

统计学漫话

陈希孺 苏 淳 编著

中国科学技术大学出版社

内 容 简 介

本书以漫谈的方式,用大量通俗的例子引入和说明了统计学的基本概念、基本思想和基本方法,使读者对统计学的全貌有所了解,并且学会用统计观点去看待现实世界中的许多事物。

本书内容深入浅出、通俗易懂,可供具有中等文化程度的读者、从事实际工作的统计工作者阅读,也可供大专院校统计专业的师生参考。

图书在版编目(CIP)数据

统计学漫话/陈希孺,苏淳编著.—2版.—合肥:中国科学技术大学出版社,2016.1

(数林外传系列:跟大学名师学中学数学)

ISBN 978-7-312-03751-1

I.统… II.①陈…②苏… III.统计学—青少年读物 IV.C8-49

中国版本图书馆CIP数据核字(2015)第302661号

出版 中国科学技术大学出版社

安徽省合肥市金寨路96号,230026

<http://press.ustc.edu.cn>

<http://shop109383220.taobao.com>

印刷 安徽省瑞隆印务有限公司

发行 中国科学技术大学出版社

经销 全国新华书店

开本 880 mm×1230 mm 1/32

印张 8.5

字数 235千

版次 1981年1月第1版 2016年1月第2版

印次 2016年1月第2次印刷

定价 28.00元

前 言

本书试图以漫谈的方式,用通俗的语言,向具有中等文化程度的读者介绍统计学的基本思想和方法.统计学在许多领域中都有广泛的应用,其重要性无需在此处强调.作者想表达一个想法:从一定程度上看,统计学的初步知识已构成一个人文化修养中必要的部分.这是因为,在现实世界中许多事物只有用正确的统计观点去看待,才能得到恰当的理解,即使在日常生活中碰到的一些事情也不例外.

目前,由国内专家编写的统计学著作已出版不少,但对具有中等文化程度的读者来说,它们大都过于专深,需要用到高等数学和概率论的知识.另一个问题是,统计学不同于纯粹数学,不能单纯从公式的数学论证中去正确理解它,而必须对其基本概念和问题提法的实际背景、方法的思想、结果的解释、使用统计方法应注意的种种问题等,作深入的思考才行.正因为如此,我们在本书的编写过程中,对以上提到的诸方面问题都作了详细论述.

本书的编写方式是,在介绍每一个主题时,先从大家都理解的一些事物入手,经过分析,提出一些想法和问题,由此逐步展开,从而引入明确的概念、方法和理论.在每一步中,我们都尽量用一些通俗的比喻,以便把一些艰深的概念形象地表达出来,但又不失科学的严谨性.这是作者定下的目标,但达到了多少,只能由广大读者来评判.

为适应具有中等文化程度的广大读者的需要,本书力图避免高深的数学知识.当然,讲统计学不能只空谈思想而不涉及具体方法.相反,只有通过介绍具体方法,才能把观点和思想讲清楚.因此,书中介绍了一些重要的统计方法,详细交代了方法的步骤及使用时要注意的地方.读本书自然不能像读一本通俗小说那么轻松,所以要求读者必须深入地思考.本书虽然力图避免枯燥,但由于作者写作能力有限,未免

力不从心,不到之处望读者谅解。

除上述读者外,我们还希望本书对学习统计学的大学生,以及具有一定统计知识和实践经验的应用工作者能多少有点用。在统计学的教学中,由于学时限制等原因,重点都放在介绍方法内容及其数学论证上,而对上面提到的若干问题则注意较少,本书希望能起一点拾遗补缺的作用。

我们还希望本书能对作者的广大同行——统计学教师有一点用处。也许作者的同行们都有这样的体会,这门课不好教,讲起来总觉得不容易使学生信服,而其难点又不在数学论证上。本书有些内容,包含了作者在教学实践中关于这些难点的若干思考和看法,虽不一定确切,但也许能引起广大同行注意这些问题。书中的叙述肯定会有不当之处,恳请广大读者不吝赐教。

在编写本书的过程中,项可风同志对第3章的写作提供了不少帮助。吴启光同志审阅了原稿,良多补益。在写作的最后阶段,陆传荣、林正炎同志安排了良好的条件,使作者得以如期完成。搁笔之际,感触颇多,除致以衷心的感谢外,特书于此,以志不忘。

作 者

目 次

前言	(I)
第 1 章 什么是统计学	(1)
1.1 统计方法和统计学	(1)
1.2 通过事物的外在数量表现考察事物的规律性	(3)
1.3 由部分推断整体、总体和样本	(5)
1.4 统计性推断的错误和误差	(9)
1.5 学一点统计学	(11)
第 2 章 获取数据(I)——抽样调查	(14)
2.1 抽样调查的意义	(14)
2.2 要注意的问题	(16)
2.3 简单随机抽样、随机数表	(18)
2.4 集团抽样	(22)
2.5 分层抽样	(25)
2.6 随机化的重要性再议	(27)
第 3 章 获取数据(II)——试验设计	(31)
3.1 引言	(31)
3.2 完全随机化设计	(32)
3.3 随机区组设计	(37)
3.4 平衡不完全随机区组设计	(41)
3.5 拉丁方设计	(45)
3.6 多因子试验	(50)

3.7	拉丁方用于多因子试验	(53)
3.8	正交拉丁方	(55)
3.9	正交表	(59)
第 4 章	平均值与比率的精度	(68)
4.1	平均值的代表性问题	(68)
4.2	总体方差	(70)
4.3	样本均值的方差	(73)
4.4	方差的估计,样本方差	(80)
4.5	均值之差的估计	(83)
附录	式(4.8)和式(4.11)的证明	(85)
第 5 章	分布与区间估计	(89)
5.1	方差的局限性	(89)
5.2	分布的概念	(91)
5.3	分布的列表形式	(93)
5.4	直方图与密度函数	(95)
5.5	标准正态分布	(100)
5.6	正态分布均值的区间估计(方差已知)	(105)
5.7	t 区间估计	(111)
5.8	大样本情况	(116)
第 6 章	概率初步知识	(120)
6.1	什么是概率	(120)
6.2	事件	(121)
6.3	古典概率	(122)
6.4	频率与统计定义	(129)
6.5	主观概率	(134)
6.6	随机变量	(135)

6.7	概率分布	(138)
6.8	均值和方差	(144)
6.9	均值的大数定律	(148)
第 7 章	假设检验	(150)
7.1	原假设和对立假设	(150)
7.2	拟合优度	(154)
7.3	检验的水平	(156)
7.4	两类错误	(162)
7.5	卡·皮尔逊的 χ^2 检验	(165)
7.6	无关联性的检验	(174)
7.7	u 检验	(178)
7.8	一样本 t 检验	(184)
7.9	与区间估计的关系	(187)
7.10	两样本 t 检验	(188)
7.11	非参数统计方法	(190)
第 8 章	相关与回归	(193)
8.1	事物的联系	(193)
8.2	相关系数	(195)
8.3	相关系数的估计和检验	(199)
8.4	偏相关和其他	(201)
8.5	“平均相关”的谬误	(203)
8.6	回归与回归方程	(205)
8.7	回归方程的估计(最小二乘法)	(208)
8.8	残差、残差平方和与偏相关	(211)
8.9	回归用于估计条件平均值	(214)
8.10	回归名称的由来	(217)
8.11	几点注意事项	(219)

附录	式(8.16)和式(8.26)的证明	(222)
第9章	方差分析法	(226)
9.1	基本思想	(226)
9.2	完全随机化设计	(229)
9.3	随机区组设计	(232)
9.4	对比试验	(237)
9.5	拉丁方设计	(240)
9.6	正交表设计	(244)
9.7	<i>F</i> 检验	(247)
9.8	估计问题和最优处理的选定	(250)
9.9	交互效应	(252)
附表	(256)

第 1 章 什么是统计学

1.1 统计方法和统计学

统计这个字眼大概对于一般人都不陌生,因为恐怕谁都听到过这类说法:明天组织去游山,需要把去的人统计一下;工厂年终发奖金,要统计一下发 1 000 元以上的有几人,500~1 000 元的有几人,等等.可见“统计”成了一个常用词.总而言之,不少单位设有统计员,国家设有各级统计机构,以收集关于经济、人口和社会等方面的资料——称为**统计数据**.对这些数据要进行整理和分析的工作,以作出种种结论和预测,所用的方法就是**统计方法**.研究这种方法的学问,就叫作**统计学**.统计学在我国也常称为**数理统计学**.其实,后者应是前者的数学理论基础部分,又称**理论统计学**.它是关于统计方法如何建立及其正确性和有效性的数学论证.概括地说,统计方法是有关收集和取得数据资料,并对之进行整理、分析,以对所研究的问题作出一定结论的那些方法,这样说基本上是正确的,但必须附加大量的补充解释,需要指出统计方法的特点何在,正是这些特点划清了统计方法和其他方法的界线.

如果使用高深的数学概念,可以用很少几行字把这个问题说清楚,但这对本书的读者未必有益.在本书的整个叙述中,我们将力求回避抽象的数学概念,而尽量用实例、大家都熟悉的事物和形象化的比喻来说明问题,因而在语言上难免失之累赘,这是要请读者见谅的.现在先列举几个例子.

例 1.1 一个生产灯泡的工厂,以往一直采用某种工艺.现在厂里的技术人员对此提出一些改进措施,以期能改善产品质量.为了验证

这个想法,取在新、老工艺下生产的灯泡各若干个去使用,记下每个参加试验的灯泡的寿命(即从开始使用到损坏所经历的时间).所得数据可以这样处理:算出使用老工艺生产的灯泡的平均寿命,例如为 420 小时;对使用新工艺生产的灯泡也这样做,结果为 440 小时.于是作出结论:新工艺的确有助于改善质量,使用时间约可增加 20 小时.

例 1.2 要调查某县的个体农户在某年每户使用化肥的平均数量.全面的普查将涉及数以十万计的农户,这是人力、物力和时间所不允许的.于是从该县的农户中抽选出若干户,比如 400 户,调查出这些户共用化肥 32 000 千克,户均 80 千克,以这个数字作为全县户均化肥用量的估计值.

例 1.3 一物件的质量 a 未知,放到天平上去称.由于天平有些误差,而对结果又有很高的精度要求,于是觉得称一次不够,就把它重复称三次,分别得出结果 2.45, 2.46 和 2.41(克),以其平均数 2.44(克)作为 a 的估计.

例 1.4 为探索吸烟与患肺癌二者之间是否有关联,调查一大批人,按是否吸烟和是否患肺癌分成四类:不吸烟也不患肺癌的、吸烟且患肺癌的、不吸烟而患肺癌的、吸烟而不患肺癌的.根据这些数据,用一定的(统计)方法,作出像“吸烟与患肺癌二者有显著关联”之类的结论(该方法的性质复杂些,待到第 7 章再作解释).一个易于理解的处理方法如下:算出在抽选出的这批人中,不吸烟者患肺癌和吸烟者患肺癌的比率,比方说分别为十万分之三与十万分之十二,于是作出结论说,吸烟者患肺癌的危险性是不吸烟者患肺癌的 4 倍.这类报道及其他性质类似的医学报道,常见于各种报刊.

以上是几个用统计方法研究问题的范例,其中都有获取、整理和分析数据的工作.由此可以总结出统计方法的哪些特点呢?这就是我们在以下几节中要讨论的问题.应当交代清楚的是,由于我们力图少用抽象的数学论证,以下几节的说明还是不全面的,有的问题将在本书以后的叙述中再行补充.

1.2 通过事物的外在数量表现 考察事物的规律性

统计方法只是从事物外在的数量表现上去研究问题,不涉及事物的质的规定性.通俗些说,统计方法可能告诉你,从试验或观察结果来看如何如何,而不能回答为什么会如何如何.如在例 1.1 中,试验结果有可能显示新工艺有助于改善质量,但其原因何在?也可能一目了然,也可能涉及专门学科领域中深奥的道理.在例 1.4 中,虽然许多统计资料都表明吸烟与患肺癌之间有关联,就是说吸烟的人看来更倾向于易得肺癌,但这种结论目前看来仍只能算是一种统计规律性——由表面上的数量关系而归纳出来的规律性.因为不仅吸烟何以引发肺癌的机制在目前尚未确切研究清楚,甚至这二者之间表面上的联系是否真正反映一种因果关系,在学者中也有分歧.有的学者认为,这二者表面上的关联,可能不过是由于它们受同一遗传基因的控制,其作用使那些易于染上吸烟嗜好的人同时也倾向于易患肺癌.若这种看法被证明为确实的,则戒烟既不会减少也不会增加患肺癌的危险性.更多的学者则认为,二者的联系是因果性的,尽管其机制目前没有被充分弄明白.

这点值得作为统计方法的一个特点(或称“性质”也可以)提出来,是因为它划清了统计学和其他专门学科的界线,如在遗传学、医学……中用了不少统计方法,但统计学绝不能代替这些专门学科,而只是有助于它们,可以说只是一个辅助性的工具而已.了解这一点,就不致对统计方法和统计学者提出过高的期望,以为他们掌握的方法是万能的,可以在许多专门领域中单枪匹马地解决种种实际问题.一个从事于实际应用问题的统计工作者,其知识面愈广,就愈易与种种专门学科领域的人员取得共同语言,因而也就愈能对他们的工作提供一些帮助.

我们说统计方法只是一个辅助性的工具,仅是就以下一点而言:

单纯的表面上的数量关系是否反映事物的本质,该本质究竟如何,必须依靠专门学科的研究才能下定论.这个提法不能理解为,统计方法的作用完全是被动性的,恰恰相反,事物的本质,其根本规律性的东西,一般都是隐藏得很深,它不时地在一些场合下有所表现.学者们注意和收集了这些资料,初看起来杂乱无章,而他如果具有一些统计的眼光,就有可能透过这些纷繁的数据而发现某种规律性的东西.这诚然还是表面上的,但可以作为专门研究的出发点,好比在一个刑事案件中,罪犯往往隐藏得很深,但他总会多少留下一些痕迹,受过训练而有经验的侦察人员,能据此对案犯作案的动机和过程提出一些设想,以作为破案工作的起点.所以我们说,统计方法在研究自然界和人类社会的规律性方面,是起着积极的、主动的作用的.科学史上有大量这样的例子.下面我们以遗传学上的一项伟大发现为例,在这问题上再多说几句.

奥地利生物学家孟德尔在 1865 年发表了一篇文章,其中事实上提出了基因的学说(“基因”一词是英国学者贝特松在 1909 年提出的),从而奠定了现代遗传学的基础.他这项伟大发现的过程很足以说明统计方法在科学研究中所起的作用.孟德尔是用豌豆做试验的,这种豌豆的果实有黄、绿两种颜色.孟德尔分别培养了一个黄色的纯系和一个绿色的纯系,其每一代所结的豆子分别全部是黄色的和绿色的,孟德尔然后将这两个纯系进行杂交,发现这种黄-绿杂交品种所结的豆子全部都是黄色的,与黄色纯系无不同.但在将这种杂交体再进行一次杂交而产生第二代时,孟德尔发现某些这种“第二代杂交豆子”呈黄色,而另一些呈绿色,其数目的比例大致接近 3 : 1.孟德尔把他的试验重复了多次,每次都得到类似的结果,到这里为止,所得到的还只是一个表面上的统计规律性,但这个表面上的规律性启发了孟德尔去发展一种理论,以解释这个现象.他假定存在一种现在称之为基因的实体以控制豆子的颜色.该实体有两个状态 y (黄)和 g (绿),共组成四种配合: yy , yg , gy , gg (称为基因型).前三种配合使豆子呈黄色,而第四种配合使豆子呈绿色(在遗传学上称 y 为显性的,而 g 为隐性的).

根据这个学说,孟德尔就容易给他的试验结果以圆满的解释:黄色纯系和绿色纯系的基因型分别是 yy 和 gg . 杂交第一代种子的基因型则只有一个可能性,即 yg , 而根据 y 为显性的假设,具有这个基因型的豆子呈黄色,在外观上与 yy 无异;但若对 yg 再进行杂交,则呈现四种可能性,即 yy , yg , gy 和 gg . 前三种呈黄色而后一种呈绿色,这解释了杂交第二代豆子中颜色黄绿之比近似为 $3:1$ 的观察结果. 为什么只是近似 $3:1$ 而非严格 $3:1$ 呢? 这好比有两个极大的盒子,每个盒子中放入为数极多的黑、白两种颜色的球,每盒中两种颜色的球的个数相同,然后你每次从两盒中各抽出一球配成一对,这样重复多次,得出许多个(个数很大,但比起盒中所有球的个数则很小)对子,在这许多个对子中,“黑-黑”对子的个数只是接近全部对子数的 $1/4$, 而不见得恰好是 $1/4$. 自然,孟德尔理论的伟大意义不是在于它给这个特殊的观察结果提供了理论解释,而是在于,用这个理论(当然是经过大大发展了的)可以解释生物体的很多遗传现象,从而形成了遗传学中的基因学派,到 20 世纪 50 年代,基因的存在已经在分子水平上获得了证实. 关于统计方法在建立孟德尔理论的过程中所起的作用,我们还可以补充一点:在从分子的水平上观察到基因的存在且完全证实这个理论以前,曾经用统计方法对依这个理论推出的大量结论进行过检验,检验的结果都证实了这个理论与观察结果符合(这个问题在第 7 章中还要讨论). 这本身就是统计方法在科学上的一项重要应用——用于客观地评价某项理论上的结论是否与观察结果相符,以作为该理论是否站得住脚的印证.

1.3 由部分推断整体、总体和样本

统计方法都具有部分推断整体这个性质. 如在例 1.2 中,整体就是全县的所有个体农户——由于我们只关心其化肥使用量,也可以说整体是由该县所有个体农户每户的化肥用量组成. 若该县有 10 万个个体农户,则整体包含 10 万个数字,所要考察的问题(化肥用量户均值)是

关系到这个整体,而不是关于其中某些户的.部分就是被抽选出的那些农户(也可以说,是抽选出的那些农户化肥用量的全部数据).我们的方法是算出这“部分”的平均值.如果停留在此处,则所得结果还只与这个“部分”有关.若再往前跨一步,而声称“以这部分的平均去估计整体的平均”,则我们工作的意义越出了这部分之外而达到整体.这一步工作称为统计推断,它是关键的一步,构成统计方法的一个重要特点.举一反三,读者不难按这种方式,对其他三例作类似的分析.

为什么要把这个强调为统计方法的一大特点呢?原因有二:一是它把统计方法与其他数学方法区别开来;二是它把大量日常工作以至生活中与数字打交道的工作和统计方法区别开来.

先说第一点.统计方法要用到许多数学工具,尽管在学者中对统计学是否可算作数学的一个分支存在分歧,但对于统计方法中使用大量数学工具、统计方法的原理依靠高深数学的论证这些事实,却不容抹杀.那么,相对其他数学方法而言,统计方法的特征何在?关键就在“部分推断整体”这一点上.举一个极简单的例子:有两块矩形木板 A 和 B ,要比较其面积谁大,大多少?量得 A 的长、宽分别为 1.52 米和 1.425 米, B 的为 1.79 米和 1.21 米.如果测量绝对准确,则根据“矩形面积 = 长 \times 宽”这个公认的数学公式,即算出 A 的面积比 B 大,大 0.0001 米².这个问题用数学方法解决了,但不是用的统计方法,因为你已掌握了与问题有关的全部资料,不存在“部分推断整体”的因素.然而,你可能觉得测量有一定误差,而二者面积测量值之差(0.0001)又很小,只测一次就下结论未必可靠.为了增加可靠性,你把 A 和 B 的长、宽各测量 100 次,算出 A 和 B 面积的 100 次测量结果的平均值之差,以此为准来定何者面积大,大多少,这就是一个典型的统计方法.为什么?就在于你只掌握了与问题有关的部分信息而非全体.因为,你既可以测量 100 次,又何尝不可再测量 200 次、300 次……直观上告诉我们,测量次数愈多,平均数愈可靠.理论上说,要“绝对”可靠,只有测量无穷多次求平均.设想你真这么做了(当然事实上不可能),你就掌握了问题的全部信息.因此在本问题中,“无穷次测量结果的全部记

录”构成一个整体,你实际做了的那 100 次测量只是这个整体中的一部分.这仍是一个由部分推断整体的格局.

再说第二点.若在例 1.2 中问题不是一个县而是一个村,则我们大可不必从其中挑出一部分农户,而可以逐户搞清楚,算出其平均值就可以了.从严格的统计学观点说,这里谈不上用到了什么统计方法,只是例行公事地作一些加法和除法的运算,就得到确实的结果(统计方法由于只用到部分资料,结果不见得确实,即有误差,这误差可能有多大,是统计学的任务,这构成统计方法的一个特点,将在下文论述).这类工作很多,虽然在习惯上也无妨承认它们用到统计方法,但这个差别却不可忽视.你把今年家里每个月所花的伙食费都记下来,到年终一平均,就得到你家今年每人每月的平均伙食费.这一切都一目了然,谈不上统计方法这个大文章.可是如果你在以往五年中都作了这个计算,分析所得结果,总结出这五年伙食费以年率 15% 的幅度增加,并进而推断在今后三年内,你家及情况类似的人家,其每人每月平均伙食支出仍将以这个幅度上升,则这整个过程就可以看作是统计方法的使用.因为所作出的结论超出了你掌握的数据资料的范围,而构成一项统计推断.

在这一节的最后,我们介绍几个在统计学上常用的专门术语,并作些补充说明.

统计方法有“部分推断整体”的特征.这个整体在统计学上常称为**总体**,也有叫**母体**的.总体依所研究的问题而定.前面已指出,如在例 1.2 中,总体由该县的全部个体农户组成.总体里的每一分子称为一个**个体**或**单元**.在例 1.2 中,每一农户构成一个个体或单元.如果你要对小学生中的近视眼比率作调查,则随着你研究规模的不同,总体可以是全国所有的小学生,或是人口在 20 万以上的城市中所有的小学生,或者是指定城市中的全部小学生.总体取得不同,研究结果适用的范围当然也就有别.

从总体中抽选出的那部分个体,统计学上称为**样本**,也有叫**子样**的.如在例 1.2 中,抽选出的那 400 户就构成样本.有时也把样本中单

个或部分个体称为样本.样本中所含的个体数,在统计上称为**样本大小**,也有叫**样本容量**的.在例 1.2 中,样本的大小为 400.从总体中抽选出样本的过程叫**抽样**,也有叫**取样的**.不论在任何问题中,由于与问题有关的往往只是个体的某项(或某几项)指标,也可把个体的指标值就说成是该个体.这时,不论总体和样本,都是由一些数字构成(这在本节开头已就例 1.2 说明过).这个看法突出了与问题有关的数量方面,便于在数学上作统一的处理.

把某一个体算作是所研究问题的总体之内,有一个明显的前提,即该个体的指标值(这项指标是问题中所关心的)必须是可知的,必要时得加以明确的规定.如在例 1.4 中,每一个个体(一个人)的指标值就是他属于四类中的哪一类(见例 1.4 开头的说明),若你没有必要的医学设备,就无法检定一个人是否患肺癌,则每一个人所属的类别就无法确定,研究工作也就无从下手了.有时,一个个体的指标值如何,需要根据一定的规则才能定下来.比方说,某甲在今日确诊患有肺癌,但他是一星期前才开始吸烟的,难于设想,这一段吸烟史与某甲患此病有关,故某甲的指标值似以定为“不吸烟,患肺癌”为合理.故这里存在一个“怎样的人算作吸烟者”的问题,这不见得是很容易解决的问题.在例 1.2 中,要求该县每户化肥用量都可知.比方说,以该户户主所报数目为准,即使有些误差也不计较,当然,若所报数目很不准,或对其误差性质有些了解,可以在问题中把这种误差考虑进来,这时总体和个体就不能像原来那样子,而是大大复杂化了,这一点在此不能细论.

在有些问题中,总体是由一些看得见、摸得着的个体构成的.例 1.2 是一个典型例子.在另一些问题中则不然,它只存在于我们的想象中.例 1.3 是这类情况的一个代表,有人可能会认为,在本例中,总体就是由这个(其质量 a 未知)物体构成,其实不然,因为在本问题中,我们关心的指标是物体的质量,而它是未知的,不符合上述“总体中的个体指标值可知”的要求.也许还有人会问:此物体的质量 a 虽未知,但可通过天平称量去了解,天平虽有误差,但在例 1.2 中,一个农户的化肥用量也不见得能准确说出,也可以有误差,二者并无不同.其实在例

1.2 中,我们要估计的对象并不是一个个农户的化肥用量,而是其整体的平均.而在例 1.3 中, a 本身就是估计的对象,我们之所以要对该物件做多次重复称量,原因正在于一次称量的精度达不到要求,若像例 1.2 那样考虑,以一个近似值去代表它,就违反了我们的初衷.

那么在例 1.3 中总体应如何定? 我们已把该物体称了三次,若还嫌精度不够,则可以再称若干次,原则上你可以无限制地称下去,每称一次就有一个数值.在想象中,我们可以有无限个数值(或者说,在想象中有无限次称量的动作,每个动作有一指标值,即该次的称量结果),它构成问题的总体.已做的那三次称量是这个总体的一部分,它构成问题中的样本.在此,抽样的过程就是把总体中的个体由想象中存在转化为具体存在的过程.

还有一个很常用的专门统计术语,叫**统计量**.它是指从样本算出的量.如在例 1.2 中,样本中包含 400 个农户,有 400 个指标值(即各户的化肥用量),其平均值 80(千克)是从这 400 个数据算出的量,它是一个统计量.在例 1.4 中,若我们调查了 8 000 个人的情况,则原始资料(样本)有 8 000 个,比率 0.000 03 和 0.000 12 是从它们计算出来的,故都是统计量.在一个问题中考虑怎样的统计量,当然要取决于所要解决的问题的性质.所提到的这两个统计量——样本值的平均与样本中带有某种属性的个体的比率,是应用上最重要的,以后还将介绍其他重要的统计量.

1.4 统计性推断的错误和误差

统计方法的另一个特点是,经由统计方法得出的结论(即统计推断),可能有错误或误差.如在例 1.1 中,由于所作结论是基于少数灯泡的试验结果,而一个好的工艺偶尔也可能生产出不好的灯泡,它就可能出错,即与以后生产实践的结果不符.在例 1.2 中,估计该县个体农户均使用化肥 80 千克,但实际也可能是 90 千克,或 70 千克等等,这样结论就包含误差.读者可能会问,统计方法作为一种认识自然和

社会的科学方法,而其结论却不能保证正确,这如何理解?事情很显然,这取决于问题的条件,在例 1.2 中,如果只允许你考察 400 户(在这个条件限制下),则不论你用什么方法,都不能作出保证无误的估计,因此,这并非统计方法的缺陷,恰恰相反,正确地使用统计方法可以最大限度地减少可能的错误或误差,并对犯错误的可能性以及误差的幅度,提供有用的估计.这些正是数理统计学这门学科讨论的主题.

统计方法作出的结论之所以可能有错误或误差,根源在于数据(样本)有误差.统计学上把数据的误差分成两大类.一类叫**系统误差**,它是由试验安排和观察工作的组织上的失误而来的.如在例 1.1 中,若在用新工艺生产灯泡时,配以优质的原材料和技术熟练的工人,而在用老工艺生产时则相反,则试验结果将包含有利于新工艺的系统误差.在例 1.2 中,如派一个懒于跋山涉水而又不负责任的人去作调查,他可能会在有交通工具可利用的河流、公路沿线调查一些农户了事,由于这些地方经济较发达,化肥用量一般也大些,这种做法可能对平均值的估计产生严重偏高的误差.尽量消除系统误差,是在安排试验和观察以取得数据时,要注意的中心问题,统计学对此有大量的讨论,本书第 2 章和第 3 章就是关于这个问题的讨论.当然,系统误差的辨识和消除,也取决于人们的知识水平与客观条件;可能有系统误差存在,但人们还没有认识到,或者虽已认识到,但需要复杂的设备才能控制.

另一类误差在统计学上叫**随机误差**,也有叫**偶然误差**的.“随机”的意思就是“随机会而定”.这种误差的产生,不是由于人们在安排试验或观察时有意的偏向或重大的失误,而是由于种种人所不能控制甚至不能察觉的偶然性因素的影响.如在例 1.3 中,尽管操作者在称量前已把天平充分校准,但称量结果还是不能摆脱大量的外界环境因素,以及操作者主观因素的影响,它们以一种偶然的方式起作用.例如,在称量时一个人偶尔从附近走过,偶尔刮一阵风把窗户震动了,室内温度偶尔有所升降,操作者瞬间心情上的波动等,这些将导致称量结果的误差.另一种随机误差以例 1.2 为代表.在例 1.2 中,即使我们在抽选农户时避免了前述的主观性,而让全县每一农户都有同等机会被选入

(这就在我们的主观意图上避免了偏向某些特定农户,也就避免了系统误差),但因为我们抽取的400户只占全县农户的很小一部分,抽取的这400户不能完全代表全县的农户,也就是说有误差存在.这种误差不是由某种偏向或失误而产生的,它是随机性的误差.

我们将以上所述总结为两点,即统计学所研究的是:① 如何安排试验和观察试验,以消除或尽量减少系统误差,使数据只受到随机误差的影响;② 如何去整理和分析这样的数据,以作出一定的结论即统计推断,对这种推断出错的可能性或误差的大小作出估计.有一个问题:既然数据带有随机误差,则基于它们所作的结论,其错误或误差也是随机会而定的,是偶然性的,那怎么可能得到它们的规律呢?这是因为,事物中包含的偶然性,在其个别实现中诚然是无秩序的、无规律的,但在该事物的大量实现中则可能呈现出某种秩序,即规律性的东西,好比投掷一个均匀骰子,它出现1~6点中的哪个点,完全是偶然的;但若将这个骰子投掷很多很多次,则会发现,每个点出现次数的比率都接近 $1/6$,投掷次数愈多,接近程度愈大,这就是在大量次投掷中体现出的秩序.在数学上,人们引进**概率**的概念来描述这个意思,称“在投掷一个均匀骰子时,每个点出现的概率都是 $1/6$ ”.对概率的研究形成数学的一个分支,叫**概率论**.它是统计方法的主要理论基础.因为在用统计方法时涉及大量数据,单个数据受到随机性的影响,显不出什么规律性,但是,与上述掷骰子的例子相似,在大量数据的集体中可能显示出规律性的东西,用概率论作工具使我们能捕捉、研究和利用这种规律性的东西,以服务于统计推断工作.由此也可以看出,统计方法的效力只能在有大量数据可利用时,才能显示出来.资料太少,统计方法也是无能为力的.当然,对数据多少的要求,依对结论精度的要求而定,并无绝对的标准.

1.5 学一点统计学

虽然近代统计学的发展可以说是起源于20世纪初,但带有统计

性质的工作却可以溯源很远.我国古时候就有所谓“结绳记事”,在浩繁的“二十四史”中有大量关于人口、钱粮、水文、天文、地震……的资料记录.在西方,“统计”(statistics)一词就是从“国家”(state)一词演化而来的,意指一种收集和整理国情资料的活动.随着近代科学技术和工农业生产方面的飞速发展,统计方法得到了愈来愈深入和广泛的应用,对人类认识和改造世界产生了重大的影响.有人把统计学列为 20 世纪几十项最重大的成就之一,如日本在战后经济恢复和高速发展中统计方法所起的作用,是人们津津乐道的话题.一些国家在国情调查中放弃普查的方法而改用抽样的方法,取得了良好的效果.用统计方法分析种种社会调查所得资料而引出的结论,往往成为有关当局决定政策的重要依据.在此我们不打算一一列举统计方法的各种应用——前面几个例子已指出了一些,它们都可以作为一类应用的代表.在本书以后的叙述中还将提到统计方法的一些应用.现在,统计学的基本知识已普遍成为高等教育的一个不可缺少的组成部分,这从高等学校中许多系科的教学计划中都包含这个内容可以看出.

除了因其广泛的应用而可以给人们在工作中助一臂之力以外,从一般的思想和文化修养的角度去看,学一点统计学也是很有益的,甚至是必要的.统计方法是认识和改造世界的重要方法之一,对这方面毫无了解,不能不说是知识结构上的一个缺陷,是可引为遗憾的.用统计的观点看待事物很重要,在许多情况下它是唯一说得通的观点.总之,随机性的思想提醒人们对事情的看法不要绝对化,习惯于这种思想的人,不会因为一些偶然落到自己头上的不愉快事件而过分耿耿于怀.另外,统计思想使人们在两可的事物中掌握适当的度.举一个例子,甲、乙两名棋手比赛 5 局,结果甲 4 胜 1 负,有人认为这个纪录肯定地说明了甲的棋艺高于乙,有人则认为还难说.就一个对统计方法略有所知的人来说,他不会陷入这种争端.他了解,两种情况都可能(一般人自然也不否定这一点),且可以算出,两种看法的“正确程度”在一定的意义下是 5 与 3 之比.这提供了看待这个问题的一个“度”,对其他许多小事以至大事,也莫不如此.因此,习惯于统计思想的人都能允执其

中国科学技术大学出版社中学数学用书

高中数学竞赛教程/严镇军 单增 苏淳 等

中外数学竞赛/李炯生 王新茂 等

第 51—76 届莫斯科数学奥林匹克/苏淳 申强

全国历届数学高考题解集/张运筹 侯立勋

中学数学潜能开发/蒋文彬

同中学生谈排列组合/苏淳

趣味的图论问题/单增

有趣的染色方法/苏淳

组合恒等式/史济怀

集合/冯惠愚

不定方程/单增 余红兵

概率与期望/单增

组合几何/单增

算两次/单增

几何不等式/单增

解析几何的技巧/单增

构造法解题/余红兵

重要不等式/蔡玉书

高等学校过渡教材读本:数学/谢盛刚

有趣的差分方程(第 2 版)/李克正 李克大

抽屉原则/常庚哲

母函数(第 2 版)/史济怀

从勾股定理谈起(第 2 版)/盛立人 严镇军

三角恒等式及其应用(第 2 版)/张运筹

三角不等式及其应用(第 2 版)/张运筹

反射与反演(第2版)/严镇军
数列与数集/朱尧辰
同中学生谈博弈/盛立人
趣味数学100题/单增
向量几何/李乔
面积关系帮你解题(第2版)/张景中
磨光变换/常庚哲
周期数列(第2版)/曹鸿德
微微对偶不等式及其应用(第2版)/张运筹
递推数列/陈泽安
根与系数的关系及其应用(第2版)/毛鸿翔
怎样证明三角恒等式(第2版)/朱尧辰
帮你学几何(第2版)/臧龙光
帮你学集合/张景中
向量、复数与质点/彭翥成
初等数论/王慧兴
漫话数学归纳法(第4版)/苏淳
从特殊性看问题(第4版)/苏淳
凸函数与琴生不等式/黄宣国
国际数学奥林匹克240真题巧解/张运筹
Fibonacci数列/肖果能
数学奥林匹克中的智巧/田廷彦
极值问题的初等解法/朱尧辰
巧用抽屉原理/冯跃峰
统计学漫话(第2版)/陈希孺 苏淳
学数学.第1卷/李潜
学数学.第2卷/李潜
学数学.第3卷/李潜