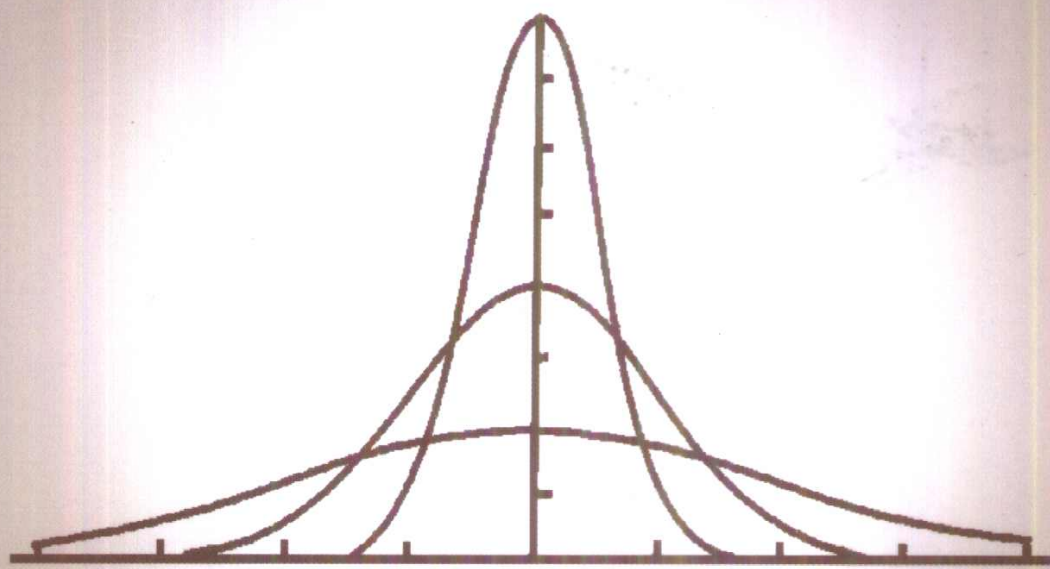




面向21世纪课程教材
Textbook Series for 21st Century

全国高等医药院校教材 供基础、预防、临床、口腔医学类专业用

生物统计学



主编 董时富



科学出版社

www.sciencep.com

40/

面向 21 世纪课程教材

Textbook Series for 21st Century

全国高等医药院校教材

(供基础、预防、临床、口腔医学类专业用)

生物统计学

主 编 董时富

副主编 方 亚 程锦泉 陈冬娥

科学出版社

2002

内 容 简 介

本教材是应新世纪形势的要求,配合医学院校教学内容和体制改革的进程而组织编写的。全书共18章,重点介绍生命科学研究中应用统计学原理进行研究设计、搜集数据、分析与解释研究结果的基本逻辑思维方式与方法。为了加强实用性,每章都提供相应的SAS分析程序。

本书可作为医学院校各专业以及其他学科如农学、林业学、管理学与心理学本科生的基础统计学教材,也可作为医学院校基础、临床各专业的教师或医务工作者的参考书。

图书在版编目(CIP)数据

生物统计学/董时富主编. —北京:科学出版社,2002.8

面向21世纪课程教材

ISBN 7-03-010691-1

I. 生… I. 董… III. 生物统计-高等学校-教材 IV. Q-332

中国版本图书馆CIP数据核字(2002)第054196号

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>

西保印刷厂 印刷

科学出版社发行 各地新华书店经销

*

2002年8月第一版 开本:850×1168 1/16

2002年8月第一次印刷 印张:20

印数:1—5 000 字数:488 000

定价:29.00元

(如有印装质量问题,我社负责调换〈新欣〉)

前 言

生物统计学的历史并不久远,但已成为当今最引人入胜的一门学科。21 世纪是生命科学的世纪,生物医学将翻开全新的一页,无疑也为生物统计学的发展与应用拓展了广阔的空间。生物统计学是以概率理论为基础研究生命科学中随机现象规律性的方法论科学。近年来国际社会对生物统计学专业人员的需求量较过去大为增加,临床医学研究者应用统计学方法评价新药或新疗法的效果;医学期刊杂志社聘请统计学专家为期刊论文把握统计处理的质量,《新英格兰医学杂志》、《柳叶刀》等著名的医学期刊还增设了负责处理统计学问题常务编委;各药业公司特聘统计学家参与新药研制研究设计与数据处理,如此等等,不胜枚举。生物统计学成为了当今的热门学科。

本教材是应新世纪形势的要求,依据教育部《关于积极推进“高等教育面向 21 世纪教学内容和课程体系改革计划”实施工作的若干意见》的精神,配合医学院校教学内容和体制改革的过程而组织编写的。本教材的特色在于:

1. 为了突出统计学应用的广泛性及学科名称的国际化,本教材采用正规学科名称——《生物统计学》。

2. 在内容编排上为了更切合生物科学研究的实际,将研究设计放于第一章,重点介绍基本的设计原则与实用方法;随后在第二、三章中介绍数据组织、表达等描述性统计方法;通过第四章给出一些最基本的概率理论,为后续章节的讨论提供基本的理论基础;第五章至第十五章介绍各种统计学推断方法;为了拓展视野,在第十六~十八章介绍了常用的多变量统计学分析方法。

3. 在写作上力求突出“三基”,强调基本概念、基本原理与基本方法,重视理论联系实际,培养学生的实际应用技能。在每一章的后面都附有该章内例题求解的 SAS 程序,让学生能通过这些程序掌握应用 SAS 系统解决自身专业实际问题的能力。

4. 本教材将简单相关与回归放于方差分析之前,破常规地增设了第十二章,使 t 检验、回归、方差分析与协方差分析等通常视为独立的参数过程,通过竞争模型假设与广义 F 检验的思维方式得到统一。

在本书编写过程中,受到了华中科技大学同济医学院教务部与公共卫生学院、武

汉大学、武汉科技大学与深圳市慢性病防治院各级领导的热情关怀和支持；受到周有尚教授、刘筱嫻教授、余松林教授及本校流行病学与卫生统计学系全体教师的关怀和帮助；付向华老师和研究生谭京广、蔡于茂、叶鹰、周海滨等同志为本书的文字处理、编排和出版付出了辛勤的劳动。在此，特向对本书给予关心和支持的各位致以衷心感谢。

由于我们的学识水平有限，加上时间紧迫，缺点和错误在所难免，恳请读者及同仁给予批评和指正。

董时富

2002年6月于武汉

目 录

生物统计学概述.....	1
第一章 研究设计中的基本统计学原则.....	7
第一节 研究设计的重要性.....	7
第二节 系统误差与控制方法.....	8
第三节 研究设计的基本类型.....	9
第四节 抽样总体与抽样方法	10
第五节 实验研究设计的基本要素	12
第六节 随机化的意义及方法	14
第七节 对照的设置与对照的均衡性	19
第八节 重复的作用与样本含量的影响因素	21
第九节 盲法及其作用	22
第二章 生物医学数据的组织与表达	26
第一节 数据与数据类型	26
第二节 频数分布表	28
第三节 统计图形表达	31
第三章 单变量综合性描述统计量	35
第一节 中心趋势度量	35
第二节 离散与变异性度量	41
第三节 率、比的均数与方差.....	45
第四章 随机变量、概率和概率分布.....	50
第一节 概率基本概念	50
第二节 随机变量及其概率分布	53
第三节 二项分布与泊松分布	55
第四节 正态分布	58
第五节 统计量的分布	62
第五章 统计学推断与单参数检验	66
第一节 样本均数与样本方差的抽样模拟	66
第二节 抽样误差与统计学推断	69
第三节 总体均数的置信估计	71
第四节 总体均数的假设检验	72
第五节 总体方差的置信估计与假设检验	75

第六节	显著性检验中的两类错误	78
第七节	参数置信区间估计与假设检验的关系	79
第六章	样本含量的估计与检验效能	82
第一节	概述	82
第二节	检验效能及其计算	83
第三节	样本含量的估计	86
第七章	总体分布的拟合优度检验	92
第一节	拟合优度检验的原理与计算步骤	92
第二节	离散型随机变量分布的拟合优度检验	93
第三节	连续型随机变量分布的拟合优度检验	97
第八章	两总体均数差异性检验	104
第一节	成组 t 检验	104
第二节	两方差间的差异性检验	107
第三节	t' 检验	108
第四节	配对 t 检验	109
第九章	一元线性相关与回归分析	114
第一节	相关与回归的概念	114
第二节	直线相关分析	115
第三节	简单线性回归分析	119
第十章	方差分析(一)	129
第一节	方差分析概述	129
第二节	单向方差分析	133
第三节	均数间的多重比较	136
第十一章	方差分析(二)	145
第一节	区组设计资料的方差分析	145
第二节	方差齐性检验	149
第三节	加权方差分析: Welch 检验	152
第四节	变量变换	154
第十二章	竞争模型假设与广义 F 检验	159
第一节	竞争模型假设检验的基本概念	159
第二节	应用于方差分析的广义 F 检验	162
第三节	应用于回归分析的广义 F 检验	164
第四节	应用于协方差分析的广义 F 检验	166
第十三章	名义分类频数表数据分析(χ^2 检验)	172
第一节	χ^2 检验的基本原理	172
第二节	χ^2 检验的基本步骤	173
第三节	两个样本率的比较	174
第四节	行 \times 列表资料的 χ^2 检验	177

第五节 行×列表 χ^2 分割分析	178
第十四章 有序列联表与配比设计方表的分析	183
第一节 有序列联表的基本分析方法	183
第二节 单向有序表数据分析	185
第三节 双向有序表数据分析	189
第四节 配比设计方形表数据分析	191
第十五章 非参数统计	199
第一节 概述	199
第二节 两独立样本检验	202
第三节 K 个独立样本检验	206
第四节 两个相关样本检验	210
第五节 K 个相关样本检验	213
第六节 等级相关与列联相关	216
第十六章 多元线性回归分析	227
第一节 多元统计分析方法概述	227
第二节 多元线性回归分析的基本原理	228
第三节 多元线性回归分析的数学模型	228
第四节 多元线性回归分析的方法步骤	229
第五节 多元线性回归分析的逐步回归法	234
第六节 多元相关分析	236
第七节 多元线性回归分析在医学中的应用	237
第十七章 Logistic 回归分析	242
第一节 Logistic 回归分析的数学模型	242
第二节 Logistic 回归模型的建立和检验	244
第三节 Logistic 回归模型系数的解释	246
第四节 配对病例-对照研究的条件 Logistic 回归分析	252
第五节 Probit 回归分析	254
第十八章 生存分析	257
第一节 生存分析中基本的概念	257
第二节 生存率的估计与生存曲线	259
第三节 生存曲线的对数秩检验	264
第四节 Cox 比例风险回归模型	266
附录 A 希腊字母表	276
附录 B 统计符号表	277
附录 C 统计检验用表	278
附表1 z 值表(标准正态分布曲线下的面积), $\Phi(-z)$ 值	278
附表2 t 界值表	279
附表3 χ^2 界值表	280

附表 4	F 值表(方差齐性检验用表)	281
附表 5.1	F 值表(方差分析检验用表)	282
附表 5.2	F 值表(方差分析检验用表)	283
附表 5.3	F 值表(方差分析检验用表)	284
附表 5.4	F 值表(方差分析检验用表)	285
附表 6	q 值表(SNK 法)	286
附表 7.1	q' 值表(Dunnett 检验)(单侧)	287
附表 7.2	q' 值表(Dunnett 检验)(双侧)	288
附表 8.1	Kolmogorov-Smirnov 拟合优度检验 D 界值表	289
附表 8.2	Kolmogorov-Smirnov 拟合优度检验 D 界值表	290
附表 8.3	Kolmogorov-Smirnov 拟合优度检验 D 界值表	291
附表 9.1	Shapiro-Wilk 的 α_{in} 系数表(I)	292
附表 9.2	Shapiro-Wilk 的 α_{in} 系数表(II)	292
附表 9.3	Shapiro-Wilk 的 α_{in} 系数表(III)	292
附表 9.4	Shapiro-Wilk 的 α_{in} 系数表(IV)	293
附表 9.5	Shapiro-Wilk 的 α_{in} 系数表(V)	293
附表 9.6	Shapiro-Wilk 拟合优度检验 W 界值表	294
附表 10	符号等级检验表(Wilcoxon 成对比较用)	294
附表 11	等级总和数临界值(双侧检验)	295
附表 12	M 值的界限值($P=0.05$)	296
附表 13.1	H 值与概率对照表	297
附表 13.2	H 值与概率对照表(续表)	298
附表 14	等级相关系数的统计学意义界限值	299
附表 15	两样本比较的 Kolmogorov-Smirnov 检验临界值表	299
附表 16	随机排列表($n=20$)	300
附录 D	英汉统计词汇对照	301
参考文献		305

生物统计学概述

一、什么是统计学

统计对于大多数人而言并不陌生,毋庸置疑的事实说明,“统计”已经进入现代人的生活。电视、广播、报刊杂志等新闻媒体每天都给人们传来很多统计信息,告诉人们有关“我国今年GDP增长率为7.8%”、“我市上周的平均气温为12.8℃”、“某一次事故中失踪和伤亡的人数”等等。它告诉我们有关某一事物或现象的事实和数值。如同其他词汇一样,不同的人对“统计”一词有不同的理解。有的人一听到这个词就想到人口出生率、死亡率、平均寿命、消费价格指数等数字、表格与表示数量的统计图形。这的确是对“统计”一词的极其生动而又恰当的表述,但这只是“统计”一词众所周知的第一种含义:即指任何用数字、表格与图形所表达的一个事实。“统计”一词的另外的含义是指统计学学科本身的含义,即统计学自身的术语、方法论及其知识总体。

概率论与数理统计(probability and statistics)简称统计学,是一门年轻而又引人入胜的学科。统计学是研究随机现象规律性的方法学。在自然界中存在着大量的随机现象,但是,“在表面上是偶然性在起作用的地方,这种偶然性始终是受内部的隐蔽着的规律支配的,而问题只是在于发现这种规律。”(《马克思恩格斯选集》第四卷,第243页,1972年)。偶然性事物的概率(即事物发生的可能性的)就是该偶然事件隐蔽着的特征。统计学就是研究这种内在特征性的一门数学方法论科学。它是为获得可靠性的推断结论,研究如何获取数据、组织与表达数据、分析数据与解释结果的科学与艺术。获取数据、组织与表达数据的过程通常称为描述统计(descriptive statistics);用观测数据做出有关总体某种结论性的决策过程称为推断统计(inferential statistics)。随着现代科学技术(特别是计算机科学技术)的迅速发展,极大地推动了统计学更广泛的应用与深化。它不仅形成了结构宏大的理论体系,而且在很多科学研究、工程技术、生物医学、社会经济与管理等领域得到了愈来愈多的应用。

二、什么是生物统计学

统计学就是把数学的语言引入具体的科学领域,把具体科学领域中要待研究的问题抽象为数学问题的过程。统计学被广泛地应用于解决自然科学和社会科学各个领域中的具体的随机现象的规律性,形成了应用于各个学科领域的统计学,即应用统计学。生物统计学(biostatistics)是以概率理论为基础研究生命科学中随机现象规律性的应用数学科学。由于生命科学特别是医学科学研究的对象的复杂性以及大量存在的各种随机因素对生物个体的复合作用,使得能用确定性

数学方法来处理的问题是极其有限的。以试验动物或人为研究对象是医学科学研究基本特征,生物个体的变异性特征决定了生物统计学在医学科学研究中的重要作用,而且这种作用已为广大生物医学科学研究者所共识。生物统计学是用概率论与数理统计的原理和方法研究生命科学领域中随机现象的数量规律的科学,包括科学研究设计、资料的搜集、整理、综合归纳、表达及分析,从而获得可信结论的应用性数学学科。

生物统计学的历史并不久远,生物学家达尔文(1809~1882)受到 Charles Lyell 的《地质学原理》(1833)一书的启发为他的进化论的产生产生了很大影响,达尔文的工作在本质上而言是属于生物统计学的;孟德尔于 1866 年所发表的有关豌豆杂交的研究也是生物统计学的问题; Karl Pearson(1857~1936)原先是一位数学物理学家,他将数学用于进化论,几乎费了半个世纪的时间从事数理统计研究,他还创始了《Biometrika》刊物及一所数理统计学校,于是生物统计学获得了推动力;K. Pearson 的学生 W. Gosset(1876~1937)关于 Student t 的著名文章于 1908 年发表在《Biometrika》上;R. A. Fisher(1890~1962)为生物统计学做出了大量而卓越的贡献,他和他的学生为数理统计方法在许多学科中的应用,尤其是在农业科学、生物学和遗传学中的应用起了很大的作用。

经历一个世纪,生物统计学获得了长足的发展,计算机与网络技术为生物统计学的应用开辟了更广阔的前景。

三、科学研究的基本程序

作为一项正规的生命科学研究,一般应包括如下 6 个过程。

1. 提出一个欲待研究的问题

“猫儿活到老,不改好奇心”,几千年的人类文明史就是人类在生产、生活中不断发现问题、解决问题,推动人类社会发展与进步的历程。可以发现,人类社会的每一次进步都伴随着一个重大“好奇”的发现过程,这种“好奇”的发现过程就是科学研究,其实质就是发现一个问题、解决一个问题的过程。作为生命科学的科学研究,就是通过回顾一些事实、总结专业实践过程中的一些经验、检索并学习相关的文献资料、结合一些公认的理论,产生一个假说,即发现一个专业问题。发现问题是任何科学研究的第一步,而且是最重要的一步,因为发现问题的性质决定了该项研究的科学价值。所谓“发现问题”就落实为人们常常所说的“选题”,它决定研究的创新性、科学价值,而且很大程度上由各自的专业内涵所主导。

2. 科学研究设计

建立一套能用科学试验或观察的过程,并能用该观察或试验结果回答所提出的假说是否成立的系统方法。研究设计包括研究的专业设计与统计学设计,前者着眼于从专业的角度对选题的创新性、科学价值、可行性等方面的考虑,后者是从统计学的原理出发对一项科学研究从获取数据到分析数据、解释结论的全过程所进行的考虑。所谓统计学设计(statistical design)是指用统计学原理对研究的全过程所作出的周密合理的统筹安排,如确定研究对象,拟定研究因素及其分配,如何执行随机、对照与重复的统计学原则,如何观察与度量效应,以及数据收集、整理与分析的方法,通过合理的、系统的安排,达到控制系统误差,以尽可能少的资源消耗(最小的人力、物力、财力和时间)获取准确可靠的信息资料及可信的结论,使效益最大化。统计设计作为医学科学

研究的重要组成部分,它是科学研究成功的基础,不容轻视。有人称:“一项完美的科学研究设计的完成,预示该项研究已经至少完成了百分之七十五。”由此可以看出研究设计的重要性。

3. 获取试验与观察的资料

该过程又称为搜集资料(collection of data),就是按照研究设计所拟定的方法与过程,通过对研究对象的观察及实验,测量并记录其结果,以形成研究的效应的原始统计数据。有很多获取原始统计学资料的渠道,即统计数据的信息源。根据信息来源可将数据分为三类:第一类为常规的工作记录。例如住院病人的病案资料、户籍与人口资料、医疗保险资料等。第二类为各种统计报表。如人口出生报告、死亡报告、居民的疾病、损伤、传染病的分月、季度与年报等资料。第三类为专题科学研究工作所获得的现场调查资料或实验研究资料。无论其资料的来源如何,研究者必须倍加关注原始资料的真实性、完整性、准确性,因为只有如此才能保证后续统计学分析结论的科学性与可靠性。

4. 数据审核与计算机录入

对于所获取的原始数据,可采用手工整理与分析。这种方法适合于较小数据含量,而且分析并不复杂的情况。当今正是计算机技术普及应用的年代,特别是对于较大的调查研究或多中心的临床随机对照试验,涉及到的变量多、记录的数据量大,分析过程复杂,一般都采用计算机程序化分析处理数据过程。无论是用手工还是用计算机过程,确保数据的准确、完整是后续分析的结论能反映客观世界真实规律的基本条件。首先应对原始数据进行仔细的核对,然后将其录入计算机,即建立一个数据文件。在录入数据时,为了避免数据在录入过程中所产生的差错,常常采用双机输入法,然后应用计算机对两个相同数据集进行逐条记录的每个变量值进行比对的方法,以减少录入误差;也可以应用计算机程序继续逻辑审核来发现原始数据中所存在的逻辑错误。应当注意的是,所谓计算机进行数据的逻辑核对过程完全是由研究者通过程序语句来操纵的,尽管计算机进行数学运算的速度是人难以比拟的,但是计算机本身并不能识别数据的真或伪。

5. 分析资料

对数据的统计学分析(analysis of data)按其分析目的可分为描述性统计分析与统计学推断分析。描述性统计(descriptive statistics)是指用统计指标、统计图、统计表等方法,对数据的特征及其分布规律进行检测与描述。统计推断(inferential statistics)是通过随机样本信息推断总体特征的过程。统计推断又包括置信区间(confidence interval)估计与统计学假设检验(hypothesis test)。统计学分析过程按变量的多寡可分为单变量分析与多重变量分析。这一部分毫无例外地是各种应用统计学教材讨论的重点。

6. 分析结果的合理解释

无论如何,统计学分析的目的就是应用观察试验的数据信息,对研究者提出的假说做出接受与否结论的推理过程;即通过统计学分析的结果对相应的专业问题做出一种理性的判断。

四、统计概念与术语

1. 同质与变异

同一总体中的每一个体都具有相同性质类别的特征称为同质(homogeneity)。同一总体中的各个个体间的差异性则称为变异(variation)。正所谓“在同一棵树上找不到完全相同的两片树叶”。

2. 总体与样本

根据研究目的所确定的具有相同性质的观察单位的集合称为总体(population)。总体又可分为:①有限总体(finitude population):指包含在一定时间、空间范围内的有限个观察单位。②无限总体(infinitude population):不宜确切划定范围的总体。从同一总体中通过随机化过程抽取的部分观察单位称为样本(sample)。

3. 样本含量

样本含量(sample size)指样本中所包含的观察或参与试验单位数,又称样本大小。

4. 概率

概率(probability)指用以度量非确定(随机)事件发生的可能性的统计学指标。所谓非确定的随机现象是指在一定的条件下,其试验的结果事件不具有惟一性与确定性特征的现象。如抛掷一枚硬币,待其落地后是哪一面朝上?一种降压药物用于10位高血压患者,每个患者对该药物的降压效应以及其他不良反应在试验之前是无法确切定论的,只有待试验结束之时方见分晓。然而,任何偶然事件都存在自身固有的规律,虽然我们不能确切地知道某一次随机试验具体的结果如何,但是我们可以用分布与概率的方式探讨其规律性,即用不确定性的概率度量与描绘随机事件与现象。对用概率度量不确定性随机现象描述得最为精彩的莫过于大文学家莎士比亚(Shakespeare):“如果你能洞穿时间的种子,并且知道哪一粒会发芽,哪一粒不会,那么就请告诉我吧!”在一个非确定性随机试验的结果尚未出现之前,虽然我们不能“洞穿”其结局,但是我们可以借助分布与概率来测度各种结局发生可能性的大小,给出一个非确定性的回答。统计学中将随机事件A的发生概率记为 $P(A)$ 。概率的取值在 $0 \sim 1$ 之间。 $P=1$ 或 $P=0$ 的事件称为必然事件, $0 < P < 1$ 的事件称为随机事件。概率接近于0(如 $P < 0.05$)的事件称为小概率事件。

5. 分布

一个随机现象的规律性通常通过随机事件及其概率来刻画。一个随机试验的所有结局事件与对应的概率的排列称为分布(distribution)。变量的值可以用来区别或描述各个个体,除了区别或描述之外,还有相对频率的概念。例如,人们对于一个新生儿的体重有6.5公斤的新闻表现出惊讶与不大相信的态度,而对一个新生儿体重有4公斤的消息被认为是极其普遍现象。实际上,在日常生活中,人们已经将一种变量的取值与该数值发生的可能性的一种度量联系起来,这就是一种事物的正常性或缺乏正常性的度量,或这样一个值出现的概率的度量,我们称这种变量具有一个分布。对应于样本数量值分布称其为频率分布(frequency distribution);对应总体数量值的分布则称其为概率分布(probability distribution)。

6. 随机变量

表现出变异性或变差特征的量称为随机变量(random variable)、机遇变量,有时简称变数或变量。随机变量这一名词是专指具有一个分布或一个概率或概率分布的变量。

7. 随机化

能使总体中每一观察单位均能以同等机会(概率)进入样本,或分配到实验组与对照组的过
程,称为随机化(randomization)。

随机化过程的特点是“机会均等”(equivalence chance)地选取或分配研究对象的“公平性”原则。我们可以列举很多执行这种“机会均等”法则的社会生活事例。如“福利彩票”、“足球彩票”、商家们主办的促销抽奖等活动;很多“众里抽一”的社会活动贯常采用“抽签”、“抓阄”、“猜拳”等

大众公认为合理的方法来获得“公平裁决”。连儿童们也会玩“剪刀、石头、布”游戏使“失败者”即使“受罚”也自认“公道”。这种游戏为何能在社会活动中流行得如此广泛,关键在于它的“公平性”,即“对于所有参与者而言,都有同等的机会接受奖励或惩罚”,并且该法则具有公认的“合理性”。这种“同等机会”在统计学中称为“等概率”,这种“合理性”的游戏规则就称为“随机化”。

8. 抽样误差

由于组成总体的各个个体间存在着差异性,因抽样过程的随机性导致样本的统计量与总体的参数不等,样本与样本的统计量之间存在差异的特性称为抽样误差(sampling error)。

9. 参数与统计量

描述总体特征的数量称为参数(parameter);常用希腊字符表示,如 μ 表示总体均数, σ 表示总体标准差, π 表示总体率。描述样本特征的数量称为统计量(statistic);常用英语字母表示,如 \bar{x} 表示样本均数, s 表示样本标准差, P 表示样本率。

五、如何学习统计学

常常有人问到这样一个问题:“生物统计学难学吗?”回答是:既难学,又不难学。因为,其一:生物统计学是生命科学与数学的交叉学科,它具有很多数学的学科特征,如抽象的数学概念、数学符号、数学公式等。这些东西要求学习者付出一定的努力,花费一定的时间与精力来掌握统计学的基本概念、基本原理、基本方法,记忆一些常用的符号与公式。学习者不应该用阅读小说那样的速度来阅读统计学教科书,平均而言,读一页统计学教科书所花的时间也许要相当于读40页小说的时间,甚至于更多。其二:生物统计学是一门实用性很强的应用性学科,重在用统计学的原理和方法来解决生命科学研究中的实际问题。学习“如何应用统计学方法来解决实际专业问题”,这是必要且相对简单的内容。但是要深入钻研“为什么要这样做”就有相当难度,而且涉及到复杂的数学理论。从实用的角度出发,没有必要对每个公式的来龙去脉弄得一清二楚。到目前为止,还没有人写出一本“不用伤脑筋的统计学”(statistics without tears);也不存在一种对任何人都适用的“新东方”式“速成掌握统计学”的秘方。尽管如此,我们可以对每一位学习者提供一个学习建议:其一是“抓住三基”,即基本概念、基本原理、基本方法;其二是“重视实际应用”。我们相信这些建议也许对于初学者而言,是目前学习统计学的最好“秘方”。实际上,本书所涉及的数学计算方法只有加、减、乘、除、平方与开方、对数与反对数和阶乘运算等常规的数学方法,所以,学习生物统计学困难的不是数学计算方法,而是学习如何用概率与分布理论来处理非确定性的事物(现象)的逻辑思维方法;学习如何在非确定性度量的基础上做出决策的过程。

生物统计学的过程就是把数学的语言引入具体的生命科学领域,把具体生命科学领域中要待研究的问题抽象为数学问题的过程。为了让大家能较好地领悟生物统计学中一些重要的思维方法,特做如下提要。让大家在进入以下章节内容的学习过程之前,建立一些粗略的概念,想必对初学者是会有帮助的。

(1) 统计学的一个重要目的是用有限的样本信息对该样本所来的总体特征进行推估,即用较小的可观察对象群(样本)信息,对不可能或不便观察的大对象群(总体)进行推估。这是统计学的精髓所在,也是其何以成为一门科学的缘由。正如法国数学家Laplace (1745~1827)所说:“由赌局引起的一门学问,居然会为人类知识最重要的研究对象,这实在令人非常惊奇。”

(2) 随机抽样(random sampling)的概念十分重要。因为,通过随机抽样所获得的样本结果可以告诉我们该样本所来的那个总体的信息。这就意味着该样本可以使我们了解到该总体中未被测量的那一部分的信息。要从样本做出关于总体特征的推断,该样本必须是总体的代表。所以,统计学中所说的样本都应该是通过随机化抽样获得的样本。要取得一个能代表总体的样本,抽样过程必须体现随机化原则,如果样本不是通过随机化方法获得,那么,由该样本所做出的统计学结论是无效的。

(3) 统计学推断的逻辑,对大多数人而言是一种新的思考方法。统计学意味着一种新的思维方法,即从非确定性或概率的角度来思考问题。统计学推断过程不是从必然推导必然,而是从偶然中发现必然。不是对现象做准确的数量刻画,而是做出近似的数学描绘和揭示各种可能出现结果的比例分布。所以,统计推断过程中没有“证明”一说。

(4) 科学的进步通常在于发现变量间的新关系。每一项科学研究的最终目的是发现变量间新的关联性。科学哲理告诉我们,除了事物的性质和数量间的关系外,没有其他办法能表达关联的意义,包含变量间的关系。变量间相关或关联性的基本特征有两点:其一是相关或关联性的强弱数量;其二是这种相关或关联性的置信度大小。可以说:所有的统计学分析都是为了探讨变量间质量上的关联性或数量上的相关性。

(5) 哪里有随机事件的问题,统计学就能在哪里派上用场。数理统计是科学研究的一种工具,不限于哪一种科学研究,无论遗传学、营养学、临床医学还是预防医学等。

(6) 计算机是统计的工具,统计离不开计算机。经历一个世纪,生物统计学获得了长足的发展,计算机与网络技术为生物统计学的应用开辟了更广阔的前景。国际通用型统计软件,如SAS (Statistical Analysis System)、SPSS (Statistical Package for the Social Science), 能为研究者提供功能齐全的统计学分析过程。学习者应在学习生物统计学的同时,也应该发展自身应用统计软件解决实际问题的能力。

(7) 统计学不是万不能,也不是万能,它只是解决具体专业问题的辅助工具,不能替代专业思考。

(8) 统计学只能帮助研究者发现规律,而不能创造规律。科学假设是根据已知的科学事实和科学理论,对所研究的自然现象及其变化的规律提出假定性推测和说明。统计方法论采用的这种科学思维方式有时被误解了,以致产生了许多关于统计和统计学家的怪论。认为“统计学家是一个从无根据的假设到预料中的必然结局之间划一条精确界限的人。”这就产生了对统计方法的滥用,这种滥用现象在杂志、报道以及电视广告中屡见不鲜。其中也不乏有说谎者利用统计学说谎;有些人对待统计学方法采用实用主义的态度,他们拿着数据来用统计方法,目的仅在于用要统计学上的一个‘ P 值’来支持他的结论。要知道统计学是帮助研究者发现规律的科学,而不是用以“创造规律”的工具。他们应用统计学是就如同醉汉应用街旁的路灯,不是用路灯为他照明,而是用路灯柱来支撑自己。我想,经过系统性地学习生物统计学的人,应该不会用不合适的眼光来看待统计学。

(董时富)

第一章 研究设计中的基本统计学原则

生物统计学是现代生命科学研究中不可缺少的工具,它在生物学、医学、农学中有着广泛的应用。从绪论中我们知道,生物统计学的工作程序主要分为四个步骤,即统计设计、收集资料、整理资料和分析资料,而设计是最为关键的一步,由于生物学研究、医学研究和农学研究在统计研究设计上的基本原理是一致的,故本章将以医学研究为主线,介绍研究设计的基本统计学原则。

第一节 研究设计的重要性

医学研究的目的在于揭示生命本质和疾病发生发展过程,认识健康与疾病的转化规律,提出有效的防治措施,增进人类生命健康。基础医学主要研究人体健康如何向疾病转化的规律;预防医学主要研究如何防止健康向疾病转化的规律;临床医学主要研究如何促进疾病向健康逆转,完成康复过程。

医学研究是一个极其复杂的认识活动,是一个由感性认识到理性认识的高度科学思维活动。1950年至1964年Doll和Hill关于吸烟与肺癌关系的研究就是一个很好的例证。他们首先提出吸烟可以导致肺癌这一假说,其根据之一是英国扫烟筒工人阴囊上皮癌多发,是由于煤炭不完全燃烧的颗粒夹杂于阴囊褶皱所致;其二是某些染料中含有苯胺易引起工人多发膀胱癌;其三是辐射可以引起某种骨肉瘤。在这些感性认识前提下,他们经过科学推理,得出一般认识,即某些化学物质可以致癌。烟草的不完全燃烧,可产生致癌物质,进而促使肺癌多发。为了验证吸烟可以引发肺癌这一假说,Doll和Hill根据演绎思维作了如下推理,并应用现代流行病学方法进行了下述调查内容的试验设计:

(1) 将肺癌患者与非肺癌患者配对,前一组吸烟者应比后一组多,吸烟史也应该长。这是由果到因的回顾性调查研究。

(2) 将吸烟医生与非吸烟医生配对随访多年后,前一组肺癌发病率应比后一组高。这是由因到果的前瞻性调查研究。

(3) 肺癌死亡率高与吸烟量多少应成正相关。

(4) 戒烟后的医生中肺癌死亡率应低于未戒烟者。

经过长期的观察,并对资料进行统计分析,得到了支持假说的结果,研究取得成功。后来的许多关于吸烟与肺癌关系的研究都基本肯定了Doll和Hill的研究成果。这说明通过科学正确的试验设计得到的结论是具有生命力的。

相反,如果一项研究缺乏良好的研究设计,则会影响到结论的真实性和可靠性。如20世纪80年代初,有研究报道孕期补充维生素(叶酸)可以减少生育神经管缺陷婴儿的危险性,即先服用维

生素后怀孕的妇女比怀孕后才开始服用维生素的妇女和拒绝参加试验的怀孕妇女所生的婴儿神经管缺陷的发生率要低得多。但遗憾的是,由于参加服用维生素试验和拒绝试验的孕妇之间存在某些生理特征上的系统差别,致使在解释试验结果时发生困难。这不能不认为是因实验设计缺乏周密考虑所造成的教训。为补救先前研究的不足,在后来的研究中采用了随机化分配受试者的方法,分叶酸补充组和安慰剂组。但在观察结束时,又因样本数过少而无法做出肯定的科学结论。直到1991年,研究者报道了一个大样本的随机化试验,获得了肯定的科学结论。在安慰剂组中的602名怀孕妇女中有21名妇女分娩出的新生儿有神经管缺陷,在叶酸补充组的592名怀孕妇女中出现新生儿神经管缺陷者只有6例,而其他维生素(不含叶酸)的补充对神经管缺陷的发生无明显影响。统计学分析证实叶酸补充组与安慰剂组之间新生儿神经管缺陷发生率存在显著性差异,说明叶酸对预防新生儿神经管缺陷确有明显的效果。

由此可见,研究设计(research design)是科研工作中的第一步基本而又至为重要的工序,是科研计划的具体实施方案,是进行实验和统计分析的先决条件,是科学研究获得预期结果的重要保证,它的好坏将直接影响着研究的质量和成功。然而,在科研工作中,有的研究因没有一个良好的实验设计为指南,出现顾此失彼,观测指标遗漏或实验数据残缺的严重情况;有的研究人员忽视统计学在科研工作中的重要作用,仅依赖现有的专业知识进行研究,只是在收集实验数据之后才开始想到运用统计学知识,可以想像其研究结果的真实性和可靠性。现代医学正在经历由经验型向科学型的变革,如果仅凭经验开展研究工作,往往花了大量的人力、物力、财力和时间,而研究的问题仍得不到解决;同样地,对于一个设计不佳的研究常常招致失败及时间和经费的浪费……。因此,在研究项目确定后,就需要进行研究设计,其意义有二:一是能够根据研究目的,规定具体的研究任务和所要采取的技术路线与方法,科学、合理、有效和周密地计划安排研究的全过程;二是能够用较为经济的人力、物力、财力和时间进行实验,最大限度地减少误差,获得可靠的结果,从而起到事半功倍的效果。

第二节 系统误差与控制方法

误差(error)即为实测值与真值之差。任何实验研究的结果都不可避免地会出现不同程度的误差,要使误差减少到最低限度,必须在实验的每个环节上加以控制。实验误差主要有随机误差(random error)和系统误差(systematic error),其中随机误差是由于一系列实验或观察条件的波动而产生的误差,这种波动是随机的,且在实验中无法避免,但可通过统计学方法进行处理。这里主要介绍系统误差及其控制方法。

一、系统误差定义

在一定实验条件下,由某种未发现或未确定的因素所引起观测值具有方向性和系统性的误差称为系统误差,又称偏倚(bias)。这种误差是不能用统计方法估计的。

二、产生系统误差的常见原因

- (1) 仪器差异 仪器未进行校正。
- (2) 方法差异 测量方法上的不同导致测得的数据存在较大差异。如用放射免疫分析、化学