本书审稿，在此一并表示衷心的感谢。

由于作者的水平所限，书中难免有不完善甚至错误的地方，尚祈读者不吝指正。

编 者

1985年8月

目 录

第一章 多元回归分析

§ 1	引言	(1)
§ 2	多元线性回归	(2)
§ 3	实例与程序	(10)
§ 4	逐步回归	(20)
§ 5	实例与程序	(25)
§ 6	非线性模型的估计	(39)
§ 7	实例与程序	(41)
§ 8	二段最小二乘法	(51)
§ 9	实例与程序	(53)

第二章 方差分析

§ 1	引言	(63)
§ 2	单因素方差分析	(63)
§ 3	实例与程序	(65)
§ 4	多因素方差分析	(70)
§ 5	实例与程序	(74)

第三章 判别分析

§ 1	引言	(82)
§ 2	贝叶斯准则及正态母体的多类判别	(83)
§ 3	分类判别效果的检验	(87)
§ 4	逐步判别分析	(90)
§ 5	实例与程序	(95)
§ 6	费歇准则二类线性判别分析	(109)
§ 7	实例与程序	(116)

第四章 聚类分析

§ 1	引言	(129)
§ 2	数据变换与相似性统计量	(130)
§ 3	逐步聚类法及谱系图的形成	(134)

§ 4 实例与程序.....	(138)
§ 5 最优分割法.....	(149)
§ 6 实例与程序.....	(154)

第五章 主成分分析

§ 1 引言.....	(165)
§ 2 主成分的确定.....	(165)
§ 3 实例与程序.....	(171)

第六章 因子分析

§ 1 引言.....	(183)
§ 2 因子模型与基本定理.....	(184)
§ 3 主因子解.....	(187)
§ 4 方差最大正交因子旋转.....	(189)
§ 5 因子得分.....	(192)
§ 6 因子分析的计算步骤.....	(193)
§ 7 实例与程序.....	(195)

第七章 典型相关分析

§ 1 引言.....	(210)
§ 2 总体典型变量与典型相关系数.....	(211)
§ 3 样本典型变量与典型相关系数.....	(214)
§ 4 典型相关系数的显著性检验.....	(215)
§ 5 实例与程序.....	(217)

第八章 对应分析

§ 1 引言.....	(229)
§ 2 对应分析的数学方法.....	(230)
§ 3 对应分析的计算步骤.....	(234)
§ 4 实例与程序.....	(235)

第九章 非线性映射

§ 1 引言.....	(247)
§ 2 简单数学原理.....	(247)
§ 3 非线性映射的计算步骤.....	(250)
§ 4 实例与程序.....	(251)

第十章 时间序列分析

§ 1	引言	(265)
§ 2	线性平稳模型	(265)
§ 3	模型的识别	(267)
§ 4	模型的参数估计	(271)
§ 5	模型的拟合检验和预测	(274)
§ 6	平稳时间序列分析的计算步骤	(276)
§ 7	实例与程序	(276)

第一章 多元回归分析

§ 1 引言

回归分析是研究随机现象中变量之间相关关系的一种数理统计方法。通过对变量实际观测值的分析、计算，建立一个变量与另一个（或一组）变量的所谓回归方程，经统计检验认为回归效果显著后，即可用于预测和控制。

回归分析在许多科学领域中都有广泛的应用，如经济中的产品销售量预测，农业中的病虫害或收成预测，自然现象中的地震、洪水预测等等。本章主要讨论多元线性回归和逐步回归。对非线性回归只作扼要的介绍。

设随机变量 y 与 N 个自变量 x_1, \dots, x_N 有线性关系，于是作 y 的 N 元线性回归模型

$$y = \beta_0 + \beta_1 x_1 + \dots + \beta_N x_N + \varepsilon, \quad (1.1.1)$$

其中 ε 是服从正态分布 $N(0, \sigma^2)$ 的随机变量（误差）， β_0, \dots, β_N 为回归系数。

回归分析的主要步骤是：

(1) 由观测值确定参数（回归系数） β_0, \dots, β_N 的估计值 b_0, \dots, b_N ，从而得 y 对 x_1, \dots, x_N 的线性回归方程

$$\hat{y} = b_0 + b_1 x_1 + \dots + b_N x_N, \quad (1.1.2)$$

其中 \hat{y} 表示对 y 的估计。

(2) 对线性关系、自变量的显著性进行统计检验。

(3) 利用回归方程进行预报或控制。

(4) 对预报精度作出估计。

§ 2 多元线性回归

(一) 回归系数的最小二乘估计

设对 y 及 x_1, \dots, x_N 作 M 次观测得

$$(y_\alpha, x_{1\alpha}, \dots, x_{N\alpha}), \alpha = 1, \dots, M.$$

将观测值代入 (1.1.1) 式，得

$$y_\alpha = \beta_0 + \beta_1 x_{1\alpha} + \dots + \beta_N x_{N\alpha} + \varepsilon_\alpha \quad (\alpha = 1, \dots, M), \quad (1.2.1)$$

其中 ε_α 相互独立且服从正态分布 $N(0, \sigma^2)$ 。

若 b_i 分别为 β_i ($i = 0, 1, \dots, N$) 的估计值，则由 (1.1.2) 式得

$$\hat{y}_\alpha = b_0 + b_1 x_{1\alpha} + \dots + b_N x_{N\alpha}. \quad (1.2.2)$$

称观测值 y_α 与回归值 \hat{y}_α 之差 $y_\alpha - \hat{y}_\alpha$ 为残差（或剩余），记作 e_α ，它是误差 ε_α 的估计值。

用最小二乘法求 b_i ，即是选取 b_i 使

$$Q = \sum_{\alpha=1}^M e_\alpha^2 = \sum_{\alpha=1}^M (y_\alpha - b_0 - b_1 x_{1\alpha} - \dots - b_N x_{N\alpha})^2 = \min.$$

据微分学知识，由 $\frac{\partial Q}{\partial b_i} = 0$ ($i = 0, 1, \dots, N$)，经整理可得

$$\left\{ \begin{array}{l} S_{11}b_1 + \dots + S_{N1}b_N = S_{y1}, \\ \dots \\ S_{1N}b_1 + \dots + S_{NN}b_N = S_{yN}, \end{array} \right. \quad (1.2.3)$$

其中

$$\left\{ \begin{array}{l} S_{ij} = \sum_{a=1}^M (x_{ia} - \bar{x}_i)(x_{ja} - \bar{x}_j) \quad (i, j = 1, \dots, N), \\ S_{yj} = \sum_{a=1}^M (y_a - \bar{y})(x_{ja} - \bar{x}_j) \quad (j = 1, \dots, N), \\ \bar{x}_i = \frac{1}{M} \sum_{a=1}^M x_{ia}, \quad \bar{y} = \frac{1}{M} \sum_{a=1}^M y_a. \end{array} \right. \quad (1.2.4)$$

(1.2.3) 式称为正规方程组, b_i 的系数组成的矩阵 $\mathbf{S} = (S_{ij})$ 称为系数矩阵. 由观测值算出离差 S_{ij} 和 S_{yj} 后, 即可解得 b_i ($i = 1, \dots, N$). 再计算 $b_0 = \bar{y} - b_1 \bar{x}_1 - \dots - b_N \bar{x}_N$ (见多元线性回归程序).

为了讨论方便起见, 记

$$\mathbf{X} = \begin{pmatrix} x_{11} - \bar{x}_1 & \cdots & x_{N1} - \bar{x}_N \\ \vdots & & \vdots \\ x_{1M} - \bar{x}_1 & \cdots & x_{NM} - \bar{x}_N \end{pmatrix},$$

$$\mathbf{b} = (b_1, \dots, b_N)',$$

$$\mathbf{Y} = (y_1 - \bar{y}, \dots, y_M - \bar{y})',$$

则 (1.2.3) 式写成矩阵式

$$\mathbf{X}' \mathbf{X} \mathbf{b} = \mathbf{X}' \mathbf{Y} \text{ 或 } \mathbf{S} \mathbf{b} = \mathbf{B},$$

其中 $\mathbf{S} = \mathbf{X}' \mathbf{X}$, $\mathbf{B} = \mathbf{X}' \mathbf{Y}$. 当 \mathbf{S} 满秩时, 有逆矩阵 $\mathbf{S}^{-1} = \mathbf{C} = (C_{ij})$, 于是

$$\left\{ \begin{array}{l} \mathbf{b} = \mathbf{S}^{-1} \mathbf{B} = (\mathbf{X}' \mathbf{X})^{-1} \mathbf{X}' \mathbf{Y}, \\ b_0 = \bar{y} - \sum_{i=1}^N b_i \bar{x}_i. \end{array} \right.$$

在逐步回归中, 常将 (1.2.3) 式的离差 S_{ij} 改为相关系数

$$r_{ij} = \frac{S_{ij}}{\sqrt{S_{ii}} \sqrt{S_{jj}}} \quad (i, j = 1, \dots, N, y).$$

显然 $r_{ii} = 1$ ($i = 1, \dots, N, y$). 此时 (1.2.3) 式写成

$$\begin{cases} r_{11}b_1^* + \cdots + r_{N1}b_N^* = r_{y1}, \\ \cdots & \cdots \\ r_{1N}b_1^* + \cdots + r_{NN}b_N^* = r_{yN}. \end{cases} \quad (1.2.5)$$

(1.2.5)式的解 b_i^* 与 (1.2.3) 式的解 b_i 有如下关系:

$$b_i = b_i^* \sqrt{S_{yy}} / \sqrt{S_{ii}} \quad (i = 1, \dots, N). \quad (1.2.6)$$

$$\text{从而 } b_0 = \bar{y} - \left(\frac{\sqrt{S_{yy}}}{\sqrt{S_{11}}} \bar{x}_1 + \cdots + \frac{\sqrt{S_{yy}}}{\sqrt{S_{NN}}} b_N^* \bar{x}_N \right). \quad (1.2.7)$$

相关系数矩阵 $R = (r_{ij})$ 的逆阵 (C_{ij}^*) 与离差矩阵 $S = (S_{ij})$ 的逆阵 $C = (C_{ij})$ 有如下关系:

$$C_{ij} = C_{ij}^* / \sqrt{S_{ii}} \sqrt{S_{jj}}.$$

在回归计算中, 常把 y 作为第 $N+1$ 个变量的 (1.2.5) 式的系数增广矩阵

$$R^{(0)} = (r_{ij}) = \begin{pmatrix} r_{11} & \cdots & r_{N1} & r_{N+11} \\ \cdots & \cdots & \cdots & \cdots \\ r_{1N} & \cdots & r_{NN} & r_{N+1N} \\ r_{1N+1} & \cdots & r_{NN+1} & r_{N+1N+1} \end{pmatrix}. \quad (1.2.8)$$

称为相关系数矩阵。

当 $N=1$ 时, (1.2.1) 式变成

$$y_\alpha = \beta_0 + \beta_1 x_\alpha + \varepsilon_\alpha \quad (\alpha = 1, \dots, M), \quad (1.2.9)$$

(1.2.2) 式变成

$$\hat{y}_\alpha = b_0 + b_1 x_\alpha, \quad (1.2.10)$$

(1.2.3) 式变成

$$\begin{cases} \frac{\partial Q}{\partial b_0} = -2 \sum_{\alpha=1}^M (y_\alpha - b_0 - b_1 x_\alpha) = 0, \\ \frac{\partial Q}{\partial b_1} = -2 \sum_{\alpha=1}^M (y_\alpha - b_0 - b_1 x_\alpha) x_\alpha = 0. \end{cases} \quad (1.2.11)$$

由 (1.2.11) 式解得

$$\begin{cases} b_1 = \frac{\sum_{a=1}^M x_a y_a - M \bar{x} \bar{y}}{\sum_{a=1}^M x_a^2 - M \bar{x}^2}, \\ b_0 = \bar{y} - b_1 \bar{x}, \end{cases} \quad (1.2.12)$$

其中 $\bar{x} = \frac{1}{M} \sum_{a=1}^M x_a, \quad \bar{y} = \frac{1}{M} \sum_{a=1}^M y_a.$

(二) 回归问题的统计检验

前面是在假定 y 与 $\mathbf{x} = (x_1, \dots, x_N)$ 具有线性关系的条件下建立线性回归方程的。但 y 与 \mathbf{x} 的关系是否为线性关系？所有 x_i ($i = 1, \dots, N$) 对 y 是否都有影响？需要作统计检验。在不同问题中，回归分析的重点各异。如果是要建立经验公式或作曲线拟合，则人们最关心的是算出回归系数估计值 b_1 ；如果是为了预报，则最关心的是求预报值及其精度，即求 \hat{y} 及方差 σ^2 ；如果是要找影响 y 的主要因素，则须检验 b_i 的显著性，从而决定 x_i 的取舍。下面我们以经济计量中的统计准则为例，介绍统计检验的具体内容，其中包括：方差分析， F 检验， t 检验，拟合优度检验及标准差。

(1) 方差分析及 F 检验

类似方差分析法，可将 y 的总离差平方和 S_{yy} 分解为

$$S_{yy} = U + Q,$$

其中

$$\begin{cases} S_{yy} = \sum_{a=1}^M (y_a - \bar{y})^2, \\ U = \sum_{a=1}^M (\hat{y}_a - \bar{y})^2 = \sum_{i=1}^N b_i S_{yi}, \\ Q = \sum_{a=1}^M (y_a - \hat{y}_a)^2 = \sum_{a=1}^M e_a^2 = S_{yy} - U. \end{cases} \quad (1.2.13)$$

U, Q 分别称为回归离差平方和和线差平方和。

它们的自由度依次是 $f_{\text{总}} = M - 1, f_U = N, f_Q = M - N - 1$ 。

对于给定的观测值, S_{yy} 也就给定了, 于是, 可用 U 和 Q 的相对大小来衡量回归效果。作统计量

$$F = \frac{U/N}{Q/(M-N-1)}. \quad (1.2.14)$$

为了检验线性回归方程是否有意义, 可作假设

$$H_0: \beta_1 = \beta_2 = \cdots = \beta_N = 0.$$

如果 H_0 成立, 就表示回归方程无意义。可以证明, 当 H_0 为真时, 统计量 F 服从自由度为 N 和 $M - N - 1$ 的 F 分布。对于给定的显著性水平 α , 若计算值 $F > F_\alpha(N, M - N - 1)$, 则在 α 显著水平上拒绝假设 H_0 , 意谓线性回归方程有显著意义; 反之, 若 $F \leq F_\alpha(N, M - N - 1)$, 则认为线性回归方程无显著意义。

(2) 各自变量的重要性—— β_i 的显著性检验

线性回归方程的显著性, 只说明回归方程中全部自变量的总体效果, 并不意味着每个自变量都有效果。为了考察各个自变量对 y 的影响是否显著, 必须逐一检验假设

$$H_{0i}: \beta_i = 0 \quad (i = 1, \dots, N).$$

为此, 可采用 t 检验。由统计理论知, 统计量

$$t_i = \frac{b_i - \beta_i}{\sqrt{C_{ii} S_y}} \quad (i = 1, \dots, N)$$

服从自由度为 $M - N - 1$ 的 t 分布。其中 C_{ii} 为逆阵 $\mathbf{S}^{-1} = \mathbf{C}$ 的对角线上的元素, S_y 为 y 的剩余标准差的无偏估计, $S_y = \sqrt{Q/(M - N - 1)}$ 。当 H_{0i} 为真时, 有

$$t_i = \frac{b_i}{\sqrt{C_{ii} S_y}} \quad (i = 1, \dots, N). \quad (1.2.15)$$

对于给定的显著性水平 α , 若计算值 $|t| > t_\alpha(M - N - 1)$, 则拒

绝假设 H_0 , 认为 β_i 显著不为 0, 即相应的 x_i 对 y 有影响; 否则 β_i 与 0 无显著差别, 即 x_i 对 y 无关重要, 应该从回归方程中剔除之. 然后再对剩下的自变量建立回归方程.

顺便指出, 在下面的逐步回归中, 检验变量 x_i 贡献的显著性是采用与 t 检验等价的 F 检验. 设从 N 个变量中去掉 x_i , 回归平方和由 U 降到 U' , 则 $W_i = U - U'$ 称为 x_i 在这 N 个变量的回归方程中的贡献. 可以证明: $W_i = b_i^2 / C_{ii}$, 且

$$F_i = \frac{W_i / 1}{Q / (M - N - 1)} = \frac{b_i^2}{C_{ii} S_y^2}$$

在假设 H_0 下服从自由度为 1 和 $M - N - 1$ 的 F 分布. 若算得 $F_i > F_a(1, M - N - 1)$, 则 x_i 的贡献显著; 反之可从回归方程中剔除 x_i .

(3) D. W 检验

在回归分析中, 有一个重要的假定, 认为 ε_i 是独立随机变量, 即在任一时刻的取值与任何前期值无关. 当 ε 取值与它的前一期的值相关时, 称这种自相关为序列相关. 在经济变量中, 常常出现这种情形, 如销售量并非只受独立变量的影响, 也部分地受到前期销售量的影响, 因此误差项 ε 存在自相关. 自相关性会影响 t 检验和 F 检验的效果, 以致整个回归计算和分析需要重新考虑. 为了检验是否存在 ε 的自相关, 常采用 Durbin-Watson 统计量

$$D = \frac{\sum_{i=1}^M (e_i - e_{i-1})^2}{\sum_{i=1}^M e_i^2} \quad (1.2.16)$$

进行所谓 $D.W$ 检验, 这里 e_i 是残差 $y_i - \hat{y}_i$, 即 ε_i 的估计值. $D.W$ 检验的步骤是:

假设 $H_0: \varepsilon_t$ 与 ε_{t-1} 无自相关。

对于给定的显著水平 α , 自变量个数 N , 样本容量 M , 查 D 分布表得上临界值 d_u 和下临界值 d_l , 由 (1.2.16) 式算得 D , 并按表 1.1 处理假设 H_0 .

表 1.1

D 值	结 论
$(4-d_l) < D < 4$	拒绝 H_0 , 有负自相关
$(4-d_u) < D < (4-d_l)$	不作结论
$d_u < D < (4-d_u)$	接受 H_0
$d_l < D < d_u$	不作结论
$0 < D < d_l$	拒绝 H_0 , 有正自相关

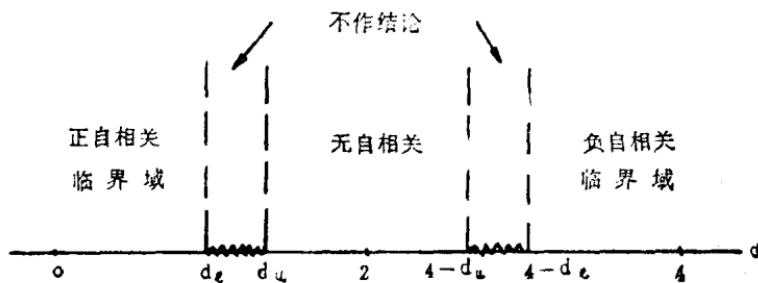


图 1.1

(4) 拟合优度检验

为了说明回归方程中全部自变量的总效果, 除前面讲的 F 检验外, 又可采用下面的无量纲值

$$R^2 = \frac{U}{S_{yy}} = 1 - \frac{Q}{S_{yy}} = 1 - \frac{\sum_{i=1}^M e_i^2}{\sum_{i=1}^M (y_i - \bar{y})^2}$$

或 $R = \sqrt{1 - \frac{Q}{S_{yy}}} \quad (1.2.17)$

实质上 U 是回归方程中全部自变量的“方差贡献”，故 R^2 即为这种贡献在总离差中所占的比例。容易看出， $0 \leq R^2 \leq 1$ ；当 $R^2 = 1$ 时， y_i 与 \hat{y}_i 完全拟合。 R^2 愈接近 1， y_i 与 \hat{y}_i 的拟合愈好，所以 R^2 是检验回归方程拟合优度的指标。 R 称为复相关系数。

(三) 利用回归方程进行预测和控制

对于给定的样本值 $x_{10}, x_{20}, \dots, x_{N0}$ ，代入回归方程算得的

$$\hat{y}_0 = b_0 + b_1 x_{10} + \dots + b_N x_{N0}$$

是 y_0 的点估计值。实用中常需给出 y_0 的预测区间。当 M 较大且 x_{10} 较接近 \bar{x}_1 时，可以近似地认为

$$y_0 - \hat{y}_0 \sim N(0, \sigma^2),$$

于是对给定显著性水平 $\alpha = 0.05$ 和 $\alpha = 0.01$ ，可分别用下面两式进行预测和控制：

$$P(\hat{y}_0 - 1.96\sigma < y_0 < \hat{y}_0 + 1.96\sigma) = 0.95$$

和 $P(\hat{y}_0 - 2.58\sigma < y_0 < \hat{y}_0 + 2.58\sigma) = 0.99.$

由于 σ 未知，用无偏估计量

$$S_y = \sqrt{\frac{Q}{M-N-1}} = \sqrt{\frac{\sum_{i=1}^M e_i^2}{M-N-1}}$$

代替，故 y_0 的 95% 和 99% 预测区间分别为

$$(\hat{y}_0 - 1.96 S_y, \hat{y}_0 + 1.96 S_y)$$

和 $(\hat{y}_0 - 2.58 S_y, \hat{y}_0 + 2.58 S_y).$

在各项经济指标的预测中，往往 M 不可能很大， x_{10} 也不一定接近 \bar{x}_1 ，这时对给定的显著性水平 α ，其预测区间为