

应用数理统计

◎ 胡良平 胡纯严 鲍晓蕾 著



中国工信出版集团



电子工业出版社
PUBLISHING HOUSE OF ELECTRONICS INDUSTRY
<http://www.phei.com.cn>

应用数理统计

◎ 胡良平 胡纯严 鲍晓蕾 著



電子工業出版社

Publishing House of Electronics Industry

北京 · BEIJING

内 容 简 介

本书的内容和写作手法介于“概率论与数理统计”与“医学统计学”或“卫生统计学”之间。可以说，本书是学习“理论统计”与“应用统计”的一座坚实的桥梁，从“待分析的数据是否值得分析”入手，阐释了“应用数理统计”与前面提及的两大类泾渭分明的统计学的区别与联系。书中从试验设计、人为定义、概率分布和抽样分布四大方面介绍了统计计算的基本原理和来龙去脉；然后紧紧抓住最小平方法和最大似然法这两大类拥有多种衍生方法的算法准则，介绍了基于这些准则构造估计方程（即求解统计模型中未知参数的过渡方程）并导出参数估计的方法。为了便于读者学习、理解和正确应用，在必要的统计推导之后，还附有许多有价值的统计应用问题与解析。

本书适合文、理、农、工、商、医、经济、法律、交通、物流等多种学科领域工作的学者、学生、教师、研究人员和医护人员学习和使用。

未经许可，不得以任何方式复制或抄袭本书之部分或全部内容。

版权所有，侵权必究。

图书在版编目(CIP)数据

应用数理统计/胡良平, 胡纯严, 鲍晓蕾著. —北京: 电子工业出版社, 2015.6

(统计分析系列)

ISBN 978-7-121-26289-0

I. ①应… II. ①胡… ②胡… ③鲍… III. ①数理统计—高等学校—教材 IV. ①O212

中国版本图书馆 CIP 数据核字(2015)第 126265 号

策划编辑：秦淑灵

责任编辑：苏颖杰

印 刷：三河市鑫金马印装有限公司

装 订：三河市鑫金马印装有限公司

出版发行：电子工业出版社

北京市海淀区万寿路 173 信箱 邮编 100036

开 本：787×1092 1/16 印张：14 字数：352 千字

版 次：2015 年 6 月第 1 版

印 次：2015 年 6 月第 1 次印刷

印 数：3000 册 定价：35.00 元

凡所购买电子工业出版社图书有缺损问题，请向购买书店调换。若书店售缺，请与本社发行部联系，联系及邮购电话：(010)88254888。

质量投诉请发邮件至 zlts@phei.com.cn, 盗版侵权举报请发邮件至 dbqq@phei.com.cn。

服务热线：(010)88258888。

前　　言

大约在 5 年前，电子工业出版社的编辑就热情地邀请笔者编写一本“应用数理统计”。为什么编辑会提出这样的要求呢？理由很简单，因为她清楚地知道，“概率论与数理统计”对于广大实际工作者来说，实在是太难懂了，简直就是看“天书”！而“医学统计学”或“卫生统计学”又颇像“傻瓜照相机”，使用者即使不懂其原理，也能得出一些分析结果来。然而，对于几乎都具有高智商的广大科研人员和医护人员来说，“傻瓜照相机”经常使他们处于迷茫状态。他们不甘心，发誓一定要学好统计学、用对统计学。可是，每当他们拿起“概率论与数理统计”书籍，就被其中的许多怪模怪样的数学符号、晦涩的定理和冗长深奥的公式证明搞得不寒而栗、望而却步；而每当他们拿起“医学统计学”或“卫生统计学”书籍，又被其中的“指挥棒”驱使着，像“双眼被黑布紧紧地包扎着”东碰西撞，生搬硬套，很多场合下得出了错误的结果却浑然不知原因何在。

当时，笔者欣然接受了编辑的邀请。然而，写作大纲却一改再改，迟迟不能开始正式的写作。若有人询问编写大纲究竟修改了多少遍，说实话，就连笔者自己也记不清了。这是怎么回事呢？道理非常简单：编辑发出的邀请，代表了几乎所有（文、理、工、商、农、医等）学科领域内广大读者的共同呼声。对于数学家和统计学家们来说，虽然他们能写出逻辑严谨、理论深奥的“概率论与数理统计”和“高等数理统计”书籍来，但在解决实际的科研课题（尤其是还没有真实的数据）时，却显得无能为力、举步维艰。有时，他们中的一些人只要看到“数据”就眼睛“放光”，根本不考虑数据是否值得分析，就忙着开始进行“统计建模”、“参数估计”和“假设检验”。而对于非数学和统计学专业的广大实际工作者来说，当他们在稀里糊涂套用了一阵子统计分析方法后，真正渴望了解的是统计学的精髓、真谛和原理。上述多领域的读者都对本书寄予了厚望，这就是为什么本书花了整整 5 年时间才得以面世的根本原因。

本书的 10 章内容是分别为上述提及的各方人士特别定制的，尽管笔者已使尽了浑身解数，仍感到离广大读者的迫切要求和殷切期待相差甚远。笔者衷心希望广大读者能群策群力，多多提出宝贵的意见和建议，不断充实、修改和完善这部“桥梁式”的统计学著作，使之为人类的进步和发展做出更大的贡献。

本书的问世，得益于无数先辈和广大同仁所付出的智慧和辛勤劳动，还有许多直接和间接为本书做出重要贡献的人。在此，笔者一并向他们表示衷心的感谢！

由于笔者水平有限，书中难免会出现这样或那样的不妥，甚至错误之处，恳请广大读者不吝赐教，以便再版时修正。

胡良平
于北京军事医学科学院
生物医学统计学咨询中心
2015 年 5 月

目 录

第1章 待分析的数据是否值得分析	(1)
1.1 哪些情形下获得的数据是不值得分析的	(1)
1.1.1 人为编造的数据是不值得分析的	(1)
1.1.2 产生于质量控制不严的试验的数据是不值得分析的	(4)
1.1.3 经过错误的方法加工整理后的数据是不值得分析的	(4)
1.1.4 不符合特定统计分析方法要求的数据是不值得分析的	(5)
1.1.5 盲目解释基于误用统计分析方法所得到的分析结果是不可取的	(7)
1.1.6 缺失值过多的数据是不值得分析的	(8)
1.2 怎样保证数据是值得分析的	(8)
1.2.1 必须制定出科学完善的科研设计方案	(8)
1.2.2 必须严格控制课题实施过程中的质量	(11)
1.2.3 必须有实时记录科研数据的规格化表格	(12)
第2章 试验设计中的统计计算	(37)
2.1 试验设计原理与方法概述	(37)
2.1.1 试验设计四个核心内容概述	(37)
2.1.2 三要素几乎不涉及统计计算问题	(37)
2.1.3 四原则中有三个原则涉及统计计算问题	(37)
2.1.4 在构造设计矩阵时有三种情形涉及复杂的统计计算问题	(38)
2.1.5 在实施质量控制时有两种情形涉及统计计算问题	(38)
2.2 设计类型构建中的统计计算	(41)
2.2.1 设计矩阵及其优良性的概念	(41)
2.2.2 与构造某种准则下最优设计方案有关的基本概念	(43)
2.2.3 依据某些数学特性来确定各种最优设计矩阵	(44)
2.2.4 依据偏差函数来确定均匀设计的设计矩阵	(48)
第3章 基于人为定义的统计计算	(51)
3.1 常用名词概念	(51)
3.1.1 一般变量与随机变量	(51)
3.1.2 一般样本统计量	(52)
3.2 相对指标的定义与计算	(53)
3.2.1 相对指标的概述	(53)
3.2.2 相对比与百分比	(53)
3.2.3 频率与概率及率的标准误	(55)
3.2.4 危险度	(56)

3.3	平均指标的定义与计算	(57)
3.3.1	四种最常用的平均指标	(57)
3.3.2	有时不存在、有时又不唯一的平均指标——众数	(60)
3.3.3	两种能消除极端值影响的稳健的平均指标	(60)
3.3.4	组合平均值	(61)
3.4	变异指标的定义与计算	(64)
3.4.1	变异指标的种类	(64)
3.4.2	两分位数间距	(64)
3.4.3	其他几个常用的变异指标	(65)
3.4.4	自由度	(66)
3.4.5	度量离散度的三种稳健尺度	(67)
3.5	相关指标的定义与计算	(67)
3.5.1	定量变量之间的相关指标	(67)
3.5.2	定性变量之间的关联指标	(70)
3.6	常用统计量的某些特性	(72)
3.6.1	算术平均值具有使方差最小的特性	(72)
3.6.2	样本方差的定义式(3-23)是总体方差的无偏估计量	(72)
第4章	离散与连续型随机变量的概率分布	(75)
4.1	随机变量的概念	(75)
4.1.1	何为随机变量	(75)
4.1.2	随机变量的种类	(75)
4.1.3	随机变量的概率分布的概念	(75)
4.2	离散型随机变量的概率分布	(76)
4.2.1	一般离散型随机变量的概率分布	(76)
4.2.2	二项分布	(76)
4.2.3	Poisson 分布	(79)
4.2.4	负二项分布	(82)
4.2.5	几何分布	(84)
4.2.6	超几何分布	(87)
4.3	连续型随机变量的概率分布	(89)
4.3.1	一般连续型随机变量的概率分布	(89)
4.3.2	正态分布	(90)
4.3.3	t 分布	(92)
4.3.4	F 分布	(93)
4.3.5	χ^2 分布	(95)
4.3.6	对数正态分布	(96)
4.3.7	指数分布	(97)
4.3.8	威布尔分布	(99)
第5章	基于概率分布的统计计算	(100)
5.1	样本含量与检验效能的估计	(100)

5.1.1	与样本含量和检验效能有关的概念问题	(100)
5.1.2	与样本含量和检验效能有关的计算问题	(101)
5.2	基于二项分布的总体率的区间估计	(103)
5.2.1	二项分布定义	(103)
5.2.2	置信区间定义	(103)
5.2.3	基于二项分布的总体率的区间估计	(103)
5.2.4	如何用 SAS 实现基于二项分布的总体率的区间估计	(103)
5.3	基于 Poisson 分布的总体均值的区间估计	(104)
5.3.1	Poisson 分布定义	(104)
5.3.2	基于 Poisson 分布的总体均值的区间估计	(104)
5.3.3	如何用 SAS 实现基于 Poisson 分布的总体均值的区间估计	(105)
5.4	基于正态分布的多种区间估计	(106)
5.4.1	正态分布定义	(106)
5.4.2	基于正态分布估计一元定量资料的参考值范围	(106)
5.4.3	基于正态分布近似估计服从 Poisson 分布随机变量的总体均值的置信区间	(107)
5.4.4	基于正态分布近似估计服从二项分布随机变量的总体率的置信区间	(107)
5.4.5	基于正态分布求总体相关系数 ρ 的置信区间	(107)
5.5	基于 t 分布的多种区间估计	(107)
5.5.1	t 分布定义	(107)
5.5.2	基于 t 分布估计单组设计一元定量资料总体均值 μ 的置信区间	(107)
5.5.3	基于 t 分布估计成组设计一元定量资料两总体均值之差 $(\mu_1 - \mu_2)$ 的置信区间	(108)
5.5.4	基于 t 分布估计单组设计一元定量资料的预测区间	(108)
5.5.5	基于 t 分布估计直线回归方程中总体截距与总体斜率的置信区间	(108)
5.5.6	基于 t 分布估计直线回归方程中与自变量 x 取特定值条件下 y 的多种区间	(109)
5.6	基于 χ^2 分布的多种区间估计	(110)
5.6.1	χ^2 分布定义	(110)
5.6.2	总体方差与总体标准差的置信区间估计	(110)
5.7	基于 χ^2 分布和正态分布估计单组设计一元定量资料的容许区间	(110)
5.7.1	几个基本概念	(110)
5.7.2	单组设计一元定量资料容许区间估计	(111)
5.8	基于参数的假设检验导出置信区间计算公式	(112)
5.8.1	关于置信区间计算公式的说明	(112)
5.8.2	总体均值置信区间公式导出方法之一	(112)
5.8.3	总体均值置信区间公式导出方法之二	(112)
5.9	基于 SAS 估计单组设计一元定量资料的三种区间	(113)
5.9.1	问题与数据结构	(113)
5.9.2	对数据结构的分析	(114)
5.9.3	统计分析的需求分析	(114)
5.9.4	用 SAS 处理该单组设计一元定量资料尽可能给出较多结果	(114)
5.9.5	SAS 输出结果及其解释	(115)

第6章	基于抽样分布的检验统计量的导出及其应用	(118)
6.1	与抽样分布有关的预备知识	(118)
6.1.1	样本算术均值 \bar{x} 服从什么分布	(118)
6.1.2	样本方差 s_n^2 服从什么分布	(119)
6.1.3	样本方差 s_{n-1}^2 服从什么分布	(119)
6.2	基于正态分布的检验统计量 Z 的导出及其应用	(119)
6.2.1	样本取自正态分布的总体且 σ^2 已知时检验统计量 Z 的导出	(119)
6.2.2	服从标准正态分布的统计量 Z 的应用场合	(120)
6.3	基于 χ^2 分布的检验统计量 χ^2 的导出及其应用	(126)
6.3.1	基于 $R \times C$ 列联表资料独立性检验统计量 χ^2 的导出	(126)
6.3.2	服从 χ^2 分布的统计量 χ^2 的应用场合	(127)
6.4	基于 t 分布的检验统计量 t 的导出及其应用	(136)
6.4.1	基于一元定量资料均值假设检验的检验统计量 t 的导出	(136)
6.4.2	服从 t 分布的统计量 t 的应用场合	(137)
6.5	基于 F 分布的各种检验统计量的导出及其应用	(137)
6.5.1	基于 F 分布的检验统计量 F 的导出	(137)
6.5.2	服从 F 分布的统计量 F 的应用场合	(138)
第7章	基于最小平方法的统计模型中参数点估计公式的导出	(142)
7.1	普通最小平方法	(142)
7.1.1	普通最小平方法定义	(142)
7.1.2	普通最小平方法(OLS)的计算原理	(142)
7.2	加权最小平方法	(144)
7.2.1	加权最小平方法定义	(144)
7.2.2	加权最小平方法的计算原理	(144)
7.3	广义最小平方法	(144)
7.3.1	广义最小平方法定义	(144)
7.3.2	广义最小平方法的计算原理	(145)
7.4	基于普通最小平方法的改进	(145)
7.4.1	普通最小平方法需要改进的场合	(145)
7.4.2	降低自变量之间多重共线性影响的改进措施	(145)
7.4.3	降低异常点影响的改进措施	(146)
7.5	偏最小平方法	(153)
7.5.1	偏最小平方法定义	(153)
7.5.2	偏最小平方法的计算原理	(153)
第8章	加权最小平方法与偏最小平方法的应用	(155)
8.1	加权最小平方法的应用	(155)
8.1.1	以自变量平方的倒数为权重进行加权最小平方估计	(155)
8.1.2	以各试验点上重复试验次数的倒数为权重进行加权最小平方估计	(160)
8.1.3	以各试验点上因变量残差平方的倒数为权重进行加权最小平方估计	(167)

8.2	偏最小平方法的应用	(174)
8.2.1	问题与数据结构	(174)
8.2.2	用两种检验方法来决定抽取几对主成分变量	(174)
8.2.3	如何获得较多统计量的计算结果	(181)
第 9 章 基于最大似然法的统计模型中参数点估计公式的导出及其应用		(185)
9.1	最大似然法	(185)
9.1.1	用日常语言表述	(185)
9.1.2	用数学语言表述	(185)
9.2	其他最大似然法	(188)
9.3	最大似然法的应用举例	(188)
9.3.1	用于概率密度函数或概率函数中参数的点估计	(188)
9.3.2	用于某些多重回归模型中回归系数的点估计	(191)
第 10 章 统计分析的关键技术		(193)
10.1	从统计计算角度考量统计分析的关键技术	(193)
10.1.1	概述	(193)
10.1.2	第一类统计分析关键技术——发现新的概率分布规律	(193)
10.1.3	第二类统计分析关键技术——构建高维空间多层次多因素多指标复杂时间 序列模型	(194)
10.1.4	第三类统计分析关键技术——发现有广泛适应性且有可扩展性的回归系数 估计方法	(194)
10.1.5	第四类统计分析关键技术——求估计方程解的各种新算法	(194)
10.2	从统计应用角度考量统计分析的关键技术	(195)
10.2.1	概述	(195)
10.2.2	第一类统计分析关键技术——为统计分析方法进行合理分类	(195)
10.2.3	第二类统计分析关键技术——合理选择统计分析方法	(201)
附录 A 胡良平统计学专著及配套软件简介		(204)

当于增加一个参数，从而增加模型的复杂性。如果增加的参数对模型的预测能力没有显著的提升，那么这个参数就是冗余的。

第1章 待分析的数据是否值得分析

经典的数理统计是分析数据的一门学科。对于给定的一些数据，数理统计会专注于研究下述内容：如何描述其概率分布规律；如何估计总体中某些参数的数值和研究估计方法的性质；如何基于样本对特定的总体参数进行假设检验和研究假设检验方法的特性；如何揭示结果变量依赖多个原因变量（即影响因素或自变量）变化而变化的规律性；如何基于多个变量的取值情况反映多个个体或样品之间的亲疏关系；如何透过现象探索多个变量之间的本质联系；如何揭示多个变量与多个样品之间的内在关联性。然而，数理统计学家在开始实施前述的任何一类数据分析之前，都是有一个最基本的假定的，这就是“待分析的数据是值得分析的”！在实践中，人们在进行数据分析之前，无论分析者是否为数理统计学家，都应首先弄清这个“假定”是否成立。因为任何拟分析的数据不会因为处理数据的人是“数理统计学家”而变得“俯首帖耳”或“值得分析”。

1.1 哪些情形下获得的数据是不值得分析的

1.1.1 人为编造的数据是不值得分析的

当你试图运用统计分析方法处理数据时，首先要问的问题不是选择什么统计分析方法更合适，而是此数据是否值得分析。

【例 1-1】 对于表 1-1 所示资料，有人曾采用三种分析策略（参见后文中张飞给出的处理方法）得出了两个相互矛盾的结论，这一命题被称为“Symposon 悖论”。事实上，对于该命题，具有聪明头脑的人（简称张飞）与具有正常心态和思维的人（简称百姓）会给出截然不同的解读方法和答案。

表 1-1 同时按年龄和性别分层后吸烟与否与是否患肺癌的调查结果

年龄(岁)	吸烟与否	例 数					
		男性:	患	未患	女性:	患	未患
≤ 40	吸烟		5	5		40	50
	不吸烟		60	55		5	5
> 40	吸烟		30	10		5	55
	不吸烟		30	5		5	35

【分析与解答】 张飞这样处理：分别按三种策略分析此资料，策略一，只考察吸烟与不吸烟患肺癌的概率之间的差别；策略二，在按性别分层的基础上，再采用策略一分析；策略三，先按性别分层，后按年龄分层，再采用策略一分析。得到两种自相矛盾的结论：结论一，吸烟有利于健康（策略一与策略三）；结论二，吸烟有害于健康（策略二）。

百姓这样处理：首先提出一个问题，此资料是否值得分析？经过审读，发现它不值得分析，因为它是造假的产物！何以见得？此资料中有两个“疑点”，由此可推测其纯属人为编

造的“调查资料”。疑点之一，调查了≤40岁男性吸烟者10人，居然有一半人患了肺癌；调查了≤40岁女性不吸烟者10人，居然也有一半人患了肺癌。这样的调查结果只有在肺癌病房中才有可能获得，即肺癌病房中有5位患者，另5位是照顾病人的家属或护工。疑点之二，随机调查400位对象，再按“年龄”、“性别”、“吸烟与否”和“是否患肺癌”将其划分成16个小组，每个小组中人数的个位数不是0就是5(由基本常识可知，一个正整数的个位数有0、1、2、…、9共10种可能性)，出现这种奇特现象的概率是极低的，大约为 $P = \left(\frac{2}{10}\right)^{16} = (0.2)^{16} = 6.5536 \times 10^{-12}$ 。以如此低的概率出现的奇特现象，在实际调查资料中几乎是不可能出现的！

另外，即便表1-1中未出现上述两个疑点，而且，调查资料是在科学完善调查设计方案指导下获得的，张飞的三个分析策略都是不正确的，它们均属于“用单因素分析方法取代多因素分析方法”的错误。正确的分析方法是直接分析三个原因（“年龄”、“性别”、“吸烟与否”）对“是否患肺癌”这个“二值结果变量”的影响是否具有统计学意义。具体地说，可以选用对数线性模型分析，也可以选用多重logistic回归模型分析。若分析目的仅仅是希望考察吸烟与不吸烟患肺癌的概率之间的差别是否具有统计学意义，还可以选用CMH校正的 χ^2 检验，此法是把性别和年龄视为两个分层因素，对每一层（研究层内吸烟与不吸烟患肺癌概率之间的差别）做一些基本计算，然后再进行合并计算，最后得出汇总分析的结果和结论（即在消除分层因素影响的前提下，分析吸烟与不吸烟患肺癌概率之间的差别是否具有统计学意义）。

事实上，基于表1-1的调查设计，即使在样本含量、样本的代表性和资料的准确性等方面都做得很好，统计分析方法选用得也正确，其结论也不一定就很科学。因为导致一个人是否患肺癌，还有很多可能的影响因素未加以考察，如是否有肺癌的家族史、生活环境（如空气污染情况、饮水的质量、饮食构成、饮食习惯、生活习惯）、锻炼身体情况、职业、劳动强度、家庭经济状况、家庭生活是否幸福美满、人际关系是否良好、心态是否平稳、睡眠质量和每晚平均深睡眠时间长短等。这一切也都与一个人“是否患肺癌”有很大关系，它们被统称为“重要的非试验因素”。在任何一项调查或试验研究中，很关键的一个问题是，是否“找准找全”了对研究目的和主要评价指标有影响的“重要非试验因素”。不仅如此，还必须在尽可能大的样本含量且对总体有极好代表性的受试对象身上准确地获得试验因素、全部重要非试验因素和主要评价指标的数值。所有这一切，都与“数据是否值得分析”密切相关。

由此可知，统计学与纯数学（这里特指数理统计）是有本质区别的！纯数学所研究的数据是有很多隐含的假定的，在这些隐含假定（即所给定的数据是值得分析的）成立的条件下，直接研究数据之间内在的规律性。而统计学则不然，在使用者分析数据之前，必须对那些“隐含假定”逐一进行考证。例如，这些数据是否来自于研究目的相符合的总体（即数据是否具有同质性、样本对于总体的代表性如何）？它们是一个指标的不同取值还是多个指标的不同取值（即是一元资料还是多元资料）？它们是定量的还是定性的（即资料的性质是什么）？它们受到多少个因素的影响（即数据是受单因素影响还是受多因素影响的结果、是什么设计类型下收集的资料）？它们是怎样被收集到的（即是随意或人为选取的还是随机获得的）？所提供的数据对要揭示的问题是否足够多（即样本含量是否充足）？不同组数据之间是否具有可比性（即是否有合理的对照组、组间所受影响是否均衡）？在实施研究或进行试验或调查过程中的质量控制是否严格（即难以控制的众多非试验因素的影响是否被控制在尽可能低的水平上；至少应有

可靠的方法确保所有非试验因素在试验因素的不同水平组中处于基本平衡状态)? 只有弄清了上述诸多疑问之后, 当答案是“此资料值得分析”时, 才能选择具体的统计分析方法(首先是进行探索性统计分析, 以了解资料的基本情况, 然后再进行正式的统计分析, 以便达到合理选择统计分析方法, 正确实现研究目的)对资料予以处理。

然而, 上述看似烦琐的“循证”过程并非是正确应用统计学的全部“法宝”, 它仅仅适用于有了数据之后。在更多情况下, 是在仅仅有了研究课题和确定了研究目标之后, 就如何能够“多快好省而又科学严谨”地实现既定的研究目标提供“强有力的保证和具有可操作性的做法”, 才是统计学的精华之所在! 也就是说, 统计学的精华是如何制订出科学完善、严谨高效、经济可靠的设计方案。在此方案的指导下, 努力做好试验或调查过程中的“质量控制”, 才有望获得值得进行统计分析的研究数据。

应当清楚, 这里所说的“数据”是广义的, 正确的表述应当是“科研资料”。有时, 科研资料是纯文本的, 更多情况下, 科研资料包括原因变量、结果变量、变量的具体取值和变量的专业含义。也就是说, 变量、变量的取值(含单位)和变量的专业含义构成了科研资料的三要素。

怎样才能确保所制订出来的设计方案质量高、可操作性强呢? 这离不开正确的统计思想的指导, 离不开“透过现象看本质”的具体措施保驾护航、离不开“基本常识”的鼎力相助, 离不开“各科专业知识”这个坚如磐石的后盾。

【例 1-2】 下面有两组试验记录, 试审查此类资料是否值得分析。第一组试验记录见表 1-2, 第二组试验记录见表 1-3。

表 1-2 某项试验的试验记录

某药浓度(?)	?	?	?
0	72	97.042	6987.024
25	72	93.431	6727.032
50	72	98.056	7060.032
100	72	96.056	6916.032
200	72	130.694	9409.968
300	72	139.403	10037.020

注: 表中“?”处为空白, 读者不知道其下方的数据代表什么含义, 也不知道它们的单位分别是什么。

表 1-3 某项试验的试验记录

动物编号	性 别	体重(g)	X	Y	Z
1	1	230	3.73	300	16
4	1	225	3.73	600	24~36
5	1	231	3.73	510	16~24
6	0	211	3.73	540	16
12	0	220	3.73	570	24~36

注: X 代表气溶胶浓度(mg/L); Y 代表翻正反射消失持续时间(min); Z 代表恢复时间(h); Z 的取值为“16”的还有 7 只动物, 即共有 9 只动物都在 16h 恢复。

【分析与解答】 在表 1-2 中, 第 3 列上有 2 个数据的小数部分都是“056”, 而第 4 列上有 3 个数据的小数部分都是“032”, 这种巧合令人难以置信! 在表 1-3 中, 有 9 只动物的恢复时间都是 16h, 还有 3 只动物的恢复时间的跨度很大, 前者过于精确, 后者过于不准, 都说明结果不真实! 如此不真实的试验结果是不值得分析的!

1.1.2 产生于质量控制不严的试验的数据是不值得分析的

有些大型科研课题，需要划分成大约 10 个分课题，每个分课题下面可能还需要再分成 5 个子课题。这样一来，子课题的总数目有 30~50 个。其中，有相当多的观测指标是相同的或在专业上是有联系的，为了探讨这些指标之间的相互关系和依赖关系，人们常常需要把它们合并在一起进行统计分析。很显然，完成这些子课题研究的研究者们的责任心、技术水平等可能相差很多，若事先没有统一的操作规程(SOP)和严格的技术培训，所收集到的试验数据必然是参差不齐甚至杂乱无章的。

【例 1-3】 笔者曾参与评审过一种新药的临床试验项目，该项目由 5 个临床试验中心(通常是 5 所规模相当的医院)共同完成。其中，有一个主要安全性评价指标是某种不良事件发生率(%)，一所医院的试验结果为 22.3%，另一所医院的试验结果为 78.2%，其他医院的试验结果介于这两个数据之间。同一个重要的评价指标在不同医院之间相差如此之大，表明此项临床试验研究的质量控制很差，其整个临床试验数据的准确性非常值得怀疑。因此，质量控制做得不好的任何科研课题(最常见的是多中心临床试验研究课题)所产生的数据都是不值得分析的。

1.1.3 经过错误的方法加工整理后的数据是不值得分析的

【例 1-4】 有些医生常需要采用两种或多种方法对疑似患有某种疾病的患者或患者的标本同时进行检测，以便诊断患者是否确实患了预先推断的某种疾病。某医生收集了同时用 B 超和 CT 检查的 94 例某病患者的检测资料，整理成表 1-4 所示的形式。收集和分析此资料的目的是希望回答可否用 B 超取代 CT。基于表 1-4 的资料，这位医生选择什么统计分析方法可以实现其分析目的呢？

表 1-4 B 超及 CT 检查结果

检查方法	例 数			
	轻 度	中 度	重 度	正 常
B 超	18	3	3	70
CT	38	7	3	46

【分析与解答】 与这位医生的分析目的对应的统计分析方法属于“一致性检验”或称“Kappa 检验”。然而，该医生将两种仪器检测 94 名患者的结果按表 1-4 的形式整理后，就无法实施一致性检验了。因为表 1-4 属于“结果变量为多值有序变量的 2×4 列联表资料”，本质上是“单因素两水平设计一元定性资料”，就是人们习惯表述的“两个独立样本资料”。针对此资料，只能采用秩和检验，其分析结果只能回答“两种仪器检测结果之间的差异是否具有统计学意义”。事实上，实际收集资料是这样操作的：用两种仪器检测每位患者，每种仪器检测结果都可能出现 4 种结果(正常、轻度异常、中度异常、重度异常)之一，故对每位患者而言，两种仪器就有 16 种可能的组合结果之一出现，即应按“配对扩大形式”来列表，见表 1-5。

基于表 1-5(注意：表中的 f_{ij} 处必须是具体的频数)采用 Kappa 检验，就可回答两种仪器检测结果之间是否具有一致性。

表 1-5 B 超及 CT 检查结果比较(例数)(表 1-4 的正确表达形式)

B 超检查结果	例 数					合 计
	CT 检查结果:正常	轻 度	中 度	重 度		
正常	f_{11}	f_{12}	f_{13}	f_{14}		70
轻度	f_{21}	f_{22}	f_{23}	f_{24}		18
中度	f_{31}	f_{32}	f_{33}	f_{34}		3
重度	f_{41}	f_{42}	f_{43}	f_{44}		3
合计	46	38	7	3		94

注: 表中的 f_{ij} 为第 i 行第 j 列上的例数, 由于原始资料整理错误而无法获知。

1.1.4 不符合特定统计分析方法要求的数据是不值得分析的

【例 1-5】 表 1-6 是某市工业部门 13 个行业 8 项指标的数据, 拟采用主成分分析方法确定 8 项指标的样本主成分(综合变量), 若要求损失信息不超过 15%。现分别基于相关系数矩阵、协方差矩阵进行计算, 所用程序及结果如下。

表 1-6 某市工业部门 13 个行业 8 项指标的数据

编 号	行 业	X_1	X_2	X_3	X_4	X_5	X_6	X_7	X_8
1	冶金	90342	52455	101091	19272	82.00	16.100	197435	0.172
2	电力	4903	1973	2035	10313	34.20	7.100	592077	0.003
3	煤炭	6735	21139	3767	1780	36.10	8.200	726396	0.003
4	化学	49454	36241	81557	22504	98.10	25.900	348226	0.985
5	机械	139190	203505	215898	10609	93.20	12.600	139572	0.628
6	建材	12215	16219	10351	6382	62.50	8.700	145818	0.066
7	森工	2372	6572	8103	12329	184.40	22.200	20921	0.152
8	食品	11062	23078	54935	23804	370.40	41.000	65486	0.263
9	纺织	17111	23907	52108	21796	221.50	21.500	63806	0.276
10	缝纫	1206	3930	6126	15586	330.40	29.500	1840	0.437
11	皮革	2150	5704	6200	10870	184.20	12.000	8913	0.274
12	造纸	5251	6155	10383	16875	146.40	27.500	78796	0.151
13	文教艺术用品	14341	13203	19396	14691	94.60	17.800	6354	1.574

注: X_1 表示年末固定资产净值(万元); X_2 表示职工人数(人); X_3 表示工业总产值(万元); X_4 表示全员劳动生产率(元/人·年); X_5 表示百元固定原资产值实现产值(元); X_6 表示资产利税率; X_7 表示标准燃料消费量; X_8 表示能源利用效果。

以下是算法 1(基于相关系数矩阵)进行主成分分析的第一部分计算结果:

Eigenvalues of the Correlation Matrix

	Eigenvalue	Difference	Proportion	Cumulative
1	3.10491252	0.20747090	0.3881	0.3881
2	2.89744162	1.96722608	0.3622	0.7503
3	0.93021555	0.28809329	0.1163	0.8666
4	0.64212226	0.33803813	0.0803	0.9468
5	0.30408413	0.21748637	0.0380	0.9848
6	0.08659776	0.05441338	0.0108	0.9957
7	0.03218438	0.02974261	0.0040	0.9997
8	0.00244178		0.0003	1.0000

从这部分分析结果可以看出，取前 3 个主成分就可以解释全部原始数据 86.66% 的信息。以下是算法 2(基于协方差阵)进行主成分分析的第一部分计算结果：

Eigenvalues of the Covariance Matrix

	Eigenvalue	Difference	Proportion	Cumulative
1	5.44387E10	4.63867E10	0.8668	0.8668
2	8052070898	7827548141	0.1282	0.9950
3	224522758	136007774	0.0036	0.9986
4	88514983.7	86405779.1	0.0014	1.0000
5	2109204.6	2105915.97	0.0000	1.0000
6	3288.62949	3274.59278	0.0000	1.0000
7	14.0367172	13.9062234	0.0000	1.0000
8	0.13049382	.	0.0000	1.0000

从这部分分析结果可以看出，只需取 1 个主成分，就可以解释原始数据 86.68% 的信息。算法 2 的分析结果看上去似乎更为合理，更好理解。

【分析与解答】 一般而言，根据上述算法 2 中的这个第 1 主成分的值，可以对 13 个行业进行排序和分类。而且，由于这个第 1 主成分主要包含的是 X_7 的信息，也就是说，本资料可以根据 X_7 的信息对 13 个行业进行排序和分类。

表面上，基于算法 2 分析本资料后得到了一个令人兴奋的结果！其实不然，在上面的分析和解释过程中，已经犯了一个不可饶恕的错误！

显然，由算法 2 所得的结果，与算法 1 所得的结果很不一样。算法 1 是基于相关系数矩阵进行的主成分分析；算法 2 是基于方差协方差矩阵进行的主成分分析。按照统计学上的基本常识，这两种算法所得结果不应相差很大，而对本资料而言，其结果却相差十分悬殊。问题出在资料不符合主成分分析的要求上！

严格地说，主成分分析适合用在单组设计多元定量资料上。所谓单组设计，即受试对象或样品来自于同一个总体，他们(或它们)对于研究目的和全部评价指标而言应具有相同的性质。此时，可以考察不同定量指标之间内在的相互和依赖关系，甚至包括它们与隐变量(即主成分变量)之间的关系。而在本例中，13 个行业之间存在着本质区别， X_7 一个指标就可明显区分出哪些行业是以“燃料消费多”为标志的行业。用这一个指标就可对这 13 个行业进行排序，这个结果似乎令人感到十分欣慰，但这是十分片面的，甚至是相当错误的，因为它仅能反映哪些行业燃料消耗量大、哪些行业燃料消耗量小，而并不能反映其他任何方面！

所以，进行主成分分析时，务必要强调受试对象或样品应具有同质性，不要盲目地滥用主成分分析。若分析本资料的目的是为了基于给定的 8 项定量指标进行综合考虑，给 13 个行业进行排序，建议采用“综合评价”方法，如秩和比法。

综上所述，什么样的数据值得进行统计分析呢？数据必须满足以下四个前提条件：其一，数据是真实的，不是造假的，也不是随意添加或删除某些数据而形成的；其二，所取数据的试验(或调查)设计无严重错误；其三，数据在收集和整理过程中没有出现任何严重的失误或偏倚；其四，数据应满足拟选用的统计分析方法的要求(或前提条件)。

以上四点都是十分重要的，对事先确定的研究目的来说：

第一点关系到研究者或数据分析师者的科研道德和科研作风问题。研究者或数据分析师者无

论出于什么目的，造假及随意添加或删除数据都是绝对不允许的！

第二点关系到对总体的定义是否清楚，所抽取的样本对于总体的代表性是否很好，样本中个体的同质性是否好，样本是否足够大，所考察的影响因素是否“找准找全”，所观测的指标是否涵盖了与研究目的对应的全部主要指标，指标的测定时间是否合适和测定结果是否准确，课题设计的质量成为所收集的数据是否值得分析的第一道防线。

第三点关系到试验(或调查)过程中可能会受到来自环境、试验条件(仪器设备、试剂等)、研究者和受试者的心理等因素或条件的改变对观测结果造成的歪曲实际的影响，“质量控制”的严格程度成为所收集的数据是否值得分析的第二道防线。

第四点应当注意的是，应根据设计类型、比较类型(如新药临床试验中，还有非劣效性检验、等效性检验、优效性检验，通常的假设检验被称为差异性检验)和统计分析目的，选择最合适的统计分析方法，而不应当盲目套用。例如，有些人针对某病患者与正常人的多项定量指标的数值，进行变量聚类分析和样品聚类分析，这是很不合适的。若研究目的是考察全部定量指标的平均值(即均值向量)在两组之间的差别是否具有统计学意义，就应当选用成组设计定量资料多元方差分析处理资料；若研究目的是希望构造出函数式，以便对新个体究竟属于某病患者还是属于正常人进行判定，就适合选用判别分析。还有人面对多项二值变量及其取值时，采用主成分分析和探索性因子分析，完全忽视了这些多元统计分析方法要求拟分析的资料应当是多元定量资料的前提条件。考察拟分析的科研数据与拟选用的统计分析方法是否吻合成为所收集的数据是否值得分析的第三道防线。

1.1.5 盲目解释基于误用统计分析方法所得到的分析结果是不可取的

有些研究者习惯于脱离实际或无中生有地解释多元统计分析的计算结果。有些多元资料并不适合所选用的某些多元统计分析方法，得出的计算结果是解释不清的。然而，分析者却凭自己对多元统计分析方法的一知半解，盲目解释计算结果，不仅牵强附会，更严重的是对未来的研究工作或政策制订会产生歪曲的指导。

【例 1-6】 原文题目为《过敏性紫癜患儿血清白三烯 B4 白介素 -5 的测定及其临床意义》，测得患儿 31 例，正常儿 27 例的血清 IL-5、LTB₄、CRP 的结果如表 1-7 所示。

表 1-7 两组血清 IL-5、LTB₄、CRP 的检验结果($\bar{x} \pm s$)

组 别	例 数	IL-5 (pg/ml)	LTB ₄	CRP
正常组	27	12.7 ± 3.2	17.6 ± 5.7	4.75 ± 2.85
患儿组急性期	31	53.8 ± 4.2	95.3 ± 12.0	36.10 ± 11.78
恢复早期	31	37.8 ± 3.9	45.7 ± 10.1	18.35 ± 6.43

原作者对此资料做了多种统计分析，其中包括直线相关分析，给出了相关系数的计算结果，见表 1-8，该计算结果可信吗？

表 1-8 血清 IL-5、LTB₄、CRP 的相关系数 $r(n=89)$

	IL-5	LTB ₄	CRP
IL-5	—	0.772	0.715
LTB ₄	0.772	—	0.735
CRP	0.715	0.735	—

【分析与解答】 原作者给出的表 1-8 中的计算结果是不可信的！因为所做的线性相关分析是错误的。道理很简单，由表 1-7 可知，任何两个定量变量的 89 对数据实际上来自 58 名儿童，“89”是“例次”而不是“例数”，因为有 31 名儿童被重复测量了两次。即使都是不同的个体，把他们的数据放在一起进行线性相关分析也是错误的，因为他们是不同质的。正常儿童与处于不同时期的患病儿童，在与所研究疾病有关的指标上的取值大小是不同的，而进行线性相关与回归分析有两个隐含的前提条件：其一，数据所测目的受试者必须具有同质性；其二，所研究的变量或指标在专业上应有联系（例如，若研究人的身高与转氨酶之间的相关关系或回归关系就缺乏专业依据了）。

1.1.6 缺失值过多的数据是不值得分析的

有些课题在实施过程中出现了不可控的严重问题，导致所收集的科研数据严重缺失。例如，在新药临床试验研究中，可能由于未采用盲法分配受试对象或盲法执行得不够严格，导致有相当比例的受试者在试验过程中要求更换治疗方案或干脆脱离此项临床试验研究，这些受试者的试验数据就不完整，较多变量或较多时间点上的变量表现为缺失值。当具有缺失值的受试者数目的比例占总例数的 20% 以上时，应宣布此临床试验失败了。又例如，在一个调查研究课题中，若发出的问卷调查表的数目为 10 万份，而收回的调查问卷为 1 万份，其中各项内容填写有效的问卷为 3 千份。若仅依据这 3 千份有效调查问卷，即使采用各种正确的统计分析方法予以处理，得出的结论很可能是无用的，甚至有可能是有严重错误的。因为有效问卷的回收率仅为 3%，这些被调查者很可能在全部 10 万名被调查者中对调查目的或其中某项项目具有很强的倾向性。换句话说，有效问卷回收率低的调查数据是不值得进行统计分析的。

1.2 怎样保证数据是值得分析的

1.2.1 必须制订出科学完善的科研设计方案

由 1.1 节介绍的内容可知，数据是否值得分析的关键有两条，其一，考察科研设计是否正确，只有在无严重错误的科研设计方案指导下收集的数据才具备值得分析的首要前提；其二，在科研工作实施（试验或调查或从文献中提取信息）过程中，考察质量控制是否严格，是否在可能出现问题的所有环节上都事先想到了，重在防范一切可能影响结果准确可靠性的非试验因素的干扰和影响，一旦出现也有切实可行的对策使其影响降到最低程度。

什么是科研设计呢？在开始一项科学研究之前，需要提出研究目标，为了很好地实现事先制订的研究目标所做的一切考虑和安排，称为该研究项目或研究课题的科研设计方案，包括课题框架设计方案和课题技术设计方案两大类。对于任何一个科研课题而言，要想圆满地完成事先确定的研究目标，要考虑的问题或方面往往是很多的，稍有不慎，一旦落下了某些重要方面或在某些方面考虑不周，就可能导致整个课题研究前功尽弃。因此，对于任何一个科研课题来说，事先制订出科学完善的科研设计方案显得尤为重要。

科研设计的涵盖面非常广，其基本内容概括起来如图 1-1 和图 1-2 所示。图 1-1 是从结构上来划分的；而图 1-2 是从功能上来划分的，它实际上就是图 1-1 中“课题技术设计方案”的一种变形。