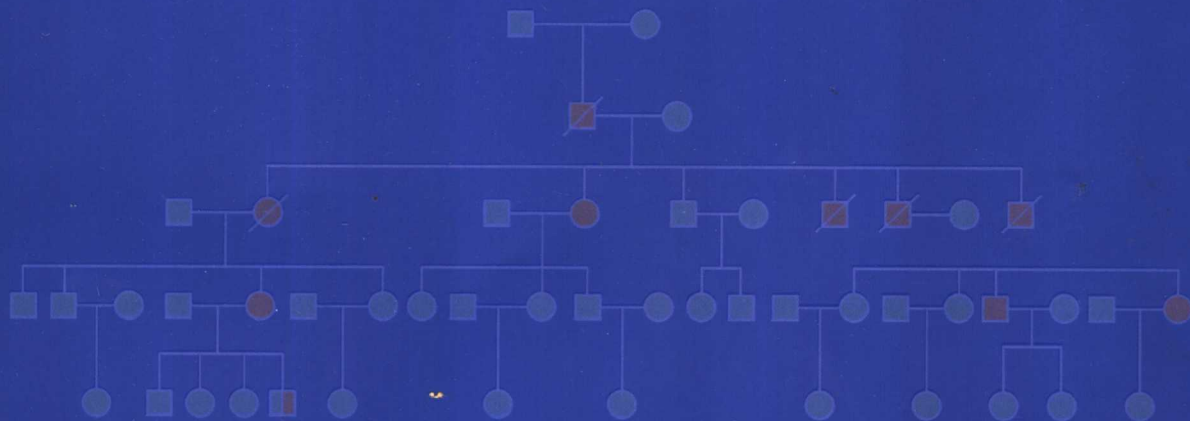


研究生教学用书

Statistical Methods
for Biomedical Research

生物医学研究的统计方法

主编 方积乾



高等教育出版社
Higher Education Press

研究生教学用书

生物医学研究的统计方法

Statistical Methods for Biomedical Research

主编 方积乾

副主编 胡良平 赵耐青 宇传华 张岩波
郝元涛 徐天和

编者(以姓氏笔画为序)

王 玖(滨州医学院)	张晋昕(中山大学)
毛宗福(武汉大学)	林爱华(中山大学)
方 亚(厦门大学)	罗艳侠(首都医科大学)
方积乾(中山大学)	周旭毓(中山大学)
刘言训(山东大学)	周诗国(军事医学科学院)
刘清海(中山大学)	赵耐青(复旦大学)
宇传华(华中科技大学)	郝元涛(中山大学)
祁爱琴(滨州医学院)	胡良平(军事医学科学院)
李 霞(哈尔滨医科大学)	施学忠(郑州大学)
李晓松(四川大学)	徐天和(滨州医学院)
李彩霞(中山大学)	凌 莉(中山大学)
余红梅(山西医科大学)	高 永(滨州医学院)
张岩波(山西医科大学)	郭秀花(首都医科大学)

秘书 余红梅 祁爱琴 吴少敏



高等教育出版社

内容简介

鉴于国内外生物医学论文普遍存在统计学缺陷的严峻局面,一批长期投身科研、热爱教学,战斗在第一线的医学统计学教授们合作编写了这本新型的教科书。依据国际学术界对生物医学论文的统计学要求精选内容,以实际问题的“原型”为中心组织统计学概念和方法的教学。全书分三篇26章:基础篇介绍统计学的思维逻辑与基本方法,应用篇进一步传授全面解决实际问题的本领,专题篇介绍生物医学研究若干热点领域常用的统计方法。每章在传授统计学知识之后,均设“结果报告”(中英文对照)、“案例辨析”、“电脑实验”、“常见疑问与小结”和“思考与练习”等5个节目。附录中有统计软件SPSS和Excel的简介。随书附送一片光盘,内有“电脑实验”的程序和输出、“案例辨析”以及“思考与练习”的参考答案。本书可以作为研究生、本科生教材,也可作为医生、护士、教师、编辑和管理者的自学用书。

图书在版编目(CIP)数据

生物医学研究的统计方法/方积乾主编. —北京:
高等教育出版社,2007.6

ISBN 978-7-04-020841-2

I. 生… II. 方… III. 生物医学工程-生物统计
IV. R318

中国版本图书馆CIP数据核字(2007)第039611号

策划编辑 安琪 责任编辑 孙葵葵 封面设计 张楠 责任绘图 朱静
版式设计 余杨 责任校对 杨雪莲 责任印制 朱学忠

出版发行 高等教育出版社
社 址 北京市西城区德外大街4号
邮政编码 100011
总 机 010-58581000
经 销 蓝色畅想图书发行有限公司
印 刷 北京佳信达艺术印刷有限公司

购书热线 010-58581118
免费咨询 800-810-0598
网 址 <http://www.hep.edu.cn>
<http://www.hep.com.cn>
网上订购 <http://www.landaco.com>
<http://www.landaco.com.cn>
畅想教育 <http://www.widedu.com>

开 本 889×1194 1/16
印 张 38.75
字 数 1 130 000

版 次 2007年6月第1版
印 次 2007年6月第1次印刷
定 价 69.00元

本书如有缺页、倒页、脱页等质量问题,请到所购图书销售部门联系调换。

版权所有 侵权必究

物料号 20841-00

序

医学统计学是当今医学各专业的必修课程之一,更是研究生不可不学的一门基础课。即使从事一线工作的医护人员,在其日常工作中,也少不了要借助统计学思维和知识阅读文献,总结经验,各类医疗卫生单位的管理干部和研究人员更是如此。

一般来说,医科类学生与理工科类学生的不同之处在于前者不习惯抽象思维,对复杂的数学公式、符号和数据缺乏兴趣,从而使统计学的“教”和“学”双方都有困难。但由于一些资深教授不断探索教学方法,积累经验,终于使以往学生感到枯燥乏味、计算繁杂的医学统计学变得生动活泼甚至是一种享受的课程。本教材的主编方积乾教授便是得到学生高度评价并多次撰文赞扬的一位教师典范。

由方教授组织多位富有教学经验的老师所编写的这本教材确实凸显了以实际问题为中心精选内容的特点,努力让读者能领会概念,掌握技能,善于表达;每一章更创新地融合了“结果报告”、“案例辨析”、“电脑实验”、“常见疑问与小结”和“思考与练习”。只有积累多年教学和解决实际问题的经验者才能做到这一步。我在阅读该教材时,感受到其文字清晰,深入浅出,通俗易懂,便于自学;计算机和统计软件的使用,使繁杂的计算问题迎刃而解,医学科研工作得以事半功倍。

这是我从事医学统计学专业教学多年以来少见的一本好教材。长江后浪推前浪,一代更比一代强,作为老一代教师,从这本书我看到了一幅喜人的前景。事物总是有创新—完善—再创新—再完善的过程,在使用此教材时,读者无疑会发现某些问题或欠缺,若能及时反馈给编者,使教材得到改进而日臻完善,必然会培育出具有中国特色的好教材。

胡孟璇

医学统计与流行病学教授

中山大学公共卫生学院

2006年11月

前 言

在中国,医学统计学或生物统计学被公认为医学和生命科学教育的重要课程。20世纪90年代以来,这类课程的教学改革进展神速,教学内容和方式发生了巨大的变化;以电脑实验展现统计理论已深入人心;传统的习题课逐渐被课堂讨论所取代。然而,继续推进医学统计学或生物统计学教学改革仍然十分必要,何以见得?让我们从国际上发生的事情说起。

Nature Medicine 2005年发表过一篇社论,题目为 *Statistically significant*,一开头就说:“去年 *Nature* 和 *Nature Medicine* 因为登载的某些文章统计分析欠佳而遭到公众批评。这些批评促使我们密切关注文章中的统计方法学”。这事起始于 *BMC Medical Research Methodology* 2004年5月发表的一篇文章,作者是西班牙 Girona 大学 Emili García-Berthou 和 Carles Alcaraz。他们查阅了2001年 *Nature* 登载的181篇研究论文,发现38%的文章至少有一处统计学错误;此后,国际性刊物出现了一系列报告,其中之一是 Robert Matthews 为 *The Financial Times* 写的文章,他分析了2000年 *Nature Medicine* 论文中的统计方法学,发现31%的作者错误地理解 *P* 值的含义,甚至有人以可笑的精度报告 *P* 值(例如, $P=0.002\ 387$)。为了弄清楚问题到底有多严重, *Nature Medicine* 请了两位哥伦比亚大学的专家对该杂志的文章进行独立的“统计学审计”,尤其要求评价2003年发表的以人为对象的21篇论文的统计学方法。这两位专家按照公认的统计学标准,运用一个清单评价这些文章,发现有的文章几乎没有定量分析,有的却使用了非常复杂的统计学和数学方法;大部分文章用了少量统计检验,但往往叙述不完整,很难评价其是否恰当。

不言而喻,统计学错误同样普遍存在于中国国内的生物医学论文中。许多人对国内具有权威性的医学杂志乃至硕士论文作过调查,发表过不少调查报告,一致认为,中国生物医学论文的统计学问题主要表现为概念错误、方法不当和表达缺陷。

面对这样的“国际、国内形势”,我们合作编写了这本新型的教科书,旨在切实帮助生物医学领域的研究生、科研人员以及期刊编辑人员学好统计学,领会概念,掌握技能,善于表达。

我们瞄准国际上生物医学科研设计和论文的要求精选内容,以实际问题为中心来组织统计概念和方法的教学,切实满足科研实践的需要。全书分基础篇、应用篇和专题篇。基础篇介绍统计学的思维逻辑与基本方法,如描述、比较、关联和回归等;应用篇在概括了研究设计的统计学原则与技术之后,进一步传授全面解决实际问题的本领,如复杂设计下资料的分析、分类变量、寿命变量的影响因素分析以及多个分类变量间关联性分析;专题篇则介绍当前生物医学研究的几个热点领域常用的统计方法,如遗传统计、生物信息和 meta 分析等。

本书凝聚了作者们长期从事教学研究的成果和为杂志审稿的心得,在教学方法上更具有如下特点:

除最后一章为“医学论文的统计学报告要求”专题外,每章在传授统计学知识之后,均设一节中英文对照的“结果报告”,为生物医学论文中统计学内容的撰写提供示范。

除“苦口婆心”地正面阐述统计学理论与方法外,每章均设一节“案例辨析”,将作者们审稿中常见的、似是而非的案例提供给读者,使读者在“辨析”中领悟正确的理论,在迎接挑战中成长。

每章都有一节“电脑实验”,用国际通用的统计软件 SPSS 实现所学的统计学计算,使读者从复杂公式和数值计算中解脱出来;用十分普及的软件 Excel 实现精心设计好的实验,生动地展现统计学概念。有条件的读者,不难亲自动手;若条件不具备,可由教师演示。

从正反两方面学了主要内容后,读者不免还有些许“疑惑”挥之不去,每章都设有“常见疑问与小结”,拾

遗补缺,梳理归纳。

每章最后的“思考与练习”以相当数量的选择题、思考题、计算题和综合题供读者实战演练、自我测试。

本书的附录中有统计软件 SPSS 和 Excel 的简介,读者借此足以入门;学习各章统计知识的同时,反复使用这两种软件,不知不觉中也就得心应手了。

随书附送一片光盘,内有各章“电脑实验”的程序和详尽的输出、“案例辨析”以及“思考与练习”的参考答案。读者要克制自己,轻易不看,必要时才去“偷看”一眼;待大彻大悟,方揭开谜底,看是否“英雄所见略同”。

本书的作者们长期投身科研,酷爱教学,常年战斗在第一线;为写好这本书,我们横向将 26 章分成 5 个模块,分别由专人牵头撰写;纵向按照正文主体、结果报告、案例辨析、电脑实验、常见疑问与小结以及思考与练习分成 6 条线,分别由专人统一协调;我们通过电子邮件群发稿件、互赐评语,彼此修改,经历了一次真正意义上的集体创作。令人难忘的是,余红梅帮助每一章修改了中英文对照的“结果报告”,字传华、高永为好几章无私奉献了自编的电脑实验 Excel 程序,祁爱琴带领其同事统一编排了全部书稿;更有许多同事以各种方式热情襄助,他们是曾芳芳、朱淑明、刘裕和蒋丽丽等。谨代表编委会全体同仁对上面提到的和未及列出的所有朋友一并鸣谢。

我们衷心希望本书有助于普及医学统计学或生物统计学知识,有助于提升生物医学领域研究生、科研人员以及期刊编辑人员的统计学素养,有助于改善生物医学研究论文的统计学水准。有人说,电影是“遗憾的艺术”,写书何尝不如此,自认周密周到,难免缺点缺陷,一旦问世,顿生“遗憾”;好在书可重印和再版,尚有挽回的机会。愿广大读者见错必纠,不吝指正,共同打造一本好的教科书。

方积乾

2006 年 11 月,广州

目 录

一、基础篇

1 绪论	3	4.4 结果报告	73
1.1 为什么要学习统计学	3	4.5 案例辨析	74
1.2 生物医学数据的来源与类型	4	4.6 电脑实验	74
1.3 常用的基本概念	7	4.7 常见疑问与小结	76
1.4 统计工作贯穿医学研究的全过程	9	思考与练习	78
1.5 结果报告	10	5 假设检验	80
1.6 案例辨析	11	5.1 假设检验的基本思想	80
1.7 电脑实验	11	5.2 假设检验的步骤	81
1.8 常见疑问与小结	13	5.3 单组样本资料的假设检验	82
思考与练习	14	5.4 假设检验的两类错误	85
2 统计描述	15	5.5 结果报告	88
2.1 定量资料的统计描述	15	5.6 案例辨析	89
2.2 定性资料的统计描述	22	5.7 电脑实验	89
2.3 常用统计图表	25	5.8 常见疑问与小结	91
2.4 结果报告	30	思考与练习	92
2.5 案例辨析	31	6 两样本定量资料的比较	94
2.6 电脑实验	32	6.1 两组独立样本的比较	94
2.7 常见疑问与小结	35	6.2 配对设计定量资料的比较	100
思考与练习	37	6.3 两组 Poisson 分布资料的比较	103
3 概率分布	39	6.4 结果报告	105
3.1 正态分布	39	6.5 案例辨析	107
3.2 二项分布	46	6.6 电脑实验	108
3.3 Poisson 分布	50	6.7 常见疑问与小结	111
3.4 结果报告	53	思考与练习	113
3.5 案例辨析	54	7 多组定量资料的比较	115
3.6 电脑实验	55	7.1 单因素方差分析	115
3.7 常见疑问与小结	59	7.2 多个样本均数的两两比较	121
思考与练习	60	7.3 Kruskal-Wallis 检验	123
4 参数估计	62	7.4 结果报告	126
4.1 抽样分布与标准误	62	7.5 案例辨析	127
4.2 Z 分布与 t 分布	66	7.6 电脑实验	128
4.3 总体参数的估计	68	7.7 常见疑问与小结	131
		思考与练习	132

8 定性资料的比较	134	10 简单线性回归分析	173
8.1 定性资料案例及比较原理	134	10.1 概述	173
8.2 两组二分类资料比较	137	10.2 简单线性回归模型	174
8.3 独立的多组二分类资料比较	142	10.3 结果报告	180
8.4 独立的多组多分类资料比较	143	10.4 案例辨析	181
8.5 结果报告	146	10.5 电脑实验	183
8.6 案例辨析	146	10.6 常见疑问与小结	188
8.7 电脑实验	147	思考与练习	190
8.8 常见疑问与小结	152	11 多重线性回归分析	192
思考与练习	153	11.1 概述	192
9 关联性分析	155	11.2 多重线性回归模型	193
9.1 概述	155	11.3 多重线性回归的应用	199
9.2 两个连续型随机变量间的相关分析	156	11.4 结果报告	204
9.3 两个分类变量间的关联分析	160	11.5 案例辨析	205
9.4 结果报告	164	11.6 电脑实验	207
9.5 案例辨析	164	11.7 常见疑问与小结	211
9.6 电脑实验	166	思考与练习	213
9.7 常见疑问与小结	170		
思考与练习	171		

二、应 用 篇

12 实验设计	217	14 调查设计	255
12.1 实验设计的概念	217	14.1 概述	255
12.2 实验设计的三要素	218	14.2 调查表的设计	257
12.3 实验设计的四原则	219	14.3 调查问卷的评价	261
12.4 实验设计类型	221	14.4 调查研究的步骤与资料收集方式	262
12.5 结果报告	227	14.5 调查研究分类	263
12.6 案例辨析	228	14.6 调查设计类型	264
12.7 电脑实验	229	14.7 调查研究的抽样方法	265
12.8 常见疑问与小结	233	14.8 调查实施中的质量控制	268
思考与练习	234	14.9 调查研究中的伦理问题	269
13 临床试验设计	237	14.10 结果报告	271
13.1 临床试验前的必要准备	237	14.11 案例辨析	271
13.2 药物 I 期临床试验	239	14.12 电脑实验	272
13.3 药物 II 期临床试验	239	14.13 常见疑问与小结	280
13.4 药物 III 期临床试验	248	思考与练习	281
13.5 结果报告	249	15 样本含量估计	283
13.6 案例辨析	249	15.1 基本概念	283
13.7 电脑实验	250	15.2 比较定量资料均值时样本含量的 估计	285
13.8 常见疑问与小结	252	15.3 比较定性资料样本频率时样本含量 估计	287
思考与练习	253		

15.4 简单线性相关分析时样本含量估计	288	思考与练习	344
15.5 抽样调查设计时样本含量估计	289		
15.6 比较定量资料样本均值时检验功效的 估计	292	18 Logistic 回归	347
15.7 比较定性资料样本频率时检验功效 估计	294	18.1 单自变量 logistic 回归	347
15.8 简单线性相关分析时检验功效估计	297	18.2 多自变量 logistic 回归	353
15.9 结果报告	297	18.3 条件 logistic 回归	360
15.10 案例辨析	298	18.4 结果报告	362
15.11 电脑实验	299	18.5 案例辨析	363
15.12 常见疑问与小结	301	18.6 电脑实验	365
思考与练习	301	18.7 常见疑问与小结	368
		思考与练习	370
16 随机区组设计和析因设计资料的分析	304	19 生存分析	372
16.1 随机区组设计资料的方差分析	304	19.1 概述	372
16.2 随机区组设计资料的多重比较	306	19.2 生存率估计	375
16.3 方差齐性检验	307	19.3 生存曲线比较	379
16.4 随机区组设计资料的秩和检验	308	19.4 Cox 比例风险回归模型	381
16.5 析因设计资料的方差分析	310	19.5 结果报告	386
16.6 结果报告	314	19.6 案例辨析	387
16.7 案例辨析	316	19.7 电脑实验	389
16.8 电脑实验	317	19.8 常见疑问与小结	393
16.9 常见疑问与小结	326	思考与练习	395
思考与练习	327	20 对数线性模型在高维列联表资料分析 中的应用	397
17 重复测量设计和交叉设计资料的 分析	330	20.1 概述	397
17.1 重复测量定量资料的分析	330	20.2 模型构建原理	398
17.2 交叉设计资料的分析	335	20.3 结果报告	404
17.3 结果报告	337	20.4 案例辨析	405
17.4 案例辨析	339	20.5 电脑实验	406
17.5 电脑实验	340	20.6 常见疑问与小结	410
17.6 常见疑问与小结	343	思考与练习	411
		三、专 题 篇	
21 多元统计方法简介	417	参考文献	437
21.1 聚类分析	417		
21.2 判别分析	420	22 时间序列分析	438
21.3 主成分分析	424	22.1 时间序列的分解	438
21.4 因子分析	428	22.2 指数平滑法	439
21.5 案例辨析	431	22.3 ARIMA 模型	442
21.6 电脑实验	431	22.4 时间序列的频域分析	447
21.7 常见疑问与小结	436	22.5 结果报告	450
思考与练习	437	22.6 案例辨析	451

22.7 电脑实验	452	25 Meta 分析	501
22.8 常见疑问与小结	454	25.1 Meta 分析的基本概念	501
思考与练习	455	25.2 Meta 分析的步骤与方法	503
参考文献	456	25.3 Meta 分析的偏倚及其控制	514
23 遗传数据基因定位的统计方法	458	25.4 Meta 分析的常用统计软件	515
23.1 基本概念	458	25.5 结果报告	519
23.2 连锁分析	461	25.6 案例辨析	520
23.3 关联分析	465	25.7 电脑实验	522
23.4 结果报告	472	25.8 常见疑问与小结	523
23.5 案例辨析	473	思考与练习	526
23.6 电脑实验	473	参考文献	527
23.7 常见疑问与小结	476	26 医学论文的统计学报告要求	528
思考与练习	477	26.1 规范医学论文统计学报告的目的和 作用	528
参考文献	477	26.2 医学研究的科学思维	528
24 基因表达谱分析的生物学信息学方法	479	26.3 医学论文统计学报告的一般要求	529
24.1 基因芯片简介	479	26.4 温哥华格式对统计学报告的要求	529
24.2 数据标准化过程中的统计学方法	480	26.5 统计学常用符号与术语	530
24.3 应用基因芯片数据进行模式分类	482	26.6 国外的 CONSORT 声明	531
24.4 特征基因挖掘新方法	486	26.7 国内 RCT 论文的统计学报告自查 清单	533
24.5 应用基因芯片数据进行聚类分析	487	26.8 其他类型论文的统计学报告要求简介	534
24.6 ArrayTools 软件应用	488	26.9 结果报告	534
24.7 结果报告	494	26.10 案例辨析	537
24.8 案例辨析	495	26.11 常见疑问与小结	538
24.9 电脑实验	495	思考与练习	539
24.10 常见疑问与小结	498	参考文献	540
思考与练习	499		
参考文献	499		
附录 A SPSS 统计软件入门			541
附录 B Excel 统计功能简介			554
附录 C 统计用表			559
表 C1 标准正态分布(Z -分布)密度曲线下的 面积($\Phi(z)$ 值)	559	表 C9 秩和检验用 T 界值表	579
表 C2 t 分布界值表	560	表 C10 Wilcoxon 符号秩检验统计量分位 数表	580
表 C3 百分率的置信区间	562	表 C11 相关系数(Pearson)检验界值表	581
表 C4 Poisson 分布 μ 的置信区间	569	表 C12 等级相关系数检验界值表	582
表 C5 χ^2 分布界值表	569	表 C13 单组样本(或配对比较)均数检验时 所需样本含量	583
表 C6 F 分布界值表	572	表 C14 单组样本率检验时所需样本含量 (单侧)	584
表 C7 秩和检验用 H 界值表	578		
表 C8 配对符号秩和检验用 T 界值表	578		

表 C15 单组样本率检验时所需样本含量 (双侧)	587	表 C19 多组样本均数检验时所需样本含量	593
表 C16 两组样本均数检验时所需样本 含量	590	表 C20 λ 值表(多组样本率检验时所需样本 含量估计用)	595
表 C17 两组样本率检验时所需样本含量 (单侧)	591	表 C21 估计单组或配对设计差值的总体平 均值时所需样本含量	596
表 C18 两组样本率检验时所需样本含量 (双侧)	592	表 C22 单组与配对设计总体概率区间估计 时所需样本含量	597
索引			598



一、基础篇

1 绪 论

统计学(statistics)是关于数据(data)的学问,是从数据中提取信息、知识的一门科学与艺术,包括研究设计、数据搜集、数据整理、数据分析和结果报告等步骤。

根据研究领域和研究对象的不同,统计学又可细分为数理统计学(mathematical statistics)、经济统计学(economic statistics)、生物统计学(biostatistics)、卫生统计学(health statistics)、医学统计学(medical statistics)等,本书的内容更接近于后三者。医学统计学侧重于介绍医学研究中的统计学原理与方法;卫生统计学与医学统计学基本相似,但更侧重于介绍社会、人群健康研究中的统计学原理与方法,如人口期望寿命、观察研究设计与统计分析等;生物统计学的研究范围更广一些,它是统计学原理与方法应用于生物学、医学的一门科学。由于方法学的通用性,目前人们更愿意采用生物统计学,而不是采用卫生统计学、医学统计学作为这门学科的名称。

1.1 为什么要学习统计学

随着科技的进步,计算机的普及,统计学显得越来越重要。为了总结经验、获得信息、发掘知识,实现科学管理与决策,几乎所有科学技术都需要统计学的帮助,生物医学研究更是如此。例如,为了检验某种新药是否对改善血液循环有帮助,涉及一系列的统计学知识:研究对象如何分组?什么对象作为药物试验组,什么对象作为试验对照组?不同组之间的疗效有无差异?为了节约成本,如何由较少对象的观察资料推断一般人群中的疗效?为了获得有用、可靠的信息,每组至少需要多少人?等等。

英国统计学家 Galton F(1822—1911)曾说过,当人类科学家在探索问题的丛林中遇到难以逾越的障碍时,唯有统计学工具可以为其开辟一条前进的道路。

统计学对于科学研究与社会管理具有相当重要的作用,具体体现在如下几个方面。

1.1.1 发现不确定现象背后隐藏的规律性

相同父母所生子女(即使是双胞胎)的身高、体重、性格等各不相同,相同老师同一教室里学习的同学考试成绩各有千秋。差异是自然界存在的普遍现象,具有可比性的对象之间的差异称为变异(variation)。变异使得观察研究或实验研究的结果具有不确定性,统计学正是发现不确定(变异)现象背后所隐藏规律的一门科学。例如,某研究者发现精神科护士与妇产科护士的出勤率各不相同,两科室出勤率真的不同吗?有两种可能:实际上两科室出勤率相同,观察得到的差异只是偶然现象;实际上两科室出勤率不同,精神科护士的出勤率确实低于妇产科护士的出勤率,观察得到的差异不是偶然的。那么哪一种情况可能性更大一些呢?借助于统计学推断可回答这一问题。

1.1.2 用统计学思维方式考虑有关研究中的问题

期刊杂志研究论文中的结果部分,常常是统计分析的结果,批判性地理解和对待这些结果十分重要。

1.1.2.1 “阳性”结果是否是虚假联系?

例如,一批感冒患者用某药治疗1周后,治愈率为90%,能否说该感冒药十分有效?仅从治愈率90%来看治疗效果较好,但感冒1周的自愈率很高,那么治愈率90%是因为感冒药的效果,还是自身免疫功能导致自愈的效果呢?要回答这一问题涉及建立统计学对照组问题。如果将足够的受试对象随机分成两组,一个组使用该感冒药,而另一组不使用,此时两组治愈率之间的差异存在统计学意义,方能说明这种感冒药有治疗效果,否则这种“阳性”结果就值得怀疑。

1.1.2.2 “阴性”结果是否真是阴性?

有人曾对发表在 *Lancet*、*JAMA*、*N Engl J Med* 等著名医学杂志上 71 篇阴性结果的论文作过分析,发现其中有 62 篇(93%)可能是由于样本含量不足造成的假阴性。回到感冒药的例子,假使感冒药的确有效,但由于存在变异,用药者效果有好有差,未用药者效果也有好有差,当接受观察的人数不多时,两总体之间的差异有可能显不出来,这时就不能认为用药和不用药效果一样。此时的统计学推断结论得出总体间无差异,主要是由于样本含量不足所引起。

1.2 生物医学数据的来源与类型

1.2.1 数据的来源

生物医学数据的几个常见来源如下:

(1) 常规保存记录 一般业务机构都有常规保存记录。例如,医院病案室长期保存有住院病人的病案首页数据,医院统计信息科保存有医疗设备利用数据,医院人力资源部门保存有职工流动情况数据。研究者可根据自己的研究兴趣,从这些数据中获取有关信息。例如,为了研究某病历年来的治疗效果,可利用住院病人病案首页数据库来分析该病的治愈率、并发症发生率及住院天数等。

获得常规保存数据相对较容易,而且相对真实可靠,因此如果能充分利用这些数据,将省时、省力、省经费。但由于此类数据不是专门为研究者设立的,并不一定能完全满足特定研究的需要。

(2) 实验记录 包括实验室记录和临床试验记录,它是生物医学研究的主要数据来源。例如,在药理实验中,将实验动物分配到不同剂量组中,观察动物的反应,然后计算出半数有效量或半数致死量。在新药临床试验中,要详细记录被观察患者的用药及病情变化,作为新药疗效评价的依据。

(3) 现场调查记录 当从常规保存记录中得不到所需数据时,可采用现场调查方法搜集数据。例如,由于有的糖尿病患者并不住院治疗,甚至有的患者尚未被发现,欲了解某地区糖尿病的患病情况,医院保存的住院病历不能满足研究需要,必须进行现场调查与观察。

(4) 其他数据 统计分析工作所需要的数据有时可取自外来资料,例如,可取自公开发表的有关报告、商业性数据库以及专题研究文献等。可供医药卫生研究参考的数据有每 10 年进行的人口普查数据、中国卫生统计年鉴、第三次国家卫生服务调查分析报告等。

1.2.2 实 例

某妇产科医生为了调查住院天数、分娩方式、妊娠结局是否与年龄、身高、体重、职业、文化程度有关,该

研究者利用常规保存的病案首页数据,在某医院搜集了 2004—2005 年共计 1 402 名妊娠分娩妇女的资料(全数据见光盘 data1-1.xls),按住院号排序后排在前面的 10 名妇女的有关数据见表 1-1。职业、文化程度、分娩方式、妊娠结局等变量的分类见表 1-2。

表 1-1 10 名妊娠分娩妇女的有关数据

住院号	年龄/y	身高/cm	体重/kg	职业	文化程度	住院天数/d	分娩方式	妊娠结局
20040001	25	162	76.0	其他	中学	9	顺产	其他
20040002	32	153	60.0	其他	小学	7	剖宫产	足月
20040003	28	158	64.0	其他	中学	10	顺产	足月
20040004	29	162	68.0	工人	大学	8	剖宫产	足月
20040005	27	158	68.0	农民	小学	6	顺产	其他
20040006	39	158	66.5	工人	中学	8	剖宫产	其他
20040007	23	162	68.0	其他	小学	11	剖宫产	其他
20040008	20	162	70.5	管理人员	大学	4	顺产	足月
20040009	27	160	71.5	其他	中学	3	顺产	其他
20040010	22	162	70.0	工人	大学	7	剖宫产	足月

表 1-2 分类变量的类别

变量	类别	变量	类别
职业	工人、农民、管理人员、知识分子、商业服务、其他	分娩方式	顺产、先兆早产、助产、剖宫产
文化程度	文盲、小学、中学、大学及以上	妊娠结局	足月、其他

原始数据一般按照表 1-1 的格式排列,即每一行代表一个研究个体(基本分析单位)的观测记录,每一列代表一个观测指标(变量);第一行为每一变量的变量名,第一列为每一研究个体的标志编号。这样陈列的数据可以认为是一种“标准格式”,采用 SAS、SPSS 等国际统计软件进行统计学分析,均无需进行数据格式变换。

1.2.3 变量的类型

只有认识了数据的特征,才能正确地选用统计学分析方法。因此,理解数据类型(types of data)对于学习统计学具有十分重要的意义。

数据由具有若干变量的观察个体所组成。所谓观察个体(observation, individual, unit, element)是指研究的基本单位,如病案首页数据库中的每个患者、动物实验中的每只动物等。所谓变量(variable)就是可以反映个体特征或属性的量,如实例中的住院号、年龄(y)、身高(cm)、体重(kg)、职业、文化程度、住院天数(d)、分娩方式、妊娠结局都是变量。不同个体结果可能有不同的取值才能称为变量,否则称为常量(constant)。例如,实例研究对象是妊娠分娩妇女,所以在此研究中“性别”不是变量而是常量。本书涉及的变量实为随机变量(random variable)。变量可分为定量变量(quantitative variable)和定性变量(qualitative variable)两大类。

(1) 定量变量 定量变量(quantitative variable)也称为数值变量(numerical variable),这是统计分析中最常见的变量。根据变量的可能取值之间有无“缝隙”(gap),常将定量变量分类为连续变量和离散变量。

可以在某一区间取任何值的变量就是连续变量(continuous variable),如年龄(y)、身高(cm)、体重

(kg)。数据之间存在“缝隙”的变量就是离散变量(discrete variable),如家庭人口数、脉搏跳动次数(次/min)等,离散变量只能取有限的几个值。

(2) 定性变量 定性变量(qualitative variable)也称为分类变量(categorical variable)。根据变量类别之间是否有顺序、等级、大小关系,常将定性变量划分为有序变量(ordinal variable)和名义变量(nominal variable)。

如果定性变量的类别之间呈现出顺序关系,则该变量就是有序变量,如文化程度(文盲、小学、中学、大学及以上),疾病严重程度(轻、中、重)等。如果定性变量的类别之间无顺序大小关系,类别只代表名称或标签含义,没有数量意义,则该变量就是名义变量,如实例中的职业、分娩方式、妊娠结局就是名义变量。

无论是有序变量还是名义变量均可根据类别数分为二项分类变量(binomial classification variable)(如性别分为男与女,考试成绩分为及格与不及格,患者随访结果分为生存与死亡)和多项分类变量(polyomial classification variable)(如职业分为工人、农民、管理人员、知识分子、商业服务、其他,血型分为 A、B、AB、O)。

定性变量只能是离散变量。

1.2.4 变量的转化与编码

(1) 定量变量转化为定性变量 根据研究的需要,有时可以将定量变量转化为定性变量。例如,一组 20~40 岁成年人的血压实际测量值为定量变量。如按舒张期分期,临床上将舒张压 <80 mmHg 定为正常血压,80~89 mmHg 定为正常高值,90~99 mmHg 定为 1 级高血压,100~109 mmHg 定为 2 级高血压, ≥ 110 mmHg 定为 3 级高血压,以此对血压实际测量值进行整理得到的结果就是有序分类变量。如果进一步按舒张压是否 ≥ 90 mmHg,将研究个体分类为正常与异常两组,则血压这一定量变量就转化为二项分类变量了。

但是,由定性变量无法再转化为原来的定量变量,因此在搜集数据阶段应尽可能搜集定量数据。定量数据所含信息比定性数据更加丰富。

(2) 定性变量的数字编码 为了对定性变量进行统计学分析,往往需要进行编码。二项分类变量可以采用 0、1 编码,如将性别男用 0 表示、女用 1 表示。因为用 0、1 分别指示了性别不同属性,所以这样获得的变量也叫指示变量(indicator variable)。通常情况下,以 1 表示研究关注的类别,以 0 表示不太关注的类别。如实例中,对“妊娠结局”变量的足月更加感兴趣,那么我们可以令足月=1、其他=0。

对于有序分类变量,可按由小到大顺序编码为 1、2、3...,也可根据实际情况给予相应的得分。例如,对于实例中的“文化程度”,可将文盲、小学、中学、大学及以上分别编码为 1、2、3、4,或按读书年数编码为 0、6、12、16。

对于名义分类变量,为了让计算机识别其分类,可以输入任何代码。每一个代码或数字只起名称或标志作用,无数值的含义。

在后面章节的多因素分析中,为了将名义分类变量代入模型,需要进行哑变量(dummy variable)编码。实例中的职业分类为工人、农民、管理人员、知识分子、商业服务、其他等 6 类,则可定义 5 个哑变量(比总的分类数 6 少 1 个),分别记为 J_1 、 J_2 、 J_3 、 J_4 、 J_5 。编码方法见表 1-3。如果某个体的职业为农民,则将 J_1 、 J_2 、 J_3 、 J_4 、 J_5 分别编码为 0、1、0、0、0;如果某个体的职业为知识分子,则将 J_1 、 J_2 、 J_3 、 J_4 、 J_5 分别编码为 0、0、0、1、0。这样, J_1 、 J_2 、 J_3 、 J_4 、 J_5 这 5 个哑变量分别代表以“其他”为参照的工人、农民、管理人员、知识分子、商业服务等职业。大多数统计软件(如 SPSS、SAS),只要说明变量属于分类变量,并告知类别数,都可以自动产生类似上述的哑变量。