



新编21世纪心理学系列教材

高级心理统计

Advanced Statistics for Psychology

刘红云 编著



新编21世纪心理学系列教材

高级心理统计

Advanced Statistics for Psychology

刘红云 编著

中国人民大学出版社
·北京·

图书在版编目 (CIP) 数据

高级心理统计/刘红云编著. —北京: 中国人民大学出版社, 2019. 3
新编 21 世纪心理学系列教材
ISBN 978-7-300-26659-6

I. ①高… II. ①刘… III. ①心理统计-高等学校-教材 IV. ①B841. 2

中国版本图书馆 CIP 数据核字 (2019) 第 016387 号

新编 21 世纪心理学系列教材

高级心理统计

刘红云 编著

Gaoji Xinli Tongji

出版发行	中国人民大学出版社	邮政编码	100080
社 址	北京中关村大街 31 号		
电 话	010-62511242 (总编室)	010-62511770 (质管部)	
	010-82501766 (邮购部)	010-62514148 (门市部)	
	010-62515195 (发行公司)	010-62515275 (盗版举报)	
网 址	http://www.crup.com.cn		
	http://www.ttrnet.com (人大教研网)		
经 销	新华书店		
印 刷	北京昌联印刷有限公司		
规 格	185 mm×260 mm 16 开本	版 次	2019 年 3 月第 1 版
印 张	27	印 次	2019 年 3 月第 1 次印刷
字 数	652 000	定 价	59.80 元

版权所有 侵权必究 印装差错 负责调换

作者简介

刘红云，北京师范大学心理学部教授，中国教育学会教育统计与测量学会秘书长、常务理事。研究领域主要有心理测量与评价、统计分析方法理论及应用，尤其关注统计与测量方法新进展及应用方面的研究。目前主要从事的应用研究有考试测评理论的发展与应用、多元统计分析技术在心理学研究中的应用。出版的主要专著、教材和译著有《追踪数据分析方法及其应用》《应用心理统计学》《行为科学统计精要》等。

内容简介

《高级心理统计》是一部系统介绍多元统计分析方法及其操作软件（SPSS、Mplus 和 HLM 等）在心理学研究中应用的著作。本书不仅包括多元回归分析、多元方差分析、判别分析和聚类分析等常用的多元统计分析方法，而且对结构方程模型（因素分析、路径分析）的基本原理、方法和应用，尤其是对研究者比较关注的一些专题如中介效应和调节效应进行了较为详细的介绍。另外，本书还包含了多层线性模型基础部分的内容。本书是一部由中级至中高程度的应用统计方法教材，可作为心理学及相关专业高年级本科生和研究生的教科书及应用工作者的参考用书。



前 言

量化研究方法的发展在一定程度上推动了社会科学实证研究的深入。近年来,随着统计软件的开发和应用,多元统计分析技术在心理学实证研究中被普遍应用,有些已经到了普及的程度。许多专业易用的操作软件的开发大大促进了这些方法在各个学科分支的应用和普及。然而,对于社会科学领域的学生而言,应用统计学的学习绝非如掌握统计软件的操作那么容易。学好多元统计分析的方法,正确选择和使用各种分析方法,呈现和报告数据分析的结果并非易事。如何从应用的角度,让读者能够通过本书的学习和操练实际掌握和应用这些方法,是我在写作本书过程中一直思考的问题。

本书从应用的角度介绍各种方法,其中不涉及复杂深奥的数学证明和估计算法的原理,具有本科统计基础的读者都可以通过学习本书的内容,练习案例数据,掌握每种方法的用途。本书应用案例^①均结合心理学的实际研究问题,采用常用的统计分析软件 SPSS 或 Mplus 演示其具体操作,并结合心理学文献中结果呈现与解释的一般规范,读者很容易模仿每种方法,并将其用于自己的研究,通过不断练习理解掌握每种方法的实际意义。以心理学研究的实际问题为切入点,采用实际研究的案例数据是本书的特色。

本书在数据清理与准备这一章的基础上,集中介绍了 12 种心理学研究中常用的统计方法及其应用,主要包括:多元方差分析(第二章)、多元回归分析(第三章)、Logistic 回归分析(第四章)、判别分析(第五章)、聚类分析(第六章)、探索性因素分析(第七章)、验证性因素分析(第八章)、路径分析(第九章)、结构方程模型(第十章)、中介效应(第十一章)、调节效应(第十二章)、多层线性模型(第十三章)及其在追踪研究中的应用(第十四章)。其中,第七章到第十二章涉及结构方程模型的主要内容。

一本教材的编写凝聚了許多人多年的努力。首先感谢选修过高级心理统计课程的所有学生,你们看似幼稚的错误和有时候有些刁钻的问题给我很多启发,让我不断调整课程的内容和进度,不断思考怎样从应用的角度理解每种方法的实质。正是你们的问题和所犯的错误帮助我不断理解每种方法的重点和难点,也促使我自己在专业上不断进步。其次感谢我所有的研究生和博士生,多年来你们不辞辛苦地给我的课程当助教、答疑、批改作业、整理教学案例,每一部分教学内容变化都凝聚着你们的努力和付出。本书书稿基于我十几年教学的讲义整理撰写而成,这里尤其要感谢帮我整理书稿的刘玥、肖悦、魏丹、李美娟、韩雨婷、张茂鑫、王璞珏和马洁同学。最后感谢中国人民大学出版社的张宏学编辑,

^① 读者可登录人大出版社官网 www.crup.com.cn, 搜索本书书名进行下载。

没有你的鼓励和督促这本书不可能完成。

本书内容对从事社会科学研究的教學人員、科研人員、研究生、高年級本科生開展實證研究將會有很大幫助，可以作為手頭必備參考書目，也可以作為高等院校研究生的教材。

由於作者水平有限和成稿過程倉促，書中錯誤和表達不妥之處在所難免，請不吝賜教。

寫給教師的話

近十幾年來，我一直在教學的第一線擔任本科生和研究生統計課程的教學，對學生在學習統計學的過程中遇到的問題和困惑深有了解。每年的教學資料都會有些許修改和變動。細細想來，無外乎以下幾個方面的變化：一是心理學研究問題不斷深入，新的統計技術和分析方法不斷出現，為了幫助學生更好地掌握新的方法，課程中不斷增加一些新的方法和技术；二是根據在教學過程中每每面對的學生的疑問和不解，甚至是學生考試、作業和論文中的誤用而做出的難點和關鍵點的調整；再者就是隨著統計軟體的不斷發展，從易用易學的角度所做的操作程序上的調整；最重要的調整可能是教學過程中應用案例和數據的選取，逐漸用實際研究的案例替代了模擬的數據。這些內容上的變化也帶來一些教學模式的變化，真正讓學生深入思考、動手不斷操練，在做中學是學習應用統計非常關鍵的一點。因此，如果您選用本書作為高年級本科生或研究生多元統計課程的教材或參考書目，我想和您分享以下幾點：

1. 注重每種方法解決的實際問題和應用條件，包括：為什麼要使用這種統計方法，或者對於心理學研究的實際問題，我們有哪些方法可以選擇？每種方法對數據有怎樣的請求？需要滿足怎樣的假設條件？這些問題是正確選擇和使用每種方法的基礎。本書幾乎每章都列出了方法使用的目的、解決的實際問題和基本假設，這部分的内容有助於學生了解每種方法所能回答的實際問題，從而選擇合適的方法應用到自己的研究中。

2. 重視結果的呈現和解釋，包括：如何在研究中規範地報告統計分析的結果？如何選用正確的圖表呈現數據？如何對統計分析結果進行合理恰當的解釋？這些都是應用統計的核心環節，也是最重要的部分。建議在講授這部分内容時能夠推薦一些採用某種方法進行實際研究的文獻給學生參考，再結合案例數據的分析，讓學生真正掌握每種方法，在研究中恰當使用和正確解釋分析的結果。

3. 通過實際案例數據，鼓勵學生多動手，勤思考。本書涉及的应用統計的軟體操作都不難，但是如何結合實際案例數據，真正掌握每種方法，尤其是在出現錯誤時找出問題所在，一定程度上依賴於不斷練習和思考過程中積累起來的對數據的感覺。建議這部分的教學除了利用本書所提供的案例數據，還應盡量讓每種方法能夠匹配一些實際的案例數據供學生練習和討論。

4. 永遠不要忽略學生學習過程中出現的問題和所犯的錯誤。講授這門課程十多年來，學生學習過程中的問題和所犯錯誤，不斷促使我反思如何更好地教學，如何站在學生角度理解每種方法，所謂教學相長說的就是這個道理。

5. 教学方式上建议采用讲授、练习和讨论相结合的方式。应用统计方法的掌握，离不开不断的练习和讨论。多年的讲授经验告诉我，讲授过的知识对学生最多就是听过了，练过做过的内容对学生才是学过了，讨论反思纠正错误后的知识才能说是掌握了。并且这一过程有时候还要反复多次，其中讲授环节的课程目前我给学生录制了线上资源，而后面两个环节是老师替代不了的，但是老师可以引导和辅导，一定要督促学生多动手练习。

帮助学生理解这些方法所能解决的问题、基本原理或思想、应用条件、软件操作和结果的解释，使学生能够将这些方法应用到自己的实际研究中，是学习应用统计的宗旨。每位老师都有独特的讲授风格，但是最终的目标往往是一致的，希望这本书和相应的线上资源对您的教学有所帮助。

写给学生的话

当你打开这本书的时候，我想你有必要先花几分钟时间浏览一下目录，了解一下这本书的侧重点，思考一下你学习这本书的目标是什么。这是一本从应用者的角度介绍各种统计方法的教材，在学习这本书的内容时，我有以下建议。

1. 学习一种方法之前，静下心来先思考一下为什么要学习和使用这种方法，这种方法能帮助你解决什么样的实际研究问题。为了帮助大家了解这些内容，书中会介绍使用每种方法的目的和所能回答的问题，并且都配有一个实际的案例帮助你了解这些内容。当然你也可以带着你的研究问题去思考对于你的研究为什么要采用这种方法而不是另外的方法。

2. 在使用每种方法之前，了解这种方法对数据的要求和基本假设。思考你想要分析的数据是否满足这些前提条件，如果不满足，可以有哪些选择和处理方法。这些条件往往是不同方法之间的区别，也是对于类似的研究问题，你选择不同方法的依据。因此，这些思考对于正确使用一种方法特别重要。

3. 数据分析是一个从简单到复杂循序渐进的过程，动手分析是你真正掌握一种方法的有效途径。数据拿来以后，不要着急采用一些复杂的分析方法，从最基本的数据清理和描述统计入手，一步一步慢慢做，从做的过程中去思考数据告诉你的结果的意义。同时，不要被复杂的原理和数理推导所吓倒，如果有一段原理你看不懂，实在无法理解，可以暂时忽略，尝试通过做的过程去理解。学习有许多途径，如果你的数理统计基础不够扎实，从做中去体会每种方法如何应用将会是一个不错的选择。

4. 实际案例数据的分析对于研究方法的掌握和迁移很有帮助。本书中提供了许多实际的案例，并提供了丰富的线上资源，这些操作的语句和结果都是已经写好供你参考学习的。因此，类似的研究问题，你完全可以通过简单的修改迁移到自己的研究中。大胆尝试，不要怕出错，尤其是软件操作的过程中，熟能生巧一点没错。

5. 结果的呈现和解释是应用统计应该关注的重点，这一环节对你真正学会应用一种方法特别重要。因此，建议在学习的过程中，能够参考一些应用这些方法的文献，与案例数据对照起来，思考得到的结果如何正确呈现和合理解释，从而真正让自己学会应用这一方法。



好的研究是智慧和耐心交互作用的结果。了解数据、分析数据和解释数据，可以说是一次你和你的数据交心交流的旅程。一定意义上讲，你从不同的角度越是了解你的数据，数据的分析结果和解释就越能为你的研究所用。希望通过本书的学习，你能有所收获，并对应用统计有新的理解。

刘红云

2018年10月

关联课程教材推荐

书号	书名	作者	定价(元)
978-7-300-24134-0	发展心理学(第3版)	雷雳	45.00
978-7-300-25588-0	变态心理学(第3版)	王建平等	59.80
978-7-300-25426-5	临床心理学(第2版)	姚树桥等	48.00
978-7-300-24309-2	实验心理学(第2版)	白学军等	45.00
978-7-300-24280-4	社会心理学(第3版)	乐国安等	52.00
978-7-300-25616-0	心理与教育科学研究方法	杨丽珠等	49.80
978-7-300-22490-9	行为科学统计精要(第8版)	弗雷德里克·J. 格雷维特	68.00
978-7-300-22245-5	心理统计学(第5版)	亚瑟·阿伦等	89.00
978-7-300-13307-2	伯克毕生发展心理学:从0岁到青少年(第4版)	劳拉·E. 伯克	79.80
978-7-300-18303-9	伯克毕生发展心理学:从青年到老年(第4版)	劳拉·E. 伯克	45.00
978-7-300-25883-6	人格心理学入门(第8版)	马修·H. 奥尔森 B. R. 赫根汉	98.00

配套教学资源支持

尊敬的老师:

衷心感谢您选择使用人大版教材!相关配套教学资源,请到人大社网站(<http://www.crup.com.cn>)下载,或是随时与我们联系,我们将向您免费提供。

欢迎您随时反馈教材使用过程中的疑问、修订建议并提供您个人制作的课件。您的课件一经入选,我们将有偿使用。让我们与教材共成长!

联系人信息:

地址:北京海淀区中关村大街31号206室 龚洪训收 邮编:100080

电子邮件: gonghx@crup.com.cn 电话:010-62515637 QQ:6130616

如有相关教材的选题计划,也欢迎您与我们联系,我们将竭诚为您服务!

选题联系人:张宏学 电子邮件: zhanghx@crup.com.cn 电话:010-62512127

人大社网站: <http://www.crup.com.cn>

心理学专业教师QQ群:259019599

欢迎您登录人大社网站浏览,了解图书信息,共享教学资源
期待您加入专业教师QQ群,开展学术讨论,交流教学心得






目 录

第一章 数据的清理与准备	1
第一节 数据清理和准备的主要目的	1
第二节 极端数据的处理	1
第三节 缺失数据的处理	5
第四节 多元分析前提假设条件的检验	11
第五节 数据清理与整理应用案例	17
第二章 多元方差分析	36
第一节 多元方差分析的一般目的和描述	36
第二节 多元方差分析主要回答的问题	37
第三节 多元方差分析的主要类型	38
第四节 多元方差分析的过程	40
第五节 多元方差分析应用案例及 SPSS 操作	52
第三章 多元回归分析	58
第一节 多元回归分析的一般目的和描述	58
第二节 多元回归分析主要回答的问题	59
第三节 多元回归分析的假设及模型	60
第四节 多元回归分析的类型	62
第五节 多元回归分析中自变量的重要性	65
第六节 多元回归分析中的统计检验	69
第七节 多元回归分析中一些值得注意的问题	71
第八节 回归分析的局限性	75
第九节 多元回归分析应用案例及 SPSS 操作	76
第四章 Logistic 回归分析	86
第一节 Logistic 回归分析的一般目的和描述	86
第二节 Logistic 回归分析主要回答的问题	87
第三节 Logistic 回归分析的前提假设与模型	87

第四节	Logistic 回归分析中一些值得注意的问题	91
第五节	Logistic 回归分析应用案例及 SPSS 操作	92
第五章	判别分析	98
第一节	判别分析的一般目的和描述	98
第二节	判别分析主要回答的问题	99
第三节	判别分析的假设条件及模型	100
第四节	判别分析的主要类型	101
第五节	判别分析的参数及解释	103
第六节	判别分析应用案例及 SPSS 操作	104
第六章	聚类分析	112
第一节	聚类分析的一般目的和描述	112
第二节	聚类分析主要回答的问题	113
第三节	聚类分析的模型及原理	113
第四节	聚类分析的主要类型	116
第五节	聚类分析中一些值得注意的问题	120
第六节	聚类分析应用案例及 SPSS 操作	121
第七章	探索性因素分析	130
第一节	探索性因素分析的一般目的和描述	130
第二节	因素分析的模型、假设及基本步骤	132
第三节	探索性因素分析前的准备	136
第四节	因素的抽取和旋转	138
第五节	探索性因素分析的应用	144
第六节	探索性因素分析中一些值得注意的问题	146
第七节	探索性因素分析应用案例及 SPSS 操作	149
第八节	探索性因素分析应用案例及 Mplus 操作	164
第八章	验证性因素分析	173
第一节	验证性因素分析与探索性因素分析的比较	173
第二节	验证性因素分析的图示、模型及基本步骤	174
第三节	验证性因素分析模型的确定和识别	176
第四节	验证性因素分析中的数据收集和参数估计	180
第五节	验证性因素分析中的模型评价与修正	183
第六节	验证性因素分析模型的应用	187
第七节	验证性因素分析应用案例及 Mplus 操作	190
第八节	等价性检验应用案例及 Mplus 操作	200

第九章 路径分析	210
第一节 路径分析的一般目的和描述	210
第二节 路径分析主要解决的问题	211
第三节 路径分析的模型和原理	214
第四节 路径分析中的模型分类和识别	217
第五节 路径分析中的效应分解及计算	223
第六节 路径分析模型的评价与修正	228
第七节 路径分析中一些值得注意的问题	232
第八节 路径分析应用案例及 Mplus 操作	233
第十章 结构方程模型	241
第一节 结构方程模型概述	241
第二节 结构方程模型中的图示、模型及假设	243
第三节 结构方程模型中的模型使用步骤	245
第四节 结构方程模型中模型比较的应用	248
第五节 结构方程模型分析中一些值得注意的问题	250
第六节 结构方程模型的局限性	253
第七节 结构方程模型估计应用案例及 Mplus 操作	253
第八节 结构方程模型多组比较应用案例及 Mplus 操作	259
第十一章 中介分析	273
第一节 中介分析的一般目的和描述	273
第二节 中介模型和中介效应	274
第三节 中介效应的检验	276
第四节 潜变量中介模型	282
第五节 中介分析中一些值得注意的问题	284
第六节 中介分析应用案例及 Mplus 操作	290
第十二章 调节效应	301
第一节 调节效应概述	301
第二节 调节效应的解释和重要性	302
第三节 显变量的调节效应	303
第四节 潜变量的调节效应	317
第五节 有调节的中介	336
第六节 有中介的调节	345
第十三章 多层线性模型简介	351
第一节 多层线性模型概述	351

第二节	多层线性分析主要回答的问题	356
第三节	多层线性分析中的模型及假设	358
第四节	多元线性分析中一些值得注意的问题	364
第五节	多层线性模型应用案例及操作	366
第十四章	多层线性模型在追踪研究中的应用	385
第一节	多层线性模型分析追踪研究数据的优势	385
第二节	追踪研究中的多层线性模型	386
第三节	一些值得注意的问题	388
第四节	应用案例一及操作	388
第五节	应用案例二及操作	407
参考文献	413



第一章

数据的清理与准备

核心要点

1. 了解异常值的概念和产生原因，掌握检验与处理异常值的方法。
2. 了解数据的缺失类型，掌握检查数据缺失程度、诊断数据缺失机制的方法，知道如何选择合适的插补方法对缺失值进行插补。
3. 熟悉多元分析的常见假设，并掌握这些假设的检验方法，知道违背多元分析假设时，该如何对变量进行校正。

第一节 数据清理和准备的主要目的

本章讨论的是在数据收集之后，进行主要的数据分析运行之前，需要解决的一系列问题。

在进行数据分析之前，花许多时间仔细检查数据是很常见的，对以下问题的考虑和处理是对数据进行诚实分析的基础：首先涉及数据录入的准确性以及异常值的处理，异常值——古怪（极端）的个案——可能会扭曲分析结果；其次，数值缺失是许多研究在数据采集中经常会出现的现象，在进行分析前必须加以评估和处理；最后，许多多元统计分析方法都是基于一定假设的，在应用这些分析方法前，必须确保数据集和前提假设之间的匹配，如果出现不匹配的情况，可以考虑变量的转换，以使它们符合分析的要求。

数据整理是每一个做量化统计的人都会遇到的问题，但也是实证工作中重要而常常不受重视的一步。在研究过程中，异常值、缺失值的出现，以及在不满足前提假设的条件下，错误地使用分析方法，都有可能造成分析结果的扭曲和错误。为保证分析结果的可靠性和准确性，必须对数据进行预处理，后面几节将分别介绍针对以上不同问题的数据清理方法。

第二节 极端数据的处理

一、异常值的定义

一般而言，异常值（outlier）可分为单变量异常值（univariate outlier）与多变量异常

值 (multivariate outlier) 两种。单变量异常值即在某个变量上明显高或者低的值。多变量异常值指在两个或多个变量上值的奇怪组合, 这使得该观测与其他观测明显不同。异常值可能影响观测结果, 也可能不影响。如图 1-1 所示, 上图的异常值会影响回归分析的观测结果, 而下图的异常值则不会产生影响。

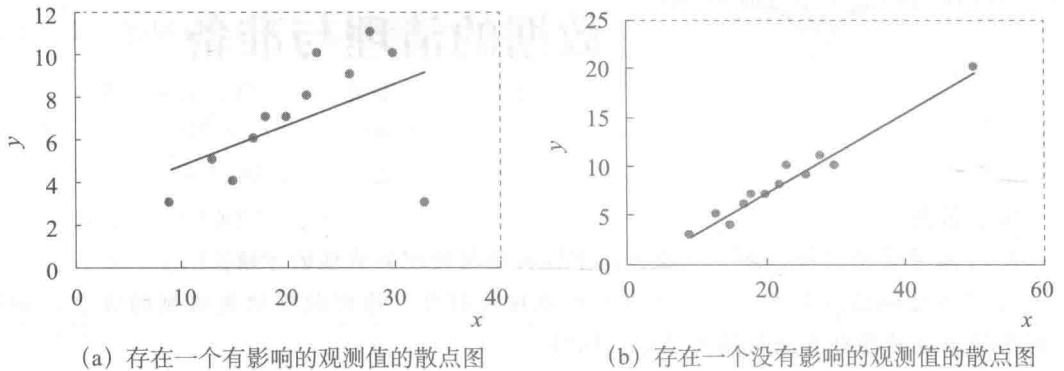


图 1-1 异常值对回归分析的影响

异常值的来源主要有以下几种:

(1) 过程性错误, 如录入、编码或缺失值定义错误。例如, 在 1 到 5 计分的李克特量表数据中出现了 6 或者其他超出计分范围的值。由这种原因造成的异常值应在数据清理阶段进行识别, 清除该值或编码为缺失。

(2) 异常事件。例如, 记录日常降水量时, 遇见台风, 使得记录值明显提升。这时应判断异常事件是否符合研究目的, 是则保留, 否则删除。

(3) 异常的观测。这是指研究者无法解释的观测值, 要考虑其是否代表了总体中的有效成分。

(4) 各变量值都正常, 但组合起来很异常。例如, 已知某人身高 165 厘米, 这高度本身不算特别高, 属于正常范围, 但如果得知该身高是测量自一位 5 岁的孩童, 则综合这两个资讯, 几乎可以肯定该身高在同龄者当中是一个与众不同的观察值组合。这种情况通常应视分析方法决定是否保留。

二、异常值的检测

下面分别介绍单变量、双变量和多变量异常值的检测方法。

(一) 单变量异常值的检测

可以通过计算标准分数或者画盒式图来检验单变量异常值。

1. 根据标准分数判定

首先将原始数据转换为标准分数, 这一操作可以通过 SPSS 的描述统计完成。在小样本 (样本量 80 及以下) 情况下, 标准分数大于等于 2.5 的观测值可被判定为异常值; 若是更大的样本, 可以提高标准分数的临界值, 最高是 4。如果不使用标准分数, 也可以根

据样本量，视 2.5 或 4 个标准差之外的观测值为异常值。

2. 盒式图

盒式图 (boxplot) 也称箱线图 (box-whisker plot), 可以直观地展示变量值的分布。可以使用 SPSS 的绘图功能 (Graphs→Legacy Dialogs→Boxplot...) 绘制盒式图。图 1-2 就是使用 SPSS 绘制的盒式图, 中间的灰色箱体为盒式图的主体, 箱体中的黑色粗线对应数据的中位数, 箱体上下边缘分别对应数据的上四分位数 (Q_3) 和下四分位数 (Q_1), 即有 75% 的数据在箱体上边缘以下, 有 25% 的数据在箱体下边缘以下。上四分位数和下四分位数之间的差值, 即四分位距 (interquartile range, IQR)。大于上四分位数加 1.5 倍四分位距 ($Q_3+1.5\times IQR$) 或小于下四分位数减去 1.5 倍四分位距 ($Q_1-1.5\times IQR$) 的值可被划分为异常值。其中, 处于 1.5~3 倍四分位距之间的异常值为温和的异常值 (mild outliers), 用空心点表示, 如图中 11 和 22 号数据点; 处于 3 倍四分位距之外的异常值为极端的异常值 (extreme outliers), 用星号表示, 如图中 18 号数据点。可在划分异常值的边界 $Q_1-1.5\times IQR$ 和 $Q_3+1.5\times IQR$ 以内, 最接近上下边界的两个值处画横线, 并将此作为箱线图的触须 (即图中大约 -4 和 10 的位置所画的横线)。

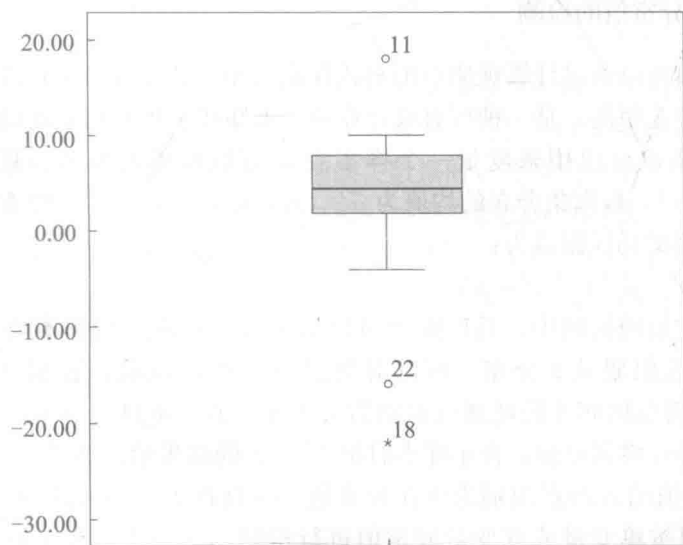


图 1-2 盒式图示例

(二) 双变量异常值的检测

双变量异常值一般通过画散点图的方式进行检测。散点图又称散点分布图, 是以一个变量为横坐标, 另一变量为纵坐标, 利用散点 (坐标点) 的分布形态反映变量统计关系的一种图形。散点图可以直观地反映: 两个变量之间是否存在相关趋势; 如果存在相关, 相关为线性还是曲线的; 是否存在双变量的异常值。

图 1-3 为依据变量 x_7 和变量 x_{10} 绘出的散点图, 其中被椭圆形圈起来的范围即期望范围, 指二元正态分布的置信区间 (α 水平通常为 90% 或 95%), 落在这个范围以外的观测值就是双变量异常值。使用 SPSS 软件可以轻易地画出双变量散点图, 但要绘制如图 1-3 中的置信椭圆则需要借助一些其他软件, 这将在后面应用案例部分详细介绍。

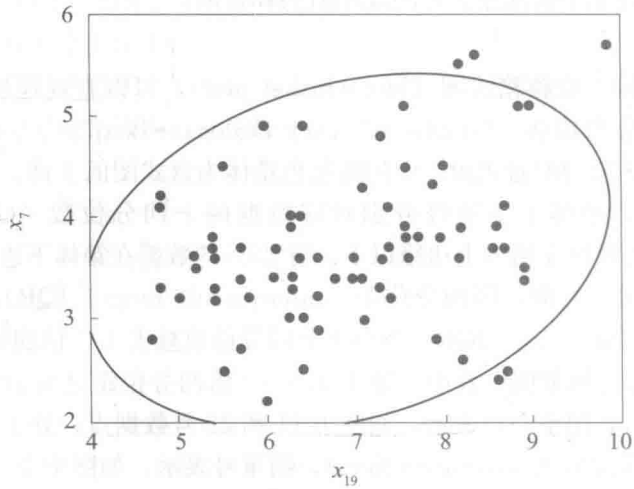


图 1-3 利用散点图检验双变量异常值

(三) 多变量异常值的检测

多变量异常值可以通过计算观测点的马氏距离 (Mahalanobis D^2) 进行检测。马氏距离代表数据的协方差距离, 是一种可有效计算两个未知样本集的相似度的方法, 具有尺度不变性。马氏距离也可以用来度量一个样本点 x 与数据集的距离。假设样本点为 $\vec{x} = (x_1, x_2, \dots, x_n)^T$, 数据集分布的均值为 $\vec{\mu} = (\mu_1, \mu_2, \dots, \mu_n)^T$, 协方差为 V , 则这个样本点 x 与数据集的马氏距离为:

$$D^2 = (\vec{x} - \vec{\mu})^T V^{-1} (\vec{x} - \vec{\mu}) \quad (1-1)$$

在多变量异常值的检测中, 马氏距离可以用来度量观测点与样本均值的多维空间距离。由于 D^2/df 近似服从 t 分布, 可以对其进行显著性检验, 置信度为 $\alpha = 0.005$ 或 0.001 , 可将落在置信区间外的观测点识别为潜在异常值。不过, 显著性检验容易受到样本量的影响。也可以根据经验, 在小样本时把 D^2/df 的临界值设为 2.5 , 大样本时设为 3 或 4 , 把超过临界值的观测点识别为潜在异常值。一旦在 D^2/df 指标上识别为潜在异常值, 该观测值就可按单变量或双变量异常值进行检测, 从而进一步了解异常的情况。在 SPSS 的回归分析中可以实现多变量马氏距离的计算, 详见后面应用案例部分。

三、异常值的处理

异常值处理的一般步骤是: 检测出数据异常后, 首先找到异常值出现的原因, 然后进行保留或者删除的处理。异常值的处理要考虑异常值是否能代表目标总体的一部分。如果确定该观测值异常, 且不属于目标总体, 则删除这个观测值; 如果可代表总体的一部分, 或不确定是否异常, 则要尽可能保留。可以通过转换变量或者改变计分, 降低异常值对分析结果的影响。