



普通高等教育“十一五”国家级规划教材

医学统计学

Medical Statistics

(第二版)

刘桂芬 ● 主编

刘玉秀 仇丽霞 ● 副主编

中国协和医科大学出版社

普通高等教育“十一五”国家级规划教材

医学统计学

Medical Statistics

(第2版)

主 编 刘桂芬
副主编 刘玉秀 仇丽霞

编委 (按姓氏笔画排)

于 浩	南京医科大学	肖 琳	中国疾病预防控制中心
仇丽霞	山西医科大学	陈长生	第四军医大学
尹 平	华中科技大学	陈平雁	南方医科大学
王立芹	河北医科大学	易 东	第三军医大学
田考聪	重庆医科大学	金水高	中国疾病预防控制中心
刘 艳	哈尔滨医科大学	赵晋芳	山西医科大学
刘玉秀	南京军区南京总医院	郜艳晖	广东药学院
刘桂芬	山西医科大学	凌 莉	中山大学
吕 桦	上海市浦东新区 CDC	宿 庄	内蒙古医学院
余红梅	山西医科大学	阎玉霞	南方医科大学
吴艳乔	华西医科大学	黄高明	广西医科大学
张岩波	山西医科大学	韩少梅	中国协和医科大学
李新华	贵阳医学院		
秘书 罗天娥	山西医科大学		

中国协和医科大学出版社

图书在版编目 (CIP) 数据

医学统计学 / 刘桂芬主编. —2 版. —北京: 中国协和医科大学出版社, 2007. 8
普通高等教育“十一五”国家级规划教材
ISBN 978 - 7 - 81072 - 946 - 8

I. 医… II. 刘… III. 医学统计 - 高等学校 - 教材 IV. R195.1

中国版本图书馆 CIP 数据核字 (2007) 第 117520 号

普通高等教育“十一五”国家级规划教材
医学统计学 (第 2 版)

主 编: 刘桂芬
责任编辑: 吴桂梅 严 楠

出版发行: 中国协和医科大学出版社
(北京东单三条九号 邮编 100730 电话 65260378)

网 址: www.pumcp.com
经 销: 新华书店总店北京发行所
印 刷: 北京丽源印刷厂

开 本: 787 × 1092 毫米 1/16 开
印 张: 30.25
字 数: 740 千字
彩 页: 2
版 次: 2007 年 8 月第二版 2007 年 8 月第一次印刷
印 数: 1—3000
定 价: 40.00 元

ISBN 978 - 7 - 81072 - 946 - 8/R · 937

(凡购本书, 如有缺页、倒页、脱页及其他质量问题, 由本社发行部调换)

再 版 前 言

为了适应新时期医学教育发展的需要，五前年我们组织编写了《卫生统计学》，汇集各位编者多年的教学科研与工作经验，于2003年8月由中国协和医科大学出版社出版，为多所高等医药院校选用，深受读者厚爱。2006年8月被正式评审确定为普通高等教育“十一五”国家级规划教材，更名为《医学统计学》。

为适应社会发展对医学生的新要求，不断完善学科知识结构，优化教材内容，我们对《卫生统计学》进行了增删修订。《医学统计学》力求做到：紧扣基础理论、基本知识和基本技能主线，反映学科新进展，保持本学科知识的系统性；强调医学科研设计与统计分析和软件应用的有机结合，注重统计结果的解释；力图结合医学应用，深入浅出，融会贯通，突出培养学生解决实际问题的能力。

全书分21章，本次再版侧重于三个方面：第一部分主要介绍医学统计基础理论与基本方法，针对《卫生统计学》多元分析方法薄弱的情况，加强了不同类型资料的回归分析方法与软件结果解释，增加了诊断试验的分析与评价。第二部分医学研究设计，在继承第一版教材内容领先优势的基础上，进一步引入了临床试验设计的最新进展。第三部分公共卫生与社区医学统计，不仅突出了公共卫生应用统计的特殊性，增加了信息系统监测数据与生物信息数据分析进展和预测评价方法，而且针对学生动手机会少，综合能力弱的问题，加强了综合分析与解决问题环节训练，补充增加了国际标准化要求的医学论文统计结果报告。书中部分章节加“*”内容，可据不同层次教学予以选择。书后附有可供分析讨论和实习应用的练习题。

本版教材由十九所院校二十七位教师参加编写，刘玉秀教授和仇丽霞教授两位副主编参与了稿件的组织与修改，为教材付梓尽心尽力；山西医科大学卫生统计教研室赵晋芳、罗天娥、萨建等老师，硕士研究生曾平、骆常好、曹红艳、王晓芳、凌建春、寇林元等同学，在教材编审的联络、编辑、复核、修改、校对等方面做了大量细致的工作。《医学统计学》的再版承蒙山西医科大学和广西医科大学领导的关怀和支持，两校卫生统计教研室为本书修订会议提供了后勤保障与服务；统计学界老一辈恩师严谨治学、无怨无悔耕耘精神的激励，第一版编委们精心编写、无私奉献和良好的铺垫，新版全体编者不计名利、夜以继日的工作。它是第一版教材的延续和提升，凝聚着两版编者的智慧和努力。值本教材出版之际，谨致以衷心感谢。

本教材除可用于本科生教学外，也可作为医药院校博士、硕士研究生和基层科研工作者的参考用书。

鉴于我们能力所限，教材中难免存有疏漏之处，诚望同仁与广大读者不吝指正。

刘桂芬

2007年7月20日

第一版前言

本教材是广大编写者集多年的教学与工作经验通力协作编写而成的，它紧扣基础理论、基本知识和基本技能主线，力图紧密结合医学应用，更好地培养学生分析、综合、解决实际问题的能力。

基于卫生统计学科的自身发展和实际工作需要，本书力求：①深入浅出，通俗易懂，可读性强；②针对目前实际工作的需要，配合专业学位研究生教学，充实了医学研究设计的新内容；③加强了卫生统计基本理论知识和技能的系统阐述，注重统计结果的正确理解和表述。全书共分18章，第一章概括介绍卫生统计基本内容，第二章至第十三章介绍基本统计方法，第十四章与第十五章介绍医学研究设计，第十六章至第十八章介绍统计方法在医学尤其是预防医学3个特殊领域中的应用。部分章节加“*”内容，可根据不同层次学习对象予以选择。另外，附有供实习使用的练习题。

本书编写过程中，承蒙山西医科大学及公共卫生学院领导的热情关怀和支持。统计学界老一辈恩师杨树勤教授、何大卫教授、苏景铭教授和刘庆欧教授悉心指导，并为本教材做了精细加工和审阅，提出许多中肯的建议。他们严谨治学、无私奉献的优秀品德，无怨无悔的耕耘精神，为我们年轻的一代树立了楷模，也为保证教材质量奠定了良好基础。

本教材由12所院校的17位教师参加编写，安徽医科大学吕桦教授、南京军区总医院刘玉秀副教授参与了稿件的组织与修改，为教材成稿付出了辛勤劳动。山西医科大学卫生统计学教研室仇丽霞副教授、肖琳讲师、赵晋芳等老师，硕士研究生孟海英、郎素萍、赵铁牛、冯志兰、孙谨芳等同学，在教材编辑、核校、修改、文印及联络等方面做了大量工作。山西医科大学科研处、南京军区总医院领导给予了支持和帮助。值本教材出版之际，谨致以衷心感谢。

本书是老、中、青教师传、帮、带、学的结晶。它可作为本科生、研究生卫生统计学教材，也是基层教学与科研工作人员不可缺少的一本非常实用的参考书。

成书过程中，虽经反复修改、核校，但由于我们水平有限，经验不足，书中肯定存在不少缺点和错误，诚恳希望广大师生和读者批评指正。

刘桂芬

2003年5月30日

目 录

第一章 绪论	(1)
第一节 医学统计学与数学和计算机	(1)
第二节 科研工作中医学统计学的作用	(2)
第三节 医学统计学中常用的几个基本概念	(4)
小结	(8)
第二章 医学资料的统计描述	(9)
第一节 频数分布表和频数分布图	(9)
第二节 定量资料集中趋势指标	(12)
第三节 定量资料的离散趋势指标	(17)
第四节 分类资料的统计描述	(19)
第五节 动态数列	(25)
小结	(27)
第三章 正态分布及其应用	(28)
第一节 正态分布的概念和特征	(28)
第二节 标准正态分布及其应用	(31)
第三节 医学参考值范围的制定	(32)
第四节 正态性判定	(36)
小结	(40)
第四章 总体均数的估计与假设检验	(41)
第一节 均数的抽样误差与标准误	(41)
第二节 t 分布	(44)
第三节 总体均数的估计	(45)
第四节 假设检验的原理和步骤	(48)
第五节 t 检验	(51)
第六节 假设检验的两型错误	(58)
第七节 假设检验时应注意的问题	(59)
小结	(61)
第五章 方差分析	(63)
第一节 完全随机设计资料的方差分析	(63)

第二节	随机区组设计资料的方差分析	(71)
第三节	析因设计资料的方差分析	(73)
第四节	重复测量资料的方差分析	(77)
小结		(82)
第六章	二项分布、Poisson 分布及其应用	(84)
第一节	二项分布的概念	(84)
第二节	二项分布的应用	(86)
第三节	Poisson 分布的概念	(89)
第四节	Poisson 分布的应用	(91)
小结		(94)
第七章	χ^2 检验	(95)
第一节	χ^2 检验的基本思想	(95)
第二节	四格表资料 χ^2 检验	(96)
第三节	行 \times 列表资料的 χ^2 检验	(99)
第四节	率的多重比较	(102)
第五节	频数分布拟合优度的 χ^2 检验	(103)
第六节	四格表资料的确切概率法	(105)
第七节	线性趋势检验	(107)
小结		(109)
第八章	基于秩的非参数检验	(110)
第一节	配对设计符号秩检验	(110)
第二节	完全随机设计两样本比较的秩和检验	(112)
第三节	完全随机设计多个样本比较的秩和检验	(115)
第四节	随机区组设计资料比较的秩和检验	(117)
第五节	多个样本资料的两两比较	(118)
小结		(120)
第九章	双变量线性回归与相关	(121)
第一节	简单线性回归	(121)
第二节	双变量相关分析	(129)
第三节	Spearman 秩相关	(133)
第四节	回归与相关分析应注意的问题	(135)
小结		(138)
第十章	观察性研究设计	(139)
第一节	观察性研究的意义与特点	(139)

第二节	观察性研究常用的方法	(139)
第三节	观察性研究设计的内容	(144)
第四节	几种常用的抽样方法及样本含量估计	(154)
第五节	调查质量的控制	(160)
小结		(164)
第十一章	实验研究设计	(165)
第一节	实验研究的基本要素	(165)
第二节	实验设计的基本原则及误差控制	(167)
第三节	常见的实验设计类型和方法	(170)
第四节	样本含量的估计	(175)
小结		(179)
第十二章	临床试验设计与分析	(180)
第一节	临床试验概述	(180)
第二节	临床试验设计与偏倚控制	(184)
第三节	临床试验数据管理与统计分析	(190)
第四节	非劣效性/等效性临床试验	(194)
小结		(202)
第十三章	诊断试验评价	(203)
第一节	评价诊断试验的常用指标	(203)
第二节	ROC 曲线及 ROC 曲线下面积的估计和检验	(209)
第三节	两样本资料诊断准确度的比较	(212)
第四节	诊断试验评价注意事项	(214)
小结		(215)
第十四章	多重线性回归	(216)
第一节	多重线性回归分析	(216)
第二节	多重线性相关	(222)
第三节	回归变量的筛选	(224)
第四节	多重线性回归应用及应注意的几个问题	(226)
第五节	回归诊断	(228)
小结		(234)
第十五章	logistic 回归	(236)
第一节	logistic 回归模型的基本概念	(236)
第二节	二分类反应变量的非条件 logistic 回归	(241)
第三节	多分类反应变量的 logistic 回归	(242)

第四节	1:1 条件 logistic 回归	(248)
第五节	剂量反应关系及半数效量估计	(250)
小结	(253)
第十六章	生存分析	(254)
第一节	生存分析的基本概念	(254)
第二节	生存曲线的估计	(257)
第三节	生存曲线的比较	(261)
第四节	Cox 回归	(263)
小结	(269)
第十七章	医学人口与疾病统计	(271)
第一节	静态医学人口统计	(271)
第二节	出生与计划生育统计	(279)
第三节	死亡统计	(284)
第四节	疾病和残疾统计	(289)
第五节	寿命表	(295)
小结	(309)
第十八章	传染病监测数据的统计分析概述	(311)
第一节	传染病监测的概念与进展	(311)
第二节	国家法定传染病报告监测数据的分析	(313)
第三节	传染病监测数据的空间分析	(317)
第四节	传染病暴发的早期探测与预警	(323)
小结	(325)
第十九章	综合评价方法	(327)
第一节	综合评价的意义和基本步骤	(327)
第二节	综合评价的基本方法	(328)
小结	(344)
第二十章	医学统计预测	(346)
第一节	统计预测概述	(346)
第二节	统计预测的基本方法	(347)
小结	(374)
第二十一章	生物信息统计分析方法	(375)
第一节	生物信息学概述	(375)
第二节	序列相似性比较方法	(376)
第三节	基因芯片的统计分析方法	(379)

第四节 基因调控网络分析方法	(383)
第五节 蛋白质序列模式和序列结构域模式	(387)
小结	(388)
第二十二章 医学论文统计结果报告	(389)
第一节 医学论文中统计学内容表达的一般要求	(389)
第二节 统计表与统计图	(392)
第三节 医学研究报告的标准化	(395)
小结	(400)
附录一 统计用表	(402)
附表 1 标准正态分布曲线下的面积	(402)
附表 2 t 界值表	(403)
附表 3 F 界值表	(404)
附表 4 q 界值表	(408)
附表 5 二项分布参数 π 的可信区间	(409)
附表 6 Poisson 分布 μ 的可信区间	(411)
附表 7 χ^2 界值表	(412)
附表 8 T 界值表 (配对比较的符号秩和检验用)	(413)
附表 9 T 界值表 (两样本比较的符号秩和检验用)	(414)
附表 10 H 界值表 (三样本比较的符号秩和检验用)	(415)
附表 11.1 随机区组设计 (Friedman) 检验统计量 M 界值表 ($\alpha=0.05$)	(416)
附表 11.2 随机区组设计 (Friedman) 检验统计量 M 界值表 ($\alpha=0.01$)	(416)
附表 12 r 界值表	(417)
附表 13 r_s 界值表	(418)
附表 14 总体均数估计时所需样本含量 ($\alpha=0.05$)	(419)
附表 15 总体率估计时所需样本含量 ($\alpha=0.05$)	(419)
附表 16 随机数字表	(420)
附表 17 随机排列表	(421)
附表 18 ψ 值表 (多个样本均数比较时所需样本例数的估计用) ($\alpha=0.05, \beta=0.1$)	(422)
附录二 练习题	(423)
附录三 常见统计学专业名词英汉对照	(464)
附录四 参考文献	(469)

第一章 绪 论

重点掌握：

1. 医学统计工作的基本步骤。
2. 常见的医学统计资料的类型。
3. 医学统计学常用的几个基本概念：总体与样本、随机误差、概率与频率。

医学统计学 (medical statistics) 是描述、归纳、探索医学数据分布特征和解释数据规律的一门学科, 也是医药卫生工作者合理运用统计学原理与方法、充分提取信息、深入揭示健康与疾病发生发展规律和预后评价的一种手段。现代统计学已不仅仅是对观察、测量和记录情况作一些整理归纳, 更重要的是利用统计方法揭示资料信息, 对研究事物现象做出科学合理的推断与决策, 以帮助我们更好地认识和掌握健康与疾病的变化规律。

第一节 医学统计学与数学和计算机

大量的循证实践在疾病的预防、治疗、健康促进与健康状况研究领域都已取得令人瞩目的成就。要循证实践就需要积累证据, 并予以严谨合理的解释。前者已把更多的人带入医学循证实践, 而后者则要求所有医药卫生工作者具有评价所做研究的能力。医学研究中大多数证据是以数据形式出现的, 医学统计学就是根据研究目的搜集、清理和分析数据, 估计和解释医学数据规律的科学, 是循证实践中起关键作用的技术。怎样从医药研究资料中提取所需信息, 帮助研究者树立科学的思维, 对研究结果给予恰如其分的评价与解释, 正是医学统计学所要解决的问题。

早在 200 多年前, 法国数学家 Laplace 就深信, 概率是发现真理的主要方法, 医学是概率应用的一个重要领域。19 世纪, 尽管 Louis、Farr 等人首先把统计方法应用于医学研究, 但并未广泛推广。直到 20 世纪中叶, 正是 Fisher、Pearson 等利用随机化实验和抽样理论, 将统计推断技术深入到各个方面, 这在充满统计推断理论和统计计算的医学杂志中是显而易见的。

有些研究者由于回避数学推导而不敢进行统计研究, 或不关心具体应用, 只把统计学作为数学的一部分去证明, 这些均已不适应医学学科发展的需要。本教材拟将基本概念的理解与计算机软件输出结果的解释列为掌握的重要内容, 把实际应用问题与概念解释有机结合起来, 不苛求呆板的数学证明。个别章节结合课堂的简单运算与公式讲解, 更好地启发和帮助研究者系统掌握解决实际问题的科研设计与资料分析能力。

随着计算机统计软件技术的深入发展, 应用统计学的研究如虎添翼, 硕果累累。统计软件作为一种专门的应用软件, 是利用计算机软件技术呈现统计数据, 进行数据分析, 模拟和

实现统计过程的一类应用软件，是应用统计的一项专门技术，是统计方法应用的重要载体，在医学统计数据分析中具有重要的地位。它不仅是一种工具，也是一种方法，在统计应用过程中起到了无可替代的积极作用，是统计学发展中不可割裂的一个重要部分。统计软件的发展也更快地推动了现代统计方法的新进展，多水平模型理论和 MLwiN 就是典型的例子。成千上万个医药研究项目的管理与分析，大都是借助统计软件来完成的。

从统计应用角度考虑，我们应注重的不再是计算问题，而是统计问题，应了解的是为什么采用这些特殊的计算，计算结果的实际意义是什么。目前研究中出现的问题不再是复杂计算的错误，而更多的是统计方法的不恰当运用，专业信息解释的不合理性。因此，掌握统计软件的应用技术和统计方法的适用范围与条件，理解研究结果并做好恰如其分的解释就显得更为重要。作为新型医学人才，学习医学统计学的原理和方法，有助于我们结合预防医学、临床医学、基础医学、药学及卫生事业管理等专业知识，更好地完善医学理论，正确地进行医学研究设计，合理地选择统计方法，恰当地解释研究结果，以独特的统计思维，不断地修正我们对生物医学现象的认识，科学地揭示数据中所蕴藏的内在规律。

本教材拟从医学研究设计（观察性研究、实验性研究和临床试验研究设计等）、统计学基本原理与方法（包括资料搜集、清理和分析的原理与方法，正确合理的推断与预测结果的解释等）和应用统计及其统计软件（包括社区医学人口统计、疾病统计及卫生领域常用的预测和综合评价方法等）几方面进行讨论，在《卫生统计学》第一版基础上，增加了生物信息、传染病监测数据分析方法等。

第二节 科研工作中医学统计学的作用

众所周知，事物的量化有助于提高对事物认识的准确性和深度。科学研究是一种探索未知的认识活动，总是要和数据打交道，使它与统计结下不解之缘。医学研究的对象是人，人是世界上最复杂的生命体，不但具有生物性，还具有社会性；不但有生理活动，还有心理活动，个体变异错综复杂。医学统计学现已成为医学研究中不可缺少的科学方法，它的作用不仅横跨医学研究的各个学科，而且纵贯医学科学研究的全过程。

一、医学研究设计

研究设计（design）是医学统计工作过程的一个重要内容，它是医学科研工作的第一步，是对医学科学研究过程、内容及具体实施方法和步骤的总设想或安排。设计就是针对具体的研究项目或问题，确定调查对象和观察单位（个体），根据是否施加干预因素，明确分析指标。针对如何获取原始数据，怎样进行资料的清理和分析，如何控制误差，预期分析结果有哪些等提出详尽的实施方案和技术路线，作好周密的考虑和安排。根据内容分为：①专业设计：反映研究者对专业知识掌握的能力和程度，主要与科研课题或项目的深度及水平有关；②统计设计：反映研究者对统计知识与技术正确应用的程度，主要与科研工作的质量有关。怎样才能以较少的人力、物力和财力，获取准确、可靠的科学结论，搞好研究设计是非常关键的一个内容，也是探索新知识、验证新理论、推广新方法等必不可少的手段。

二、搜集资料

搜集资料 (collection of data) 是统计工作的基础, 它直接关系到科研工作的质量。其任务是研究人员按科研设计要求, 获取准确、可靠、有用的原始数据, 并做好质量监控与评估。若所搜集资料不准确、不完善, 使用再“科学”的加工方法所得结果也只能给人以假象, 失去研究的意义与价值。数据的准确性是指观察、测量、记录或计算的数据, 均无虚假差错之处, 尽可能做到界限明确, 真实可靠, 不造成混淆。数据的完整性是指用来研究分析的项目没有遗漏、重复和缺失。数据搜集的及时性是指资料在一定条件下, 按规定的时间完成填写与登记等。

医学统计资料的来源主要有:

1. 常规保存的工作记录 如人口登记、出勤记录、职工流动、工伤、出生、职业病报告、恶性肿瘤报告卡、体检、手术记录、化验、病理报告等, 它是医学研究资料的一个重要来源。但由于这些数据不一定是为某专项研究所收集, 有时会给分析带来诸多不便, 或由于工作人员责任心不强, 出现漏报、重复和缺失等。资料搜集过程中除应加强技术督查外, 尚应注重不断积累。

2. 卫生服务信息监测系统及统计报表 主要指由医疗卫生机构根据国家有关部门规定统一管理, 监测哨点逐级上报的内容。如法定传染病预警系统, 慢性病、死因监测系统, 出生缺陷监测与医院管理系统, 社区居民健康档案和医院工作报表等。它可为了解居民健康状况, 拟订医疗卫生健康服务计划, 合理配置医疗卫生资源等提供科学依据。

3. 专项调查与实验研究资料 专项调查或实验研究一般指为解决某个(些)问题或验证某个(些)假说等所进行的专门研究。如全国7~10岁儿童龋患率现场调查、某地2型糖尿病普查、不同中草药配方治疗高血压病临床试验研究记录等。

4. 外源资料 指为工作需要取自公开发表的报告、专业参考文献、基因数据库、商业数据库、人口、公安、保险等与人类健康相关的统计资料等。

三、清理资料

清理资料 (cleaning data) 是按设计要求将一些分散的、表现个体特征的原始数据系统化、条理化, 以便更好地揭示所研究事物的内在规律, 它往往结合数据督查、数据库构建和具体研究过程而进行。应考虑:

1. 资料核查 除搜集资料时调查员自查、调查员间互查和专业人员核查外, 清理资料尚应对全部分析资料进行逻辑检查和数字核准等。包括对原始调查表中项目的审核, 缺失数据的检查, 误填、漏填项目的核准、修正, 数据类型以及编码等问题的考虑, 它是保证和提高数据分析质量的前提。

2. 数据结构与特点 任何观察或实验研究获得的结果, 都必须结合专业知识转变为数据后, 才能进行分析。无论是字符型, 还是数值型变量, 均可用二维数据结构矩阵来表述。一般情况下, 行表示观察单位, 列表示分析指标或变量。Excel、SAS、SPSS 等软件均可以此形式作为数据录入的主要格式。其中观察单位 (observation unit) 亦称个体 (subject), 随研究目的而异, 它可大可小, 可以是一个人、一组群、一份样品、一个采样点、一毫升水或

一个病室等。一次对观察单位的测定,可获得一个测量记录,也可对同一个体进行多次观测,获取几个测量记录,但每个观察单位所得数据,一般放在同一行上。主要研究指标可以是研究项目、混杂因素,也可以是研究对象的基本特征,可以是分析变量,也可以是分组变量或协变量等。它们可由测量结果直接录入,也可以经数据转换生成新的数据。

3. 数据编码 数据汇总时,应由专业人员根据专业要求进行编码。编码技术包括:①设计编码:如调查问卷设计时数值变量值的位数、取值范围控制,不详数据的编码,定性数据的数量化,连续变量的离散化等;②过录与检查编码:无论是调查员还是专业技术人员,搜集到数据后,在核查项目内容和数据的基础上,将编码过录到调查表或电子档案确定的位置上;③录入编码:根据调查内容核准无误后,选择两名以上责任心强、数据录入人员,进行双份录入和一致性检查,对发现的错误进行纠正。

4. 设计分组 一般应据研究目的、资料性质、观察单位数多少和习惯用法来考虑分组。主要有质量分组和数量分组。质量分组即按研究事物的属性特征进行归组,如性别:男性15人,女性13人。数量分组即按研究指标观测值数量大小来归组,详见第二章第一节。

5. 预分析 根据不同的应用软件,建立数据库时应注意的问题略有不同,但都可进行变量分布、数据特征描述、探索性研究等预分析,以利于更好地揭示研究事物内部的共性和对比组之间的差异性。

四、分析资料

资料分析 (analysis of data) 亦称统计分析,包括统计描述和统计推断两方面。统计描述是按研究设计要求,计算相应的统计指标,选用适当的统计表或统计图来概括数据特征,阐明事物现象的水平及内在联系;统计推断是根据抽样原理,在概括样本信息的基础上,对所研究总体的特征进行推断,主要包括估计、预测、假设检验和结果的正确解释,它是本教材内容的主体。

总之,以上四个步骤是紧密联系、环环相扣、不可分割的整体,任一阶段有缺陷,都将造成一定的损失,甚至导致科学研究工作失败。

第三节 医学统计学中常用的几个基本概念

一、同质与变异

同质性 (homogeneity) 指研究事物现象存在的共性,它是统计研究的基础,是资料清理和分析的前提。任何源于事实的数据,皆应以组内尽可能相同或相近,对比组间具有均衡可比性为依据。

尽管在同质总体中,不同个体某研究指标观测值间经常也存在不确定性。这种同质群体中自然状态下个体值之间千差万别、参差不齐的情况称为变异 (variation)。变异是客观存在的,是统计研究的对象,而同质是相对的。统计学就是用来描述某总体同质性,处理不确定性 (变异) 的科学和艺术。

二、医学统计资料的类型

要进行统计分析就需要有足够量的反映不确定性的数据,无论用何种方式搜集数据,都应根据研究目的,划清同质总体的范围,确定研究对象和观察单位。随研究目的不同,其内容可有不同。如欲了解某市 10 岁健康男童身高水平,凡在该市居住两年以上,年龄满 10 周岁,排除了患有影响身高疾患的男童,均可作为本次研究的调查对象。每个男童就是一个观察单位。观察单位的研究特征称作变量(variable),变量的观察结果称为变量值或观测值。如本次研究中,身高是研究特征,对每个男童身高测量的结果称为身高变量值,简称身高值或变量值。对变量取值的过程就是测量,而取值所需的标准称为测量尺度,它是获取正确、稳定、一致测量结果的条件。

(一) 常用的测量尺度

1. 名义尺度(nominal scale) 指变量的结果是按某事物属性分类来进行测量的,如性别变量:男、女;血型变量:A型、B型、O型、AB型,所用符号与属性一一对应,同一符号内各变量值的本质相同。

2. 顺序尺度(ordinal scale) 指变量值不但可以分类,而且各类之间有某种特征程度上的不同,可用数学上大于或小于来表达它们之间的关系,如治疗结果中的无效、好转、有效、显效、痊愈;工作面污染程度的轻度、中度和重度等。显然,评价尺度可改变,但它们的顺序或等级不变。有时用序数或秩次 1, 2, …, R 来表示,尽管形式上看序次间级差相等,但其相邻两序次间的数量级差却不一定相同,且难以精确量化。

3. 区间尺度(interval scale) 指用数量大小来度量某种特征,一般回答“是多少”。它不仅表示顺序,而且可把两次测量值之间的相差表达出来。无论在何时、何种情况下再做测量,两次测量值不变,其差值可表达为某一常数。依此尺度可用仪器、工具或其他定量方法测定个体某项指标的结果。其变量值 x 可以是实数轴上的一个连续区间,任意两个取值之间可有无穷多个数值,表现为连续型随机变量,如身高(cm)、体重(kg)、血压(kPa)、呼吸次数(次/分)等。变量值 x 也可以是整数范围内的随机变量,如育龄妇女生育子女数、龋齿数等。

4. 比数尺度(ratio scale) 指以比值、比率等来度量某种特征,如中性粒细胞占白细胞数总数的百分比、体重指数、某指标治疗后占治疗前百分比等。

(二) 变量分类

根据变量的测量结果不同,将变量分为定量变量和分类变量两类。

1. 定量变量(quantitative variable) 其变量值表现为大小不等的数值,一般带有度量衡单位,可按区间和比数尺度测得。应用时以计数或测量的多少将定量变量分为离散型(discrete)和连续型(continuous)两种:离散型定量变量是指测量值只取整数的情况,如育龄妇女生育的子女数、患龋齿数等;若测量值是区间内任意值,则称为连续型定量变量,如身高、体重等。由一组同质的定量变量值所组成的资料称为定量资料(quantitative data)。

2. 分类变量(categorical variable) 其变量值表现为事物的属性、特征或类别。若按属性特征分类,也称定性变量(qualitative variable),一般按名义尺度测得,如性别变量只具有相互对立的两种情况,即二分类变量(binary variable);若变量的观察结果表现为相互

对立的多种情况,称作多分类变量,如血型变量。将按照某事物属性分组,清点各组的观察单位数而得到的资料称为定性资料(qualitative data),如血型变量:A型18例、B型19例、O型35例、AB型23例。若按事物的等级或顺序分类,也称作顺序变量(ordinal variable)即有序变量或等级变量,常由顺序尺度测得。其顺序变量值表现为类别,其取值不但可以分类,且可表现出分类值程度或等级间的差别,如化验结果-、±、+、++、+++、++++。若将变量结果的等级或程度称作水平分组,清点各组观察单位数而得到的资料,称为等级资料(ordinal data)或半定量资料。其变量值之间不仅有类别的不同,且不同水平间也有顺序存在,但它又无法用精确的数值表示出相差的大小,如职工体检眼底动脉硬化级别检查结果:正常326例、轻度硬化18例、中度硬化13例、重度硬化3例。

根据不同的资料类型,可以选用不同的统计分析方法。有时分析条件不同,也可根据研究目的将资料进行由定量-等级-定性资料的转化。如某医师测得10名3岁儿童血红蛋白含量(g/L)结果如下:108、110、116、95、109、87、92、113、120、116,可知该资料为定量资料。若只考虑检测儿童有无贫血,可按临床参考值整理为无贫血者 ≥ 110 (g/L)5例,有贫血者 < 110 (g/L)5例,即转化为定性资料。欲分析贫血发生的严重程度,也可将该资料整理为无贫血者 ≥ 110 (g/L)5例,轻度贫血者90~110(g/L)4例,中度贫血者 ≤ 90 (g/L)1例,则又转化为等级资料。

三、总体与样本

总体(population)是根据研究目的确定的同质的所有观察单位某种变量值的集合,如某地所有15岁健康男孩身高值总体。据总体中观察单位数(N)是否可知,分为有限总体和无限总体。有限总体指总体中观察单位数是有限的或可知的总体,如某地区II期高血压病人总体,无论我们是否已了解其总体特征与观察单位总数,但总可以调查得到某特定时间或空间内的观察单位数,则可认为该总体属有限总体。无限总体是指总体中的观察单位数是无限的或不可知的总体,如空气中 SO_2 含量浓度总体。反映总体特征的指标称作参数(parameter)。根据研究目的从研究总体中随机抽取反映总体特征的部分观察单位,其实测值组成样本(sample)。样本中的观察单位数称作样本例数(n)或样本含量(sample size)。把从研究总体中按一定的概率规则,抽取部分观察单位进行研究的方法,称作随机抽样研究。随机抽样(random sampling)不是随意选择(purpose selection)。所谓随机是指研究总体中每个观察单位按其在总体中的分布情况,被抽到样本中的机会均等且互不影响,只有这样,才有可能保证抽到的样本有代表性,它是统计推断的基础,统计学中将反映样本特征的指标称作统计量(statistic),也称作参数估计值。

四、误差

任何研究总是期望对总体做出客观、可靠、真实的评价。但在实际工作中,调查结果可能会受到机遇和偏倚的干扰与影响而偏离真实情况。统计研究中,将实测值与真实值之差称为误差(error),据其产生的原因分为随机误差和非随机误差两类。

1. 随机误差 包括抽样误差和随机测量误差等。①抽样误差指随机抽样研究中,由于抽样而引起的样本指标与总体参数间的相差,其大小随样本不同而改变,它也是一个随机变

量；②随机测量误差 (random measurement error) 指对同一观察单位某项指标在同一条件下进行反复测量所产生的误差。即使严密控制研究条件，但由于一些偶然因素或就目前医疗技术水平尚无法控制的因素，也可使实测值与真实值之间产生一定的相差。实际工作中，即便测量仪器或方法多么成熟，测量误差也无法避免，但应控制在容许误差范围之内。

随机误差的出现是随机的，分布是有规律的，其值可正可负、可大可小，当观察次数足够大时，随机误差服从正态分布。经典统计方法就是依其特点和规律由样本对总体做出推断。

2. 非随机误差 非随机误差指所得资料偏离研究的真实情况，致使推断、预测出现的偏差，包括系统误差和过失误差。系统误差 (systematic error) 或偏倚 (bias) 可产生于设计人员、调查者或调查对象，也可由于设计考虑不当，资料搜集不准，汇总、计算有误等造成。一般带有倾向性，如有恒向、恒量、周期性或有特定的变化规律。其产生原因复杂，贯穿于研究全过程并对研究结果有影响，又很难用统计方法评价其影响的大小，必须依靠科学的研究设计，正确的资料搜集，完善的分析与计算，严谨的工作态度与作风，将其减小或控制在最小容许范围内。过失误差指由于工作人员责任心不强，检查核对制度不严，或故意修改等而造成的检查、记录、观察、录入数据错误等而产生。过失误差是错误，一般应杜绝。一旦发生应彻底纠正，否则有可能得出谬论。

五、频率与概率

频率 (frequency) 是指在相同条件下，进行有限 n 次重复试验，某随机事件 A 发生次数 X 与 n 次试验的比值，其值介于 $0 \sim 1$ 之间。如某地段产院 2000 年记录在册的出生人数 120 名，其中男性 70 名，得该产院 2000 年出生的男婴占出生总人数的频率为 58.33%，此结果称作该产院男婴的出生频率。频率是个变数，随样本变化而改变。

概率 (probability) 是描述某随机事件 A 发生可能性大小的度量，常记作 P ，可用小数或百分数表示，其值为 $0 < P < 1$ 。医学研究现象中，绝大多数属随机现象。若采用同一种方法治疗某病患者，只知道治疗转归可能有治愈、好转、无效、死亡四种结果之一，但对一个刚入院治疗的患者，治疗后究竟出现哪一种结果是不确定的。这里，每一种可能的结果都是一个随机事件，亦称偶然事件，简称事件。若将患者转归为“治愈”记作事件 A ，则治愈的概率记为 $P(A)$ ，简记为 P 。这是一个很有意义的、研究者颇为关心的未知数值 (概率是个定值)。若该事件必然发生， $P = 1$ ；若该事件不可能发生，则 $P = 0$ ，它们是随机事件的特例。当随机事件 A 发生的可能性 $0 < P < 1$ 时， P 值越大，该事件发生的可能性就越大。统计研究中，很多结论都是带有概率性的，习惯上把 $P \leq 0.05$ 的随机事件称作小概率事件，表明该随机事件发生的可能性很小，对于某一次试验，可认为该事件几乎不可能发生。由此引出“小概率原理”，它是统计推断的一条重要原理。

虽然随机事件 A 在一次试验中可能出现或不出现，但在多次重复试验中，它呈现出明显的统计规律性。假设在相同条件下，独立地进行 n 次重复试验，随 n 充分增大，频率摆动的幅度越来越小，则称该事件 A 为随机事件，其频率可作为概率的估计值。但在 n 较小时，频率波动性很大，用以估计概率是不可靠的。

总之，医学统计学是医学研究领域非常活跃的一门学科。它以其独特的思维方式深深地