



混合效应模型 在林业建模中的应用

姜立春 李凤日◎著



科学出版社

混合效应模型 在林业建模中的应用

姜立春 李凤日 著

科学出版社

北京

内 容 简 介

本书是关于线性和非线性混合效应模型理论、方法及其在林业建模上应用的专著。内容主要介绍单水平和多水平统计模型。将单水平混合效应模型如何应用于树高曲线、单木生长、树干削度、枝条特征及木材微纤丝角和管胞长度模型的构建；介绍如何使用多水平混合效应模型对木材密度、树皮厚度进行统计建模。本书采用国际通用的著名统计软件 S-Plus 和 SAS 来进行各种模型的分析。最后，结合具体的实例，由浅入深地逐步介绍如何使用 S-Plus 软件的 lme 和 nlme 模块来进行各种模型的分析。

本书可供从事森林经理、林业建模工作者和高校相关专业的师生参考使用。

图书在版编目 (CIP) 数据

混合效应模型在林业建模中的应用 / 姜立春, 李凤日著. —北京: 科学出版社, 2014.6

ISBN 978-7-03-040928-7

I. ①混… II. ①姜… ②李… III. ①林业—建立模型 IV. ①S711

中国版本图书馆 CIP 数据核字(2014)第 120572 号

责任编辑: 王海光 王 好 / 责任校对: 张怡君

责任印制: 徐晓晨 / 封面设计: 北京铭轩堂广告设计有限公司

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>

北京科印技术咨询有限公司印刷

科学出版社发行 各地新华书店经销

*

2014 年 6 月第 一 版 开本: B5 (720×1000)

2014 年 6 月第一次印刷 印张: 9 1/2

字数: 210 000

定价: 68.00 元

(如有印装质量问题, 我社负责调换)

前 言

混合效应模型 (mixed effects model) 是近代发展起来的新的统计方法, 其中, 多水平混合效应模型的应用分析是最重要的发展内容之一。混合效应模型已广泛应用于教育学、心理学、社会学、经济学、生物统计学、医学、药物动力学、农业、林业、工业研究等领域。

混合效应模型是既包含固定效应又包含随机效应的一类模型。主要应用于分组数据, 如纵向数据 (longitudinal data)、重复测量数据 (repeated measures data)、分层数据 (hierarchical data) 及多变量多层数据 (multivariate multilevel data)。与普通最小二乘法 (ordinary least square, OLS) 相比, 混合效应模型能满足独立同分布的假设, 能得到渐进无偏参数估计, 并用方差-协方差结构来反映数据间的相关性及异质性。混合效应模型包括总体和个体预测, 总体预测只能反应平均水平而不能体现组间的差异, 而个体预测通过二次抽样校正随机参数来反应组间的差异。

落叶松 (*Larix spp.*) 是我国东北地区主要速生用材林树种之一。目前, 东北三省落叶松人工林面积超过 $400 \times 10^4 \text{ hm}^2$, 约占东北地区人工林面积的 55% (雷加富, 2005)。六十年来, 森林估算方法已经从一个简单的图形化的形式过渡到复杂的回归模型形式, 通常这种回归模型具有随机效应成分。近年来随着计算混合效应模型统计软件的发展, 混合效应模型已经成为林分生长和收获模型分析的重要工具。

林业领域的许多数据都具有层次结构, 如在单木树高生长模型中, 树高生长嵌套与单木中, 单木嵌套于样地中, 样地嵌套于林分中, 形成了一个具有三个层次 (单木-样地-林分) 的数据结构。同一层次的数据不是相互独立的, 不同层次的数据具有较强的异质性。传统的回归分析不再满足这类数据的基本假设。从统计分析角度上讲, 传统分析方法在分析这类数据时所遇到的问题可通过混合效应模型得到解决。由于混合效应模型可对组水平和个体水平的数据进行同时分析, 在一个模型中可同时检验组变量和个体变量的效应。另外, 混合效应模型分析不需要假设数据中的观察值相互独立, 因而可以修正因观察数据的非独立性引起的参

数标准误差估计。

本书包含了单水平模型、两水平模型和三水平模型，重点介绍混合模型的方法及其应用。本书共分为 10 章，第 1 章介绍单水平和多水平线性和非线性混合效应模型的一些基本知识。第 2 章至第 7 章利用单水平模型进行不同类型的落叶松数据分析及模型构建，包括树高曲线、单木生长、树干削度、枝条特征及木材微纤丝角和管胞长度模型。第 8 章介绍了落叶松木材基本密度的株间变异、径向变异及沿树干方向的变异规律。同时考虑单木和树干高度效应，利用两水平模型构建落叶松木材密度模型，并与单水平模型进行了对比分析。第 9 章同时考虑区组、样地和单木，利用三水平模型构建落叶松树皮因子模型，并与单水平和两水平模型进行了对比分析。第 10 章结合具体的实例，由浅入深地逐步介绍如何使用 S-Plus 软件的 lme 和 nlme 模块来进行各种模型的分析。

本书内容是国家自然科学基金项目（30972363，31170591）和“十二五”国家科技支撑计划项目（2012BAD22B02）共同资助的研究成果。此外，本书还获得中央高校基本科研业务费专项资金（DL12DA01）的资助。在此一并表示感谢！

由于作者水平有限，书中难免有疏漏和不足之处，敬请读者批评指正。

著 者
2014 年 4 月

目 录

前言

1 绪论	1
1.1 线性模型.....	2
1.2 广义线性模型.....	2
1.3 单水平线性混合模型.....	3
1.3.1 模型定义.....	3
1.3.2 模型的参数估计.....	4
1.4 多水平线性混合模型.....	6
1.5 非线性混合效应模型.....	6
1.5.1 单水平非线性混合模型.....	7
1.5.2 多水平非线性混合模型.....	7
1.6 模型拟合和比较.....	8
1.7 混合效应模型构建的步骤.....	8
1.7.1 运行空模型(empty model)	9
1.7.2 加入解释变量(explanatory variable)	9
1.7.3 确定参数效应.....	9
1.7.4 确定随机效应参数的协变量(covariate variable)	9
1.7.5 确定组内方差-协方差结构.....	10
1.7.6 确定随机效应的方差-协方差结构.....	12
1.7.7 模型检验.....	12
2 基于混合效应的落叶松树高曲线模拟	14
2.1 研究地区概况.....	14
2.2 研究方法.....	15
2.2.1 数据.....	15
2.2.2 方法.....	15
2.3 结果与分析.....	16
2.3.1 基础混合树高曲线模型的拟合.....	16
2.3.2 含有林分变量的树高曲线模型构建.....	17
2.3.3 含有林分变量的树高曲线模型拟合.....	19

2.3.4	随机效应的方差-协方差结构	19
2.3.5	模型检验	21
2.4	小结	21
3	落叶松单木生长混合效应模型	23
3.1	研究方法	24
3.1.1	样本准备	24
3.1.2	基础模型	24
3.1.3	单木生长混合模型构建	25
3.2	结果与分析	25
3.2.1	基于单木效应混合模型模拟	25
3.2.2	基于样地效应混合模型模拟	27
3.2.3	方差-协方差结构	31
3.2.4	模型评价	33
3.2.5	模型检验	33
3.3	小结	35
4	落叶松树干削度模型	39
4.1	研究方法	41
4.1.1	数据	41
4.1.2	方法	41
4.2	结果与分析	45
4.2.1	树干削度和材积相容模型	45
4.2.2	树干削度非线性混合模型	48
4.2.3	冠长率和林分密度对树干干形的影响	51
4.3	小结	54
5	基于线性混合模型的落叶松枝条特征模型	56
5.1	研究地区概况	57
5.2	研究方法	57
5.2.1	数据	57
5.2.2	线性混合效应模型	57
5.2.3	模型评价和检验指标	58
5.3	结果与分析	59
5.3.1	基础模型构建	59
5.3.2	枝条长度混合模型拟合	59
5.3.3	枝条角度混合模型拟合	60
5.3.4	枝条基径混合模型拟合	60

5.3.5	随机效应的方差-协方差结构	60
5.3.6	误差的自相关性和异质性	61
5.3.7	模型评价	62
5.3.8	模型检验	64
5.3.9	混合模型应用	66
5.4	小结	67
6	落叶松微纤丝角混合效应模型	68
6.1	研究地区概况	68
6.2	材料与方法	69
6.2.1	数据采集和处理	69
6.2.2	方法	69
6.3	结果与分析	70
6.3.1	基础模型拟合	70
6.3.2	混合模型拟合	70
6.3.3	方差协-方差结构	72
6.3.4	模型评价	73
6.4	小结	74
7	落叶松早晚材管胞长度混合效应模型	76
7.1	材料与方法	76
7.1.1	数据采集和处理	76
7.1.2	方法	77
7.2	结果与分析	78
7.2.1	基础模型拟合	78
7.2.2	早材管胞长度混合模型拟合	78
7.2.3	晚材管胞长度混合模型拟合	79
7.2.4	方差-协方差结构	80
7.2.5	模型评价	82
7.3	小结	83
8	基于多水平混合效应模型的落叶松木材密度模拟	84
8.1	材料与方法	84
8.1.1	数据采集和处理	84
8.1.2	方法	85
8.2	结果与分析	86
8.2.1	落叶松木材基本密度的变异及早期选择	86
8.2.2	固定密度模型的确定	90

8.2.3	混合效应模型拟合	90
8.2.4	模型评价	94
8.3	小结	97
9	基于多水平混合效应模型的落叶松树皮因子模拟	99
9.1	研究方法	99
9.1.1	数据	99
9.1.2	方法	100
9.2	结果与分析	101
9.2.1	固定模型的确定	101
9.2.2	三水平混合效应模型	101
9.2.3	单水水平混合效应模型拟合	102
9.2.4	2水平混合效应模型拟合	103
9.2.5	3水平混合效应模型拟合	103
9.2.6	误差的异方差性	104
9.2.7	模型评价	104
9.2.8	模型应用	109
9.3	小结	109
10	混合效应模型在 S-Plus 软件中的实现	111
10.1	S-Plus 软件介绍	111
10.2	数据	111
10.3	线性混合效应模型在 S-Plus 软件 lme 程序中的运行	111
10.3.1	空模型	111
10.3.2	加入协变量解释组间变异	114
10.3.3	在模型中纳入组内解释变量	116
10.3.4	组内和组间变量交互作用评估	118
10.3.5	增加随机效应评估	120
10.3.6	不同方差-协方差结构比较	124
10.3.7	ML 和 REML 方法比较	125
10.3.8	多水平线性混合效应模型	129
10.4	非线性混合效应模型在 S-Plus 软件 nlme 程序中的运行	132
10.4.1	单水平非线性混合效应模型	132
10.4.2	多水平非线性混合效应模型	134
	参考文献	136

1 绪 论

森林系统的动态特性促使我们不断采集有关森林变量的信息。为了有效地利用、保护和管理森林资源,我们必须利用这些信息进行林分生长和收获的预测。一般情况下,林分生长和收获模拟数据展现一个纵向结构或重复测量特征,在本书中这类数据统一称为分组数据。因此,针对这样的分组数据,在建模过程中应考虑随机变异的三种成分:①随机效应,即分组之间的异质性;②序列相关性,即在时间或空间上残差的相关性;③测量误差。

分组数据中存在的非独立观察可以用组内相关系数(*intra-class correlation coefficient*, ICC)来测量。研究显示,即便是很小的 ICC 在统计检验中也可能导致很大的第 I 类错误(*type I error*),从而错误地拒绝真的统计假设。一些传统的统计分析方法,如 *t* 检验、线性回归、方差分析(*analysis of variance*, ANOVA)等,都要求观察相互独立(*independent observations*)、方差齐性(*homoscedasticity*)及正态分布(*normal distribution*)。而分组数据中的多次测量数据之间可能存在某种相关性,用常规的统计方法就不能充分揭示出其内在的特点,有时甚至会得出错误的结论。

近二十年来,混合效应模型已广泛地应用于社会科学、生物统计学、计量经济学和统计学等领域,如教育学、心理学、社会学、经济学、生物学、医学、药物动力学、农业、林业、工业研究等学科。在文献中,混合模型被冠以各种不同的名称,如混合效应模型(Goldstein, 1995; Pinheiro and Bates, 2000); 随机效应模型(*random effects model*)(Laird and Ware, 1982; Longford, 1987); 随机系数模型(*random-coefficient model*)(Longford, 1993); 协方差成分模型(*covariance components model*)(Dempster et al., 1981); 分层线性模型(*hierarchial linear model*)(Bryk and Raudenbush, 1992)。

混合效应模型是既包含固定效应又包含随机效应的一类模型。主要应用于分组数据,如纵向数据(*longitudinal data*)、重复测量数据(*repeated measures data*)、分层数据(*hierarchical data*)及多变量多层数据(*multivariate multilevel data*)。从统计分析角度上讲,传统分析方法在分析分组数据时所遇到的问题可通过混合效应模型得到解决。由于混合效应模型可对组水平和个体水平的数据进行同时分析,在一个模型中可同时检验组变量和个体变量的效应。另外,混合效应模型分析不需要假设数据中的观察值相互独立,因而可以修正因观察数据

的非独立性引起的参数标准误差估计。混合效应模型可以分为线性混合效应模型和非线性效应混合模型。

1.1 线性模型

假定有 n 个观测值 y_1, \dots, y_n ，同时有一组解释变量 $x_{11}, \dots, x_{1p}, x_{21}, \dots, x_{2p}, \dots, x_{n1}, \dots, x_{np}$ 。构造的一般线性模型：

$$y_i = \sum_{j=1}^p x_{ij} \beta_j + e_i \quad i=1, \dots, n \quad (1-1)$$

式中， β_1, \dots, β_p 为未知的固定效应参数； e_i 为未知独立一致正态分布随机误差向量， $e_i \sim N(0, \sigma^2)$ ； σ^2 为方差。

将上述模型写成详细的表达式：

$$\begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{bmatrix} = \begin{bmatrix} x_{11} & x_{12} & \cdots & x_{1p} \\ x_{21} & x_{22} & \cdots & x_{2p} \\ \vdots & \vdots & \cdots & \vdots \\ x_{n1} & x_{n2} & \cdots & x_{np} \end{bmatrix} \begin{bmatrix} \beta_1 \\ \beta_2 \\ \vdots \\ \beta_n \end{bmatrix} + \begin{bmatrix} e_1 \\ e_2 \\ \vdots \\ e_n \end{bmatrix} \quad (1-2)$$

通常将上式可写成矩阵形式：

$$\begin{cases} y = X\beta + e \\ e \sim N(0, \sigma^2) \end{cases} \quad (1-3)$$

1.2 广义线性模型

在线性模型中，我们总是要求模型误差 e 的各个分量是不相关的，并且具有相同的方差 σ^2 ，即协方差满足条件 $Cov(e) = \sigma^2 I$ 。但不等方差的情况是经常存在的，在这种情况下协方差应满足条件 $Cov(e) = \sigma^2 R$ 。

为了和线性混合模型对比，广义线性模型写成矩阵形式：

$$\begin{cases} y = X\beta + e \\ e \sim N(0, \sigma^2 R) \end{cases} \quad (1-4)$$

式中， y 为 $n \times 1$ 维观测向量； X 为 $n \times p$ 维设计矩阵； β 为 $p \times 1$ 维向量； e 为 $n \times 1$ 维误差向量； N 为正态分布； σ^2 为方差； R 为相关矩阵，当 $R = I_n$ 时，模型为一般线性模型； I_n 为单位阵。

1.3 单水平线性混合模型

1.3.1 模型定义

线性混合效应模型是广义线性模型的进一步深化, 在广义线性模型中, 所有的参数都是固定效应参数, 而在线性混合效应模型中, 含有随机效应参数, 同时误差 e 有更为灵活的结构, 包括相关性和方差不齐性。单水平线性混合效应模型 (single level linear mixed effects model, SLMEM) 的形式:

$$\begin{cases} y_i = X_i\beta + Z_i b_i + e_i, & i = 1, \dots, n \\ b_i \sim N(0, G) \\ e_i \sim N(0, \sigma^2 R_i) \end{cases} \quad (1-5)$$

式中, y_i 为第 i 组中 $n_i \times 1$ 维观测向量; X_i 为 $n_i \times p$ 维已知固定效应设计矩阵; β 为 $p \times 1$ 维未知的固定效应参数向量; Z_i 为 $n_i \times q$ 维已知随机效应设计矩阵; b_i 为 $q \times 1$ 维未知的随机效应参数向量, b_i 的期望值 $E(b_i) = 0$, b_i 的方差为 $V(b_i) = \sigma^2 G_1$, $Cov(b_i, b_{i'}) = 0$; e_i 为 $n_i \times 1$ 维误差向量, e_i 的期望值 $E(e_i) = 0$, e_i 的方差为 $V(e_i) = \sigma^2 R_i$, $Cov(e_i, e_{i'}) = 0$ (当 $i \neq i'$), 这表示组间的观察是彼此独立的, 但 e_i 的组内数据间可能存在相关性, 即组内的观察可能彼此不独立; b_i 与 e_i 的协方差 $Cov(b_i, e_i) = 0$; G 为随机效应 $q \times q$ 维协方差矩阵; $\sigma^2 R_i$ 为 $n_i \times n_i$ 维协方差矩阵。模型的假定:

$$E \begin{pmatrix} b \\ e \end{pmatrix} = \begin{bmatrix} 0 \\ 0 \end{bmatrix}, \quad Var \begin{pmatrix} b \\ e \end{pmatrix} = \begin{bmatrix} G & 0 \\ 0 & R \end{bmatrix} \quad (1-6)$$

上述模型可用向量和矩阵形式:

$$\begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{bmatrix} = \begin{bmatrix} X_1 \\ X_2 \\ \vdots \\ X_n \end{bmatrix} [\beta] + \begin{bmatrix} Z_1 & 0 & \cdots & 0 \\ 0 & Z_2 & \cdots & 0 \\ \vdots & \vdots & \cdots & \vdots \\ 0 & 0 & \cdots & Z_n \end{bmatrix} \begin{bmatrix} b_1 \\ b_2 \\ \vdots \\ b_n \end{bmatrix} + \begin{bmatrix} e_1 \\ e_2 \\ \vdots \\ e_n \end{bmatrix} \quad (1-7)$$

模型的假定可以写成详细的表达式:

$$E(b) = 0, \quad E(e) = 0 \quad (1-8)$$

$$Var(b) = \sigma^2 \begin{bmatrix} G_1 & 0 & \cdots & 0 \\ 0 & G_1 & \cdots & 0 \\ \vdots & \vdots & \cdots & \vdots \\ 0 & 0 & \cdots & G_1 \end{bmatrix} = \sigma^2 G \quad (1-9)$$

$$\text{Var}(e) = \sigma^2 \begin{bmatrix} R_1 & 0 & \cdots & 0 \\ 0 & R_2 & \cdots & 0 \\ \vdots & \vdots & \cdots & \vdots \\ 0 & 0 & \cdots & R_n \end{bmatrix} = \sigma^2 R \quad (1-10)$$

y 的方差:

$$V(y) = ZV(b)Z' + V(e) = \sigma^2(ZGZ') + \sigma^2 R = \sigma^2(ZGZ' + R) = \sigma^2 V$$

式中, $V = ZGZ' + R$ 。

当构造随机效应设计矩阵 Z 且给定 G 和 R 的协方差结构后, 就可以得到 V 。当 $Z=0$ 和 $R=I_n$ 时, 线性混合效应模型就转化为一般线性模型; 当 $Z=0$ 和 $R \neq I_n$ 时, 线性混合效应模型就转化为广义线性模型。

1.3.2 模型的参数估计

混合效应模型中由于存在组内误差结构和随机效应的方差-协方差结构, 使模型的参数估计更为复杂。在混合效应模型中有两个设计矩阵: 固定效应设计矩阵 X 和随机效应设计矩阵 Z , 这两个矩阵都需要构造。此外需要同时估计 4 种参数, 即固定效应参数 β 、随机效应参数 b 、随机效应的方差-协方差矩阵 G 和随机误差的方差-协方差矩阵 R 。下面举例介绍线性混合模型设计矩阵 X 、 G 、 Z 、 R 的构造形式。

假定有 m 株树, 每株树有连续 3 年的每年 1 次树高 y 的测量值。研究 y 随时间的线性趋势模型。所构造的 X 矩阵为一个 $3m \times 2$ 维矩阵, 其结构:

$$X = \begin{bmatrix} 1 & 1 \\ 1 & 2 \\ 1 & 3 \\ \vdots & \vdots \\ 1 & 1 \\ 1 & 2 \\ 1 & 3 \end{bmatrix}$$

X 矩阵中的第一列是截距, 第二列是斜率(时间)。

G 、 Z 和 R 可分别写成如下形式:

$$G = \begin{bmatrix} \sigma_1^2 & 0 & \cdots & 0 \\ 0 & \sigma_1^2 & \cdots & 0 \\ \vdots & \vdots & \cdots & \vdots \\ 0 & 0 & \cdots & \sigma_1^2 \end{bmatrix}$$

不同的模型时,模型估计需采用极大似然法;当比较随机效应不同,而固定效应相同的模型时,模型估计需采用限制极大似然法。(Raudenbush and Bryk, 2002)。一般来说,极大似然法比较灵活,当模型的固定效应不同,随机效应也不同,但两个模型的随机效应是嵌套时,也可以采用极大似然法进行模型比较。注意,在 S-Plus 软件中的 lme 模块默认的是限制极大似然法,而在 nlme 模块默认的是极大似然法。因为有时许多非线性混合效应嵌套模型的固定效应和随机效应都是不同的,采用限制极大似然法无法进行某些似然比检验。限制极大似然法有时会过低估计模型的方差,许多分析倾向于使用限制极大似然法。因此当模型比较完成后,需应用限制极大似然法运行最优模型,得到模型的估计结果。

1.4 多水平线性混合模型

单水平线性混合模型可以非常容易地扩展到多水平线性混合模型(multilevel linear mixed effects model, MLMEM)。如利用3水平分级数据(树木嵌套于样地,样地又嵌套于小班,小班又嵌套于林班)可以建立一个3水平多层模型。模型表达式:

$$\begin{cases} y_{ijk} = X_{ijk}\beta + Z_{i,jk}b_i + Z_{ij,k}b_{ij} + Z_{ijk}b_{ijk} + e_{ijk}, \\ \quad i = 1, \dots, m, j = 1, \dots, m_i, k = 1, \dots, m_{ij} \\ b_i \sim N(0, D_1), b_{ij} \sim N(0, D_2), b_{ijk} \sim N(0, D_3) \\ e_{ijk} \sim N(0, \sigma^2 R_i) \end{cases} \quad (1-13)$$

式中, y_{ijk} 为第 i 个 1 水平中的第 j 个 2 水平中的第 k 个 3 水平内的观察值; m 为第 1 水平的分组数; m_i 为对应于第 1 水平的第 2 水平的分组数; m_{ij} 为对应于第 1 和 2 水平的第 3 水平的分组数; X_{ijk} 为已知设计矩阵; β 为固定参数向量; $Z_{i,jk}$, $Z_{ij,k}$ 和 Z_{ijk} 分别为 1 水平、2 水平和 3 水平的随机效应设计矩阵; b_i , b_{ij} 和 b_{ijk} 分别为 1 水平、2 水平和 3 水平的随机参数向量; D_1 , D_2 和 D_3 分别为 1 水平、2 水平和 3 水平的随机参数的方差-协方差矩阵; R_i 为方差-协方差结构, e_{ijk} 为模型的误差项。

1.5 非线性混合效应模型

非线性混合效应模型可以看作是线性混合效应模型或者非线性模型的扩展(Bates and Watts, 1988),随机效应参数加入到非线性模型的固定参数中,并且随不同的区组(group)变化而变化。

1.5.1 单水平非线性混合模型

单水平非线性混合效应模型(single level nonlinear mixed effects model, SNLMEM)的形式如下:

$$\begin{cases} y_{ij} = f(\varphi_{ij}, v_{ij}) + e_{ij}, & i=1, \dots, M, j=1, \dots, n_i \\ \varphi_{ij} = A_{ij}\beta + B_{ij}b_i \\ b_i \sim N(0, D) \\ e_{ij} \sim N(0, \sigma^2 R_i) \end{cases} \quad (1-14)$$

式中, y_{ij} 为第 i 个区组中的第 j 次观测值; M 为区组的数量; n_i 为在第 i 个区组上观测的次数; f 为含有参数向量 φ_{ij} 和协变量向量 v_{ij} 的函数; β 为 $(p \times 1)$ 维固定效应向量, b_i 是带有方差-协方差矩阵 D 的 $(q \times 1)$ 维随机效应向量; A_{ij} 和 B_{ij} 为相应的设计矩阵; e_{ij} 为服从正态分布的误差项; σ^2 为方差; R_i 为区组 i 的方差-协方差矩阵; D 为随机效应的方差-协方差矩阵。

1.5.2 多水平非线性混合模型

单水平非线性混合模型可以非常容易地扩展到多水平非线性混合模型(multilevel nonlinear mixed effects model, MNLMEM)。以两水平模型为例, 模型形式如下:

$$\begin{cases} y_{ijk} = f(\varphi_{ijk}, v_{ijk}) + e_{ijk}, & i=1, \dots, M, j=1, \dots, M_i, k=1, \dots, n_{ij} \\ \varphi_{ijk} = A_{ijk}\beta + B_{i,jk}b_i + B_{ijk}b_{ij} \\ b_i \sim N(0, D_1) \\ b_{ij} \sim N(0, D_2) \\ e_{ij} \sim N(0, \sigma^2 R_{ij}) \end{cases} \quad (1-15)$$

式中, y_{ijk} 为第 i 个 1 水平中的第 j 个 2 水平内的第 k 次观察值; M 为第 1 水平的分组数量; M_i 为第 1 水平内第 2 水平的分组数量; n_{ij} 为第 i 个 1 水平中的第 j 个 2 水平内的观测次数; f 为含有参数向量 φ_{ijk} 和协变量向量 v_{ijk} 的函数; β 为 $(p \times 1)$ 维固定效应向量; A_{ijk} 为设计矩阵; b_i 为第 1 水平带有方差-协方差矩阵 D_1 的 $(q_1 \times 1)$ 维随机效应向量, b_{ij} 为第 2 水平带有方差-协方差矩阵 D_2 的 $(q_2 \times 1)$ 维随机效应向量, b_i 和 b_{ij} 不相关; $B_{i,jk}$ 和 B_{ijk} 为随机效应的设计矩阵; e_{ij} 为服从正态分布的误差项; σ^2 为方差, R_{ij} 为第 i 个 1 水平中的第 j 个 2 水平内的方差-协方差矩阵。

1.6 模型拟合和比较

混合模型拟合评价提供了多种信息标准(information criteria)。常用的三种信息标准为: 赤池信息准则(akaike information criteria, AIC), 贝叶斯信息准则(bayesian information criteria, BIC)和对数似然值(log likelihood, LL)或-2 倍的对数似然值(-2LL)。在这些标准中, AIC, BIC, -2LL 的值越小和 LL 的值越大, 表明模型拟合数据越好。AIC 和 BIC 的计算式如下:

$$AIC = -2LL + 2d \quad (1-16)$$

$$BIC = -2LL + d \ln(n) \quad (1-17)$$

式中, LL 为最大似然函数的对数值; d 为模型中估计参数个数; n 为有效观察个数; $\ln(n)$ 为 n 的自然对数。

当用不同模型拟合同一个数据集时, 模型之间的比较可以是嵌套模型(nested model), 也可以是非嵌套模型(non-nested model)。对于嵌套模型, 即一个模型是另一个模型的亚模型(sub-model)时, 可用似然比检验(likelihood ratio test, LRT)进行比较。假定现有一个一般模型(general model), 另一个是限制模型(restricted model)。 L_2 和 L_1 分别是一般模型和限制模型的似然值(likelihood value)

$$LRT = 2 \log \left(\frac{L_2}{L_1} \right) = 2 [\log(L_2) - \log(L_1)] \quad (1-18)$$

应用式(1-18), 必须有 $L_2 > L_1$, 则 $\log(L_2) > \log(L_1)$, 此时 LRT 为正值。LRT 的分布近似卡方分布 (χ^2 distribution), 其自由度为两个模型的参数个数之差 ($df = k_2 - k_1$)。如果有 $LRT > \chi^2(k_2 - k_1, \alpha)$, 则对应于 L_2 的模型是备选模型。

对于非嵌套模型, LRT 检验不再适用。在这种情况下, 可用 AIC, BIC, -2LL 值或 LL 值进行模型比较。通常以 AIC 和 BIC 为主要判断指标。它们既可以用于嵌套模型的比较, 也可以用于非嵌套模型的比较。如果这些统计值很近似, 则选取含参数个数最少的模型。

1.7 混合效应模型构建的步骤

模型的构建是一个复杂的系统过程, 没有一个特定的系统适合于所有类型的数据。一般来说, 模型的建立通常是一个既基于统计学考虑, 又基于理论考虑的探索过程, 同时模型又能提供生物学意义的解释。混合模型的建模步骤在许多研究中都有过讨论(Hox, 1995; Singer, 1998; Pinheiro and Bates, 1998; 王济川等, 2008; 石磊等, 2013)。结合前人所推荐的方法, 我们总结了混合效应模型构建的7个步骤。