

实用数据 分析方法

吴国富 安万福 刘景海 编



中国统计出版社

实用数据分析方法

吴国富 安万福 刘景海 编

中国统计出版社

(京)新登字041号

内 容 提 要

本书作为《微型机对话型数据统计分析方法程序集》的姊妹篇，全书共十章，内容有基础统计分析方法，各种随机数的产生，多元数据图表示的各种方法，非参数检验的各种方法，多元分析中的各种方法及数量化理论。该书注重介绍各种方法的实施步骤和应用范围，并扼要介绍其主要数学理论依据，但不作繁琐的推导证明。

本书可供从事应用数据分析方法的工程技术人员和管理人员学习及解决实际问题使用，也可供大专院校学生学习和实习参考。

实用数据分析方法

SHIYONG SHUJU FENXI FANGFA

吴国富 安万福 刘景海 编

*

中国统计出版社出版

(北京三里河月坛南街38号 100826)

新华书店北京发行所发行

北京市通县永乐印刷厂印刷

*

850×1168毫米 32开本 11.125印张 27万字

1992年10月第1版 1992年10月北京第1次印刷

印数：1—4000

ISBN 7-5037-0882-4/C·568

定价：11.00元

前 言

众所周知，无论是实验得到的数据，还是通过抽样或调查得到的数据，只有针对不同的数据，选用某些统计分析方法进行处理和分析，才能够揭示这些数据所反应的自然规律，进而获得解决问题的方法。因此，数据的各种统计分析方法，在各个学科中得到了越来越广泛的应用，特别是在地质、医学、气象、管理及各种社会调查的数据分析、处理等诸多方面的应用，成绩更为显著。

本书的内容与《微型机对话型数据统计分析方法程序集》一书相对应，那里侧重于程序的编写及调用，而这里则侧重于方法的介绍和应用，推导和证明较少。

本书是在1987年完成第一稿的基础上，经过中国科学院应用数学研究所多年来的教学与实践，几次修改后定稿。另外，我们还编制了与本书涉及的全部统计方法对应的应用软件（共2张盘，现已是第三版），以供社会各部门、各单位需要。

本书主要包括三部分内容：

第一部分（1~5章）为基础方法部分，主要介绍各种检验，包括非参数检验、各种随机数的产生、各种图分析法等。

第二部分（6~9章）为多元统计分析方法部分，主要介绍多元回归分析、主成份分析、因子分析、判别分析、聚类分析的各种方法。

第三部分（10章）为数量化理论部分，包括数量化方法Ⅰ~Ⅳ，这些方法在管理及社会科学研究方面有着广泛的应用。

本书的目的在于给广大工程技术人员和管理人员在利用《微型机对话型数据统计分析方法程序集》解决实际问题时，能带来

一些方便。

限于我们的水平，错误和不当之处在所难免，欢迎广大读者批评指正。

编者

1992.5

目 录

第一章 描述统计

§ 1.1 特征值、方差、斜度、峰度.....	(1)
1.1.1 总和与均值.....	(2)
1.1.2 最大值、最小值、中位值、众数.....	(3)
1.1.3 方差、标准差、均差.....	(3)
1.1.4 极差、四分位极差.....	(4)
1.1.5 斜度、峰度.....	(4)
1.1.6 箱尾图.....	(5)
§ 1.2 频数分布.....	(6)
1.2.1 频数分布与直方图.....	(6)
1.2.2 树干和叶.....	(9)
§ 1.3 相关系数.....	(11)
1.3.1 相关系数.....	(11)
1.3.2 回归直线.....	(12)
§ 1.4 顺序相关系数.....	(15)
§ 1.5 数据平滑.....	(18)
1.5.1 滑动平均法.....	(18)
1.5.2 中位数扫描法.....	(19)

第二章 估计、检验

§ 2.1 区间估计.....	(21)
2.1.1 总体均值的区间估计.....	(21)

2.1.2	总体方差的置信区间	(23)
§ 2.2	正态总体均值和方差的检验	(24)
2.2.1	均值检验	(24)
2.2.2	方差检验	(25)
§ 2.3	两个总体均值之差的检验	(27)
§ 2.4	拟合优度检验	(31)
§ 2.5	独立性检验	(32)
2.5.1	2×2 列联表	(32)
2.5.2	$k \times l$ 列联表	(34)

第三章 非参数法

§ 3.1	两种处理效果之差的检验	(38)
3.1.1	Wilcoxon 检验	(39)
3.1.2	Mann—Whitney 检验	(40)
3.1.3	中位数检验	(42)
3.1.4	Van·der Waerden 检验	(43)
§ 3.2	配对检体两种处理效果之差的检验	(45)
3.2.1	符号检验	(45)
3.2.2	带符号的 Wilcoxon 检验	(46)
§ 3.3	配对检体两种处理效果的独立性 (一致性) 检验	(48)
3.3.1	Spearman 顺序相关系数	(48)
3.3.2	Kendall 顺序相关系数	(50)
§ 3.4	三个以上处理效果之差的检验	(51)
3.4.1	Kruskal-Wallis 检验	(51)
3.4.2	Jonckheere 检验	(53)
3.4.3	k -重图	(54)
§ 3.5	配对情形 m (matching) 种处 理效果之差的检验	(57)

3.5.1	Friedman 检验	(58)
3.5.2	Page 检验	(59)
§ 3.6	各处理间的一致性 (独立性) 检验	(60)

第四章 随机数

§ 4.1	一维随机数的产生	(63)
4.1.1	二项随机数	(63)
4.1.2	普阿松随机数	(64)
4.1.3	指数随机数	(65)
4.1.4	正态随机数	(66)
4.1.5	Weikull 随机数	(69)
4.1.6	Γ -随机数	(70)
§ 4.2	多维随机数的产生	(71)
4.2.1	多维正态随机数的产生	(71)
4.2.2	单位球面上的随机数	(74)
4.2.3	超平面上的随机点	(75)
4.2.4	随机排列	(76)
§ 4.3	随机数的检验	(77)
4.3.1	频数检验	(78)
4.3.2	链的检验	(80)
4.3.3	组合检验	(83)

第五章 多元数据的图分析法

§ 5.1	雷达图	(87)
§ 5.2	脸谱图	(89)
§ 5.3	三角多项式图表示法	(97)
§ 5.4	星座图	(100)
§ 5.5	连接向量图	(102)
§ 5.6	单变量的正态性检验	(106)

§ 5.7	多元数据的正态检验·····	(111)
-------	----------------	---------

第六章 多元回归分析

§ 6.1	什么是多元回归分析·····	(116)
§ 6.2	多元线性回归分析·····	(117)
6.2.1	多元线性回归模型·····	(117)
6.2.2	回归系数的最小二乘估计·····	(118)
6.2.3	标准回归系数·····	(121)
6.2.4	方差分析表与回归方程的检验·····	(121)
6.2.5	偏相关系数·····	(124)
6.2.6	回归系数的检验与置信区间·····	(125)
6.2.7	利用回归方程进行预测和控制·····	(127)
§ 6.3	逐步回归分析·····	(130)
6.3.1	何为逐步回归分析·····	(130)
6.3.2	变量选择的方法·····	(131)
6.3.3	引入变量和剔除变量的依据·····	(133)
6.3.4	引入和剔除变量的消去变换·····	(135)
6.3.5	选择最佳模型的准则·····	(136)
§ 6.4	回归诊断·····	(150)
6.4.1	何为回归诊断·····	(150)
6.4.2	残差及标准化残差·····	(151)
6.4.3	各种图分析·····	(153)
6.4.4	残差的序列相关性检验·····	(154)
6.4.5	对回归方程影响特别大的样品的发现·····	(155)
6.4.6	多重共线性的发现·····	(158)

第七章 判别分析

§ 7.1	何为判别分析·····	(166)
§ 7.2	二群线性判别·····	(167)

7.2.1	线性判别函数的导出 (I)	(167)
7.2.2	线性判别函数的导出 (II)	(171)
7.2.3	误判概率	(172)
7.2.4	判别系数的检验	(173)
§ 7.3	二群非线性判别	(183)
7.3.1	二次判别函数的导出	(183)
7.3.2	关于误判概率	(184)
7.3.3	等方差性检验	(184)
§ 7.4	多群线性判别 (逐步判别)	(191)
7.4.1	多群线性判别函数	(191)
7.4.2	Wilks的 λ 统计量	(193)
7.4.3	各变量贡献率的检验	(194)
7.4.4	变量选择	(195)
§ 7.5	典型判别分析	(207)

第八章 主成份分析与因子分析

§ 8.1	主成份分析	(218)
8.1.1	何为主成份分析	(218)
8.1.2	主成份的导出	(219)
8.1.3	用雅可比法求特征值和特征向量	(221)
8.1.4	贡献率与累积贡献率	(224)
8.1.5	数据的标准化的	(225)
§ 8.2	典型相关分析	(228)
8.2.1	何为典型相关分析	(228)
8.2.2	典型相关分析的计算步骤	(230)
§ 8.3	因子分析	(239)
8.3.1	何为因子分析	(239)
8.3.2	因子载荷矩阵的估计	(242)
8.3.3	因子的解释	(242)

8.3.4 因子载荷的正交旋转.....	(243)
----------------------	---------

第九章 聚类分析

§ 9.1 聚类分析研究的内容.....	(255)
§ 9.2 变量的类型.....	(256)
§ 9.3 距离与相似系数.....	(257)
9.3.1 距离.....	(258)
9.3.2 相似系数.....	(258)
§ 9.4 聚类的方法.....	(261)
§ 9.5 加入法.....	(297)

第十章 数量化理论

§ 10.1 数量化理论——I类.....	(300)
10.1.1 数量化理论——I类的数据.....	(301)
10.1.2 数学模型.....	(304)
10.1.3 预测精度及各项目的贡献.....	(305)
§ 10.2 数量化理论——II类.....	(309)
10.2.1 数量化理论——II类的数据.....	(310)
10.2.2 判别函数的给出(I).....	(310)
10.2.3 判别函数的给出(II).....	(317)
10.2.4 相关比、范围和偏相关系数.....	(320)
10.2.5 判据的确定.....	(320)
§ 10.3 数量化理论——III类.....	(327)
§ 10.4 数量化理论——IV类.....	(335)
10.4.1 亲近度.....	(335)
10.4.2 数量化理论——IV类的给出.....	(337)

第一章 描述统计

统计的方法，通常有描述性统计和推测性统计。描述统计就是把数据本身包含的信息，加以总结、概括、浓缩、简化，使问题变得更加清晰、简单，易于理解，便于处理。例如，表、图、数字特征等的计算就属于描述统计。

§ 1.1 特征值、方差、斜度、峰度

实验数据有如下三种形式：

(a) n 个数据 x_1, x_2, \dots, x_n

(b) k 个不同的数 x_1, x_2, \dots, x_k 以及每个数的频数 f_1, f_2, \dots, f_k ，列成表如下：

数据 x_i	x_1	x_2	...	x_k	合 计
频数 f_i	f_1	f_2	...	f_k	n

(c) 给出分组和频数

分 组	组 中 值	频 数
$(a_0, a_1]$	m_1	f_1
$(a_1, a_2]$	m_2	f_2
\vdots	\vdots	\vdots

续表

分 组	组 中 值	频 数
$(a_{i-1}, a_i]$	m_i	f_i
\vdots	\vdots	\vdots
(a_{k-1}, a_k)	m_k	f_k
合 计		n

其中 $a_0 < a_1 < a_2 < \dots < a_k$, $m_i = (a_{i-1} + a_i)/2$, f_i 表示数据落入第 i 组 (a_{i-1}, a_i) 内的频数。

1.1.1 总和与均值

对如上三种形式的数据, 其总和、均值的表达式分别为:

总和:

$$(a) \quad T = x_1 + x_2 + \dots + x_n = \sum_{i=1}^n x_i$$

$$(b) \quad T = f_1 x_1 + f_2 x_2 + \dots + f_k x_k = \sum_{i=1}^k f_i x_i$$

$$(c) \quad T = f_1 m_1 + f_2 m_2 + \dots + f_k m_k = \sum_{i=1}^k f_i m_i$$

均值:

$$(a) \quad \bar{x} = \frac{1}{n} \sum_{i=1}^n x_i$$

$$(b) \quad \bar{x} = \frac{1}{n} \sum_{i=1}^k f_i x_i$$

$$(c) \quad \bar{x} = \frac{1}{n} \sum_{i=1}^k f_i m_i$$

1.1.2 最大值、最小值、中位数、众数

设 n 个数据按从小到大的顺序排列为

$$x_{(1)} \leq x_{(2)} \leq \cdots \leq x_{(n)}$$

则最大值即为 $x_{(n)}$ ，最小值为 $x_{(1)}$ ，而中位数

$$M = \begin{cases} x_{(\frac{n+1}{2})} & \text{当 } n \text{ 为奇数时} \\ \frac{1}{2} [x_{(\frac{n}{2})} + x_{(\frac{n}{2}+1)}] & \text{当 } n \text{ 为偶数时} \end{cases}$$

众数：在 (b)、(c) 情况下，使频数最大的 x_i 或 m_i 。

1.1.3 方差、标准差、均差

方差：

$$(a) \quad S^2 = \frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})^2 = \frac{1}{n} \sum_{i=1}^n x_i^2 - \bar{x}^2$$

$$(b) \quad S^2 = \frac{1}{n} \sum_{i=1}^k f_i (x_i - \bar{x})^2 = \frac{1}{n} \sum_{i=1}^k f_i x_i^2 - \bar{x}^2$$

$$(c) \quad S^2 = \frac{1}{n} \sum_{i=1}^k f_i (m_i - \bar{x})^2 = \frac{1}{n} \sum_{i=1}^k f_i m_i^2 - \bar{x}^2$$

标准差：

$$S = \sqrt{S^2}$$

均差：

$$(a) \quad D = \frac{1}{n} \sum_{i=1}^n |x_i - \bar{x}|$$

$$(b) \quad D = \frac{1}{n} \sum_{i=1}^k f_i |x_i - \bar{x}|$$

$$(c) D = \frac{1}{n} \sum_{i=1}^k f_i |m_i - \bar{x}|$$

1.1.4 极差、四分位极差

极差:

$$R = \text{最大值} - \text{最小值} = x_{(n)} - x_{(1)}$$

四分位极差:

$$QR = Q_3 - Q_1$$

其中: Q_3 ——第三个四分位数 (75%的点)。

Q_1 ——第一个四分位数 (25%的点)。

四分位偏差:

$$QD = (Q_3 - Q_1) / 2$$

四分位变差:

$$QV = 100(Q_3 - Q_1) / (Q_3 + Q_1)$$

四分位斜度:

$$QMS = |(Q_3 - M) - (M - Q_1)| / (Q_3 - Q_1)$$

其中 M 为中位数。

变异系数:

$$CV = S / \bar{x} \quad (\text{标准差/均值})$$

1.1.5 斜度、峰度

斜度:

$$\gamma_1 = -\mu_3 / S^3$$

峰度:

$$\gamma_2 = (\mu_4 / S^4) - 3$$

其中 S 为标准差, 而 μ_j 对三种数据的定义如下:

$$(a) \mu_j = \frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})^j$$

$$(b) \mu_j = \frac{1}{n} \sum_{i=1}^k f_i (x_i - \bar{x})^j \quad (j=3, 4)$$

$$(c) \mu_j = \frac{1}{n} \sum_{i=1}^k f_i (m_i - \bar{x})^j$$

由于正态分布的斜度、峰度皆为零，所以常把 γ_1 、 γ_2 做为偏离正态分布的尺度。

1.1.6 箱尾图

Tukey 提出的箱尾图由箱子和其上引出的两个尾组成（见图1.1），故有箱尾图之称。这种图用来表示在一定时间内一个班成绩的变化、物体位置的变化、原料的变化、产品标准的变化等。

例1.1 如下表1.1表示高一年级10门课程5次考试成绩，对此数据，取横坐标为测验次数，画出箱尾图。

我们对每次考试成绩，分别求出均值、标准差、最大值、最小值，便可画出箱尾图如下：

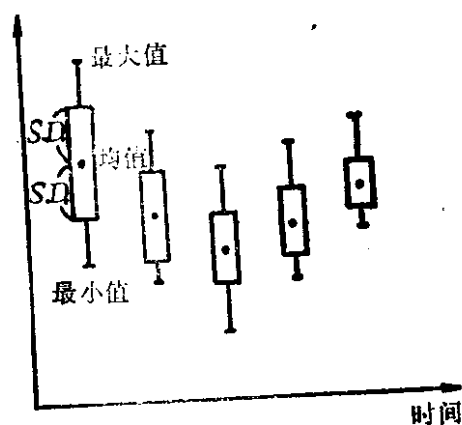


图1.1 箱尾图

表1.1 考试成绩

课 目 次 数	课 目									
	A	B	C	D	E	F	G	H	I	J
1	77	56	61	54	60	68	90	39	68	75
2	63	69	77	66	65	52	58	43	66	64
3	78	58	82	65	86	90	60	41	45	80
4	66	58	60	73	66	63	64	51	74	53
5	51	42	75	68	46	40	44	38	34	58

§ 1.2 频数分布

在这一节里我们将介绍两种描述频数分布的方法：直方图法和树干叶法。

1.2.1 频数分布与直方图

实验数据 x_1, x_2, \dots, x_n 的均值 \bar{x} 、方差 S^2 等数字特征，概括地描述了实验数据的一些基本统计特征，但要进一步研究数据取值的分布情况，则需计算它们的经验分布，做出直方图就是一种最常用的方法，具体做法如下：

先用 a, b 分别表示该数据中的最大、最小值，即

$$a = \max_{1 \leq i \leq n} \{x_i\}, \quad b = \min_{1 \leq i \leq n} \{x_i\}$$

再把区间 $[a, b]$ 分成 k 个等间隔的小区间：

$$[a_{i-1}, a_i], \quad i=1, 2, \dots, k$$

这里 $a = a_0 < a_1 < a_2 < \dots < a_k = b$ 。这样的 k 个小区间把 n 个数据分成了 k 个组，且称每组的中间值 $m_i = (a_{i-1} + a_i)/2, i=1, 2, \dots, k$ 为该组的组中值，往往就以它来代表该组的数据。

然后统计落入各组内的数据个数 f_i ，它就是第 i 组 $[a_{i-1}, a_i]$ 的频数，显然应该有

$$f_1 + f_2 + \dots + f_k = \sum_{i=1}^k f_i = n$$

若用 n 去除各组的频数，又可得各组的频率

$$f_i/n \quad (i=1, 2, \dots, k)$$

