

SPSS 应用系列丛书 (2)

世界优秀统计工具

SPSS 11

统计分析教程

高级篇

张文彤 主编



北京希望电子出版社
Beijing Hope Electronic Press
www.bhp.com.cn

SPSS 应用系列丛书 (2)

世界优秀统计工具

SPSS 11

统计分析教程

高级篇

张文彤 主编



北京希望电子出版社
Beijing Hope Electronic Press
www.bhp.com.cn

内 容 简 介

SPSS 是世界最为优秀的统计工具之一, 深受各行业用户的青睐, SPSS 11.0 是其最新版本。

本书为《SPSS 11.0 统计分析教程》的高级篇, 由 4 部分 15 章及 1 个附录组成。主要包括: 一般线性模型、混合线性模型、多元线性回归与曲线拟合、分类资料的回归分析、非线性回归及其他回归过程、对数线性模型、聚类分析与判别分析、因子分析与对应分析、信度分析与多维尺度分析、生存分析、缺失值分析等。

本书作者从统计专业用户的角度出发, 结合自身多年的 SPSS 使用经验, 在以风趣、明快的笔触介绍软件操作的同时, 注意将相应的统计学知识融入其中。书中既有深入浅出的软件功能介绍, 又有针对实际问题的解决办法, 更侧重于对统计新方法、新观点的讲解。

本书不仅是 SPSS 10~11 版的通用入门教材, 也是各行业数据开发、应用的广大从业人员的重要指导书, 同时也可作为大专院校相关专业的参考书。

本版 CD 为 SPSS11.0 相关材料和书中有关数据, 并赠送 SPSS11.0 试用版软件。

系列盘书名 : SPSS应用系列丛书(2)

盘 书 名 : 世界优秀统计工具SPSS 11.0统计分析教程(高级篇)

文 本 著 者 : 张文彤 等

责 任 编 辑 : 郭淑珍

C D 制 作 者 : 希望多媒体开发中心

C D 测 试 者 : 希望多媒体测试部

出 版、发 行 者 : 北京希望电子出版社

地 址 : 北京市海淀区知春路63号卫星大厦三层 100080

网 址 : www.bhp.com.cn

E-mail: lwm@bhp.com.cn

电 话 : 010-62520290,62521724,62528991,62630301,62524940,62521921,82610344

(发行) 010-62613322-215 (门市) 010-82675588-501,82675588-201 (编辑部)

经 销 : 各地新华书店、软件连锁店

排 版 : 希望图书输出中心 全卫

C D 生 产 者 : 北京中新联光盘有限责任公司

文 本 印 刷 者 : 北京媛明印刷厂

开 本 / 规 格 : 787 毫米×1092 毫米 1/16 23.00 印张 528 千字

版 次 / 印 次 : 2002 年 6 月 第 1 版 2002 年 6 月 第 1 次 印刷

印 数 : 0001-5000 册

本 版 号 : ISBN 7-900101-23-3

定 价 : 40.00 元 (本版 CD)

说明: 凡我社产品如有残缺, 可执相关凭证与本社调换。

前 言

计算机永远是属于年轻人的，统计软件也是如此。作为一个极具活力和开拓精神的软件，SPSS (Statistical Product and Service Solutions) 自身的进步非常迅速，而它近几年来更是以令人惊讶的速度在国内得到了迅速普及，这无疑是使用者对该软件本身的肯定。但是，在拥有这一优秀软件的同时，广大用户希望得到一本优秀 SPSS 参考书的呼声也越来越高。现在随着 11.0 版的正式推出，这一问题变量更加明显。近年来，国内也出现了数种非官方的 SPSS 教材，总体看来，各有优势，但不足之处也比较明显：

- ◇ SPSS 是一个非常权威而严肃的统计软件，可现在的许多应用型教材都存在着各种各样的常识性错误，统计理论也似是而非，用户完全无法从中体会到 SPSS 的强大功能，最后以为该软件只能作一些简单的分析，华而不实。这对 SPSS 的形象造成了极坏的影响。
- ◇ 现有几种统计专业人士编写的 SPSS 书籍在统计理论上是非常严肃的，但大多以编程为主线，或者仍然以老版本界面操作为主，内容从 6.0 一直到 10.0 都通用，完全没有体现出 SPSS 10~11 版许多出色的新功能。同时这些教材没有照顾到非统计专业人员的特点，写得过于专业、晦涩难懂，并不适合初学者入门。
- ◇ 据笔者不完全统计，所有的教材都只涉及到 SPSS 约 2/3 的常用功能，中、高级统计分析要么完全不涉及，要么走马观花的一笔带过，完全没有实用价值。特别是在市场研究领域被广泛应用的几个 SPSS 模块始终没有相应教材可用，这不能不说是一个遗憾。
- ◇ 没有真正易学易用的教材。理想的教材应当是深入浅出，同时幽默风趣，以激发读者自学的兴趣。非常可惜，完全满足这些特色的教材还没有出现。

针对以上问题，本书的写作目的是尽快提高广大使用者的水平，真正掌握最新的 SPSS 11.0 版的强大功能。按照读者的不同层次，本书分为基础篇和高级篇两册，共 31 章，内容以 SPSS 11 为准，包括了全部 10 个模块的所有主要功能。在写法上充分注意了通俗易懂：基础篇涵盖了常用的统计分析方法，入门部分强调文字轻松愉快，同时体现出 SPSS 操作中最具特色的功能和操作技巧。其中第 1 章专为初学者而准备，如果希望快速入门，读者可以在学习完第 1 章后直接跳到基础统计部分继续学习；基础统计部分充分考虑到了非统计专业人员的特点，将统计理论融入软件介绍之中，力求深入浅出；高级篇中的统计模块介绍则以统计学理论为准绳，立足于应用实例将统计方法、界面操作与结果解释结合讲述，使读者学后能真正掌握相应方法，而不是只明白了对话框的中文含义。为保证质量，其中比较重要的一般线性（方差分析）模型、多元线性回归模型、Logistic 模型、时间序列模型、生存分析均不惜篇幅详细介绍，并多由专人负责编写。

所谓术业有专攻，笔者主要从事的是医学统计和市场研究统计分析，因此书中的大部分数据实例都来自这两个领域，但本书更多的是作为 11.0 版的通用入门教材而编写，因此适用于各行业的初/中级用户。由于 SPSS 的功能极为强大，全部学习完毕需要相当长的时间，对于希望快速入门的朋友，这里针对不同专业给出参考阅读顺序如下：

2007

- ◇ 临床科研工作者：基础篇 1、7、11~15 章，高级篇 1、4、5、13 章。但如果需要从事新药临床试验工作，则请务必阅读基础篇第 4 章，以补充编程知识。
- ◇ 市场研究工作者：基础篇 1、7、9~16 章，高级篇 1、4、8~11 章，其中高级篇 8~11 章涉及到了许多市场研究专用方法，是学习的重点。
- ◇ 社会学工作者：基础篇 1、7、11~15 章，高级篇 1、4、8~10 章。

但是，为了保证学习的系统性和连贯性，只要时间允许，笔者仍然强烈建议朋友们按章节顺序渐进学习，这样才能真正体会到 SPSS 软件的强大实力。

本书的雏形来自笔者在医学统计之星网站 (<http://www.MesStatStar.com>) 上连载的 SPSS 10 教程，但更多的心得体会则来自于长期从事 SPSS、SAS 教学的积累，以及数年来的 SPSS 使用经验。从 2000 年 4 月创作网上教程算起，全书的写作一共用了 20 个月的时间，虽然网络教程一面世就立刻受到了广大网友的热烈欢迎，但为了保证质量，真正向大家奉献一本精品，我们没有急于求成，将其匆匆结册出版，而是三易其稿，力求文字浅显易读，统计理论正确无误，并结合最新的 11.0 版的情况，对内容进行了长达三个月的详细修订。一言以蔽之，我们无愧于心。

本书同时也是作为复旦大学研究生用教材，编委基本上都是复旦大学卫生统计与社会医学教研室的年轻统计教师。其中田晓燕负责编写统计绘图部分，刘晓云负责分类资料的回归分析部分，罗剑峰负责时间序列部分，董伟负责生存分析部分，张文彤负责其余各章节和全书的统稿工作。

非常感谢苏立民总经理为我提供了这样一个机会，使我能够将自己的使用经验和大家一起分享。在成书过程中我还得到了卫生统计学、数理统计学和社会学界许多前辈、师长和朋友们的热心指导和帮助，在此一并致谢。

为节省篇幅，全书大部分的实例数据并未在书中列出，大家可以在书后所附光盘上找到它们，也可以到我的医学统计之星网站上下载。书后光盘中同时包括了 SPSS 公司主要产品的一月试用版，包括 SPSS 11 Base 一月试用版，从而大家就可以按照书中的叙述进行操作，让结果在屏幕上真实再现，相信这样对朋友们的学习更有帮助。

教材的编写有各种各样的风格，本书风格的最大的特点就是幽默风趣，易于理解。如果将 SPSS 用户手册比作辞海的话，我们的编写目的就是提供一本大家学习 SPSS 时得心应手、图文并茂的新华字典，成为初学者快速成长为 SPSS 专家的桥梁。SPSS 是一款非常出色的统计软件，祝大家在本书的帮助下能够将它用得开心，玩得愉快。

限于作者水平，书中错缪之处难免，还请同行专家和广大读者不吝赐教。为便于交流，现列出各编委的电子信箱如下，欢迎大家就任何问题与我们联系。

张文彤：wtzhang@spss.com.cn

田晓燕：xytian@epscn.com

刘晓云：liuxiaoyun@hotmail.com

罗剑峰：jfluo@shmu.edu.cn

董伟：dwshmu@sina.com

张文彤

2002 年元旦于复旦公卫学院

目 录

第一部分 一般线性与混合线性模型

第 1 章 征服一般线性模型

—General Linear Model

菜单详解(上)	2
1.1 方差分析模型简介	3
1.1.1 模型入门	3
1.1.2 常用术语	4
1.1.3 方差分析模型的适用条件	5
1.2 Univariate 过程入门	6
1.2.1 引例	6
1.2.2 界面说明	7
1.2.3 结果解释	12
1.2.4 对引例的进一步分析	13
1.3 常用试验设计及分析方法详解	15
1.3.1 完全随机设计 (Completely Random Design)	15
1.3.2 配伍设计 (Randomized Block Design)	15
1.3.3 交叉设计 (Cross-over Design)	16
1.3.4 析因设计 (Factorial Design)	18
1.3.5 拉丁方设计 (Latin Square Design)	19
1.3.6 正交设计 (Orthogonal Design)	21
1.3.7 星点设计 (Central Composite Design)	24
1.3.8 嵌套设计(Nested Design) 与裂区设计 (Split-plot Design)	24
1.4、协方差分析	27

1.4.1 概述	27
1.4.2 预分析:线性趋势的判断	28
1.4.3 预分析:检验各组总体 斜率是否相等	28
1.4.4 正式分析:比较修正均数 有无差异	29

第 2 章 征服一般线性模型

—General Linear Model

菜单详解(下)	32
2.1 Multivariate 过程	32
2.1.1 引例与界面说明	33
2.1.2 结果解释	33
2.1.3 对引例的进一步分析	35
2.2 Repeated Measures 过程	36
2.2.1 引例	37
2.2.2 界面说明	38
2.2.3 结果解释	39
2.2.4 对引例的进一步分析	42
2.3 Variance Components 过程	43
2.3.1 引例	43
2.3.2 界面说明	44
2.3.3 结果解释	45

第 3 章 混合线性模型入门

—Mixed Model 菜单详解

3.1 模型简介	46
3.1.1 模型入门	47
3.1.2 混合效应模型的用途	49
3.2 Linear 过程	49
3.2.1 引例与界面说明	49
3.2.2 结果解释	54
3.2.3 对引例的进一步分析	55
3.3 混合线性模型分析实例	58
3.3.1 家庭聚集性数据	59
3.3.2 重复测量数据	59
3.3.3 嵌套设计数据	60

第二部分 回归分析

第4章 多元线性回归与曲线拟合

——Regression 菜单详解(上) . . . 64

4.1	Linear 过程	65
4.1.1	线性回归模型简介	65
4.1.2	引例与界面说明	68
4.1.3	结果解释	73
4.1.4	对引例的进一步分析	74
4.1.5	一个多元回归实例	77
4.2	关于线性回归的高级话题	79
4.2.1	衡量多元线性 回归方程的标准	79
4.2.2	强影响点的诊断及对策	81
4.2.3	多重共线性问题及对策	82
4.2.4	分类自变量的设置与 哑变量的使用	84
4.2.5	趋势面分析	86
4.2.6	通径分析(Path Analysis)	86
4.3	Curve Estimation 过程	87
4.3.1	引例	87
4.3.2	界面说明	88
4.3.3	结果解释	89

第5章 分类资料的回归分析

——Regression 菜单详解(中) . . . 91

5.1	Binary Logistic 过程	91
5.1.1	模型简介	91
5.1.2	引例	92
5.1.3	界面说明	92
5.1.4	结果解释	97
5.1.5	对引例的进一步分析	99
5.2	关于 Logistic 模型的高级话题	100
5.2.1	模型中的假设检验方法	100
5.2.2	模型的自变量设置方法	101
5.2.3	模型诊断	104
5.2.4	配对 Logistic 回归模型	107
5.3	Multinomial Logistic 过程	109
5.3.1	引例	110
5.3.2	界面说明	110

5.3.3 结果解释 113

5.4	Ordinal 过程	115
5.4.1	引例	115
5.4.2	界面说明	116
5.4.3	结果解释	118
5.5	Probit 过程	119
5.5.1	引例	119
5.5.2	界面说明	120
5.5.3	结果解释	121

第6章 非线性回归及其他回归过程

——Regression 菜单详解(下) . . . 124

6.1	Nonlinear Regression 过程	124
6.1.1	引例与界面说明	125
6.1.2	结果解释	129
6.1.3	对非线性模型的深入探讨	131
6.2	Weight Estimation 过程	132
6.2.1	引例与界面说明	132
6.2.2	结果解释	133
6.2.3	对引例的进一步分析	135
6.3	Two-Stage Least-Squares 过程	135
6.3.1	引例与界面说明	136
6.3.2	结果解释	138
6.4	Optimal Scaling 过程	139
6.4.1	引例与界面说明	140
6.4.2	结果解释	144
6.4.3	对引例的进一步分析	145

第三部分 多元统计分析方法

第7章 对数线性模型

——Loglinear 菜单详解 148

7.1	模型简介	148
7.1.1	原理	148
7.1.2	模型选择	149
7.2	General 过程	150
7.2.1	引例	150
7.2.2	界面说明	151
7.2.3	结果解释	152
7.2.4	对引例的进一步分析	156
7.3	Logit 过程	156

7.3.1 引例与界面说明.....	157	10.1 Reliability Analysis 过程.....	213
7.3.2 结果解释.....	158	10.1.1 引例与界面说明.....	214
7.3.3 对引例的进一步分析.....	161	10.1.2 结果解释.....	216
7.4 Model Selection 过程.....	161	10.2 Multidimensional Scaling 过程.....	217
7.4.1 引例.....	161	10.2.1 引例与界面说明.....	218
7.4.2 界面说明.....	162	10.2.2 结果解释.....	221
7.4.3 结果解释.....	163	10.3 Multidimensional Scaling (PROXSCAL)过程.....	224
第 8 章 聚类分析与判别分析		10.3.1 引例.....	225
——Classify 菜单详解.....	166	10.3.2 界面说明.....	225
8.1 K-means Cluster 过程.....	166	10.3.3 结果解释.....	230
8.1.1 引例与界面说明.....	167	10.3.4 对引例的进一步分析.....	232
8.1.2 结果解释.....	169	第 11 章 结合分析.....	234
8.1.3 对引例的进一步分析.....	170	11.1 模型简介.....	234
8.2 Hierarchical Cluster 过程.....	171	11.1.1 为什么使用结合分析.....	234
8.2.1 引例.....	171	11.1.2 常用术语.....	235
8.2.2 界面说明.....	172	11.1.3 结合分析的基本步骤.....	236
8.2.3 结果解释.....	175	11.1.4 SPSS 中的相应过程.....	236
8.3 Discriminant 过程.....	177	11.2 Orthogonal Design 子菜单.....	237
8.3.1 模型简介.....	177	11.2.1 Generate 项.....	237
8.3.2 引例.....	180	11.2.2 Display 项.....	240
8.3.3 界面说明.....	181	11.3 CONJOINT 过程.....	241
8.3.4 结果解释.....	183	11.3.1 引例及语法说明.....	241
8.3.5 对引例的进一步分析.....	186	11.3.2 结果解释.....	243
第 9 章 因子分析与对应分析		11.3.3 对引例的进一步分析.....	246
——Data Reduction 菜单详解.....	190	第四部分 其他高级统计分析方法	
9.1 Factor Analysis 过程.....	190	第 12 章 岁月如歌	
9.1.1 模型简介.....	191	——Time Series 菜单详解.....	250
9.1.2 引例.....	193	12.1 时间序列的建立和平稳化.....	251
9.1.3 界面说明.....	194	12.1.1 缺失值的填补	
9.1.4 结果解释.....	197	——Replace Missing Values	
9.1.5 对引例的进一步分析.....	200	过程.....	251
9.2 Correspondence Analysis 过程.....	202	12.1.2 时间变量的定义	
9.2.1 引例与界面说明.....	203	——Define dates 过程.....	252
9.2.2 结果解释.....	205	12.1.3 时间序列的平稳化	
9.3 Optimal Scaling 过程.....	208	——Create Time Series	
9.3.1 引例与界面说明.....	208	过程.....	254
9.3.2 结果解释.....	210	12.2 时间序列的图形化观察.....	258
第 10 章 信度分析与多维尺度分析			
——Scale 菜单详解.....	213		

12.2.1 序列图 (Sequence Chart)	258	13.3.1 引例与界面说明	297
12.2.2 自相关图 (Autocorrelation Chart)	260	13.3.2 结果解释	300
12.2.3 互相关图 (Cross-correlation Chart)	264	13.3.3 对引例的进一步分析	302
12.2.4 谱密度图 (Spectral Chart)	266	13.3.4 Life Tables 过程 与 Kaplan-Meier 过程的比较	303
12.2.5 交叉谱图 (The Cross-Spectrum)	268	13.4 Cox Regression 过程	304
12.3 Exponential Smoothing 过程	269	13.4.1 模型简介	304
12.3.1 模型简介	269	13.4.2 引例及界面说明	305
12.3.2 引例与界面说明	270	13.4.3 结果解释	309
12.3.3 结果解释	272	13.4.4 对引例的进一步分析	311
12.4 Autoregression 过程	273	13.5 关于 Cox 模型的高级话题	312
12.4.1 模型简介	273	13.5.1 分类自变量的定义 与比较方法	312
12.4.2 引例与界面说明	274	13.5.2 Cox 模型中的分层分析	312
12.5 ARIMA 过程	276	13.5.3 配对 Logistic 回归	313
12.5.1 ARMA 模型简介	277	13.5.4 竞争风险 (Competing risks) 的 Cox 模型	315
12.5.2 标准建模步骤	278	13.5.5 复发性疾病的 Cox 模型	315
12.5.3 界面说明	279	13.6 Cox w/Time-Dep Cov 过程	316
12.5.4 综合分析实例	280	13.6.1 模型简介	316
12.6 季节解构 ——Seasonal Decomposition 过程	286	13.6.2 引例与界面说明	317
12.6.1 引例与界面说明	286	13.6.3 结果解释	319
12.6.2 结果解释	287	13.6.4 分析时依 Cox 模型时的 注意事项	320
第 13 章 生存分析——Survival 菜单详解	290	第 14 章 缺失值分析 ——Missing Value Analysis 过程详解	321
13.1 生存分析简介	290	14.1 缺失值理论简介	321
13.1.1 应用背景	290	14.1.1 数据的缺失方式	321
13.1.2 基本术语	291	14.1.2 SPSS 中可用的缺失值 处理方法	322
13.1.3 SPSS 中相应模块简介	292	14.2 界面说明	323
13.2 Life Tables 过程	292	14.3 分析实例	326
13.2.1 引例	293	14.3.1 缺失值的生成及分析操作	326
13.2.2 界面说明	294	14.3.2 结果解释	327
13.2.3 结果解释	295	14.3.3 对引例的进一步分析	328
13.3 Kaplan-Meier 过程	297		

第 15 章 其他统计分析功能	
— 不得不说的故事	331
15.1 典型相关分析	331
15.1.1 方法简介	331
15.1.2 引例及语法说明	331
15.1.3 结果解释	332
15.2 岭回归分析	335
15.2.1 方法简介	335
15.2.2 引例及语法说明	336
15.2.3 结果解释	336
15.3 广义线性模型简介	337
附录 SPSS 公司部分软件介绍	340
参考文献	345

第一部分 一般线性与混合线性模型

第 1 章 征服一般线性模型

——General Linear Model 菜单详解 (上)

第 2 章 征服一般线性模型

——General Linear Model 菜单详解 (下)

第 3 章 混合线性模型入门

——Mixed Model 菜单详解


第1章 征服一般线性模型


——General Linear Model 菜单详解（上）

在我看来，学习统计软件时不去了解所用统计原理和模型，就像是学英语时只背单词而不记语法一样，是永远不可能说出一口流利的英文的。

——张文彤

请注意，本章的标题用了一些修辞手法，一般线性模型可不是三言两语就能说清楚的，因为它包括的内容实在太多了。那么，究竟什么是一般线性模型呢？用最通俗的话讲：方差分析模型和线性回归模型都是分析一个/多个自变量对一个/多个连续性应变量的影响，并且都假设应变量和自变量是线性数量关系，因此它们都算是般线性模型的特例，凡是和它们沾边的都可以用 GLM 来做。包括成组设计的方差分析（即单因素方差分析）、配伍设计的方差分析（即两因素方差分析）、多因素方差分析、多元方差分析、重复测量的方差分析、协方差分析、多元线性回归分析等等。因此，能真正掌握 GLM 菜单的用法，会使大家的统计分析能力有极大的提高。

 什么？GLM 还可以完成线性回归分析！是的，它的确可以。只不过它作回归分析的功能没有专用的回归分析模块强，所以该功能较少被用到。

 有些书中将 General Linear Models 翻译成广义线性模型，实际上 GLM 模块的功能只限于一般线性模型。广义线性模型对应的英文是 Generalized Linear Models，指的是范围更广、功能更强大的另一大类模型，在 SPSS 中可以编程拟合，但并非 GLM 菜单提供的对话框就可以完成的，在本书的最后一章中对其有基本的介绍。

好了，本章要讲述的菜单名为 General Linear Model，既然一般线性模型的能力如此强大，那么下属的四个过程各自的功能是什么呢？请看：

- ◇ Univariate 过程：四个菜单中的大哥大，当应变量只有一个时，我们所进行的分析就要用它来完成。显然，它是用的最多的一个。
- ◇ Multivariate 过程：当结果变量（应变量）不止一个时，当然要用它来分析啦！
- ◇ Repeated Measures 过程：顾名思义，重复测量的数据就要用它来分析，这一点可能要强调一下，用前两个菜单似乎都可以分析出来结果，但在许多情况下该结果是不正确的。在相应章节我会详细讲述。
- ◇ Variance Components 过程：用于对层次数据拟合方差成份模型，它是普通线性模型向随机效应的进一步扩展，是一种可以考察各个层次因素的变异大小，从而为哪些层次上可能存在组内聚集性、如何可能减小数据变异提供信息的统计方法，也就是现在非常热门的多水平模型的最原始形式。

由于 Univariate 过程是最常用的一个过程，本章我们就围绕它来进行讲解，另外三个过程将放在下一章学习。

1.1 方差分析模型简介

现实世界中变量间的联系是错综复杂的,比如要研究性别对身高的影响,显然就要考虑到年龄、遗传、营养状况等因素的作用。这时 t 检验作为单因素分析方法就无能为力了,而方差分析可以在控制其他因素影响的同时研究两者之间的关系,分析的效率更高,适用范围更广。

同时,许多时候各自变量之间还会存在交互作用,如研究催化剂对化学反应的催化能力,如果该催化剂只在某个温度范围内效果最佳,则只单独研究该催化剂的催化作用是没有实际意义的,此时这种交互作用也成为了我们研究的重点,即必须要研究在什么温度条件下该催化剂的催化能力最佳。对交互作用的分析也是方差分析模型的特长。

1.1.1 模型入门

在均数的比较一章中我们实际上已经接触到了方差分析的基本思想,这里首先来复习一下单因素方差分析中变异的分解公式:

总变异=处理因素导致的变异+随机变异

实际上,该公式只是为了便于大家理解,标准的单因素方差分析模型如下:

$$X_{ij} = \mu + \alpha_i + \varepsilon_{ij}$$

其中 X_{ij} 表示第 i 组的第 j 个观察值; μ 表示总体的平均水平; α_i 表示影响因素在 i 水平下对应变量的附加效应,并假设所有 α_i 之和应当为 0; ε_{ij} 为一个服从正态分布 $N(0, \sigma^2)$ 的随机变量,代表随机误差。一般情况下,我们做假设检验实际上就是检验各个 α_i 是否均为 0,如都为 0,即各组总体均数都相等,则 X_{ij} 就会成为服从正态分布 $N(\mu, \sigma^2)$ 的一个变量。



有的时候以上模型也被写为 $X_{ij} = \mu_i + \varepsilon_{ij}$, 此时检验的含义就变成了各组均数 μ_i 是否相同。显然,它和上面的公式是等价的,但这种写法应当更加容易理解。比如在上一章中我们比较三组石棉矿工的用力肺活量有无差别,那么石棉肺患者组中某一位观察对象的用力肺活量数值就等于石棉肺患者组的平均水平再加上一个随机误差项。而如果三组总体均数相同,则它就应当等于总体均数(平均水平)再加上一个随机误差项,实际上就变成了同一个变量分布中的某一点。

多因素方差分析模型只是对以上公式的进一步扩展,以两因素方差分析模型为例,公式如下:

$$X_{ijk} = \mu + \alpha_i + \beta_j + \varepsilon_{ijk}$$

其中 α_i 、 β_j 分别表示 A 因素 i 水平和 B 因素 j 水平的附加效应, ε_{ijk} 仍为随机误差变量。此时如果要说明因素 A 有无影响,就是检验如下假设:

$H_0: \alpha_i = 0, H_1: \text{至少有一个 } \alpha_i \neq 0$

如果要说明因素 B 有无影响,就是检验如下假设:

$H_0: \beta_i = 0, H_1: \text{至少有一个 } \beta_i \neq 0$

而所说的模型无显著性就是指上面两个 H_0 同时成立（均不能被拒绝）。

更为复杂的是如考虑交互作用的情形，模型如下：

$$X_{ijk} = \mu + \alpha_i + \beta_j + \alpha_i\beta_j + \varepsilon_{ijk}$$

其中 α_i 、 β_j 分别表示 A 因素 i 水平和 B 因素 j 水平的附加效应。 $\alpha_i\beta_j$ 则为两者的交互效应。

可能有的朋友已经看得头痛了，但无论怎样，只要记住方差分析的原理即可：根据资料类型以及研究目的，可将总变异分解为两个或多个部分，除一部分代表随机误差的作用外，每个部分的变异可由某因素的作用（或交互作用）来解释，通过比较可能由某因素所致的变异与随机误差的大小，借助 F 分布做出推断，即可了解该因素对结果变量的影响是否存在。

1.1.2 常用术语

方差分析中的常用术语有：

1. 因素 (Factor)：因素是可能对应变量有影响的变量，一般来说，因素会有不止一个水平，而分析的目的就是考察或比较各个水平对应变量的影响是否相同。例如影响农作物产量的因素有气温、降雨量、日照时间等。在方差分析中，因素的取值范围不能无限，只能有若干个水平，即应当为分类变量。

2. 水平 (Level)：因素的不同取值等级称作水平，例如性别有男、女两个水平。需要注意的是有时候水平是人为划分出来的，比如身高被分为高、中、低三个水平。

3. 单元 (Cell)：指各因素水平之间的组合，我们所说的方差齐就是指的各个单元间的方差齐。注意有的试验设计并不会给出所有可能的水平组合，如拉丁方设计。

4. 元素 (Element)：指用于测量应变量值的最小单位，比如研究石棉矿工用力肺活量，则肺活量是从每一位矿工身上测得，矿工就是试验的元素。一个单元格内可以有多个元素，也可以只有一个，甚至于没有元素。

5. 均衡 (Balance)：如果一个实验设计中任一因素各水平在所有单元格中出现的次数相同，且每个单元格内的元素数相同，则该试验是均衡的，否则，就被称为不均衡。不均衡的实验设计在分析时较为复杂。

6. 固定因素 (Fixed Factor) 与随机因素 (Random Factor)：两者都是因素的不同种类，固定因素指的是该因素在样本中所有可能的水平都出现了。换言之，该因素的所有可能水平仅此几种，针对该因素而言，从样本的分析结果中就可以得知所有水平的状况，无需进行外推。比如要研究糖尿病、IGT（糖耐量异常）和正常人的血糖有无差别，则按照糖尿病有无可将所有人分为糖尿病、IGT 和正常人三种，此时该因素就被认为是固定因素。另外，有些人设定的因素，比如下面引例中的 group，它被分为工前、工中、工后这三个值，如果我们所做的检验就是想弄清楚这三个水平对应变量有无影响，不需要外推到其他水平（如工后半小时），则它也可以被认为是固定因素。

和固定因素相对应的是随机因素，它指的是该因素所有可能的取值在样本中没有都出现，或不可能都出现。如下面引例中的 worker，它表示的是工人这个配伍因素，实际上总体中当然不可能只有这 10 个工人，他们只是所有工人的代表而已。因此要用

样本中 10 个工人对应变量的影响情况来推论总体中全体工人对应变量的影响情况, 包括未出现的那些工人, 这不可避免的存在误差(即随机效应), 需要估计该误差的大小, 因此被称为随机因素。又如我们要研究什么温度下催化剂的效果最好, 样本中取了 30、40、50°C 三个水平, 如果我们在分析结果中能同时外推 35°C、45°C 这些水平的情况, 此时温度也是随机因素。

一般来说固定因素和随机因素在分析时应分别指定, 如果将随机因素按固定因素来分析, 则可能得出错误的分析结果。但是如果所有单元格内都至多只有一个元素, 则随机效应无法被估计出来, 此时两种做法的统计分析结论完全相同。

7. 交互作用(Interaction): 如果一个因素的效应大小在另一个因素不同水平下明显不同, 则称为两因素间存在交互作用。当存在交互作用时, 单纯研究某个因素的作用是没有意义的, 必须分另一个因素的不同水平研究该因素的作用大小。

如果所有单元格内都至多只有一个元素, 则交互作用无法测量, 只能不予考虑, 最典型的例子就是配伍设计的方差分析。

1.1.3 方差分析模型的适用条件

方差分析并不是万金油, 它也有自己的适用条件(以 H_0 假设成立为前提), 具体说有以下几点:

- ◇ 各样本的独立性: 只有各样本为相互独立的随机样本, 才能保证变异的可加性(可分解性)。
- ◇ 正态性: 即所有观察值系从正态总体中抽样得出。
- ◇ 方差齐: 这里所说的方差齐是指假设总的模型无意义时方差齐, 亦即每一个单元格中的方差齐。

道理大家其实都懂, 但做起来就不一定都明白了, 为什么多因素方差分析一般都不提这些条件呢? 首先在以上条件中, 对独立性的要求是最严的, 但它一般都可以保证。其次为正态性和方差齐性, 具体来说就是:

- ◇ 单因素方差分析: 适用条件是必须要考虑的问题, 尤其是正态性和方差齐性一般都需要进行考察。
- ◇ 配伍设计的方差分析: 不考虑这两个问题, 这是由于正态性和方差齐性的考察是以单元格为基本单位的, 此时每个格子中只有一个元素, 当然没法分析了! 不止是配伍设计的方差分析, 只要是无重复数据的方差分析, 如交叉设计、正交设计等, 都不考虑这两个问题。
- ◇ 有重复数据的多因素方差分析: 由于方差齐性的考察是以单元格为基本单位, 此时单元格数目众多, 真正分到每个格子中的样本例数一般都只有 3~5 例, 此时很难检验出差别; 或者出现另一种极端情况, 因为极个别格子方差参差不齐而导致检验不能通过, 这种情况实际上对分析结果影响并不太严重。因此在多因素方差分析中, 方差齐性往往只限于理论探讨的程度。真正在应用时, 只要数据分布不是明显偏态, 不存在极端值即可, 这两种情况对结果的影响要比方差不齐严重得多。



根据 Box 的研究结果, 如果各组的例数相同 (即均衡), 或总体呈正态分布, 则方差分析模型对方差略微不齐有一定的耐受性, 只要最大/最小方差之比小于 3, 分析结果都是稳定的。



在实际操作中, 对数据正态性的考察有一个办法, 就是拟合完毕后作出残差分布图, 如果残差呈随机分布, 则可知 (单元格内) 原始数据满足正态条件。

1.2 Univariate 过程入门

如前所述, Univariate 过程是最为常用的一个过程, 几乎所有的实验设计都可以利用它来进行分析, 它也是后面三个过程的基础。下面, 我们就以配伍组设计的方差分析为例来看一下该过程的用法。

1.2.1 引例

例 1.1 某厂医务室测定了 10 名氟作业工人工前、工中及工后 4 小时的尿氟浓度 ($\mu\text{mol/L}$)。问氟作业工人在这三个不同时间的尿氟浓度有无差别? (杨树勤, 《卫生统计学》第三版 P46)

工人编号	工前	工中	工后
1	90.53	142.12	87.38
2	88.43	163.17	65.27
3	47.37	63.16	68.43
4	175.80	166.33	210.54
5	100.01	144.75	194.75
6	46.32	126.33	65.27
7	73.69	138.96	200.02
8	105.27	126.33	100.01
9	86.32	121.06	105.27
10	60.01	73.69	58.95

	group	worker	fu
1	1.00	1.00	90.5
2	1.00	2.00	88.4
3	1.00	3.00	47.4
4	1.00	4.00	175.8
5	1.00	5.00	100.0
6	1.00	6.00	46.3

图 1.1 数据格式示意

解: 显然, 该数据的特点和以前学习过的配对数据非常相似, 也属于同一受试对象不同处理间的比较。但是, 这里的处理组成为了三组 (工前、工中和工后), 如果采用两两求出差值的方法来做配对 t 检验, 显然会犯和多个均数比较时用两两 t 检验一样的错误, 会扩大一类错误。这里我们应采用两因素的方差分析方法来分析。

根据统计分析的要求, 我们建立了三个变量来包括上述信息, time 表示测定时间为工前、工中或工后, worker 代表区组, weight 表示最终的体重增量, 如图 1.1 所示。建好的数据集存为文件 twoway.sav。

由于在配伍设计的数据中, 每个格子中只有一个元素 (无重复数据), 因此交互作用和方差齐性都无法检验, 下面只分析主效应即可。

Analyze → General Linear Model → Univariate

Dependent List 框: fu

要分析的应变量为 fu

Fixed Factor 框: group、worker

固定效应变量为 group、worker

Model:

Custom

要求自定义方差分析模型

Build Terms 下拉列表: Main effects	模型中准备纳入主效应
Model 框: group、worker	模型中只纳入主效应 group、worker
Continue	
Post Hoc:	
Post Hoc Test For 框: group	要求对 group 作两两比较
<input checked="" type="checkbox"/> S-N-K	具体的两两比较方法为 SNK 法
Continue	
OK	

这里 worker 本应被选入随机因素框，此处将其也选入固定因素框是为了让输出结果和原书结果一致，同时也是基于下面的事实：在无重复数据的多因素方差分析中，随机效应无法被估计出来，两种做法在统计分析结论上完全相同。

1.2.2 界面说明

【主对话框】(见图 1.2)

1. **Dependent Variable** 框：选入需要分析的变量（应变量），只能选入一个。
2. **Fixed Factor(s)** 框：选入自变量中的固定因素，对固定因素和随机因素的解释请参见前面的“常用术语”。
3. **Random Factors** 框：用于选入自变量中的随机因素变量，如果你看完前面的解释后还弄不明白，假装没看见这个框框就是了。

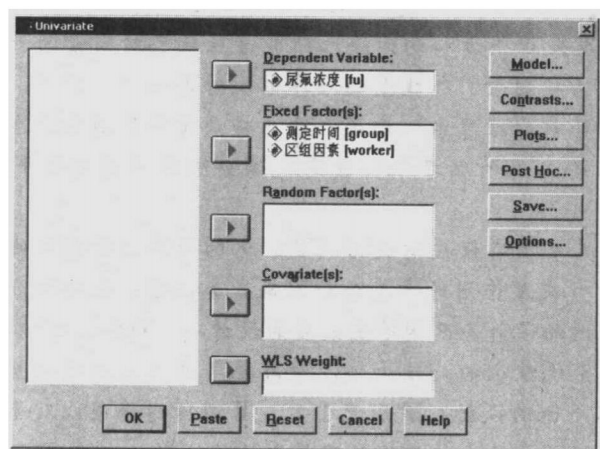


图 1.2 主对话框

4. **Covariate** 框：用于选入协变量，所谓协变量就是指一些与应变量、自变量可能都有关系的连续性变量，它们的存在可能会影响分析结果的正确性，从而不得不在分析中加以控制。这种控制了协变量的分析被称为协方差分析，后面有专门的一节介绍。

5. **WLS Weight** 框：即用于选入加权最小二乘法的权重系数。

此处进行的加权最小二乘法分析实际上和 Liner 过程主对话框中 WLS 框提供的功能相同，只不过在方差分析中该方法用得比较少，因此在这里我们不