



全国统计教材编审委员会“十一五”规划教材

统计学： 从数据到结论

第二版

★ 吴喜之 编著



中国统计出版社
China Statistics Press



全国统计教材编审委员会“十一五”规划教材

统计学： 从数据到结论

第二版

★ 吴喜之 编著



中国统计出版社
China Statistics Press

(京)新登字 041 号

图书在版编目(CIP)数据

统计学:从数据到结论/吴喜之编著. —2版.

—北京:中国统计出版社,2006.9

ISBN 7-5037-4996-2

I. 统…

II. 吴…

III. 统计学

IV. C8

中国版本图书馆 CIP 数据核字(2006)第 104408 号

统计学:从数据到结论

作 者/吴喜之

责任编辑/吕 军

装帧设计/艺编广告

出版发行/中国统计出版社

通信地址/北京市西城区月坛南街 57 号 邮政编码/100826

办公地址/北京市丰台区西三环南路甲 6 号

电 话/邮购(010)63376907 书店(010)68783172

印 刷/河北天普润印刷厂

经 销/新华书店

开 本/787×1092mm 1/18

字 数/290 千字

印 张/20.25

印 数/1—5000 册

版 别/2006 年 10 月第 2 版

版 次/2006 年 10 月第 1 次印刷

书 号/ISBN 7-5037-4996-2/C·2110

定 价/38.00 元

中国统计版图书,版权所有。侵权必究。

中国统计版图书,如有印装错误,本社发行部负责调换。

出版说明

“十一五”时期是继续深化教育改革,加强素质教育,努力建设有利于创新型科技人才生长的教育培训体系的关键时期。为了更好地培育统计创新型科技人才,适应统计教育培训的新形势,全国统计教材编审委员会制定了《“十一五”全国统计教材建设规划》(以下简称规划)。规划坚持“以人为本”的科学发展观,坚持统计教育与实践相结合,坚持统计教育同国际接轨,坚持培养创新型的统计人才的指导思想,编写符合国民经济发展需要和统计事业发展需要的统计教材。

这批教材是在深入分析统计教育形势和统计教材建设发展状况,总结多年来统计教材建设经验的基础上,本着以建设本科统计教材为主的方针,积极探索研究生层次的统计教材,力争使规划统计教材的编写做到层次分明,有针对性和实用性。建设精品教材,是编委会自成立以来就孜孜以求的目标。考虑到统计教材建设的实际情况,“十一五”期间,本科教材主要以修订为主,对以往规划统计教材中使用面广,得到广大教师和学生普遍认可的教材组织了修订。修订后的教材,淘汰了过时的内容和例子,增加了计算机操作和大量的案例,编写手法也做了一定的调整,在实用性、可操作性等方面有了较大的改进。

近年来,我国现代化建设快速发展,高等教育规模持续扩大,尤其是研究生教育规模的扩大,使得高等学校研究生统计教学工作面临着许多新情况、新问题,任务艰巨。因此,必须坚持科学发展观,在规模持续发展的同时,把提高研究生统计教学质量放在突出的位置,培养全面发展的创新型的统计人才。教材是统计教学的载体,建设高质量

的研究生层次的统计教材是统计教育发展的需要。因此，编委会在“十一五”期间对研究生的统计基础课教材做了些有益的探索。根据《规划》的要求，这批教材主要采取招标和邀请的方式组织有关院校的专家、学者编写。

值得特别提出的是，在这批教材中，有《非参数统计》、《概率论与数理统计》、《经济计量学教程》、《医学统计》、《应用时间序列分析》、《多元统计分析》、《统计学》7种教材入选国家教育部组织编写的“普通高等教育‘十一五’国家级规划教材”，更加充实和完善了“十一五”期间统计教材的建设。

为了便于教学和学习，这批教材里面包含了与之相配套的《学习指导与习题》，使得这批教材在编辑出版上形成了比较完整的体系。我们相信，这批教材的出版和发行，对于推动我国统计教育改革，加快我国统计教材体系和教材内容更新、改造的步伐，打造精品教材，都将起到积极的促进作用。

限于水平和经验，这批教材的编审、出版工作还会有缺点和不足，诚恳欢迎教材的使用单位、广大教师 and 同学们提出批评和建议。

全国统计教材编审委员会

2006年6月

前言

有什么在本书中等待着你们去发现,去探讨,去欣赏呢?当然不是数学公式和定理定义的堆砌,也不是和枯燥的公文报表相关的政府工作的培训。这是一门充满了哲学韵味的认识世界的学问。

不知读者们是否意识到,统计已经渗入到人们的社会、生活、工作等各个领域。每天新闻媒介报道的各个方面都离不开各种统计数据和各种分析与预测。人们可能对于这些统计内容觉得习以为常,也可能会有一些好奇或神秘感。由于国情不同,统计的地位与人们对统计的看法也不同。在发达国家,一般民众觉得统计学和数学类似,是一门高不可攀但极易找到满意工作的学问。在中国,又有一些人认为统计就是处理政府报表的职业。但自从中国向世界开放之后,越来越明确的一点是,没有什么学科或领域能够真正离开统计。

以应用为目标学习统计,究竟是为了什么?是为了流利地背诵一大堆定义、概念和抽象的名词和术语吗?是为了学习如何进行推导和证明一些复杂的定理和公式吗?这些问题不仅学生会思考,更重要的是统计教师要思考。本书的目的是希望读者在学习之后,能够知道实际中哪些是统计问题,最好能够自己解决一部分统计问题,即使不能解决也知道能够在哪里查到答案和向谁请教。知识固然重要,更重要的是通过学习获得解决和处理问题的能力。

学习并不总是一个令人生畏或至少成为某种负担的过程。人们学会走路、说话、骑车、下棋、打球等大都是在一种乐趣中进行的。为什么涉及到日常生活的每一个方面的统计就不能和看侦探小说那么引人入胜呢?其实任

何一门科学,都有其趣味性;而只有把科学研究当成游戏的人才会真正成为大师。这门课并不想使读者都成为统计学家,而仅仅想让读者如同学会使用电脑、手机,学会辩论、上网或讨价还价那样愉快地认识或理解在人生中无法躲开的统计。

本书由浅入深地把统计最基本和最有用的部分在这么一本不厚的教科书中完整地介绍给读者;而且让读者可以边学习,边着手用统计软件处理数据。篇幅大、语言啰嗦的教材对读者是个负担,不但浪费了资源,也抓不住要领。因此,作者力图惜墨如金;既节省篇幅,又要将该解释的全部说清。希望读者慢慢咀嚼,不必图快。

很少有一本统计教材包括像本书那么多的统计内容。我觉得,这些内容本来并不深奥;只是其貌似复杂的数学工具把它搞成阳春白雪;再加上强调数学推导的教学方式,使得统计显得高不可攀。本教材要还这些统计应用以其本来面目。使得统计变成人人都能够基本上理解和掌握的有用工具。多数使用计算机的人都不是计算机专业的;多数开汽车的都不会修汽车;但这对他们毫无妨碍。难道不会推导或背诵与统计有关的数学公式就不能应用统计这个工具了吗?

本书每一章的主要部分是用日常语言来引进和解释一些概念,如果可能就通过例子来说明。如果不涉及应用,这部分就足够了。涉及应用的各章后面的小结中,有一部分是说明如何通过统计软件来处理本章的数值例子;这会给多数想要自己动手分析数据的读者以方便。小结的最后还展示了与概念及计算有关的一些数学公式,使那些精力充沛的读者能更深刻地理解内容。这种安排使得本教材能够适用于各种不同水平、不同要求的读者群体。本教材不仅可供没有学过概率论和数理统计的非统计专业的本科生和研究生使用,也可以供统计专业的本科生作为理解统计本来含义的教材使用(以代替不能满足需要的“描述统计学”等类课程);它还可以为各领域的广大实际工作者作为应用各种统计方法的参考书。为读者可以使

用各种软件来进行分析,本书所涉及的所有电子版数据都有文本格式、SPSS 格式及部分 EXCEL 格式。

软件方面,本书主要使用 SPSS 和部分地区应用 Excel。我们不可能介绍所有软件,也不可能介绍某个软件的所有细节。经验证明,只要把某个方法在某个软件的基本选项指出,学生就可以通过自己的经验(最多借助于帮助)来得到所需要的结果。在课本中罗列使用软件时的出现的各种对话(选项)框的做法对本课程完全没有必要。

在前计算机时代,几乎所有的统计教科书都给出了各种与分布有关的表格。但随着计算机的普及,所有统计软件(无论是商业的还是免费的)都给出了和各种分布有关的各种函数,把人们从繁琐而又不精确的查表中解放出来。目前很多国外的统计教科书都不再提供了既占用篇幅又比较粗糙的分布表。本书不准备提供任何和分布有关的表格。本书第四章会介绍如何使用软件来进行与概率分布有关的计算。

这个教材的全部内容曾作为非统计专业硕士和博士的课程分别在北京大学光华管理学院及中国人民大学讲授过,受到普遍欢迎。实践证明,这个课程前 15 章的内容完全能够轻轻松松地在一个学期(每周三个学时)中全部讲完。一些热心而又好奇的非统计背景的人士也曾读过本教材的全部内容,没有任何理解上的问题。当然,根据不同的教学对象和需要,有些章节可以完全不讲或少讲。

本书前面的章节,是对统计基本概念的介绍。而后面的部分则是更有针对性的一些统计模型和方法。一般传统统计学的课程包括前六章,或最多前九章的内容;而第十章到第十四章一般属于多元统计分析的课程内容;第十五章一般属于时间序列课程包含的内容;第十六章一般属于非参数统计课程的内容;第十七章介绍了生存分析;第十八章对指数进行了必要的介绍。目前大多数流行的统计应用都已包含在本教材内。

本书的编写是在国家统计局教育中心的建议和鼓励下产生,并得到其大力支持。本书还受到北京大学、中国

人民大学以及各兄弟院校老师和学生的鼓励和帮助。中国统计出版社一直关心着本书的写作和出版。SPSS北京办事处的专家也一直积极对写作过程中出现的有关计算问题予以帮助。特别要指出的是敬爱的汪仁官老师又一次为我所写的统计教材进行了非常认真的审校,使我重新感受到做学生的幸福;中国统计界的老前辈茆诗松老师也热心地对本书提出了许多宝贵而又中肯的建议。他们的审校和建议使本书避免了许多错误和不妥之处。没有这些支持和帮助,本书是不可能面世的。谨在此对所有各方面表示衷心的感谢。

吴喜之
2003年6月

第二版说明

本书的第一版发行不到一年,已经作为参考书或教科书在许多学校使用。各个学校的师生对本书提出许多宝贵的意见,并且指出了很多错误和不妥之处。他们的支持和鼓励,促成了本书第二版的诞生。

和第一版相比较,第二版对许多内容完全重新写过,还进行了一些调整,同时加强了对概念和方法的解释,使得该书更加容易理解。第二版还对例题和习题都做了很多修订和增减。此外,还增加了一些内容;除了基于 SPSS 的操作和输出结果的分析 and 解释之外,还增加了 SAS 软件的使用;特别是增加了关于如何通过免费的,功能强大的,需要自己动手写程序的 R 软件来理解概念及处理数据。R 软件的代码公开及透明的优势是一些“黑匣子”式的傻瓜软件所无法比拟的。R 软件是使用 S 语言来编程的(和 S-plus 的编程语言一样);在其问世的不到 10 年的时间,已经成为国外统计研究生的首选软件。它有强大的网上支持系统。多数最新的统计计算方法,在进入商业软件之前,就已经以 R 语言的形式在 R 网站上免费提供。建议使用本书的师生最好也使用 R 语言。在掌握 R 软件之后,对其他统计方向的学习和研究都会有很大的帮助,甚至会有一种到了自由天地的感觉。在 R 软件的帮助中有完整的入门材料和各种命令的意义和使用例子;因此本书没有自己编写关于 R 语言的附录来增加篇幅。

本书选择 SPSS 软件选项和 SAS 软件语句(或选项)的原则是容易理解和掌握;当然,由于编者知识有限,对于有些方法,没有找到(因此也无法提供)最合适、最方便或者最新的软件选项或模块;希望读者提出建议,使得再版

时予以弥补。

由于使用软件比查表更加方便和可靠,可能会有人说,你自己都不查表,为什么要教学生去查表呢?的确,编者除了在最初等的统计课教学过程中曾经涉及到少数统计表之外,从来都是使用软件。“己所不欲,勿施于人”,本书不提供任何统计分布表。希望有条件的读者尽量使用计算机,而不去查表。实际上,如果没有计算机的支持,很难对有一定规模的数据在任何统计方向进行较深入的分析。

有许多人(比如各层管理人员)不一定要进行第一线的实际数据计算,但为了理解手中关于本单位及有关方面信息的意义,为了更好地进行明白的决策,他们必须理解各种统计推断结果的意义。对这些人,不一定要求能够熟练使用软件,更不需要理解数学推导,但必须明白各种统计概念和方法以及输出结果的意义。相信本书能够对帮助他们有所帮助。

作为教科书,本书内容对于每周两学时的课程似乎太多。但是,什么讲或者什么不讲可以根据学生的需要由老师自己安排。实际上,对于任何课程,最好是由任课教师来决定讲哪些内容以及如何讲。因为他们最了解他们所面对的学生。教科书编者的思维方式不见得和讲课老师的一致,而老师最好按照自己的理解来讲述。一本好的教科书,应该给教师以较大的余地和自由。

希望读者继续对本书予以宝贵的支持和批评指正。

吴喜之

2006年3月

第一章 一些基本概念

- 1.1 统计是什么? 1
- 1.2 现实中的随机性和规律性, 概率和机会 2
- 1.3 变量和数据 3
- 1.4 变量之间的关系 4
 - 1.4.1 定量变量间的关系 5
 - 1.4.2 定性变量间的关系 7
 - 1.4.3 定性和定量变量间的混和关系 8
- 1.5 统计、计算机与统计软件 8
- 1.6 小结 11
- 1.7 习题 11

第二章 数据的收集

- 2.1 数据是怎样得到的? 13
- 2.2 个体、总体和样本 13
- 2.3 收集数据时的误差 15
- 2.4 抽样调查和一些常用的方法 16
- 2.5 计算机中常用的数据形式 18
- 2.6 小结 20
- 2.7 习题 22

第三章 数据的描述

- 3.1 如何用图来表示数据? 23
 - 3.1.1 定量变量的图表示: 直方图、盒形图、茎叶图和散点图 23
 - 3.1.2 定性变量的图表示: 饼图和条形图 28
- 3.2 如何用少量数字来概括数据? 32
 - 3.2.1 数据的“位置” 33
 - 3.2.2 数据的“尺度” 34
 - 3.2.3 数据的标准得分 37
- 3.3 小结 38
 - 3.3.1 本章软件使用说明 38
 - 3.3.2 本章的概括和公式 40

3.4	习题	41
第四章 机会的度量:概率和分布		
4.1	得到概率的几种途径	42
4.2	概率的运算	43
4.3	变量的分布	47
4.3.1	离散随机变量的分布	47
4.3.2	二项分布	47
4.3.3	多项分布	50
4.3.4	Poisson 分布	50
4.3.5	超几何分布	51
4.3.6	连续随机变量的分布	52
4.3.7	正态分布	53
4.3.8	χ^2 -分布	56
4.3.9	t -分布	57
4.3.10	F -分布	59
4.3.11	均匀分布	60
4.3.12	累积分布函数	60
4.4	抽样分布、中心极限定理	61
4.5	用小概率事件进行判断	64
4.6	小结	64
4.6.1	本章例题和软件使用说明	64
4.6.2	本章的概括和公式	68
4.7	习题	74
第五章 简单统计推断:总体参数的估计		
5.1	用估计量估计总体参数	76
5.2	点估计	77
5.3	区间估计	78
5.3.1	一个正态总体均值 μ 的区间估计	79
5.3.2	两个正态总体均值之差 $\mu_1 - \mu_2$ 的区间估计	80
5.3.3	总体比例 (Bernoulli 试验成功概率) p 的区间估计	82

目 录

5.3.4	总体比例(Bernoulli 试验成功概率)之差 $p_1 - p_2$ 的区间估计	83
5.4	关于置信区间的注意点	83
5.5	小结	84
5.5.1	使用软件解本章例题的说明	84
5.5.2	本章的概括和公式	88
5.6	习题	91
第六章 简单统计推断:总体参数的假设检验		
6.1	假设检验的过程和逻辑	94
6.2	对于正态总体均值的检验	98
6.2.1	根据一个样本对其总体均值大小进行检验	98
6.2.2	根据来自两个总体的独立样本对其总体均值的检验	101
6.2.3	成对样本的问题	103
6.3	对于比例的检验	104
6.3.1	对于总体比例的检验	104
6.3.2	对于连续变量比例的检验	106
6.4	从一个例子说明“接受零假设”的说法不妥	108
6.5	小结	110
6.5.1	使用软件解本章例题的说明	110
6.5.2	本章的概括和公式	115
6.6	习题	117
第七章 相关和回归分析		
7.1	问题的提出	119
7.2	定量变量的相关	122
7.3	定量变量的线性回归分析	127
7.4	自变量中有定性变量的回归	132
7.5	Logistic 回归	134
7.6	小结	137
7.6.1	使用软件解本章例题的说明	137
7.6.2	本章的概括和公式	140

7.7	习题	142
第八章 列联表、χ^2 检验和对数线性模型		
8.1	列联表数据	143
8.2	二维列联表的检验	144
8.3	高维列联表和(多项分布)对数线性模型	145
8.4	Poisson 对数线性模型	149
8.5	小结	150
8.5.1	使用软件解本章例题的说明	150
8.5.2	本章的概括和公式	154
8.6	习题	156
第九章 方差分析		
9.1	方差分析(只考虑主效应,不考虑交互效应及协变量)	158
9.2	方差分析(考虑交互效应但不考虑协变量)	161
9.3	方差分析(考虑协变量)	163
9.4	小结	164
9.4.1	使用软件解本章例题的说明	164
9.4.2	本章的概括和公式	167
9.5	习题	169
第十章 寻找多个变量的代表:主成分分析和因子分析		
10.1	主成分分析	171
10.2	因子分析	176
10.3	因子分析和主成分分析的一些注意事项	179
10.4	小结	179
10.4.1	使用软件解本章例题的说明	179
10.4.2	本章的概括和公式	182
10.5	习题	184
第十一章 把对象分类:聚类分析		
11.1	如何度量距离远近?	186
11.2	事先要确定分多少类: k -均值聚类	187
11.3	事先不用确定分多少类:分层聚类	188

目 录

11.4	处理连续和分类变量混合的大数据集:两步聚类	190
11.5	聚类要注意的问题	191
11.6	小结	192
11.6.1	使用软件解本章例题的说明	192
11.6.2	本章的概括和公式	194
11.7	习题	195
第十二章 把对象归到已知的类中:判别分析		
12.1	几种判别分析方法	196
12.2	判别分析要注意什么	205
12.3	小结	205
12.3.1	使用软件解本章例题的说明	205
12.3.2	本章的概括和公式	207
12.4	习题	209
第十三章 两组变量之间的相关:典型相关分析		
13.1	两组变量的相关问题	210
13.2	典型相关分析	211
13.3	小结	215
13.3.1	使用软件解本章例题的说明	215
13.3.2	本章的概括和公式	216
13.4	习题	217
第十四章 行变量和列变量的关系:对应分析		
14.1	对应分析方法	219
14.2	小结	222
14.2.1	使用软件解本章例题的说明	222
14.2.2	本章的概括和公式	223
14.3	习题	225
第十五章 随时间变化的对象:时间序列分析		
15.1	时间序列的组成部分	227
15.2	指数平滑	228
15.3	Box-Jenkins 方法:ARIMA 模型	231

15.3.1	ARIMA 模型介绍	231
15.3.2	ARMA 模型的识别和估计	232
15.3.3	用 ARIMA 模型拟合例 15.1	236
15.3.4	用 ARIMA 模型拟合带有独立变量的时间序列	239
15.4	小结	241
15.4.1	使用软件解本章例题的说明	241
15.4.2	本章的概括和公式	245
15.5	习题	248
第十六章 总体分布未知时的检验:非参数检验方法		
16.1	关于非参数检验的一些常识	249
16.2	单样本检验	251
16.2.1	关于单样本中位数(α -分位数)的符号检验	251
16.2.2	关于单样本位置的 Wilcoxon 符号秩检验	254
16.2.3	关于单样本的 Kolmogorov-Smirnov 检验	257
16.2.4	关于随机性的游程检验(runs test)	263
16.3	两独立样本检验	266
16.3.1	比较两独立总体中位数的非参数检验: Wilcoxon (Mann-Whitney)秩和检验	266
16.3.2	关于两样本分布的 Kolmogorov-Smirnov 检验	269
16.3.3	两样本 Wald-Wolfowitz 游程检验	271
16.4	关于多个独立样本的检验	272
16.4.1	Kruskal-Wallis 关于多个样本的秩和检验	272
16.4.2	Jonckheere-Terpstra 关于多个样本的秩检验	274
16.4.3	Brown-Mood 中位数检验	276
16.5	多个相关样本的检验	278