



# 若干混合效应模型的 统计推断研究

叶仁道 罗堃 著



科学出版社

# 若干混合效应模型的统计推断研究

叶仁道 罗堃 著



科学出版社  
北京

## 内 容 简 介

本书系统阐述若干混合效应模型的基本理论、方法和应用。全书共 8 章。第 1 章通过实例引入平衡随机效应模型、非平衡随机效应模型、面板数据模型等混合效应模型。第 2 章讨论矩阵方面的补充知识, 以及多元正态分布、多元偏正态分布及广义推断方法等。第 3~第 8 章详细讨论几类混合效应模型的基本理论与方法, 包括一般平衡随机效应模型、非平衡单向分类随机效应模型、非平衡双向分类随机效应模型、面板数据模型、混合效应模型、偏正态混合效应模型等。本书凸显混合效应模型小样本推断研究特色, 并将其研究结论广泛应用于实际案例。

本书可以作为高等学校统计系、生物统计系、数学科学系、数量经济系等相关专业的高年级本科生及研究生的学位课或选修课教材, 以及数学、生物、医学、经济、金融、工程等领域的教师或者科技工作者的参考书。

### 图书在版编目(CIP)数据

若干混合效应模型的统计推断研究 / 叶仁道, 罗堃著.  
—北京: 科学出版社, 2016

ISBN 978-7-03-046685-3

I. ①若… II. ①叶…②罗… III. ①统计推断—研究 IV.  
①0212

中国版本图书馆 CIP 数据核字(2015)第 306732 号

责任编辑: 魏如萍 / 责任校对: 张海燕  
责任印制: 徐晓晨 / 封面设计: 蓝正设计

科学出版社 出版

北京东黄城根北街 16 号

邮政编码: 100717

<http://www.sciencep.com>

北京中石油彩色印刷有限责任公司 印刷

科学出版社发行 各地新华书店经销

\*

2016 年 3 月第 一 版 开本: 720×1000 1/16

2016 年 3 月第一次印刷 印张: 13 1/4

字数: 281 000

定价: 72.00 元

(如有印装质量问题, 我社负责调换)

# 前 言

混合效应模型是现代统计学中一类重要模型，广泛应用于经济、金融、生物、医学、地质、气象、农业、工业、工程技术等众多领域。对于这类模型，传统的基于精确分布的统计方法几乎是不可能得到的，而基于大样本理论的近似分布又往往与真实分布偏离较远。因此，现有关于混合效应模型参数估计、检验和预测等方面的基本理论和传统方法，仍难以对现实中许多复杂数据和模型做出准确分析。鉴于此，本书对混合效应模型的统计推断理论做出深入探讨。

针对一般平衡随机效应模型、非平衡单向分类随机效应模型、非平衡两向分类随机效应模型、面板数据模型、混合效应模型、偏正态混合效应模型等，本书运用矩法、谱分解方法、矩阵技术、结构方法、Monte Carlo 方法等多种统计方法，研究感兴趣参数的优良估计、精确检验等统计推断问题，建立一系列新的有效的统计推断方法，并将其应用到经济、金融、生物、医学等领域的实际数据分析中，从而实现混合效应模型的应用价值。

借本书出版之际，我们要感谢曾经给予过我们帮助的导师、朋友及学生。首先，我们要特别感谢导师王松桂教授，感谢他对我们多年来研究工作所给予的关心、指导和帮助。其次，要感谢美国新墨西哥州立大学数学科学系的王通会教授，感谢他在我访问该校期间所提供的热情指导与关照。最后，要感谢我的研究生徐立军同学，协助我们制作本书中所涉及的图形与表格，并负责部分章节的打字工作。

本书得到国家自然科学基金(11401148)、国家社会科学基金(12CJY012)、国家统计局重点项目(2015LZ14)、教育部人文社科研究项目(14YJC910005)、浙江省自然科学基金(LY14A010030、Y6110017)、浙江省信息化与经济社会发展研究中心及杭州电子科技大学校级重点学科“应用统计学”专著出版基金的资助，在此表示诚挚的谢意。

由于编者水平所限，书中若有不足之处还望国内同行及广大读者不吝赐教。

联系方式: yerendao@hdu.edu.cn。

编 者

2015年10月20日

# 目 录

第 1 章 绪论	1
1.1 模型概论	1
1.2 国内外研究现状	4
1.3 研究方法与思路	7
1.4 研究特色与创新之处	9
1.5 研究内容与框架	10
第 2 章 预备知识	12
2.1 矩阵知识	12
2.2 多元正态分布	19
2.3 多元偏正态分布	23
2.4 矩阵维偏正态分布	32
2.5 广义 $p$ -值和广义置信区间	40
第 3 章 一般平衡随机效应模型	43
3.1 方差分量的非负估计	44
3.2 方差分量广义推断	53
3.3 过程能力指数广义推断	59
第 4 章 非平衡单向分类随机效应模型	71
4.1 组间方差分量的区间估计	71
4.2 暴露水平评价	84
4.3 可靠性推断	89
第 5 章 非平衡两向分类随机效应模型	98
5.1 暴露水平评价	98
5.2 组内相关系数推断	109
第 6 章 面板数据模型	124
6.1 回归系数广义推断	125
6.2 方差分量广义推断	131
第 7 章 混合效应模型	136
7.1 方差分量广义推断	136

7.2 协方差阵估计 .....	143
<b>第 8 章 偏正态混合效应模型</b> .....	<b>158</b>
8.1 含两个随机效应的偏正态混合效应模型 .....	158
8.2 含 $k+1$ 个随机效应的偏正态混合效应模型 .....	178
参考文献.....	189
符号表.....	198
索引.....	200

# 第1章 绪 论

混合效应模型在经济、金融、生物、医学、地质、气象、农业、工业、工程技术等领域具有非常广泛的应用背景，上述领域的许多现象或问题都可以借助混合效应模型做出较好的解释。与传统模型相比，混合效应模型引入了随机效应，能够充分描述数据之间的相关信息，从而大大提高模型的精度。因此，这类模型能够较为有效地处理结构复杂的数据，如纵向(longitudinal)数据、成组(clustered)数据、空间(spatial)数据等。为此，混合效应模型已成为现代统计学中应用最为广泛的模型之一。本著作将系统探讨若干混合效应模型的基本理论、方法和应用。

本章首先通过实例引入各类混合效应模型，进而介绍其国内外研究现状、研究方法思路、研究特色与创新之处、研究内容与框架等，使读者对这类模型的丰富实际背景及相关研究有一些了解。

## 1.1 模型概论

混合效应模型一般表达式为

$$y = X\beta + Z\epsilon + \epsilon_0 \quad (1.1)$$

其中， $y$  为  $n \times 1$  观测向量； $X$  和  $Z$  分别为  $n \times p$  和  $n \times q$  已知设计矩阵； $\beta$  为  $p \times 1$  未知参数向量，称为固定效应； $\epsilon$  为  $q \times 1$  随机向量，称为随机效应； $\epsilon_0$  为  $n \times 1$  随机误差向量。一般我们假定  $E(\epsilon) = 0$ ， $\text{Cov}(\epsilon) = D$ ， $E(\epsilon_0) = 0$ ， $\text{Cov}(\epsilon_0) = D_0$ ，且  $\epsilon$  和  $\epsilon_0$  互不相关。其中， $D$  为非负定矩阵， $D_0$  为正定矩阵。则

$$E(y) = X\beta, \text{Cov}(y) = ZDZ' + D_0$$

对  $D$  和  $D_0$  的不同假定，则可得到不同类型的混合效应模型。假设它们依赖于一个未知参数向量，则模型(1.1)的随机部分  $Z\epsilon + \epsilon_0$  可以分解为

$$Z\epsilon + \epsilon_0 = \sum_{i=1}^k Z_i \epsilon_i + \epsilon_0$$

其中,  $Z_i$  为  $n \times q_i$  的已知设计矩阵;  $\epsilon_i$  为  $q_i \times 1$  的随机效应, 且  $\epsilon_i$  和  $\epsilon_j$  互不相关 ( $i \neq j$ )。此时, 我们一般假定

$$E(\epsilon_i) = 0, \text{Cov}(\epsilon_i) = \sigma_i^2 I_{q_i}, \text{Cov}(\epsilon_i, \epsilon_j) = 0, i \neq j$$

其中,  $I_a$  为  $a$  阶单位阵, 在不致混淆情况下, 可以略去下标  $a$ 。于是, 模型(1.1)可进一步表示为

$$y = X\beta + \sum_{i=1}^k Z_i \epsilon_i + \epsilon_0 \quad (1.2)$$

进而有

$$E(y) = X\beta, \text{Cov}(y) = \sum_{i=0}^k \sigma_i^2 Z_i Z_i' \triangleq \Sigma(\sigma^2)$$

其中,  $\sigma^2 = (\sigma_0^2, \dots, \sigma_k^2)'$ ,  $\sigma_i^2$  称为方差分量, 因此模型(1.2)也称为方差分量模型。在模型(1.2)中, 若  $X$  和  $Z_i$  为由 0 和 1 构成的指示矩阵, 则模型(1.2)称为混合方差分析模型; 若固定效应只有常数项, 即  $X\beta = 1_n \mu$ , 其中  $1_n$  表示分量均为 1 的  $n$  维列向量,  $\mu$  为总体均值, 则模型(1.2)称为一般随机效应模型。模型(1.2)还包含多类具有广泛应用背景的混合效应模型, 如面板数据模型、单向(两向、多向)分类混合效应模型、套分类混合效应模型等。

### 例 1.1 两向分类随机效应模型。

研究人的平均血压在一天内的变化规律。在一天内选择  $a$  个时间点测量被观测者的血压, 假定观测了  $b$  个人, 用  $y_{ij}$  表示第  $i$  个时间点的第  $j$  个人的血压, 则  $y_{ij}$  可以表示为

$$y_{ij} = \mu + \alpha_i + \beta_j + e_{ij}, i = 1, 2, \dots, a; j = 1, 2, \dots, b \quad (1.3)$$

其中,  $\alpha_i$  为第  $i$  个时间点对血压的影响, 称为第  $i$  个时间点的效应;  $\beta_j$  为第  $j$  个人对血压的影响, 称为第  $j$  个人的效应。由于我们的目的是研究人在一天内的平均血压, 所以,  $a$  个观测时间点和  $b$  个被观测的人都是随机抽取的, 即  $\alpha_i$  和  $\beta_j$  为随机效应, 则模型(1.3)为两向分类随机效应模型。此外, 如果我们关心的问题是人的血压在一天内的变化规律, 则  $a$  个观测时间点可以是固定的, 此时模型(1.3)为两向分类混合效应模型。进一步, 如果我们感兴趣的是特定的  $b$  个人, 则  $\beta_j$  也是非随机的, 从而模型(1.3)就是两向分类固定效应模型。

引进适当的矩阵记号, 模型(1.3)可以写成模型(1.2)的形式。记

$$y = (y_{11}, \dots, y_{1b}, \dots, y_{a1}, \dots, y_{ab})', \epsilon_0 = (e_{11}, \dots, e_{1b}, \dots, e_{a1}, \dots, e_{ab})', \\ Z_2 = I_a \otimes 1_b, Z_1 = 1_a \otimes I_b, \epsilon_2 = (\alpha_1, \alpha_2, \dots, \alpha_a)', \epsilon_1 = (\beta_1, \beta_2, \dots, \beta_b)'$$



其中,  $\otimes$  表示矩阵的 Kronecker 乘积。此时模型(1.3)变形为

$$y = 1_{ab}\mu + Z_2\varepsilon_2 + Z_1\varepsilon_1 + \varepsilon_0 \quad (1.4)$$

一般我们总是假定所有随机效应都是互不相关的,  $\text{Var}(\alpha_i) = \sigma_\alpha^2$ ,  $\text{Var}(\beta_j) = \sigma_\beta^2$ ,  $\text{Var}(e_{ij}) = \sigma_e^2$ 。则  $y$  的协方差矩阵为

$$\text{Cov}(y) = \sigma_\alpha^2 Z_2 Z_2' + \sigma_\beta^2 Z_1 Z_1' + \sigma_e^2 I_{ab} = \sigma_\alpha^2 (I_a \otimes J_b) + \sigma_\beta^2 (J_a \otimes I_b) + \sigma_e^2 I_{ab}$$

其中,  $J_n = 1_n 1_n'$ 。

### 例 1.2 带交互效应的非平衡两向分类随机效应模型。

从某工厂中随机抽取  $a$  个工人和  $b$  个区域, 每个工人轮流在这些区域工作不等数量个班次, 我们利用带交互效应的非平衡两向分类随机效应模型来拟合经过对数变换的暴露数据。令  $x_{ijk}$  为第  $i$  个工人在第  $j$  个区域第  $k$  个班次工作时间的暴露量,  $i=1, 2, \dots, a$ ;  $j=1, 2, \dots, b$ ;  $k=1, 2, \dots, n_{ij} > 0$ 。假定  $x_{ijk}$  服从对数正态分布, 即  $y_{ijk} = \ln(x_{ijk})$  服从正态分布, 则

$$y_{ijk} = \mu + \alpha_i + \beta_j + \gamma_{ij} + e_{ijk}, \quad i=1, 2, \dots, a; \quad j=1, 2, \dots, b; \quad k=1, 2, \dots, n_{ij} \quad (1.5)$$

其中,  $\mu$  表示总体均值,  $\alpha_i$  表示第  $i$  个工人的随机效应,  $\beta_j$  表示第  $j$  个区域的随机效应,  $\gamma_{ij}$  表示第  $i$  个工人在第  $j$  个区域的交互效应,  $e_{ijk}$  表示随机误差项。记

$$y = (y_{111}, \dots, y_{11n_{11}}, \dots, y_{ab1}, \dots, y_{abn_{ab}})' , \quad \varepsilon_0 = (e_{111}, \dots, e_{11n_{11}}, \dots, e_{ab1}, \dots, e_{abn_{ab}})'$$

$$Z_3 = \text{diag}(1_{n_{1.}}, 1_{n_{2.}}, \dots, 1_{n_{a.}}), \quad Z_2 = (z_1', z_2', \dots, z_b')', \quad Z_1 = \text{diag}(z_1, z_2, \dots, z_b)$$

$$\varepsilon_3 = (\alpha_1, \alpha_2, \dots, \alpha_a)', \quad \varepsilon_2 = (\beta_1, \beta_2, \dots, \beta_b)', \quad \varepsilon_1 = (\gamma_{11}, \dots, \gamma_{1b}, \dots, \gamma_{a1}, \dots, \gamma_{ab})'$$

其中,  $z_i = \text{diag}(1_{n_{i1}}, 1_{n_{i2}}, \dots, 1_{n_{ib}})$ ;  $n_{i.} = \sum_{j=1}^b n_{ij}$ ,  $i=1, 2, \dots, a$ 。记  $n_{..} =$

$\sum_{i=1}^a \sum_{j=1}^b n_{ij}$ , 则模型(1.5)可以写成矩阵形式, 即

$$y = 1_{n..} \mu + Z_3 \varepsilon_3 + Z_2 \varepsilon_2 + Z_1 \varepsilon_1 + \varepsilon_0 \quad (1.6)$$

如果假定  $\text{Var}(\alpha_i) = \sigma_\alpha^2$ ,  $\text{Var}(\beta_j) = \sigma_\beta^2$ ,  $\text{Var}(\gamma_{ij}) = \sigma_\gamma^2$ ,  $\text{Var}(e_{ijk}) = \sigma_e^2$ , 且  $\alpha_i$ 、 $\beta_j$ 、 $\gamma_{ij}$ 、 $e_{ijk}$  互不相关, 则  $y$  的协方差矩阵为

$$\text{Cov}(y) = \sigma_\alpha^2 Z_3 Z_3' + \sigma_\beta^2 Z_2 Z_2' + \sigma_\gamma^2 Z_1 Z_1' + \sigma_e^2 I_{n..}$$

### 例 1.3 面板数据模型。

面板数据模型是计量经济学中应用最为广泛的一类混合效应模型<sup>[1~3]</sup>。假设我们对  $N$  个个体(如个人、家庭、公司、城市、国家或区域等)进行  $T$  个时刻的

观测, 则观测数据可表示为

$$y_{it} = x_{it}'\beta + \mu_i + e_{it}, \quad i=1, 2, \dots, N; \quad t=1, 2, \dots, T \quad (1.7)$$

其中,  $y_{it}$  表示第  $i$  个个体在时刻  $t$  的某项经济指标;  $x_{it}$  为  $k \times 1$  已知向量, 它刻画了第  $i$  个个体在时刻  $t$  的一些自然特征;  $\mu_i$  为第  $i$  个个体的个体效应,  $e_{it}$  为随机误差项。如果我们的目的是研究整个市场的运行规律, 而不是关心这特定的  $N$  个个体, 这  $N$  个个体只不过是从总体中抽取的随机样本, 这时个体效应就是随机的。记

$$\begin{aligned} y &= (y_{11}, \dots, y_{1T}, \dots, y_{N1}, \dots, y_{NT})' \\ X &= (x_{11}, \dots, x_{1T}, \dots, x_{N1}, \dots, x_{NT})' \end{aligned}$$

$Z_2 = I_N \otimes 1_T$ ,  $\epsilon_2 = (\mu_1, \mu_2, \dots, \mu_N)'$ ,  $\epsilon_0 = (e_{11}, \dots, e_{1T}, \dots, e_{N1}, \dots, e_{NT})'$   
则模型(1.7)可写成矩阵形式, 即

$$y = X\beta + Z_2\epsilon_2 + \epsilon_0 \quad (1.8)$$

若假定  $\text{Var}(\mu_i) = \sigma_\mu^2$ ,  $\text{Var}(e_{it}) = \sigma_e^2$ , 且  $\mu_i$ 、 $e_{it}$  互不相关, 则

$$\text{Cov}(y) = \sigma_\mu^2 Z_2 Z_2' + \sigma_e^2 I_{NT} = \sigma_\mu^2 (I_N \otimes J_T) + \sigma_e^2 I_{NT}$$

在模型(1.7)中, 若把时间效应也考虑进来, 则模型(1.7)可表示为

$$y_{it} = x_{it}'\beta + \mu_i + v_t + e_{it}, \quad i=1, 2, \dots, N; \quad t=1, 2, \dots, T \quad (1.9)$$

其中,  $v_t$  为第  $t$  个时刻的时间效应。假定  $T$  个观测时间点是随机选取的, 即  $v_t$  为随机效应,  $\text{Var}(v_t) = \sigma_v^2$ ,  $v_t$  与  $\mu_i$ 、 $e_{it}$  互不相关。记  $Z_1 = 1_N \otimes I_T$ ,  $\epsilon_1 = (v_1, v_2, \dots, v_T)'$ , 则得到如下模型:

$$y = X\beta + Z_2\epsilon_2 + Z_1\epsilon_1 + \epsilon_0 \quad (1.10)$$

其协方差矩阵为

$$\text{Cov}(y) = \sigma_\mu^2 Z_2 Z_2' + \sigma_v^2 Z_1 Z_1' + \sigma_e^2 I_{NT} = \sigma_\mu^2 (I_N \otimes J_T) + \sigma_v^2 (J_N \otimes I_T) + \sigma_e^2 I_{NT}$$

## 1.2 国内外研究现状

### 1.2.1 方差分量的估计

在混合效应模型的参数估计中, 方差分量的估计无疑是最为丰富多彩的一个篇章。自 Airy<sup>[4]</sup> 的研究至今, 统计学家提出了若干种估计方法<sup>[5~8]</sup>, 如方差分析估计 (analysis of variance estimate, ANOVAE)、极大似然估计 (maximum likeli-

hood estimate, MLE)、限制极大似然估计(restricted maximum likelihood estimate, REMLE)、最小范数二次无偏估计(minimum norm quadratic unbiased estimate, MINQUE)及谱分解估计(spectral decomposition estimate, SDE)等。史建红<sup>[9]</sup>将这些估计方法归为三类,即方差分析估计和谱分解估计基于矩法、极大似然估计和限制极大似然估计基于分布、最小范数二次无偏估计基于准则。

针对含有两个方差分量的混合效应模型, Wu 等<sup>[10]</sup>对方差分量的谱分解估计和方差分析估计进行比较,获得它们相等的充要条件以及彼此优于对方的充分条件。Ye 和 Wang<sup>[11]</sup>则构造两种协方差矩阵估计类,并在二次损失函数下证明新估计量优于协方差矩阵的谱分解估计。在此基础上, Ye 等<sup>[12]</sup>基于协方差矩阵的方差分析估计,给出一类改进估计,并在两种二次损失函数下证明这个新估计量优于方差分析估计。而针对含有两个方差分量的多元混合效应模型,马铁丰和王松桂<sup>[13]</sup>建立了方差分量矩阵的谱分解估计,并在二次损失函数下证明其优于方差分析估计以及两者等价性条件。Kubokawa 和 Tsai<sup>[14]</sup>则利用 Stein-Haff 等式(参见文献[15]和文献[16]),构造方差分量矩阵的截断估计,并在 Stein 损失函数下证明其优于尺度同变估计(scale-equivariant estimate)。进一步地, Sun 等<sup>[17]</sup>建立了方差分量矩阵的非负估计类,并在二次损失函数下证明其优于最小二乘估计和方差分析估计。马铁丰等<sup>[18]</sup>则利用谱分解估计方法,分别给出平衡和非平衡情况下方差分量矩阵的非负估计类,并证明其统计优良性。这些研究成果为进一步统计推断奠定了理论基础,但对于结构更为复杂的混合效应模型,如一般平衡随机效应模型、非平衡随机效应模型、混合效应模型等,其相关研究成果仍然较少。

### 1.2.2 参数的广义检验

在混合效应模型的假设检验中,未知参数可被分为感兴趣参数(interest parameter)和多余参数(nuisance parameter)。通常情况下,多余参数的数量多于感兴趣参数,这导致对感兴趣参数检验的研究变得较为复杂。此时,基于精确分布的传统方法往往难以奏效。对此, Tsui 和 Weerahandi<sup>[19]</sup>另辟新径,创造性地提出广义  $p$ -值(generalized  $p$ -value)的概念,为研究复杂数据和模型的检验问题开辟了一种全新的思路。Weerahandi<sup>[20]</sup>推广置信区间的定义,进一步提出广义置信区间(generalized confidence interval)的概念。基于广义  $p$ -值和广义置信区间所建立的广义方法,具有稳健性、计算简便以及易应用于小样本问题等优良特点。

近年来,有关广义方法的研究受到高度重视,并取得长足发展。例如,

试读结束: 需要全本请在线购买: [www.ertongbook.com](http://www.ertongbook.com)

Tian<sup>[21]</sup>、Park<sup>[22]</sup>、Krishnamoorthy 和 Mathew<sup>[23]</sup>、Krishnamoorthy 等<sup>[24]</sup>、Gamage 等<sup>[25]</sup>针对正态总体、对数正态总体和逆高斯总体,研究总体均值的推断问题。类似的, Krshnamoorthy 和 Lu<sup>[26]</sup>、Lin 和 Lee<sup>[27]</sup>、Ye 等<sup>[28]</sup>针对正态总体和逆高斯总体,考虑共同均值的推断问题。又如, Weerahandi<sup>[29]</sup>、Arendacka<sup>[30]</sup>、Mathew 和 Webb<sup>[31]</sup>针对各种简单混合效应模型,讨论方差分量的推断问题。Weerahandi 和 Berger<sup>[32]</sup>、Lin 和 Lee<sup>[33]</sup>、Chi 和 Weerahandi<sup>[34]</sup>则针对各种简单生长曲线模型,探讨固定处理效应的检验问题。此外, Tian<sup>[35]</sup>、Ye 和 Wang<sup>[36]</sup>、Gilder 等<sup>[37]</sup>针对各种随机效应模型,研究组内相关系数的推断问题。再如, Weerahandi 和 Johnson<sup>[38]</sup>、Roy 和 Mathew<sup>[39]</sup>、Krishnamoorthy 和 Lin<sup>[40]</sup>针对正态总体、指数总体和威布尔总体,考虑可靠性参数的推断问题。此外, Mathew 等<sup>[41]</sup>、Hsu 等<sup>[42]</sup>针对正态总体,讨论几种常见过程能力指数的区间估计问题;李新民等<sup>[43]</sup>、Hannig 等<sup>[44]</sup>利用信仰推断理论,分别研究广义  $p$ -值和广义置信区间的构建问题;徐兴忠和刘芳<sup>[45]</sup>针对混合模型,考虑混合比的推断问题;Xiong 等<sup>[46]</sup>针对异方差线性模型,讨论位置参数和尺度参数的检验问题。这些研究成果可为实际数据分析提供切实可行的方法,但也产生了一些新的值得深入研究的问题。其中,广义方法在一般平衡随机效应模型、非平衡随机效应模型、一般混合效应模型等复杂统计模型下的应用及其理论性质,显得尤为重要。

### 1.2.3 偏正态分布研究

在经济、金融、生物、医学等众多领域,存在的大量纵向数据既可揭示研究对象的动态变化,又可反映个体异质性。近年来针对此类数据统计建模的研究取得了长足发展。较早的研究如 Diggle 等<sup>[47]</sup>、Fitzmaurice 等<sup>[48]</sup>较为系统地介绍了纵向数据几类常见模型的基本理论与方法。在这些常见模型中,混合效应模型可谓最重要的一类,亦是本领域众多研究的重点。例如, Verbeke 和 Molenberghs<sup>[49]</sup>、Demidenko<sup>[50]</sup>研究了混合效应模型的估计、检验、影响诊断等问题。针对混合效应模型的参数估计问题,王松桂和尹素菊<sup>[8]</sup>提出了回归系数和方差分量的谱分解估计方法, Wu 等<sup>[51]</sup>则建立了回归系数和方差分量的估计量同时成为最小方差估计的充要条件。此外, Wu 和 Zhang<sup>[52]</sup>、梁华和师义民<sup>[53]</sup>针对纵向数据建立了非参数混合效应模型,并将其应用于荷尔蒙、艾滋病、恐慌焦虑症等实际案例的数据分析中; Baltagi<sup>[1]</sup>、Frees<sup>[54]</sup>则针对经济社会领域的纵向数据,利用混合效应模型建立纵向数据分析方法,并将其应用到长期收支动态、劳动力市场状况等经济问题的研究中。

然而,经典理论与方法大多基于正态分布假设,但实际纵向数据却更为常见偏态分布,这就导致统计推断的精度大幅下降。当实际纵向数据不满足正态分布假设时,部分研究者如 Wang 和 Chow<sup>[55]</sup>、王松桂等<sup>[56]</sup>尝试利用 Box-Cox 变换,对纵向数据进行“综合治理”,以改善它们的正态性、对称性和方差相等性。然而,这种纯粹的数据变换往往会导致因变量缺乏实际意义,尤其是在利用不同的变换函数对多维向量分量进行变换的情况下。对此,近期研究者,如 Genton<sup>[57]</sup>、Gupta 等<sup>[58]</sup>、Jara 等<sup>[59]</sup>、Zhou 和 He<sup>[60]</sup>,开始寻找一类比正态分布更加“灵活”的参数分布族,作为纵向数据统计建模的分布假设。在众多偏态分布中,偏正态分布严格地包含正态分布且是偏态的,不仅具有正态分布的优良统计性质,而且不失偏态分布的基本特征,成为偏态分布统计推断领域的一大研究热点。

而关于偏正态分布的研究成果虽不鲜见,但早期研究大多强调分布本身,此类分布下的统计建模问题仍需系统性地深入分析。例如, Azzalini<sup>[61]</sup>、Henze<sup>[62]</sup>首次提出偏正态分布的概念,并利用正态随机变量和截断正态随机变量构造出偏正态分布的概率表示。Azzalini<sup>[63]</sup>、Azzalini 等<sup>[64,65]</sup>、Arellano-Valle 和 Genton<sup>[66]</sup>、Balakrishnan 和 Scarpa<sup>[67]</sup>则进一步提出多维偏正态分布,讨论矩生成函数、边缘分布、条件分布、偏度度量方式等统计性质,并将其应用到判别分析、回归分析、图模型等多元分析中。此外, Loperfido<sup>[68]</sup>、Gupta 和 Huang<sup>[69]</sup>、Genton 等<sup>[70]</sup>、Wang 等<sup>[71]</sup>基于偏正态分布假设,研究随机向量二次型的分布形式、矩生成函数、独立性等相关问题。但上述研究往往更多关注偏正态分布本身的特征、性质等问题,仍较少涉及此类分布下的统计建模理论与方法。对此,近期研究如 Lachos 等<sup>[72]</sup>针对三种常见的偏正态多元回归模型,建立极大似然估计的 EM 算法,并将其应用于牙菌斑和胆固醇水平等实际纵向数据分析。Lin 和 Lee<sup>[73]</sup>则针对偏正态线性混合效应模型,构造了未知参数极大似然估计的数值算法。尽管如此,上述研究并未给出未知参数估计的显示解,亦未深入探讨其理论性质,而这对于提高统计推断精度、促进实际数据分析,恰恰具有更为重要的意义。

### 1.3 研究方法思路

本著作针对一般平衡随机效应模型、非平衡单向分类随机效应模型、非平衡双向分类随机效应模型、面板数据模型、混合效应模型、偏正态混合效应模型

等,运用矩法、谱分解方法、矩阵技术、结构方法、Monte Carlo方法等多种统计研究方法,研究未知参数的优良估计、精确检验等统计推断问题,建立一系列新的有效的统计推断方法,并将其应用到经济、金融、生物、医学等领域的实际数据分析中。具体研究方法如下。

(1)谱分解方法。先对混合效应模型的协方差矩阵进行谱分解,若其特征值均为方差分量的线性组合,则通过对原模型进行适当的线性变换,即可得到若干个奇异线性模型。在此基础上,运用最小二乘统一理论,给出模型参数的谱分解估计。

(2)矩阵技术。在研究方差分量、相关系数、协方差矩阵等未知参数的估计问题时,矩阵技术中的高度技巧,如矩阵分解、矩阵偏序、矩阵广义逆和矩阵不等式等,是一种有效的研究工具。本著作巧妙地将这些矩阵技术与基于矩法的估计方法(方差分析估计和谱分解估计)相结合,以构造未知参数的可行估计,并证明其统计优良性。

(3)结构方法。在研究未知参数(包括暴露水平、可靠性参数和过程能力指数)的假设检验和区间估计问题时,利用 Hannig 等<sup>[44]</sup>提出的结构方法,构造相应的广义  $p$ -值和广义置信区间。在此基础上,研究其理论性质,如不变性、理论功效和理论置信水平等。

(4)Monte Carlo方法。针对本著作所建立的估计和检验方法,利用 Monte Carlo方法模拟风险函数、第一类错误、功效函数、覆盖概率和区间长度等评价指标,并将其与已有结果进行比较,以判断所给新方法的合理性和有效性。

基于上述多种统计研究方法,本著作的研究思路如下。

其一,针对结构较为复杂的混合效应模型,利用基于谱分解方法和矩法,并结合矩阵分解、矩阵偏序、矩阵广义逆、矩阵不等式等矩阵技术中高度技巧,研究方差分量、相关系数、协方差矩阵等未知参数的估计问题,构造相应的可行估计,并证明其统计优良性。

其二,针对一般平衡随机效应模型、非平衡随机效应模型、面板数据模型、一般混合效应模型等多种复杂模型,利用结构方法,研究未知参数(包括暴露水平、可靠性参数和过程能力指数)的假设检验和区间估计问题,建立相应的广义推断方法。进而,研究其理论性质,如不变性、理论功效和理论置信水平等。

其三,基于非中心偏  $F$  分布,构造感兴趣参数的精确检验,建立偏正态混合效应模型的假设检验理论与方法。

其四,针对上述研究方法,利用 Monte Carlo方法模拟风险函数、第一类错误、功效函数、覆盖概率和区间长度等评价指标,以判断这些方法的合理性和有

效性。

其五, 将本著作所建立的方法应用于经济、金融、生物、医学等领域的实际数据分析, 为解决实际问题提供一种简单方便、效果理想的方法。

## 1.4 研究特色与创新之处

本著作研究可深化现有对混合效应模型参数估计和检验的理论探索, 亦有助于提高统计推断的精度, 改善实际纵向数据分析的效果。具体而言, 本著作的研究特色与创新之处如下。

(1) 利用矩阵技术的高度技巧, 构造未知参数的优良估计。

参数估计的传统方法基于一些常见的总体, 难以应用于结构较为复杂的混合效应模型, 故而构造未知参数的优良估计较为棘手。对此, 本著作巧妙地利用矩阵分解、矩阵偏序、矩阵广义逆和矩阵不等式等矩阵技术中的高度技巧, 并将其与谱分解方法和矩法有机结合, 从而大大简化未知参数的估计问题, 并构造未知参数的优良估计。这一尝试不仅有助于改进现有参数估计方法, 亦有助于深化混合效应模型的统计推断研究。

(2) 扩大广义方法的应用范围, 进一步揭示其理论性质与本质特征。

基于广义  $p$ -值和广义置信区间所建立的广义方法目前多应用于相对简单的模型, 且有关其理论性质的结果较少, 这制约了广义方法的进一步应用与完善。鉴于此, 本著作不仅探讨广义方法在一般平衡随机效应模型、非平衡随机效应模型、面板数据模型、一般混合效应模型等多种复杂统计模型中的应用, 更进一步研究其理论性质, 如不变性、理论功效和理论置信水平等, 提出一些更加深刻的理论结果, 以刻画这种方法的本质特征。

(3) 将可靠性参数和过程能力指数的研究从常见总体推广到复杂混合效应模型。

现有关于可靠性参数和过程能力指数的研究往往基于一些比较常见的总体, 如正态总体、指数总体和 Wishart 总体等, 相应的理论与方法难以应用于结构较为复杂的混合效应模型。本著作针对单向分类随机效应模型和一般平衡随机效应模型, 系统地研究可靠性参数和过程能力指数的统计推断问题。这是一次有益的尝试, 亦有助于深化现有关于混合效应模型的理论探索。

(4) 基于非中心偏  $F$  分布, 构造感兴趣参数的精确检验, 建立偏正态混合效

应模型的假设检验理论与方法。

前期研究表明,在偏正态分布假设下,对于回归系数、方差分量等感兴趣参数,基于两个观测向量二次型比值所构造的检验统计量,不服从通常意义下的非中心  $F$  分布,故而无法建立相应的检验方法。对此,本著作定义一种新的分布即非中心偏  $F$  分布,以此作为偏正态分布假设下观测向量二次型比值的分布,在此基础上,构造回归系数、方差分量等感兴趣参数的精确检验。这是一个难度较大但理论与实际意义均很强的工作,是本著作的重要创新性探索之一。

本著作系统研究一般平衡随机效应模型、非平衡单向分类随机效应模型、非平衡两向分类随机效应模型、面板数据模型等若干混合效应模型的基本理论、方法和应用。其学术价值和科学意义如下。

(1)构造混合效应模型的优良估计和精确检验,可深化与推广现有关于混合效应模型的统计推断理论。

(2)本著作所建立的统计推断理论与方法,可提高混合效应模型对纵向数据的拟合精度,改善实际纵向数据分析的效果,更好地实现混合效应模型的应用价值。

(3)本著作研究成果可在一定程度上揭示偏正态纵向数据统计推断的一般规律,对偏  $t$ 、偏椭圆等等其他形式的偏态纵向数据分析,具有借鉴意义。

## 1.5 研究内容与框架

本著作针对一般平衡随机效应模型、非平衡单向分类随机效应模型、非平衡两向分类随机效应模型、面板数据模、混合效应模型、偏正态混合效应模型等多种复杂模型,系统探讨感兴趣参数的统计推断理论、方法及应用。具体研究内容与框架如下。

第1章为绪论。本章通过实例引入各类混合效应模型,并介绍国内外研究现状、研究方法思路、研究特色与创新之处、研究内容与框架等,使读者对这类模型的丰富实际背景及相关研究有一些了解,有助于对后续章节内容的理解。

第2章为预备知识。为构建混合效应模型统计推断理论与方法,本章讨论矩阵方面的补充知识、多元正态分布、多元偏正态分布及广义推断方法等。具体包括矩阵分解、幂等阵、正交投影阵、Kronecker 乘积、矩阵微商等矩阵技术,随机向量的均值、协方差阵、二次型分布、线性型分布等统计性质,广义检验变



量、广义  $p$ -值、广义枢轴量、广义置信区间等基本概念。

第3章为一般平衡随机效应模型。本章针对一般平衡随机效应模型，探讨方差分量及过程能力指数的估计与检验问题，构建方差分量的非负估计和广义推断方法，在此基础上，给出过程能力指数的广义置信区间，并从数值角度研究其统计优良性。

第4章为非平衡单向分类随机效应模型。本章针对非平衡单向分类随机效应模型，研究组间方差分量、暴露水平和可靠性参数的区间估计与假设检验问题。基于广义  $p$ -值和广义置信区间的概念，建立一系列新的有效的广义推断方法。

第5章为非平衡两向分类随机效应模型。本章针对暴露水平评价问题，将非平衡单向分类随机效应模型推广到非平衡两向分类随机效应模型，建立其广义推断理论与方法。在此基础上，讨论两种组内相关系数的统计推断问题，建立广义推断方法和大样本推断方法。

第6章为面板数据模型。本章针对单个回归系数，构建其检验方法和置信区间。进而，针对方差分量线性组合，建立其检验方法和置信区间。在此基础上，从分析角度探讨上述方法的统计优良性。最后，基于 Monte Carlo 模拟方法，从数值上对所给检验方法和置信区间进行比较。

第7章为混合效应模型。本章针对混合效应模型，研究方差分量和协方差阵的估计或检验问题，建立一系列新的统计推断方法，并通过 Monte Carlo 模拟研究，进一步论证所给方法的合理性和有效性。最后，将上述方法应用于武器发射精度比较问题。

第8章为偏正态混合效应模型。本章针对偏正态混合效应模型，系统探讨密度函数、矩生成函数、均值向量、协方差矩阵、独立性条件等统计性质。在此基础上，证明偏正态混合效应模型观测向量二次型服从非中心偏  $\chi^2$  分布的充要条件，并给出该模型下 Cochran 定理。最后，基于非中心偏  $F$  分布，构造回归系数和方差分量的精确检验。