

线性模型分析方法

——适用于动物科学和动物医学

王继华 李旭东 著



科学出版社

线性模型分析方法

——适用于动物科学和动物医学

王继华 李旭东 著

科学出版社

北京

内 容 简 介

本书通过大量实例详细介绍线性模型的建立方法和统计分析的基本原理、方法与常见问题,包括回归分析模型及其应用、方差-协方差分析模型及其应用、多元线性模型及其应用、线性混合模型及其应用、线性混合模型参数估计方法等,统计分析方法包括最小二乘(LS)法、最小范数二次无偏估计(MINQUE)、最大似然(ML)法和约束最大似然(REML)法。

本书介绍的线性模型分析方法和技术,只需要理工科大学本科的数学基础。书中给出的算法规律性很强,便于理解、使用和记忆。本书还对线性模型建立和分析中常见的问题及注意事项做了详细介绍。

本书适合动物科学、动物医学和农学等非数学专业的科技人员、高年级本科生和研究生阅读,对于相关专业的中青年大学教师,也有参考价值。

图书在版编目(CIP)数据

线性模型分析方法:适用于动物科学和动物医学/王继华,李旭东著. —北京:科学出版社,2015.6

ISBN 978-7-03-044957-3

I. ①线… II. ①王… ②李… III. ①线性模型-分析方法 IV. ①O212

中国版本图书馆 CIP 数据核字(2015)第 129330 号

责任编辑:李 欣 / 责任校对:钟 洋

责任印制:徐晓晨 / 封面设计:陈 敬

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>

北京数图印刷有限公司印刷

科学出版社发行 各地新华书店经销

*

2015 年 9 月第 一 版 开本:720×1000 B5

2015 年 9 月第一次印刷 印张:16 3/4

字数:338 000

定价:98.00 元

(如有印装质量问题,我社负责调换〈印科〉)

作者简介



王继华(1957年11月~),男,河北省大名县人,教授,硕士生导师,河北省省级优秀教师.主要研究动物科学,包括动物育种原理与方法、动物营养与饲料等.已发表学术论文85篇.已出版学术著作9部:《家畜育种学导论》《鸡饲料配方设计技术》《动物科学研究方法》《仔猪饲料配方设计实用新技术》《仔猪饲料配方设计高级技术》《肥育猪饲料配方设计技术》《猪的生长发育及营养调控技术》《饲用矿物元素配合物的研究及应用》《蛋鸡饲料配方设计技术》等.



李旭东(1982.1.26-),河海大学博士,现在四川省农田水利局工作.

序

目前国内外已有一些线性混合模型类出版物可读,但有的并不适合非数学类工程科技人员,尤其是动物、农业和医学等生命科学领域的科技人员,而这是使用线性混合模型方法最多的人群.此书通过大量实例,深入浅出地介绍线性模型分析方法的基本技术,包括基本概念、原理、方法和应用的完整技术.

此书系统总结国内外对线性模型分析方法的研究成果和王继华教授多年来在动物科学研究和生产上应用线性模型分析方法的研究心得.此书介绍的这些分析线性模型的统计方法,只涉及理工科大学本科的数学知识,便于非数学专业的科技人员理解和使用.

此书特点是理论联系实际,利用线性模型分析方法处理动物饲养试验的实际问题,例如,动物饲养试验分析方法、动物营养需要量的研究方法、饲料利用率研究方法等,资料丰富、数据翔实,读者可以按照书中实例模仿使用.

承蒙作者信任和重托,本人有幸率先阅读了这部书稿,看到线性模型分析的新方法和新应用,蕴涵了作者的大量心血,是他多年来从事动物科学研究和实践的智慧结晶.本人乐为此书作序,并热切希望本书的出版能够提升我国动物科学研究和生产实践水平,促进我国养殖业的可持续性健康发展.

中国畜牧兽医学会动物遗传育种学分会 理事长
中国农业大学 动物科技学院教授、博士生导师

张 勤

2015年3月1日

前 言

多年来,方差-协方差分析和回归分析一直是统计分析模型的基础,这些技术有个基本假定,那就是模型残差或误差项是独立同分布的.线性混合模型方法放松了这一基本假定,使我们可以处理更加复杂的数据结构.

线性混合模型在理工农医等各个领域得到迅速推广和应用,原因是在实际问题中,参与试验的个体是随机抽取的,并且我们的研究目标常不是这些个体本身的特征,而是它们所在总体的特征,这时把个体效应视为随机效应引入模型,可以巨大提高模型精度.有时候一些固定效应也属于干扰因素,例如,在不同猪场测定某些添加剂的使用效果,猪场效应就是干扰因素,可以在设计试验时把一些固定效应引入模型.设计试验时在模型中引入干扰因素,无论引入随机效应还是引入固定效应,都可巨大提高模型精度,巨大提高模型参数的推断质量.所以统计学家一直非常重视线性混合模型的研究,使得线性混合模型理论迅速发展.

使用混合模型有很多好处,有时候可以提高估计量的精度,得到更广泛的推断;有时候使用更加合适的数学模型使我们能够更加深入地洞察数据结构,启迪我们的研究思路和研究方法.尤其是,常见通用软件,如 SAS 和 SPSS 都增加了线性混合模型分析程序,使得线性混合模型技术更便于普及.但是,如何合理使用这些软件包,还需要正确理解线性混合模型分析方法.本书没有介绍有关软件包的使用法,是出于两种考虑:一是,网上可以下载到有关程序使用说明,例如,SAS 的程序 Proc Mixed 或 GLIMMIX,和 SPSS 的程序 MIXED 等;二是,本书使用 OFFICE 的 EXCEL 软件计算全部实例,有助于读者理解线性混合模型的算法和过程.

本书介绍的线性模型分析技术适用于许多应用领域,希望为非数学专业的工程科技人员提供有用的工具,包括动物科学、动物医学和农学专业的研究人员,或想从事量化研究的非数学专业的研究生,以及相关专业的中青年大学教师.

本书介绍的线性模型分析方法和技术,是作者多年研究与教学实践的结晶,这些方法只需要理工科大学本科的数学基础(包括微积分、线性代数和数理统计或生物统计学).本书初稿作为动物科学专业的高年级本科生和研究生讲课的讲稿使用,已在本校内部印刷了 2 次.教学实践表明,动物科学专业高年级本科生完全可以理解和正确使用我们给出的线性模型分析方法.本书给出的添加约束的方法不仅便于理解和使用,而且,由于模型约束条件本身的特性,使得这些方法具有很强的规律性,与传统的线性混合模型正规方程一样,非常便于写出,便于记忆.

中国畜牧兽医学会动物遗传育种学会理事长张勤教授百忙中仔细审阅了本

书全稿,提出不少修改意见和建议,使本书质量大大提高.先生的学风和境界,学为人师,德为世范,我们再次对先生的厚爱和支持表示衷心的感谢和敬意.

本书的初始形态为讲稿,十几年来,每遇见简单易懂便于学生理解的精辟论述,就会毫不犹豫地插入讲稿,所以本书实际上参考了很多文献,在这里向原作者致以衷心的感谢.

限于本人知识和技术水平,书中难免有不足之处,恳望读者不吝赐教,电子邮箱是 hdwangjihua@126.com.

王继华

2015年5月1日

于河北工程大学动物科学系

目 录

第 1 章 引论	1
1.1 动物科学中的线性数学模型	1
1.2 数学模型化是现代研究方法的核心	4
1.3 线性模型学习方法	4
参考文献	5
第 2 章 线性模型基础知识	6
2.1 概述	6
2.2 线性模型的种类	7
2.2.1 线性回归模型	7
2.2.2 方差分析模型	11
2.2.3 协方差分析模型	15
2.2.4 固定效应模型与随机效应模型	17
2.2.5 混合效应模型	19
2.3 固定因子与随机因子	22
2.3.1 固定因子与随机因子的概念	23
2.3.2 固定因子与随机因子的辨识	23
2.4 真模型、选模型与等价模型	25
2.4.1 真模型与选模型	25
2.4.2 等价模型	27
2.4.3 随机误差与模型残差的区别	27
2.5 线性模型实例	29
2.5.1 机理分析模型与试验模拟模型	29
2.5.2 线性模型建模实例	30
2.6 线性模型的技术含量——交叉设计实例	40
2.6.1 平衡数据	40
2.6.2 不平衡数据	43
参考文献	45
第 3 章 最小二乘法	46
3.1 线性模型的基本假定	46
3.1.1 一元线性模型的概念	46

3.1.2	一元线性模型的基本假定	47
3.1.3	线性模型的常见形式及其关系	49
3.2	线性模型参数的可估性	50
3.2.1	线性模型的误差平方和	50
3.2.2	模型参数的可估性	50
3.2.3	观测数据的结构平衡性	51
3.3	模型参数的最小二乘估计	52
3.3.1	估计原则	52
3.3.2	模型参数的最小二乘估计方法	53
3.3.3	最小二乘估计的优良性质	54
3.3.4	多因子的最小二乘正规方程	55
3.3.5	最小二乘估计量的剩余误差方差及其估计	56
3.3.6	最小二乘参数估计量的统计特征	57
3.4	关联矩阵列不满秩的处理	57
3.4.1	正规方程系数矩阵不满秩的例子	58
3.4.2	正规方程系数矩阵不满秩的处理	60
3.4.3	实际例子的处理	62
3.5	多因子模型系数矩阵列不满秩的参数估计	63
3.6	假设检验	67
3.6.1	假设检验的原理	68
3.6.2	假设检验的类型	69
3.6.3	假设检验的方法	69
3.6.4	假设检验的简约方法——子模型法	69
3.6.5	假设检验:关联矩阵和检验条件矩阵都满秩	70
3.6.6	假设检验:一般情况	72
3.6.7	假设检验的两类错误	73
3.6.8	统计检验的效力	75
3.7	置信区间	75
3.7.1	单个样本平均数的置信区间	75
3.7.2	多个平均数的线性组合的置信区间	76
3.7.3	正态向量总体平均数的检验	77
3.7.4	正态向量总体平均数的置信区间	78
3.8	多重比较	80
	参考文献	83
第4章	线性回归模型及其应用	84

4.1 线性回归模型的基本假定	84
4.1.1 回归分析模型的基本概念	84
4.1.2 线性回归模型的基本假设	85
4.1.3 线性回归模型的模型诊断	86
4.2 线性回归分析原理与方法	87
4.2.1 线性回归分析概述	87
4.2.2 线性回归方程的显著性检验	89
4.2.3 回归系数的显著性检验	91
4.2.4 依变量的预测	93
4.2.5 广义回归模型	95
4.2.6 线性回归分析方法要点	96
4.3 回归分析实例	96
4.3.1 模型参数估计	96
4.3.2 舍入误差与算法	99
4.3.3 回归系数可靠性检验	101
4.3.4 模型的决定系数	102
4.3.5 模型参数的置信区间	103
4.4 线性回归模型的同—性检验	103
4.4.1 一般原理	103
4.4.2 比较不同回归模型的实例	105
4.5 回归模型的建立	106
4.5.1 全部自变量的可能组合	106
4.5.2 逐步搜索选择变量	107
4.5.3 简单向前搜索选择变量	108
4.5.4 向后剔除法	108
4.6 线性回归模型的使用技术	109
4.6.1 数据缺失与子模型分析方法	109
4.6.2 数据规模与模型中有效参数的个数	109
4.6.3 多重共线性	110
4.6.4 定量预测的准确性	111
4.6.5 回归模型建模策略	113
4.7 回归模型在动物营养与饲料研究中的应用	114
4.7.1 用抛物线模型估计动物营养需要量	114
4.7.2 动物营养需要量估计实例	115
4.7.3 用斜率比法或平行线法估计养分的生物学效率	116

4.7.4	影响生物学效率评估结果的因素	118
4.7.5	动态饲料数据库及饲料配方的可靠性	119
4.8	分段回归模型在动物营养与饲料研究中的应用	121
4.8.1	动物营养剂量反应规律	121
4.8.2	用折线回归模型估计动物营养需要量	124
4.8.3	多折点回归模型	125
4.8.4	分段回归模型	127
	参考文献	129
第5章	方差与协方差分析模型及其应用	132
5.1	单向分类的方差分析模型	132
5.2	双向分类无互作效应的方差分析模型	135
5.2.1	双向分类模型	135
5.2.2	广义最小二乘解	136
5.2.3	解的特性	138
5.2.4	可估函数	138
5.2.5	最小二乘均数	139
5.2.6	可估函数的方差	140
5.3	方差分析模型的假设检验	140
5.3.1	方差分析	140
5.3.2	假设检验	141
5.3.3	子模型分析方法	142
5.4	有互作效应的方差分析模型	143
5.5	数据缺失	146
5.5.1	数据缺失的原因	146
5.5.2	缺失整组数据	147
5.5.3	观测效应的关联性	148
5.6	协方差分析模型	149
5.6.1	协方差分析模型及分析要点	150
5.6.2	动物饲养试验实例	153
5.6.3	更复杂的协方差分析模型	156
	参考文献	158
第6章	一般线性模型及其扩展应用	159
6.1	一般线性模型	159
6.1.1	与一元线性模型的关系	159
6.1.2	模型参数的估计	160

6.1.3 关联矩阵列不满秩时的参数估计	161
6.2 误差结构矩阵之逆的简化计算	164
6.3 一般线性模型的统计推断	164
6.3.1 预估问题	164
6.3.2 参数的统计检验	165
6.3.3 设计矩阵列不满秩时的参数检验	165
6.4 多元线性模型	166
6.4.1 多元线性模型概述	167
6.4.2 多元线性模型的参数估计	168
6.4.3 多元线性模型的假设检验	169
6.4.4 多元线性模型的预估及其精度	170
6.5 多元线性模型分析示例	171
6.5.1 建立数学模型	172
6.5.2 估计模型参数和协方差矩阵	174
6.5.3 模型的预测	174
6.5.4 预测值的误差限	175
6.5.5 假设检验	175
6.5.6 多元线性模型假设检验的非典型情况	176
参考文献	178
第7章 线性混合模型及其在动物科学中的应用	179
7.1 一般线性混合模型及其参数估计	179
7.2 最大似然(ML)估计	182
7.2.1 数据分布	182
7.2.2 由密度函数到似然函数	183
7.2.3 未知参数的最大似然估计	184
7.2.4 线性固定模型参数的最大似然估计	185
7.2.5 线性固定模型参数 ML 估计量的假设检验	185
7.3 线性混合模型参数的最大似然估计	186
7.3.1 最大似然估计的方法原理	186
7.3.2 最大似然法的计算方法——EM 算法	189
7.3.3 最大似然法计算实例	189
7.4 随机效应间有相关的情况	190
7.4.1 随机效应间有相关时的迭代算法	190
7.4.2 随机效应相关矩阵分析	191
7.5 模型参数估计量的统计检验	192

7.5.1	模型参数估计的质量	192
7.5.2	G 和 R 已知	194
7.5.3	G 和 R 未知	194
7.6	动物饲养试验分析实例	195
7.7	固定效应估计与显著性检验实例	200
7.8	随机效应的BLUP与显著性检验示例	206
	参考文献	211
第8章	线性混合模型的参数估计及应用	213
8.1	方差分量估计的概念与意义	213
8.1.1	为什么估计方差分量	213
8.1.2	方差分量的概念	214
8.1.3	方差分量分析方法	215
8.2	MINQUE法	216
8.2.1	MINQUE法的数学原理	216
8.2.2	通过混合模型方程组求MINQUE	218
8.2.3	MINQUE计算实例	219
8.2.4	方差分量的统计检验与区间估计	221
8.3	约束最大似然法	222
8.3.1	约束最大似然法的要点	222
8.3.2	方差分量REML估计的计算	222
8.3.3	REML估计方程的导出	224
8.3.4	REML法计算实例	227
8.3.5	随机变量间有相关时的REML	229
8.4	ML和REML估计量的可靠性检验	230
8.4.1	ML和REML方差估计量的可靠性检验	230
8.4.2	固定效应和随机效应估计量的显著性检验	232
8.4.3	ML和REML估计量的置信区间	233
8.5	线性混合模型参数估计方法的应用问题	233
8.5.1	负的方差分量估计值	233
8.5.2	方差参数的精确性	235
8.5.3	固定效应与随机效应标准误差的偏差	235
8.5.4	ML和REML的比较	235
8.5.5	MINQUE, I-MINQUE和ML, REML比较	236
	参考文献	237
第9章	多元线性混合模型	239

9.1	两个依变量时方差协方差组分的 REML	239
9.2	多性状动物模型	241
9.3	扩展的 Canonical 转换	244
9.3.1	只含一个随机效应时的数据转换	244
9.3.2	含多个随机效应时的数据转换	246
9.4	方法示例	247
9.4.1	直接多元分析	248
9.4.2	通过转换作间接多元分析	249
9.5	多性状动物模型的 REML	251
9.6	REML 方法总结	252
	参考文献	253

第 1 章 引 论

在数量遗传学和数量生态学带动下,20 世纪 70 年代开始兴起用数学模型方法研究动物现象.科学家已由真理发现者转换为模型建立者(model builder),构造假设、建立模型、发展理论是科学家的首要任务.科学研究过程由归纳-推理过程转变为假设-求证过程.科学理论的判断标准也不再是看它能否被证明绝对正确无误,而是看它能否被重复检验.

1.1 动物科学中的线性数学模型

(1) 生命存在的方式——质和量.任何事物都有质和量两方面,一个研究对象,量的表现常可用一个系统来描述,可把影响系统运动的各种因素表示为不同因子或变量(常量可看做是只取一个数值的变量),用数学模型描述变量间及变量与系统间的关系,由此研究系统的运动规律.

家畜系统也表现为质和量两方面.影响一个家畜系统的各种因子或变量间的关系,以及这些因子(或变量)与家畜系统之间的关系,同样可以而且应该用数学模型描述.只有把家畜系统中各变量间以及各变量与系统间的复杂关系用数学模型表述,这门科学才算发展到高级阶段.

(2) 动物科学发展趋势.过去的动物科学,没有数学也取得了辉煌成绩.现在和今后的研究是否可以不用数学?科学发展史表明,任何一门科学,初级阶段总是以定性描述开始,深入到一定阶段,积累大量研究结果,对本门科学有一定程度了解后,便开始总结对研究对象的认识,形成知识.在地球生物圈中,除人类本身外,家畜是人类接触最早、认识最深的一个生命系统,每种家畜都构成一个子系统.人类对家畜系统进行了千万年观察研究,积累了大量感性认识,已总结出不同侧面的知识,形成了动物科学的不同分支,例如生理学、生物化学、营养学、遗传学、繁殖学等.

模型化方法就是把研究对象(原型)的一些次要的细节、非本质的联系舍去,从而以简化和理想化的形式去再现原型的各种复杂结构、功能和联系.作为一种现代科学认识手段和思维方法,模型具有两方面的含义,即抽象化和具体化.

数学模型技术正成为现代动物科学研究的核心技术.世界权威杂志 J. Anim. Sci. 在 1970 年以前,没发表“模型化”论文,20 世纪 70 年代发表的有 3.2%,80 年代发表的有 17%,90 年代发表的有 31.5%,2000 年后发表的论文有 50%涉及数

学模型,这会是一个增长最快的领域.再者,“模型化”内涵也变为机理建模,而不停留在经验建模水平.动物系统的机理模型化已为理解动物科学原理和生产过程贡献很多,数学模型化方法正大步进入动物科学所有学科.

(3) 数学模型的种类.数学模型可按不同方式分类:

按应用领域(或所属学科)可分为遗传模型、营养模型、生态模型、数学生理学模型、人口模型、环境模型等,范畴更大则形成许多边缘学科,如生物数学、医学数学、数量经济学等.

按建模目的分为描述模型、分析模型、预报模型、优化模型、决策模型、控制模型等.

按照对模型结构的了解程度分:有白箱、灰箱、黑箱模型,把研究对象比喻成一只暗箱,要通过模型化来揭示它的奥妙.

按模型的表现特性分:有确定性模型和随机性模型、静态模型和动态模型(是否考虑时间因素)、离散模型和连续模型(指模型中的变量为离散型还是连续型的)等.虽然从本质上讲大多数实际问题是随机、动态、非线性的,但由于确定性、静态、线性模型容易处理,并常可作为初步近似来解决问题,所以模型化时,常先考虑确定性、静态、线性模型.而且实际上,很多非线性问题在小范围内可用线性模型逼近.连续模型便于利用微积分方法求解,作理论分析,而离散模型便于计算机计算,所以用哪种模型要具体问题具体分析.在具体模型化过程中将连续模型离散化,或将离散变量视作连续,也是常用方法.

按变量间的关系可分为线性模型和非线性模型.按变量的形式可分为实变量模型和虚变量模型.实变量模型是指模型中的变量为连续型变量,可取任意实数值,虚变量模型是指模型中的变量为离散型变量,有些变量只取值 0 或 1.虚变量模型中,按模型中变量的性质,又可将模型分为固定模型、随机模型和混合模型.根据研究对象的观测值来源可分为单向分类模型、双向分类模型、多向分类模型和系统分类模型.

按模型对实际问题的配合程度看有 3 类模型:真模型(real model)可精确描述数据,没有不可解释的变异.真模型多数不能精确知道,也不一定是线性模型.理想模型(ideal model)是真模型的简化或理想化.可操作模型(apply model)是理想模型的简化版,能用于数据分析.模型不能无限简化,不能简化到没分析价值.要掌握理想模型简化为可操作模型的过程所做的假定,才能判断可操作模型的质量.

数学模型化的优点是能够创建科学理论,进行科学推断、解释和预测、得到观测范围外的结果.缺点是,如果模型过于简化,为在概念上可控,把原型简化成概念的骨架,那么,由此得到的结论在回用于原型时,符合程度就很差.原型越复杂,过分简化的缺点就越大.

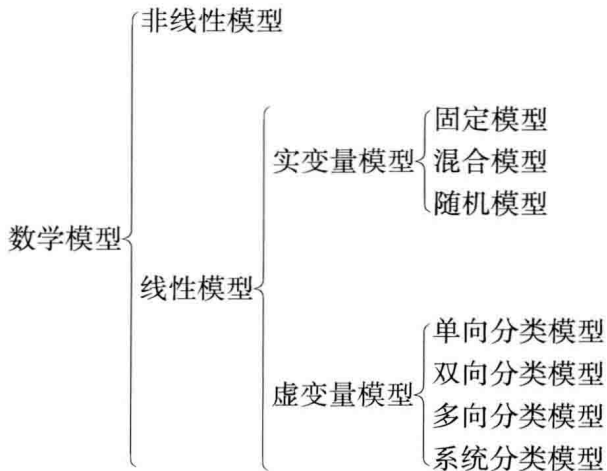
(4) 线性数学模型.在数理统计模型中,数学模型是指描述观测值与影响观测

值变异的各因素间关系的数学方程. 所有统计分析都是基于一定数学模型进行的. 线性统计模型简称线性模型(linear model), Rao(1973)最早给出系统介绍. 这种统计模型统一了许多传统的线性统计问题, 例如, 平均数估计、线性回归、方差-协方差分析和数量化方法 I, 所以是生物学研究和应用领域最重要、最基本的一种数学模型. 多数参数方法都可归为线性模型.

由线性模型发展来的各种方法和技术, 可在一定范围内推广到非线性模型. 实际问题多可用线性模型描述, 或作为初级近似描述. 模型参数的估计方法有多种, 这些方法多可用线性模型理论解决. 大学本科生物统计教材上介绍的经典问题, 多可用线性模型表述.

在应用领域, 一般从三方面对线性模型分类: 一是按因子数目, 分为单因子、二因子、三因子和多因子模型等; 二是按因子性质分为固定模型、随机模型、混合模型等; 三是按模型功能分为回归模型、方差分析模型、协方差分析模型、方差分量模型、混合模型等. 上述分类方法只是考虑到模型某一方面的特征, 实践中常要考虑各种特征.

从应用角度, 一般把数学模型分为如下种类:



统计模型化有两个基本思想. 模型化一般基于某些假设, 所以建立模型前验证这些假设非常重要. 建好模型后还要从两方面评价模型: 第一是模型拟合度. 由模型拟合的数据是否接近观测到的样本数据? 拟合数据与样本数据间的差是否呈随机分布? 第二是评价把模型用于预测更广范围的数据时的可靠性. 在一个范围内成立的数学模型, 在另一个范围内未必成立, 例如, 元明粉(无水硫酸钠)对生长猪的生理作用, 当日粮中添加量在 $0 \sim 0.3\%$ 时, 有健胃、促生长作用; 在 $0.4\% \sim 0.6\%$ 时, 有健胃、促生长和软化粪便作用; 在 0.6% 以上时, 有倾泻作用.