

世纪统计学系列教材

# 多元统计分析

何晓群 编著

中国人民大学出版社

21 世纪统计学系列教材

# 多元统计分析

何晓群 编著

中国人民大学出版社

## 图书在版编目 (CIP) 数据

多元统计分析/何晓群编著.  
北京:中国人民大学出版社,2004  
(21世纪统计学系列教材)

ISBN 7-300-05402-1/F·1683

I. 多…

II. 何…

III. 多元分析:统计分析-教材

IV. 0212.4

中国版本图书馆 CIP 数据核字 (2004) 第 015170 号

21 世纪统计学系列教材

**多元统计分析**

何晓群 编著

---

出版发行	中国人民大学出版社		
社 址	北京中关村大街 31 号	邮政编码	100080
电 话	010-62511242 (总编室)	010-62511239 (出版部)	
	010-82501766 (邮购部)	010-62514148 (门市部)	
网 址	<a href="http://www.crup.com.cn">http://www.crup.com.cn</a> <a href="http://www.ttrnet.com">http://www.ttrnet.com</a> (人大教研网)		
经 销	新华书店		
印 刷	北京鑫丰华彩印有限公司		
开 本	787×965 毫米 1/16	版 次	2004 年 4 月第 1 版
印 张	24.5	印 次	2004 年 4 月第 1 次印刷
字 数	448 000	定 价	28.00 元

---

**版权所有 侵权必究 印装差错 负责调换**

## 《21 世纪统计学系列教材》编委会

编委会主任 易丹辉

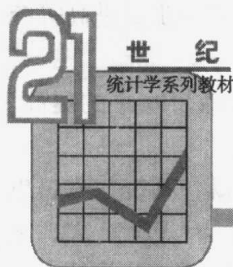
编委会委员 (按姓氏笔画排序)

尹德光 冯士雍 张尧庭

陈希孺 吴喜之 赵彦云

柯惠新 袁 卫 倪加勋

顾 岚 袁寿庄 耿 直



## 总 序

改革开放以来，高等统计教育有了很大的发展。随着课程设置的不断调整，有不少教材出版，同时也翻译引进了一些国外优秀教材。作为培养我国统计专门人才的摇篮，中国人民大学统计学系自 1952 年创建以来，走过了风风雨雨，一直坚持着理论与应用相结合的办学方向，培养能够理论联系实际、解决实际问题的高层次人才。随着新知识经济和网络时代的到来，我们在教学科研的实践中，深切地感受到，无论是自然科学领域、社会科学领域的研究，还是国家宏观管理和企业生产经营管理，甚至人们的日常生活，信息需求量日益增多，信息处理技术更加复杂，作为信息技术支柱的统计方法，越来越广泛地应用于各个领域。

面对新的形势，我们一直在思索，课程设置、教材选择、教学方式等怎样才能使学生适应社会经济发展的客观需要。在反复酝酿、不断尝试的基础上，我们决定与统计学界的同仁，共同编写、出版一套面向 21 世纪的统计学系列教材。

这套系列教材聘请了中科院院士、中国科技大学陈希孺教授，上海财经大学数量经济研究院张尧庭教授，中国科学院数学与系统科学研究所冯士雍研究员等作为编委。他们长期任中国人民大学的兼职教授，一直关心、支持着统计学系的学科建设和应用统计的发展。中国人民大学应用统计科学研究中心 2000 年已成为国家级研究基地，这些专家是首批专职或兼职研究人员。这一开放性研究基地

的运作，将有利于提升我国应用统计科学研究的水平，也必将进一步促进高等统计教育的发展。

这套教材是我们奉献给新世纪的，希望它能促进应用统计教育水平的提高。这套教材力求体现以下特点：

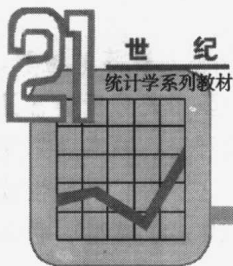
第一，在教材选择上，主要面向经济类统计学专业。选材既包括统计教材也包括风险管理与精算方面的教材。尽管名为统计学系列教材，但并不求大、求全，而是力求精选。对于目前已有的内容较为成熟、适合教学需要、公认的较好的教材，并未列入本次出版计划。

第二，每部教材的内容和写作，注意广泛吸收国内外优秀教材的成果。教材力求简明易懂、内容系统和实用，注重对统计方法思想的阐述，并结合大量实际数据和实例说明统计方法的特点及应用条件。

第三，强调与计算机的结合。为着力提高学生运用统计方法分析解决问题的能力，教材所涉及的统计计算，要求运用目前已有的统计软件。根据教材内容，选择使用 SAS、SPSS、TSP、STATISTICA、EViews、MINITAB、Excel 等。

感谢中国人民大学出版社的同志们，他们怀着发展我国应用统计科学的热情和提高统计教育水平的愿望，经过反复论证，使这套教材得以出版。感谢参与教材编写的同行专家、统计学系的教师。愿大家的辛勤劳动能够结出丰硕的果实。我们期待着与统计学界的同仁，共同创造应用统计辉煌的明天。

易丹辉  
2000年8月  
于中国人民大学



## 前 言

多元统计分析是统计学中一个非常重要的分支,在国外,从20世纪30年代起,已开始在自然科学、管理科学和社会、经济等领域广泛应用。我国自20世纪80年代起在许多领域拉开了多元统计分析应用的帷幕,本书正是为了适应新的需求形势而编著的。

本书写作的指导思想是:在不失严谨的前提下,明显不同于纯数理类教材,努力突出实际案例的应用和统计思想的渗透,结合统计软件较全面地系统介绍多元分析的实用方法。为了贯彻这一思想,本书参考了国内外大量书籍及文献,在系统介绍多元分析基本理论和方法的同时,尽力结合社会、经济、自然科学等领域的研究实例,把多元分析的方法与实际应用结合起来,注意定性分析与定量分析的紧密结合,努力把同行们以及我们在实践中应用多元分析的经验 and 体会融入其中。几乎每种方法都强调它们各自的优缺点和实际运用中应注意的问题。为使读者掌握本书内容,又考虑到这门课程的应用性和实践性,每章后面给出一些简单的思考与练习。我们鼓励读者自己利用一些实际数据去实现这些方法。多元分析的应用离不开计算机,本书的案例主要运用在我国广泛流行的SPSS软件实现,部分方法用SAS软件完成。本书一个显著的特点是在每种方法后结合实例概要介绍了SPSS或SAS软件的实际操作实现过程。在每章后面还都注明了参考

文献，有兴趣的读者可进一步阅读。

全书共分 14 章。主要内容有多元正态分布、均值向量和协方差阵的检验、聚类分析、判别分析、主成分分析、因子分析、对应分析、典型相关分析等常见的主流方法，还参考国内外大量文献系统介绍了这些年在市场研究、顾客满意度研究、金融研究、环境研究等领域应用颇广的一些较新方法。这些内容有定性数据的建模分析、对数线性模型、Logistic 回归、路径分析、结构方程模型、联合分析、多变量的图表示法、多维标度法等。

本书可作为统计学专业本科生的多元分析课程教材。由于本书的内容较多，教师在选用此书为教材时可以灵活选讲。本书还可作为非统计专业研究生量化分析教材。根据我们多年的教学实践，本书讲授 72 课时较为合适，若有计算机和投影设备的配合，教学将会更为方便和有效。

在本书的写作过程中，始终得到中国人民大学 21 世纪统计学系列教材编审委员会和中国人民大学出版社的支持。编写大纲经过教材编审委员会的认真讨论，教材初稿得到吴喜之教授的认真审阅，提出不少中肯意见。在此基础上本人对教材做了认真修改，若现在书中仍有不妥之处，当属笔者自负。本书的大部分案例是我们多年教学和科研工作的积累，有部分案例为体现其典型性引用他人著作。我还要特别感谢长期鼓励我进行应用研究的几位导师方开泰、陈希孺、张尧庭先生。在此，谨向对本书出版有过帮助的师长和朋友表示衷心的感谢。

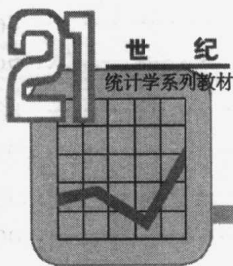
本书的完成可以说是我们师生合作的共同成果。我系博士研究生王作成，硕士研究生陈少杰、李强对本书的计算设计及输入做了大量工作，还有李书争、高玉兰、崔迎等也为本书的写作做过一些具体工作。在我们的合作中，我不仅仅是他们的老师，因为我常常从他们的研究和提问中得到重要启发，教学相长在合作过程中得到真正体现。由于我们水平所限，书中难免有不足之处，尤其是在一些应用研究的体会性讨论中，恐有偏颇之处，恳切希望读者批评指正。对于读者提出的意见和建议，笔者将在中国人民大学六西格玛质量管理研究中心的网站及时给予反馈。网址为 <http://www.ruc-6sigma.com>，欢迎大家登录。

何晓群

2004 年 3 月

于中国人民大学应用统计科学研究中心  
中国人民大学六西格玛质量管理研究中心





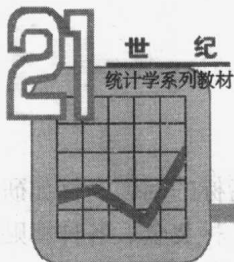
# 目 录

<b>第1章 多元正态分布</b> .....	1
1.1 多元分布的基本概念 .....	2
1.2 统计距离和马氏距离 .....	6
1.3 多元正态分布.....	10
1.4 均值向量和协方差阵的估计.....	16
1.5 维希特 (Wishart) 分布 .....	17
<b>第2章 均值向量和协方差阵的检验</b> .....	20
2.1 均值向量的检验.....	21
2.2 协方差阵的检验.....	29
2.3 形象分析.....	31
2.4 有关检验的上机实现.....	36
<b>第3章 聚类分析</b> .....	54
3.1 聚类分析的基本思想.....	55
3.2 相似性度量.....	58
3.3 类和类的特征.....	63
3.4 系统聚类法.....	66

3.5	模糊聚类分析	76
3.6	K-均值聚类和有序样品的聚类	79
3.7	计算步骤与上机实现	83
3.8	社会经济案例研究	92
<b>第4章</b>	<b>判别分析</b>	<b>98</b>
4.1	判别分析的基本思想	99
4.2	距离判别	100
4.3	Bayes 判别	103
4.4	Fisher 判别	104
4.5	逐步判别	105
4.6	判别分析方法步骤及框图	106
4.7	判别分析的上机实现	113
4.8	判别分析的一个案例	126
<b>第5章</b>	<b>主成分分析</b>	<b>135</b>
5.1	主成分分析的基本思想与理论	136
5.2	主成分分析的几何意义	137
5.3	总体主成分及其性质	140
5.4	样本主成分的导出	148
5.5	有关问题的讨论	150
5.6	主成分分析步骤及框图	156
5.7	主成分分析的上机实现	157
<b>第6章</b>	<b>因子分析</b>	<b>167</b>
6.1	因子分析的基本理论	168
6.2	因子载荷的求解	173
6.3	因子分析的步骤与逻辑框图	180
6.4	因子分析的上机实现	181
<b>第7章</b>	<b>对应分析</b>	<b>195</b>
7.1	列联表及列联表分析	196
7.2	对应分析的基本理论	199
7.3	对应分析的步骤及逻辑框图	206
7.4	对应分析的上机实现	207
<b>第8章</b>	<b>典型相关分析</b>	<b>220</b>
8.1	典型相关分析的基本理论及方法	221

8.2	典型相关分析的步骤及逻辑框图 .....	228
8.3	典型相关分析的上机实现 .....	234
8.4	社会经济案例研究 .....	239
<b>第9章</b>	<b>定性数据的建模分析</b> .....	<b>244</b>
9.1	对数线性模型基本理论和方法 .....	245
9.2	对数线性模型的上机实现 .....	247
9.3	Logistic 回归基本理论和方法 .....	252
9.4	Logistic 回归的方法及步骤 .....	260
9.5	Logistic 回归的上机实现 .....	262
<b>第10章</b>	<b>路径分析</b> .....	<b>267</b>
10.1	基本概念和理论 .....	268
10.2	分解相关系数 .....	273
10.3	路径模型的调试和检验 .....	277
10.4	路径分析流程图及 SPSS 指令 .....	280
10.5	案例分析 .....	282
<b>第11章</b>	<b>结构方程模型</b> .....	<b>287</b>
11.1	结构方程的基本思想及模型设定 .....	288
11.2	结构方程模型的识别和估计 .....	291
11.3	结构方程模型的评价和修改 .....	292
11.4	结构方程模型的上机实现 .....	293
11.5	一个实例 .....	299
<b>第12章</b>	<b>联合分析</b> .....	<b>304</b>
12.1	联合分析的基本理论和方法 .....	304
12.2	联合分析的方法步骤及框图 .....	311
12.3	联合分析的上机实现 .....	317
<b>第13章</b>	<b>多变量的图表示法</b> .....	<b>323</b>
13.1	散点图矩阵 .....	324
13.2	脸谱图 .....	326
13.3	雷达图 .....	328
13.4	星座图 .....	331
<b>第14章</b>	<b>多维标度法</b> .....	<b>334</b>
14.1	MDS 的基本理论和方法 .....	335
14.2	MDS 的古典解 .....	336

14.3	古典解的优良性	343
14.4	非度量方法	345
14.5	多维标度法的上机实现	347
14.6	社会经济案例研究	351
附录		355
附表 1	$T^2(p, n)$ 表	355
附表 2	$\Lambda_a(p, m_1, m_2)$ 表	361
附表 3	$L(p, v)$ 表	377
附表 4	$M(p, v_0, r)$ 表	379



## 第 1 章

# 多元正态分布

众所周知,一元正态分布在统计学的理论和实际应用方面都有重要的地位。同样,在多元统计学中,多元正态分布占有相当重要的位置。原因是:(1)许多实际问题研究中的随机向量确实遵从正态分布,或近似遵从正态分布;(2)对于多元正态分布,已有一整套统计推断方法,并且得到了许多完整的结果。

多元正态分布是最常用的一种多元概率分布。除此之外,还有多元对数正态分布、多项式分布、多元超几何分布、多元  $\beta$  分布、多元  $\chi^2$  分布、多元指数分布等。本章从多维变量及多元分布的基本概念开始,着重介绍多元正态分布的定义及一些重要性质。

### 本章目标

1. 掌握多元分布的有关概念。
2. 掌握统计距离及马氏距离的概念。
3. 理解多元正态分布的定义及其有关性质。
4. 了解 Wishart 分布的定义及其基本性质。

## 1.1 多元分布的基本概念

在研究社会、经济现象和许多实际问题时，经常遇到多指标的问题。例如研究职工工资构成情况时，计时工资、基础工资与职务工资、各种奖金、各种津贴等都是同时需要考察的指标；又如在研究公司的运营情况时，要涉及公司的资金周转能力、偿债能力、获利能力及竞争能力等财务指标，这些都是多指标研究的问题。显然，仅研究某个指标或是将这些指标割裂开来分别研究，都不能从整体上把握所研究问题的实质。一般地，假设我们所研究的问题涉及  $p$  个指标， $n$  次观测，这将会得到  $np$  个数据，我们的目的就是观测对象进行分组、分类，或分析、考察这  $p$  个变量之间的相互关联程度，或找出内在规律等等。下面我们简要介绍多元分析中涉及的一些基本概念。

### 一、随机向量

我们所讨论的是多个变量的总体，所研究的数据是同时观测  $p$  个指标（即变量），又进行了  $n$  次观测得到的，我们把这  $p$  个指标表示为  $X_1, X_2, \dots, X_p$ ，常用向量

$$\mathbf{X} = (X_1, X_2, \dots, X_p)'$$

表示对同一个体观测的  $p$  个变量。若观测了  $n$  个个体，则可得到如表 1—1 的数据，称每一个个体的  $p$  个变量为一个样品，而全体  $n$  个样品形成一个样本。

表 1—1

序号 \ 变量	$X_1$	$X_2$	...	$X_p$
1	$x_{11}$	$x_{12}$	...	$x_{1p}$
2	$x_{21}$	$x_{22}$	...	$x_{2p}$
⋮	⋮	⋮	...	⋮
$n$	$x_{n1}$	$x_{n2}$	...	$x_{np}$

横看表 1—1，记

$$\mathbf{X}_{(\alpha)} = (x_{\alpha 1}, x_{\alpha 2}, \dots, x_{\alpha p})', \alpha = 1, 2, \dots, n$$

表示第  $\alpha$  个样品的观测值。竖看表 1—1，第  $j$  列的元素

$$\mathbf{X}_j = (x_{1j}, x_{2j}, \dots, x_{nj})', j = 1, 2, \dots, p$$

表示对第  $j$  个变量  $X_j$  的  $n$  次观测数值。

因此，样本资料矩阵可用矩阵语言表示为：

$$\mathbf{X} = \begin{bmatrix} x_{11} & x_{12} & \cdots & x_{1p} \\ x_{21} & x_{22} & \cdots & x_{2p} \\ \vdots & \vdots & & \vdots \\ x_{n1} & x_{n2} & \cdots & x_{np} \end{bmatrix} = (\mathbf{X}_1, \mathbf{X}_2, \dots, \mathbf{X}_p) = \begin{bmatrix} \mathbf{X}'^{(1)} \\ \mathbf{X}'^{(2)} \\ \vdots \\ \mathbf{X}'^{(n)} \end{bmatrix}$$

若无特别说明，本书所称向量均指列向量。

**定义 1.1** 设  $X_1, X_2, \dots, X_p$  为  $p$  个随机变量，由它们组成的向量  $\mathbf{X} = (X_1, X_1, \dots, X_p)'$  称做随机向量。

## 二、分布函数与密度函数

描述随机变量的最基本工具是分布函数，类似地描述随机向量的最基本工具还是分布函数。

**定义 1.2** 设  $\mathbf{X} = (X_1, X_2, \dots, X_p)'$  是一随机向量，它的多元分布函数是

$$F(\mathbf{x}) = F(x_1, x_2, \dots, x_p) = P(X_1 \leq x_1, \dots, X_p \leq x_p) \quad (1.1)$$

式中， $\mathbf{x} = (x_1, x_2, \dots, x_p) \in R^p$ ，并记成  $\mathbf{X} \sim F$ 。

多元分布函数的有关性质此处从略。

**定义 1.3** 设  $\mathbf{X} \sim F(\mathbf{x}) = F(x_1, x_2, \dots, x_p)$ ，若存在一个非负的函数  $f(\cdot)$ ，使得

$$F(\mathbf{x}) = \int_{-\infty}^{x_1} \cdots \int_{-\infty}^{x_p} f(t_1, \dots, t_p) dt_1, \dots, dt_p \quad (1.2)$$

对一切  $\mathbf{x} \in R^p$  成立，则称  $\mathbf{X}$  (或  $F(\mathbf{x})$ ) 有分布密度  $f(\cdot)$ ，并称  $\mathbf{X}$  为连续型随机向量。

一个  $p$  维变量的函数  $f(\cdot)$  能作为  $R^p$  中某个随机向量的分布密度，当且仅当

$$(i) f(\mathbf{x}) \geq 0, \forall \mathbf{x} \in R^p$$

$$(ii) \int_{R^p} f(\mathbf{x}) d\mathbf{x} = 1$$

**【例 1.1】** 若随机向量  $X = (X_1, X_2, X_3)'$  有密度函数

$$f(x_1, x_2, x_3) = x_1^2 + 6x_3^2 + \frac{1}{3}x_1x_2$$

$$0 < x_1 < 1, \quad 0 < x_2 < 2, \quad 0 < x_3 < \frac{1}{2}$$

容易验证它符合分布密度函数的两个条件 (i) 和 (ii)。

最重要的连续型多元分布——多元正态分布将留在 1.3 节讨论。

### 三、多元变量的独立性

**定义 1.4** 两个随机向量  $X$  和  $Y$  称为是相互独立的, 若

$$P(X \leq x, Y \leq y) = P(X \leq x)P(Y \leq y) \quad (1.3)$$

对一切  $x, y$  成立。若  $F(x, y)$  为  $(X, Y)'$  的联合分布函数;  $G(x)$  和  $H(y)$  分别为  $X$  和  $Y$  的分布函数, 则  $X$  与  $Y$  独立当且仅当

$$F(x, y) = G(x)H(y) \quad (1.4)$$

若  $(X, Y)'$  有密度函数  $f(x, y)$ , 用  $g(x)$  和  $h(y)$  分别表示  $X$  和  $Y$  的分布密度, 则  $X$  和  $Y$  独立当且仅当

$$f(x, y) = g(x)h(y) \quad (1.5)$$

注意在上述定义中,  $X$  和  $Y$  的维数一般是不同的。

类似地, 若它们的联合分布等于各自分布的乘积, 称  $p$  个随机向量  $X_1, \dots, X_p$  相互独立。由  $X_1, \dots, X_p$  相互独立可以推知任何  $X_i$  与  $X_j$  ( $i \neq j$ ) 独立, 但是, 若已知任何  $X_i$  与  $X_j$  ( $i \neq j$ ) 独立, 并不能推出  $X_1, \dots, X_p$  相互独立。

### 四、随机向量的数字特征

#### 1. 随机向量 $X$ 的均值

设  $X = (X_1, \dots, X_p)'$  有  $p$  个分量。若  $E(X_i) = \mu_i$  ( $i = 1, 2, \dots, p$ ) 存在, 我们定义随机向量  $X$  的均值为:

$$E(X) = \begin{bmatrix} E(X_1) \\ E(X_2) \\ \vdots \\ E(X_p) \end{bmatrix} = \begin{bmatrix} \mu_1 \\ \mu_2 \\ \vdots \\ \mu_p \end{bmatrix} = \mu \quad (1.6)$$



$\mu$  是一个  $p$  维向量, 称为均值向量。

当  $A, B$  为常数矩阵时, 由定义可立即推出如下性质:

$$(1) E(\mathbf{AX}) = \mathbf{AE}(X) \quad (1.7)$$

$$(2) E(\mathbf{AXB}) = \mathbf{AE}(X)\mathbf{B} \quad (1.8)$$

## 2. 随机向量 $X$ 自协方差阵

$$\begin{aligned} \Sigma &= \text{cov}(\mathbf{X}, \mathbf{X}) = E(\mathbf{X} - E\mathbf{X})(\mathbf{X} - E\mathbf{X})' = D(\mathbf{X}) \\ &= \begin{bmatrix} D(X_1) & \text{cov}(X_1, X_2) & \cdots & \text{cov}(X_1, X_p) \\ \text{cov}(X_2, X_1) & D(X_2) & \cdots & \text{cov}(X_2, X_p) \\ \vdots & \vdots & & \vdots \\ \text{cov}(X_p, X_1) & \text{cov}(X_p, X_2) & \cdots & D(X_p) \end{bmatrix} \\ &= (\sigma_{ij}) \end{aligned} \quad (1.9)$$

称它为  $p$  维随机向量  $X$  的协方差阵, 简称为  $X$  的协方差阵。

称  $|\text{cov}(\mathbf{X}, \mathbf{X})|$  为  $X$  的广义方差, 它是协方差阵的行列式之值。

## 3. 随机向量 $X$ 和 $Y$ 的协方差阵

设  $\mathbf{X} = (X_1, \dots, X_n)'$  和  $\mathbf{Y} = (Y_1, \dots, Y_p)'$  分别为  $n$  维和  $p$  维随机向量, 它们之间的协方差阵定义为一个  $n \times p$  矩阵, 其元素是  $\text{cov}(X_i, Y_j)$ , 即

$$\text{cov}(\mathbf{X}, \mathbf{Y}) = (\text{cov}(X_i, Y_j), i=1, \dots, n; j=1, \dots, p) \quad (1.10)$$

若  $\text{cov}(\mathbf{X}, \mathbf{Y}) = 0$ , 称  $X$  和  $Y$  是不相关的。

当  $A, B$  为常数矩阵时, 由定义可推出协方差阵有如下性质:

$$(1) D(\mathbf{AX}) = \mathbf{AD}(\mathbf{X})\mathbf{A}' = \mathbf{A}\Sigma\mathbf{A}'$$

$$(2) \text{cov}(\mathbf{AX}, \mathbf{BY}) = \mathbf{A}\text{cov}(\mathbf{X}, \mathbf{Y})\mathbf{B}'$$

(3) 设  $\mathbf{X}$  为  $n$  维随机向量, 期望和协方差存在, 记  $\mu = E(\mathbf{X}), \Sigma = D(\mathbf{X})$ ,  $\mathbf{A}$  为  $n \times n$  常数阵, 则

$$E(\mathbf{X}'\mathbf{AX}) = \text{tr}(\mathbf{A}\Sigma) + \mu'\mathbf{A}\mu$$

对于任何随机向量  $\mathbf{X} = (X_1, \dots, X_p)'$  来说, 其协方差阵  $\Sigma$  都是对称阵, 同时总是非负定 (也称半正定) 的。大多数情形下是正定的。

## 4. 随机向量 $X$ 的相关阵

若随机向量  $\mathbf{X} = (X_1, \dots, X_p)'$  的协方差阵存在, 且每个分量的方差大于零, 则  $X$  的相关阵定义为: