

河北省省级精品课教材

应用多元统计分析

李春林 陈旭红 编著

清华大学出版社

应用多元统计分析

李春林 陈旭红 编著

清华大学出版社
北京

内 容 简 介

本书是在河北省精品课“多元统计分析”课程建设的基础上，贴近省属院校实际，以学生应用分析技能为主要培养目标，以方法、案例引导，对学生开展方法学习、案例分析、数据处理、结果讨论、文献阅读和论文撰写全方位的应用分析技能训练，是一本主要面向省属院校统计学各专业和其他相关专业的高年级本科生或研究生的应用型教材。

版权所有，侵权必究。侵权举报电话：010-62782989 13701121933

图书在版编目(CIP)数据

应用多元统计分析 / 李春林, 陈旭红编著. -北京: 清华大学出版社, 2013

ISBN 978-7-302-33296-1

I. ①应… II. ①李… ②陈… III. ①多元分析—统计分析—高等学校—教材
IV. ①O212.4

中国版本图书馆 CIP 数据核字(2013)第 168807 号

责任编辑：陈 明 洪 英

封面设计：傅瑞学

责任校对：刘玉霞

责任印制：杨 艳

出版发行：清华大学出版社

网 址：<http://www.tup.com.cn>, <http://www.wqbook.com>

地 址：北京清华大学学研大厦 A 座 邮 编：100084

社 总 机：010-62770175 邮 购：010-62786544

投稿与读者服务：010-62776969, c-service@tup.tsinghua.edu.cn

质量反馈：010-62772015, zhiliang@tup.tsinghua.edu.cn

印 装 者：三河市金元印装有限公司

经 销：全国新华书店

开 本：170mm×230mm 印 张：14.5 字 数：268 千字

版 次：2013 年 8 月第 1 版 印 次：2013 年 8 月第 1 次印刷

印 数：1~2000

定 价：29.00 元

产品编号：051546-01

前言

多元统计分析是统计学科中的一个重要分支，在自然科学、社会科学等领域具有广泛的应用，是探索多元世界强有力的工具。河北经贸大学的“多元统计分析”课程是统计学各专业的主干课程，是河北省的省级精品课程。在精品课程建设的过程中，我们结合丰富的教学、科研实践和大量鲜活的案例，贴近省属院校实际，以学生的应用分析技能为主要培养目标，以方法、案例引导进行多元统计分析方法的学习，对学生开展方法学习、案例分析、数据处理、结果讨论、文献阅读和论文撰写全方位的应用分析技能训练。

作为省属院校，我们切身体会到应用分析能力的培养对学生未来发展的重要性，也切实感受到国内纯应用性专业教材匮乏的无奈。因此，我们在建设省级精品课程的同时，结合科研和教学经验，紧贴应用分析技能培养这条省属院校学生培养与就业的生命线，编写了这本以应用为主线、以方法与软件相结合更好地解决实际问题为核心的《应用多元统计分析》教材。

本书用浅显的语言阐明各种多元统计方法的功能和原理，针对具体的案例，通过在国内广泛使用的统计分析软件 SPSS，讲授方法的上机实现和应用，尽可能详尽地介绍统计软件的各种操作选项和提供数据处理结果的解释，结合文献阅读和论文撰写对学生进行应用分析技能的培养。

本书涵盖了常用的多元统计分析方法，是一本主要面向省属院校统计学和经济学、管理学、生物医学统计等有关专业的高年级本科生或研究生的应用型教材和教学参考书，也可作为社会统计工作者和数据分析人员的实用参考书。

应用多元统计分析

本书在编写过程中,研究生孟杰、刘扬、冯丽红、李圣瑜、俱翠、胡一帆、王洪彪做了大量的基础性工作,清华大学出版社对教材的编写和出版给予了大力支持,陈明编辑为本书做了大量的组织工作,在此一并表示感谢!由于作者水平有限,书中难免出现疏漏和错误,希望广大读者提出宝贵意见,以便进一步修改。

李春林

2013年7月于石家庄

目 录

第 1 章 多元统计分析的理论基础	1
1.1 多元分布	1
1.1.1 随机向量	1
1.1.2 多元分布函数与密度函数	1
1.1.3 随机向量的数字特征	3
1.2 多元正态分布	4
1.2.1 多元正态分布的定义和性质	4
1.2.2 多元正态分布均值向量和协方差阵的估计	4
1.3 多元正态分布均值向量和协方差阵的检验	5
1.3.1 单总体均值向量的检验	5
1.3.2 多总体均值向量的检验	8
1.3.3 协方差阵的检验	11
1.3.4 多元正态分布均值向量和协方差阵检验的上机实现	16
习题	27
第 2 章 多元数据图	29
2.1 矩阵散点图	29
2.2 多维箱线图	31
2.3 雷达图	33
2.4 星形图	36
2.5 脸谱图	37
习题	39
第 3 章 数据预处理	41
3.1 数据集成与数据审核	41
3.1.1 数据集成	41

应用多元统计分析

3.1.2 数据审核	42
3.2 数据清理.....	43
3.2.1 缺失值数据	43
3.2.2 异常值数据	45
3.3 数据转换.....	48
3.3.1 数据标准化	49
3.3.2 数据的代数运算	50
3.3.3 数据的离散化	51
习题	52
第4章 因子分析	53
4.1 因子分析的基本理论.....	53
4.1.1 主成分分析的基本思想与模型	53
4.1.2 因子分析的基本思想与模型	55
4.1.3 因子分析的主要步骤	57
4.1.4 因子分析与主成分分析的区别与联系	58
4.2 因子分析的上机实现.....	59
4.2.1 因子分析的适用性检验	59
4.2.2 主因子个数的确定	60
4.2.3 因子旋转	61
4.2.4 因子得分	63
4.3 因子分析的案例分析.....	64
4.3.1 我国各地区社会发展状况的因子分析	64
4.3.2 我国制造业产业竞争力的因子分析	71
习题	76
第5章 聚类分析	77
5.1 聚类分析的基本理论.....	77
5.1.1 聚类分析的概念和基本思想	77
5.1.2 点与点之间的相似性度量方法	79
5.1.3 类与类的相似性度量方法	81
5.1.4 聚类的方法	83
5.2 聚类分析的上机实现.....	85
5.2.1 系统聚类方法	85
5.2.2 K 均值聚类	95

5.3 聚类分析应用实例	100
习题.....	103
第 6 章 判别分析.....	105
6.1 判别分析的基本理论	105
6.1.1 判别分析的概念和基本思想.....	105
6.1.2 距离判别法.....	106
6.1.3 费歇尔判别法.....	107
6.1.4 贝叶斯判别法.....	107
6.1.5 逐步判别法.....	108
6.2 判别分析的上机实现	108
6.3 判别分析的案例分析	124
习题.....	136
第 7 章 对应分析.....	139
7.1 对应分析的基本知识	139
7.2 对应分析的上机实现	144
7.3 多元对应分析的案例分析	150
习题.....	163
第 8 章 典型相关分析.....	167
8.1 典型相关分析的基本理论	168
8.1.1 典型相关分析的概念及基本思想.....	168
8.1.2 典型相关分析的数学描述.....	168
8.1.3 典型相关系数和典型变量的求解.....	169
8.1.4 典型相关系数的显著性检验.....	170
8.2 典型相关分析的上机实现	171
8.3 典型相关分析的案例分析	173
习题.....	181
第 9 章 回归分析.....	183
9.1 多元回归分析	183
9.1.1 多元回归分析概述.....	183
9.1.2 多元回归参数的估计与检验.....	185
9.1.3 多元回归案例的上机实现.....	188
9.2 多重多元回归分析	193
9.2.1 多重多元回归的数学模型.....	193

应用多元统计分析

9.2.2 多重多元回归分析的基本思想	195
9.2.3 多重多元回归分析的上机实现	197
习题	204
第 10 章 logistic 回归	209
10.1 logistic 回归模型的建立	209
10.1.1 logistic 函数及其性质	209
10.1.2 logistic 回归系数的意义	210
10.1.3 logistic 方程的检验	211
10.2 二项 logistic 回归模型的上机实现	213
10.2.1 上机操作	213
10.2.2 模型结果解释	217
10.3 多项 logistic 回归模型	219
10.3.1 多项 logistic 回归模型的建立及上机实现	219
10.3.2 多项 logistic 回归模型与二项 logistic 回归模型的转换	222
习题	223
参考文献	224
社会学类	224
教育学类	224
心理学类	224
医学类	224
其他类	224
统计学类	224
管理学类	224
经济学类	224
法学类	224
文学类	224
历史学类	224
哲学类	224
理学类	224
工学类	224
农学类	224
军事学类	224
其他类	224

CHAPTER



第1章 多元统计分析的理论基础

1.1 多元分布

1.1.1 随机向量

多元统计分析讨论的是多变量总体。把 p 个随机变量放在一起得 $\mathbf{X} = (X_1, X_2, \dots, X_p)$ 是一个 p 维随机向量, 如果同时对 p 个变量作一次观测, 则得到一个样本 $(x_{11}, x_{12}, \dots, x_{1p})$, n 次观测得到的样本排成一个 $n \times p$ 矩阵, 称为样本数据矩阵(或样本资料矩阵), 记为

$$\mathbf{X} = \begin{bmatrix} x_{11} & x_{12} & \cdots & x_{1p} \\ x_{21} & x_{22} & \cdots & x_{2p} \\ \vdots & \vdots & \ddots & \vdots \\ x_{n1} & x_{n2} & \cdots & x_{np} \end{bmatrix} \stackrel{\text{def}}{=} \begin{bmatrix} \mathbf{X}_{(1)} \\ \mathbf{X}_{(2)} \\ \vdots \\ \mathbf{X}_{(n)} \end{bmatrix} \stackrel{\text{def}}{=} (\mathbf{X}_1, \mathbf{X}_2, \dots, \mathbf{X}_p) \quad (1-1)$$

矩阵 \mathbf{X} 的第 i 行: $\mathbf{X}_{(i)} = (x_{i1}, x_{i2}, \dots, x_{ip})$ ($i=1, 2, \dots, n$) 表示对第 i 个样本的观

测值, 是一个 p 维的随机向量。矩阵 \mathbf{X} 的第 j 列: $\mathbf{X}_j = \begin{bmatrix} x_{1j} \\ x_{2j} \\ \vdots \\ x_{nj} \end{bmatrix}$ ($j=1, 2, \dots, p$) 表示对第 j 个变量的 n 次观测, 是一个 n 维随机向量。

1.1.2 多元分布函数与密度函数

1. 联合分布

设 $\mathbf{X} = (X_1, X_2, \dots, X_p)$ 是 p 维随机向量, 称 p 元函数

应用多元统计分析

$$F(x_1, x_2, \dots, x_p) = P\{X_1 \leq x_1, X_2 \leq x_2, \dots, X_p \leq x_p\} \quad (1-2)$$

为 \mathbf{X} 的联合分布函数。

若存在非负函数 $f(x_1, x_2, \dots, x_p)$, 使得随机向量 \mathbf{X} 的联合分布函数对一切 $(x_1, x_2, \dots, x_p) \in \mathbb{R}^p$ 均可表示为

$$F(x_1, x_2, \dots, x_p) = \int_{-\infty}^{x_1} \int_{-\infty}^{x_2} \cdots \int_{-\infty}^{x_p} f(t_1, t_2, \dots, t_p) dt_1 dt_2 \cdots dt_p \quad (1-3)$$

则称 $f(x_1, x_2, \dots, x_p)$ 为连续型随机向量 \mathbf{X} 的联合概率密度函数, 简称为多元密度函数或密度函数。

多元密度函数 $f(x_1, x_2, \dots, x_p)$ 满足以下两条性质:

(1) 对一切实数 x_1, x_2, \dots, x_p , $f(x_1, x_2, \dots, x_p) \geq 0$;

(2) $\int_{-\infty}^{\infty} \int_{-\infty}^{\infty} \cdots \int_{-\infty}^{\infty} f(x_1, x_2, \dots, x_p) dx_1 dx_2 \cdots dx_p = 1$.

2. 边缘分布

称随机向量 \mathbf{X} 的部分分量 $(X_{i_1}, \dots, X_{i_m})$ ($1 \leq m < p$) 的分布为边缘分布。

设 $\mathbf{X}^{(1)} = (X_{i_1}, \dots, X_{i_r})$ 为 r 维随机向量, $\mathbf{X}^{(2)} = (X_{i_{r+1}}, \dots, X_{i_p})$ 为 $p-r$ 维随机向量。若 $\mathbf{X} = (\mathbf{X}^{(1)}, \mathbf{X}^{(2)})$, 则 $\mathbf{X}^{(1)}$ 的边缘分布密度函数为

$$\begin{aligned} f_1(x^{(1)}) &= f_1(x_{i_1}, \dots, x_{i_r}) \\ &= \int_{-\infty}^{\infty} \cdots \int_{-\infty}^{\infty} f(x_1, x_2, \dots, x_p) dx_{i_{r+1}} \cdots dx_{i_p} \end{aligned} \quad (1-4)$$

$\mathbf{X}^{(2)}$ 的边缘分布密度函数为

$$\begin{aligned} f_2(x^{(2)}) &= f_2(x_{i_{r+1}}, \dots, x_{i_p}) \\ &= \int_{-\infty}^{\infty} \cdots \int_{-\infty}^{\infty} f(x_1, x_2, \dots, x_p) dx_{i_1} \cdots dx_{i_r} \end{aligned} \quad (1-5)$$

3. 条件分布

称给定 $\mathbf{X}^{(2)}$ 时 $\mathbf{X}^{(1)}$ 的分布为条件分布。当 \mathbf{X} 的密度函数为 $f(x^{(1)}, x^{(2)})$ 、 $\mathbf{X}^{(2)}$ 的密度函数为 $f_2(x^{(2)})$ 时, 则给定 $\mathbf{X}^{(2)}$ 时 $\mathbf{X}^{(1)}$ 的条件密度函数为

$$f_1(x^{(1)} | x^{(2)}) = f(x^{(1)}, x^{(2)}) / f_2(x^{(2)}) \quad (1-6)$$

4. 独立性

设 X_1, \dots, X_p 是 p 个随机变量, X_i 的分布函数记为 $F_i(x_i)$ ($i=1, \dots, p$),

$F(x_1, \dots, x_p)$ 是 $\mathbf{X} = (X_1, X_2, \dots, X_p)$ 的联合分布函数。若对一切实数 x_1, x_2, \dots, x_p ,

$$F(x_1, \dots, x_p) = F_1(x_1) \cdots F_p(x_p) \quad (1-7)$$

均成立, 则称 X_1, \dots, X_p 相互独立。

在连续型随机变量的情况下, X_1, \dots, X_p 相互独立, 当且仅当 $\mathbf{X} = (X_1, X_2, \dots, X_p)$ 的联合密度函数 $f(x_1, x_2, \dots, x_p)$ 满足: 对一切实数 x_1, x_2, \dots, x_p ,

$$f(x_1, x_2, \dots, x_p) = f_1(x_1)f_2(x_2) \cdots f_p(x_p) \quad (1-8)$$

均成立, 其中 $f_i(x_i)$ 是 X_i 的密度函数 ($i=1, \dots, p$)。

1.1.3 随机向量的数字特征

设 $\mathbf{X} = (X_1, X_2, \dots, X_p), \mathbf{Y} = (Y_1, Y_2, \dots, Y_q)$ 是两个随机向量。

1. 均值向量

若 $E(X_i) = \mu_i$ 存在, 则称 $E(\mathbf{X}) = (E(X_1), E(X_2), \dots, E(X_p)) = (\mu_1, \mu_2, \dots, \mu_p)$ 为随机向量 \mathbf{X} 的均值向量。

2. 协方差阵

若 X_i 和 Y_j 的协方差 $\text{cov}(X_i, Y_j)$ 存在 ($i=1, 2, \dots, p; j=1, 2, \dots, q$), 则称

$$\text{cov}(\mathbf{X}, \mathbf{Y}) = E[(\mathbf{X} - E\mathbf{X})(\mathbf{Y} - E\mathbf{Y})]$$

$$= \begin{bmatrix} \text{cov}(X_1, Y_1) & \text{cov}(X_1, Y_2) & \cdots & \text{cov}(X_1, Y_q) \\ \text{cov}(X_2, Y_1) & \text{cov}(X_2, Y_2) & \cdots & \text{cov}(X_2, Y_q) \\ \vdots & \vdots & & \vdots \\ \text{cov}(X_p, Y_1) & \text{cov}(X_p, Y_2) & \cdots & \text{cov}(X_p, Y_q) \end{bmatrix} \quad (1-9)$$

为随机向量 \mathbf{X} 和 \mathbf{Y} 的协方差阵。若 $\text{cov}(\mathbf{X}, \mathbf{Y}) = \mathbf{O}$ (其中 \mathbf{O} 表示零矩阵), 则称 \mathbf{X} 和 \mathbf{Y} 不相关。若 $\mathbf{Y} = \mathbf{X}$, 则 $\text{cov}(\mathbf{X}, \mathbf{X}) = D(\mathbf{X})$ 为 \mathbf{X} 的(协)方差阵。

3. 相关阵

记 $r_{ij} = \frac{\text{cov}(X_i, X_j)}{\sqrt{\text{var}(X_i)} \sqrt{\text{var}(X_j)}} = \frac{\sigma_{ij}}{\sqrt{\sigma_{ii}\sigma_{jj}}} \quad (i, j = 1, 2, \dots, p)$, 称 $\mathbf{R} = (r_{ij})_{p \times p}$

为 \mathbf{X} 的相关阵, r_{ij} 也称为分量 X_i 与 X_j 之间的(线性)相关系数。

1.2 多元正态分布

1.2.1 多元正态分布的定义和性质

1. 多元正态分布的定义

定义 1(密度定义) 若 p 元随机向量 $\mathbf{X} = (X_1, X_2, \dots, X_p)$ 的概率密度函数为

$$f(x_1, x_2, \dots, x_p) = \frac{1}{(2\pi)^{p/2} |\Sigma|^{1/2}} \exp \left\{ -\frac{1}{2} (\mathbf{X} - \boldsymbol{\mu})^T \boldsymbol{\Sigma}^{-1} (\mathbf{X} - \boldsymbol{\mu}) \right\} \quad (1-10)$$

其中, $\boldsymbol{\mu}$ 为 $p \times 1$ 的常数向量, $\boldsymbol{\Sigma}$ 为正定阵 ($\boldsymbol{\Sigma} > \mathbf{O}$), $|\boldsymbol{\Sigma}|$ 为协差阵 $\boldsymbol{\Sigma}$ 的行列式, 则称 \mathbf{X} 服从 p 元正态分布, 也称 \mathbf{X} 为 p 元正态变量。记为 $\mathbf{X} \sim N_p(\boldsymbol{\mu}, \boldsymbol{\Sigma})$ 。

当 $p=2$ 时, 可以得到二元正态分布的密度公式。

2. 多元正态分布的性质

性质 1 若 $\mathbf{X} \sim N_p(\boldsymbol{\mu}, \boldsymbol{\Sigma})$, 则 $E(\mathbf{X}) = \boldsymbol{\mu}$, $D(\mathbf{X}) = \boldsymbol{\Sigma}$ 。

本性质将正态分布的参数 $\boldsymbol{\mu}$ 和 $\boldsymbol{\Sigma}$ 赋予了明确的统计意义。

性质 2 设 Y_1, Y_2, \dots, Y_k 相互独立, $Y_i \sim N_p(\boldsymbol{\mu}_i, \boldsymbol{\Sigma}_i)$ ($\forall i$), 则

$$\sum_{i=1}^k Y_i \sim N_p \left(\sum_{i=1}^k \boldsymbol{\mu}_i, \sum_{i=1}^k \boldsymbol{\Sigma}_i \right) \quad (1-11)$$

性质 3 如果正态随机向量 $\mathbf{X} = (X_1, X_2, \dots, X_p)$ 的协方差阵 $\boldsymbol{\Sigma}$ 是对角阵, 则 \mathbf{X} 的各分量是相互独立的随机变量。

性质 4 多元正态向量 $\mathbf{X} = (X_1, X_2, \dots, X_p)$ 的任意线性变换仍然服从多元正态分布。

1.2.2 多元正态分布均值向量和协方差阵的估计

1. 均值向量的估计

设样本 $\mathbf{X}_{(1)}, \mathbf{X}_{(2)}, \dots, \mathbf{X}_{(n)}$ 间相互独立, 且服从于 p 元正态分布 $N_p(\boldsymbol{\mu}, \boldsymbol{\Sigma})$, $n > p$, $\boldsymbol{\Sigma} > \mathbf{O}$, 则总体参数均值 $\boldsymbol{\mu}$ 的估计量是

$$\begin{aligned}\hat{\mu} &= \bar{\mathbf{X}} = \frac{1}{n} \sum_{i=1}^n \mathbf{X}_{(i)} = \frac{1}{n} \left(\sum_{i=1}^n X_{i1}, \sum_{i=1}^n X_{i2}, \dots, \sum_{i=1}^n X_{ip} \right) \mathbf{X} \\ &= (\bar{X}_1, \bar{X}_2, \dots, \bar{X}_p)\end{aligned}\quad (1-12)$$

均值向量 μ 的估计量就是样本均值向量,且 $\hat{\mu}$ 是 μ 的无偏估计。

2. 协方差阵的估计

总体参数协方差阵 Σ 的极大似然估计是

$$\begin{aligned}\hat{\Sigma}_p &= \frac{1}{n} \mathbf{L} = \frac{1}{n} \sum_{i=1}^n (\mathbf{X}_{(i)} - \bar{\mathbf{X}})(\mathbf{X}_{(i)} - \bar{\mathbf{X}})^T \\ &= \frac{1}{n} \begin{bmatrix} \sum_{i=1}^n (X_{i1} - \bar{X}_1)^2 & \sum_{i=1}^n (X_{i1} - \bar{X}_1)(X_{i2} - \bar{X}_2) & \cdots & \sum_{i=1}^n (X_{i1} - \bar{X}_1)(X_{ip} - \bar{X}_p) \\ \vdots & \ddots & \ddots & \vdots \\ \sum_{i=1}^n (X_{i2} - \bar{X}_2)^2 & \cdots & \sum_{i=1}^n (X_{i2} - \bar{X}_2)(X_{ip} - \bar{X}_p) \\ \vdots & \ddots & \ddots & \vdots \\ \sum_{i=1}^n (X_{ip} - \bar{X}_p)^2 & & & \end{bmatrix}\end{aligned}\quad (1-13)$$

其中 \mathbf{L} 是离差阵,它是每一个样本(向量)与样本均值(向量)的离差积形成的 n 个 p 阶对称阵的和。同一元相似, $\hat{\Sigma}_p$ 不是 Σ 的无偏估计,为了得到无偏估计我们常用样本协差阵 $\hat{\Sigma} = \frac{1}{n-1} \mathbf{L}$ 作为总体协差阵的估计。

1.3 多元正态分布均值向量和协方差阵的检验

1.3.1 单总体均值向量的检验

设总体 $\mathbf{X} \sim N_p(\mu, \Sigma)$, 随机样本 $\mathbf{X}_{(\alpha)}$ ($\alpha = 1, \dots, n$)。检验

$$H_0: \mu = \mu_0 \quad (\mu_0 \text{ 为已知向量}), \quad H_1: \mu \neq \mu_0$$

1. 当 $\Sigma = \Sigma_0$ 已知时均值向量的检验

取检验统计量为

$$T_0^2 = n(\bar{\mathbf{X}} - \mu_0)^T \Sigma_0^{-1} (\bar{\mathbf{X}} - \mu_0) \stackrel{H_0}{\sim} \chi^2(p) \quad (1-14)$$

应用多元统计分析

按传统的检验方法,对给定的显著性水平 α ,查 χ^2 分布临界值表得 λ_α ,使 $P\{T_0^2 > \lambda_\alpha\} = \alpha$,则否定域为 $\{T_0^2 > \lambda_\alpha\}$ 。

由样本值 $x_{(a)} (a=1, \dots, n)$,计算 \bar{X} 及 T_0^2 值,若 $T_0^2 > \lambda_\alpha$,则否定 H_0 ,否则 H_0 相容。

利用统计软件还可以通过计算显著性概率值(p 值)给出检验结果,且由此得出的结果更丰富。

假设在 H_0 成立情况下,随机变量 $T_0^2 \sim \chi^2(p)$,由样本值计算得到 T_0^2 的值为 d ,同时可以计算以下概率值

$$p = P\{T_0^2 \geq d\} \quad (1-15)$$

常称此概率值为显著性概率值,或简称为 p 值。

对给定的显著性水平 α ,当 $p < \alpha$ 时,则在显著性水平 α 下否定假设 H_0 。在这种情况下,可能犯“以真当假”的第一类错误,且 α 就是犯第一类错误的概率。

当 $p \geq \alpha$ 时,在显著性水平 α 下 H_0 相容。在这种情况下,可能犯“以假当真”的第二类错误,且犯下第二类错误的概率 β 为

$$\beta = P\{T_0^2 \leq \lambda_\alpha \mid \mu = \mu_1 \neq \mu_0\} \quad (1-16)$$

其中检验统计量 $T_0^2 \sim \chi^2(p, \delta)$,非中心参数

$$\delta = n(\mu_1 - \mu_0)^T \Sigma_0^{-1} (\mu_1 - \mu_0) \quad (1-17)$$

p 值的直观含义可以这样看,检验统计量 T_0^2 的大小反映 \bar{X} 与 μ_0 的偏差大小,当 H_0 成立时 T_0^2 的值应该较小。现由观测数据计算 T_0^2 值为 d ;当 H_0 成立时统计量 $T_0^2 \sim \chi^2(p)$,由 χ^2 分布可以计算该统计量 $\geq d$ 的概率值(即 p 值)。比如 $p = 0.02 < \alpha = 0.05$,这时出现一个比较小的概率标准($\alpha = 0.05$)还要小的事件 $\{T_0^2 \geq d\}$ 。也就是说,在 $\mu = \mu_0$ 假设下,观测数据中极少情况会出现 T_0^2 的值大于等于 d 值,故在 0.05 显著性水平下有足够的证据否定原假设,即认为 μ 与 μ_0 有显著的差异。

又比如当 $p = 0.22 \geq \alpha = 0.05$ 时,表示在 $\mu = \mu_0$ 假设下,观测数据中经常会出现 T_0^2 的值大于等于 d 值的情况,故在 0.05 显著性水平下没有足够的证据否定原假设,即认为 μ 与 μ_0 没有显著的差异。

2. 当 Σ 未知时均值向量的检验

考虑统计量

$$\begin{aligned} T^2 &= (n-1)[\sqrt{n}(\bar{X} - \mu_0)]^T A^{-1} [\sqrt{n}(\bar{X} - \mu_0)] \\ &= (n-1)n(\bar{X} - \mu_0)^T A^{-1} (\bar{X} - \mu_0) \sim T^2(p, n-1) \end{aligned} \quad (1-18)$$

再利用 T^2 与 F 分布的关系,统计量取为

$$F = \frac{(n-1)-p+1}{(n-1)p} T^2 \sim F(p, (n-1)-p+1) \\ = F(p, n-p) \quad (1-19)$$

进行检验。

3. 似然比统计量

在数理统计中关于总体参数的假设检验,通常是利用最大似然原理导出似然比统计量来进行检验。在多元统计分析中几乎所有重要的检验都是利用最大似然比原理给出的。下面我们回顾一下最大似然比原理。

设 p 元总体的密度函数为 $f(x, \theta)$, 其中 θ 是未知参数,且 $\theta \in \Theta$ (参数空间),又设 Θ_0 是 Θ 的子集,我们希望对下列假设:

$$H_0: \theta \in \Theta_0, \quad H_1: \theta \notin \Theta_0$$

作出判断,这就是假设检验问题。称 H_0 为原假设(或零假设), H_1 为对立假设(或备择假设)。

从总体 X 抽取容量为 n 的样本 $X_{(t)} (t=1, \dots, n)$ 。把样本的联合密度函数

$$L(x_{(1)}, \dots, x_{(n)}; \theta) = \prod_{t=1}^n f(x_{(t)}; \theta) \quad (1-20)$$

记为 $L(X; \theta)$, 并称它为样本的似然函数。

引入统计量

$$\lambda = \max_{\theta \in \Theta_0} L(X; \theta) / \max_{\theta \in \Theta} L(X; \theta) \quad (1-21)$$

它是样本 $X_{(t)} (t=1, \dots, n)$ 的函数,常称 λ 为似然比统计量。由于 $\Theta_0 \subset \Theta$, 从而 $0 \leq \lambda \leq 1$ 。

由最大似然比原理知,如果 λ 取值太小,说明 H_0 为真时观测到此样本 $X_{(t)} (t=1, \dots, n)$ 的概率比 H_0 为不真时观测到此样本 $X_{(t)} (t=1, \dots, n)$ 的概率要小得多,故有理由认为假设 H_0 不成立,所以从似然比出发,以上检验问题的否定域为

$$\{\lambda(X_{(1)}, \dots, X_{(n)}) < \lambda_\alpha\} \quad (1-22)$$

按传统的检验方法, λ_α 是由显著性水平 α 确定的临界值, 它满足当 H_0 成立时使得

$$P\{\lambda(X_{(1)}, \dots, X_{(n)}) < \lambda_\alpha\} = \alpha \quad (1-23)$$

为了得到 λ_α , 必须研究似然比统计量 λ 的抽样分布。在一些特殊的情况下,可以得到 λ 的精确分布;但在很多情况下是得不到 λ 的精确分布的。当样本量很大且满足一定正则条件时, $-2\ln\lambda$ 的抽样分布与 χ^2 分布十分接近。下面不加证明地给出一条很有用的结论。

定理 1 当样本容量 n 很大时,

$$-2\ln\lambda = -2\ln \left[\frac{\max L(\mathbf{X}; \theta_0)}{\max_{\theta \in \Theta} L(\mathbf{X}; \theta)} \right] \quad (1-24)$$

近似服从自由度为 f 的 χ^2 分布, 其中 $f = \Theta_0$ 的维数 - Θ_0 的维数。

本章将讨论的一些检验问题, 就是利用似然比统计量的近似分布进行检验的方法。

当 Σ 未知时检验均值向量 $\mu = \mu_0$ 的似然比统计量及其分布。

设样本的似然函数为 $L(\mu, \Sigma)$ 。检验均值向量 $\mu = \mu_0$ 的似然比统计量为

$$\lambda = \frac{\max_{\mu=\mu_0, \Sigma > 0} L(\mu_0, \Sigma)}{\max_{\mu, \Sigma > 0} L(\mu, \Sigma)} \quad (1-25)$$

否定域为

$$\{\lambda < \lambda_\alpha\} \Leftrightarrow \{T^2 > T_\alpha^2\} \Leftrightarrow \{F > F_\alpha\}$$

其中

$$F = \frac{n-p}{p} \frac{T^2}{n-1} \sim F(p, n-p) \quad (1-26)$$

1.3.2 多总体均值向量的检验

1. 两正态总体均值向量的检验

1) 两总体协方差阵相等(但未知)时均值向量的检验

设 $\mathbf{X}_{(\alpha)} (\alpha = 1, \dots, n)$ 为来自总体 $\mathbf{X} \sim N_{(p)}(\boldsymbol{\mu}^{(1)}, \Sigma)$ 的随机样本; $\mathbf{Y}_{(\alpha)} (\alpha = 1, \dots, m)$ 为来自样本总体 $\mathbf{Y} \sim N_{(p)}(\boldsymbol{\mu}^{(2)}, \Sigma)$ 的随机样本, 且相互独立, Σ 未知。检验

$$H_0: \boldsymbol{\mu}^{(1)} = \boldsymbol{\mu}^{(2)}, \quad H_1: \boldsymbol{\mu}^{(1)} \neq \boldsymbol{\mu}^{(2)}$$

当 $p=1$ 时, 取检验统计量为

$$t = \frac{(\bar{\mathbf{X}} - \bar{\mathbf{Y}}) \left/ \sqrt{\frac{1}{n} + \frac{1}{m}} \right.}{\sqrt{\left[\sum_{i=1}^n (\mathbf{X}_{(i)} - \bar{\mathbf{X}})^2 + \sum_{j=1}^m (\mathbf{Y}_{(j)} - \bar{\mathbf{Y}})^2 \right] / (n+m-2)}} \sim t(n+m-2) \quad (1-27)$$

即

$$t^2 = \frac{nm}{m+n} (\bar{\mathbf{X}} - \bar{\mathbf{Y}})^T \left[\frac{\sum_{i=1}^n (\mathbf{X}_{(i)} - \bar{\mathbf{X}})^2 + \sum_{j=1}^m (\mathbf{Y}_{(j)} - \bar{\mathbf{Y}})^2}{n+m-2} \right]^{-1} (\bar{\mathbf{X}} - \bar{\mathbf{Y}})$$