

抽样调查

导论

INTRODUCTION
TO
SURVEY
SAMPLING

11
115

Graham Kalton

格雷厄姆·卡尔顿 (美)著
郝虹生 翻译 王文颖 校译

中国统计出版社
China Statistics Press



抽样调查 导论

**INTRODUCTION
TO
SURVEY
SAMPLING**

中国统计出版社
China Statistics Press



INTRODUCTION TO SURVEY SAMPLING

GRAHAM KALTON

ISBN: 0-8039-2126-8

Copyright © 1983 by sage publications.Inc.

Simplified Chinese translation edition jointly published by China Statistics Press

本书中文简体字翻译版由中国统计出版社出版。未经出版者
预先书面许可，不得以任何方式复制或抄袭本书的任何部分。

北京市版权局著作权合同登记号： 01-2003-2714

(京) 新登字 041 号

图书在版编目 (CIP) 数据

抽样调查导论 INTRODUCTION TO SURVEY SAMPLING

/ (美) 卡尔顿(Kalton, G.)著；郝虹生等译。

- 北京：中国统计出版社，2003. 6

ISBN 7-5037-4054-X

I 抽…

II ①卡… ②郝…

III 抽样调查

IV C811

中国版本图书馆 CIP 数据核字 (2003) 第 012701 号

责任编辑 / 吕 军

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 75 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 / (010) 63459084, 63266600-22500 (发行部)

印 刷 / 科伦克三莱印务(北京)有限公司

经 销 / 新华书店

开 本 / 850 × 1168 1/32

字 数 / 99 千字

印 张 / 4

印 数 / 1-4000 册

版 别 / 2003 年 5 月第 1 版

版 次 / 2003 年 5 月北京第 1 次印刷

书 号 / ISBN 7-5037-4054-X/C 2029

定 价 / 12.00 元

版权所有。未经许可，本书的任何部分不准以任何方式在世界任何地区以任何文字翻印、拷贝、仿制或转载。中国统计版图书，如有印装错误，本社发行部负责调换。

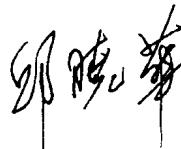
格雷厄姆·卡尔顿（美）著
郝虹生 翻译 王文颖 校译

序

满页数学公式的抽样调查书籍往往使实际部门的统计工作者望而生畏,而 Kalton 教授的《抽样调查导论》则令人耳目一新。该书通过通俗易懂的描述和简单实用的举例,既巧妙地避开了复杂的数学公式,又系统、完整地阐明了抽样调查的基本思想和方法。该书不足八万字,但内容非常丰富,不仅包括了抽样调查技术的基本内容,而且包括抽样调查实践中经常会遇到的许多具体问题。对于具备一定抽样理论的人或是初学者,尤其是实际部门的统计工作者,该书是一本非常有价值的工具书。

随着我国统计制度方法改革的深化,抽样调查技术的应用将会越来越广泛,成为我国政府统计调查的主体模式。从这个意义上讲,我们的统计工作者都应该认真学习和掌握抽样调查的基本思想和方法,而 Kalton 教授的《抽样调查导论》则是一本非常好的教课书。

Kalton 教授是国际上著名的抽样专家,在抽样理论和实践方面都造诣颇深。Kalton 教授作为亚洲开发银行和中国国家统计局的合作项目《建立中国企业抽样调查体系》聘请的专家,为国家统计局设计了《工业企业抽样调查方案》并指导在江苏省进行了试点,取得了良好效果。借此书出版之际,我向亚洲开发银行、Kalton 教授和 Dalisay S. Maligalig 女士,对他们为我国企业抽样调查工作所作的贡献表示诚挚的感谢。



2003 年 3 月 27 日

写给中文版本的序言

我很高兴此书被译为中文并由中国统计出版社出版。抽样调查这门学科关注于如何在调查研究项目中抽选调查元素(如住户、企业等)。从选中元素采集数据的目的是推算总体,因此抽选样本的方法是至关重要的。本书概述的复杂概率抽样设计重点放在从大总体中抽取大样本。中国是一个非常大的国家,显然会经常使用复杂抽样设计,用于对国家、省以及城市等范围的总体进行估计。我希望此书能在促进中国有效推广运用概率抽样方法方面贡献微薄之力。

中文版与英文版内容相同,但我借此机会增加了一些非常有用的参考书目,这些参考书都出版于本书英文版出版之后。

我十分感激郝虹生博士翻译了此书。郝博士在密歇根大学学习时曾修过我讲授的抽样课程,而且在中国承担过许多调查项目的抽样设计工作,因此尤为胜任本书的翻译。我同样要感谢国家统计局企业调查总队副总队长王文颖先生联系并关注出版本书,此外,他还对本书的翻译和校审作了颇有价值的贡献。我还要感谢国家统计局企业调查总队总队长宋跃征先生和中国统计出版社副主编严建辉先生对本书出版予以的支持。没有这些朋友和同事们的努力,也就没有本书的出版。我真诚地感激他们。

A handwritten signature in black ink, appearing to read "Graham Kalton". The signature is fluid and cursive, with the name written in a single continuous line.

丛书编辑导言

现在,无论在学术界还是在更为广泛的领域,调查研究都是非常重要的。在每一门社会科学以及相关学科中它都是一个基本工具,也是诸如全国选举研究、全国舆论研究中心的一般社会调查、以及密歇根调查研究中心范围广泛的消费者调查这些重要研究的基础。它被频繁地用于各种各样的应用性研究,既用于描述的目的,也作为评估的基础。另外,由于在政治竞选活动中的应用,它已获得了非同寻常的公众认知。

当然,调查研究的基础在于抽样调查过程。如果抽样本身的设计和实施欠佳,那么无论所提的问题有多么好,也无论分析有多么细致,我们也获取不到多少知识。尽管这些道理显而易见,有人却可能会认为抽样不过是个技术问题,最好留给统计学家去考虑。而我的看法则完全不同。对于许多项目来说,虽然抽样统计学家确实起着举足轻重的作用,但是,指导调查和为实质性目的而使用调查资料的研究人员至少应该对抽样的原理有适当程度的了解。

卡尔顿的这本书在解释抽样方法上采取了一种巧妙的中间路线。它并不是以抽样统计学家为对象写作的。与此相反,本书有很强的可读性,只要掌握了一定程度的基础统计学就能够理解。书中所有的概念都以举例来详细说明,从而为读者理解调查设计的主要组成部分提供了一个扎实的基础。本书的一个重要特色是它还包括了实践中需要考虑的问题。举例来说,若不是在实践中遇到抽样框和无回答这类问题,有关这些问题的章节可能就会被忽视。

卡尔顿的这本书内容涵盖广泛,既包括了抽样理论方面的内容,也涉及了诸多的实践问题,因此,无论对初学者还是曾学过一些抽样理论的人来说,它都是一本很有价值的书。

里查德 G. 尼米
丛书编辑

目 录

1. 引言	(1)
2. 简单随机抽样	(5)
3. 等距抽样	(14)
4. 分层	(18)
比例分层	(20)
非比例分层	(23)
层的选择	(25)
5. 整群与多阶段抽样	(28)
6. 与规模成比例的概率抽样	(38)
7. 其它概率抽样设计	(48)
两相抽样	(48)
重复抽样	(50)
固定样本设计	(54)
8. 抽样框	(59)
缺矢元素	(60)
群元素	(61)
空缺与异类元素	(65)
重复登记	(66)
9. 无回答	(67)
10. 调查数据分析	(74)
权数	(74)
抽样误差	(80)
11. 样本量	(88)

12. 两个实例.....	(92)
一项全国范围的面对面访谈调查.....	(92)
一项电话访谈调查.....	(94)
13. 非概率抽样.....	(99)
14. 结束语	(104)
参考文献.....	(106)
附录：本书中有关术语英汉对照	(113)

1

引言

如今,无论是为了研究目的还是为了管理目的,抽样调查作为一种为种类繁多的主题提供统计数据的手段已被广为接受。在社会学、社会心理学、人口学、政治学、经济学、教育以及公共卫生等领域,抽样调查都已得到广泛应用,这些调查被用于形成、检验和改进研究假设。一些国家的中央政府大量采用抽样调查来了解其人口在就业与失业、收入与支出、住房状况、教育、营养、健康、旅行方式以及许多其它方面的情况。此外,中央政府还对企业和机构进行抽样调查,如制造商、零售商、农场、学校和医院等。同样,地方政府也为了地方性规划的目的而使用抽样调查。市场研究人员通过抽样调查研究产品市场,发现其产品是如何被使用的以及使用性能情况,并确定消费者的反应。民意测验则被用于随时了解政治领袖及其政党的声望,并就各种时事问题观测公众舆论。

抽样调查目前已被广泛使用,而我们所熟知的抽样调查的历史却并不长,这多少有些令人惊讶。抽样调查的历史主要是限于本世纪,而且其实际应

用自 1930 年代以来才开始明显增多。本世纪，调查方法论在各个方面都取得了相当可观的进展，特别是本文的主题——抽样方法方面。本世纪初，尽管已知抽样调查原则上是可行的 (O'Muircheartaigh and Wong, 1981)，统计学家们还是就非全面调查是否足以推断总体而进行过辩论。此后，抽样调查逐渐被广泛接受，人们发明了一系列令人印象深刻的抽样方法，用来在各种不同的实际条件下抽出既有效又经济的样本。

一项调查的设计需要就若干因素作出许多相互联系的决策，如资料收集的方式（是使用面对面访谈、电话访谈，还是自填表格）、所提问题的设计、资料处理的方法以及抽样设计等（参见 Moser and Kalton, 1971；Warwick and Lininger, 1975）。虽然本书所涉及的仅是抽样设计，但需要认识到，在实践中抽样设计只是整个调查设计的一个组成部分。特别是资料收集过程中所涉及的经济性问题，对抽样设计的选择会产生相当大的影响。

调查设计的首要步骤之一是界定所要研究的总体。在此处，“总体”这一术语是在技术意义上使用的，它指的是全部被研究的元素。“元素”在此处是分析单位，可以是人，也可以是户、农场、学校或任何其它单位。对于总体，需要根据调查目的进行明确而严密的定义，因为调查的结果将取决于所采用的定义。例如，考虑在一个城市进行一项调查，了解公众对于采用一种新的公共汽车系统的支持程度。是否应该仅限于对居住在城市边界之内的人进行调查？调查总体的最低年龄为多少岁？是否应包括在城市选举中没有投票资格的居民？是否应排除城市中的暂住者？如果应该排除，怎样对他们进行定义？在对大多数总体进行定义时都会出现许多诸如此类的问题，使得对总体的定义并不像初看起来那样简单明了。

定义总体的一种有益作法，是先按照调查目的的要求定义一个理想的总体，即目标总体。然后，这种定义常常要根据实际情况的限制而加以修正，形成调查总体。例如，在理想的情况

下,美国的许多全国性调查应包括驻守在海外基地的军人、居住在夏威夷和阿拉斯加的人口、住在医院、旅馆、监狱、军营以及其他机构中的人。然而,由于很难对这些人进行调查,使得这些人常常被排除在调查总体之外。先定义理想的目标总体的优点是可以明确地识别被排除的部分,使我们能够对推断的局限性的程度和后果做出估计。

对总体做出定义之后,就可以提出如何从中抽选样本的问题。当然,一种可能是对总体中的所有元素进行全面的调查,但这种做法常常不适用。显然,只要有足够的估计精确度,仅从总体的一部分收集资料能减少费用,因此抽样比全面调查更为经济。抽样调查可以更快地实施,并更快地处理资料,因而可以更加及时地报告调查结果。此外,将资源集中用于总体的一部分,有可能使所收集资料的质量优于全面调查。其结果是抽样调查事实上可能得到更准确的结果。由于这些原因,除非总体很小,我们几乎总是采用抽样的方法进行调查。

抽样设计所关注的主题是如何从总体中抽选出所要调查的那部分样本。对样本的一个基本区分是其抽选是概率抽样,还是非概率抽样。对于概率样本来说,每个元素都有一个已知且非零的机会入选样本,这样可以避免抽选偏差,并可以利用统计学理论得到调查估计量的性质。非概率抽样涉及许多方法,包括使用志愿参加者,以及根据对总体的“代表性”有意识地选择部分元素作为样本。所有非概率抽样方法的弱点是其主观性,缺乏一个为之提供支持的理论框架。一个志愿参加者样本,或由专家挑选的代表性样本,只能依靠主观估计来取得,而不是根据不依赖于假定的统计方法。鉴于非概率抽样的这一弱点,本书将主要讨论概率抽样,但是在第13章中将包括一些有关非概率抽样的讨论。

对于任何形式的概率样本来说,一个基本要求是要有一个可以从中抽选样本元素的抽样框。简单情况下,当总体所有元

素的名单可以获得时, 抽样框可以是该名单。若没有名单, 抽样框则是某种等效的用以识别总体元素的程序。区域抽样是使用这种抽样框的一个很好实例。在使用这种技术时, 总体的每一个元素都与一个特定的地理区域相联系(例如, 人或户与他们居住的区域相联系, 若他们有一个以上的居住地, 则与主要的居住地相联系)。然后, 抽选一个由区域组成的样本, 在被抽选出的区域中, 或者调查所有的元素, 或者从中再抽选部分元素(见第12章)。抽样框的一般构成情况以及它所包含的关于总体元素的信息常常对抽样设计的选择有很大的影响。抽样框中的缺陷, 例如没能包括调查总体中的所有元素, 可以给样本带来有害的作用。在第8章中将对抽样框问题进行更详细的讨论。

人们已开发出各种概率抽样技术以提供有效、实用的抽样设计。其中最广泛应用的是等距抽样、分层抽样、多阶段(整群)抽样, 以及与规模成比例的概率抽样。为了便于说明, 以下各章将分别讨论这些技术, 但在实践中这些技术常常共同使用而形成复合设计。在第12章中给出了两个实例, 以说明如何使这些技术结合起来。我们将首先讨论适用于对小而紧凑的总体进行抽样的相对简单的技术, 然后进而讨论对更大的且更分散的总体进行抽样所需的更为复杂的技术。

2

简单随机抽样

简单随机抽样(SRS)为讨论概率抽样方法提供了一个自然的出发点,这并不是由于它被广泛应用(其应用并不广泛),而是由于它是最简单的方法,而且是许多更复杂方法的基础。在给简单随机抽样下定义之前,我们引入如下符号:样本量由 n 表示,总体规模由 N 来表示。然后我们将简单随机抽样正式地定义为具有如下性质的一种抽样方法:来自总体 N 个元素的由任意可能的 n 个不完全相同元素组成的子集,都有同样的可能性被选为样本。这个定义意味着,总体中的每一个元素都有相同的概率被选入样本,但上述定义要比这更严格。我们以后将会看到,更为复杂的抽样也经常是等概率抽选方法(epsem),但是在这些设计中,被抽中元素集合的联合概率却是不相等的,而在简单随机抽样时则是相等的。

以下我们将以一个具体例子来讨论简单随机抽样。假定要在一所中学进行一项调查,了解学生的闲暇爱好。从学校可以得到该校 1872 名学生的按学号排列的名单。这些号码从 0001 排到

1917, 其中的若干间断是由于有一些已具有学号的学生后来离开了学校。假定该项调查要求抽一个 $n = 250$ 的简单随机样本。(关于 n 的选择将在第 11 章中讨论)。

抽出所要求的简单随机样本的一种方式是使用抽彩方法, 即准备 1872 个完全相同的碟片, 将每个学生的姓名或学号写在其中一个碟片上。将这些碟片放在一个坛子里并进行彻底的混合, 然后任意地从中选出 250 个。如果这些操作程序得以完满地执行, 则被选出的碟片即能够确定出一个由 250 个学生组成的简单随机样本。这种方法虽然在概念上很简单, 但执行起来却很麻烦, 而且还依赖于碟片得到彻底混合的假设, 因此这种方法极少被采用。

抽出简单随机样本的另一种方式是使用随机数字表。这些表经过精心编制和验证, 以保证每一位数、每一对数、以及由此类推的更多位数最终以同样的频率出现在表中。表 1 给出了由坎代尔和史密斯(1939)编制的随机数字表的一个摘录。

表 1 随机数字表

67 28	96 25	68 36	24 72	03 85	49 24
85 86	94 78	32 59	51 82	86 43	73 84
40 10	60 09	05 88	78 44	63 13	58 25
94 55	89 48	90 80	77 80	26 89	87 44
11 63	77 77	23 20	33 62	62 19	29 03

资料来源: Kendall, M. G. and B. B. Smith, Tables of Random Sampling Numbers. Copyright(c)1939 by Cambridge University Press. Reprinted by permission.

由于学生的识别号码包括四位数, 我们需要以四个数字为一组来抽选随机数。在实践中, 应从表中随意选择的某一个点开始, 但在此处为简单起见, 我们从表的左上角开始。然后, 沿着第一组的四列数字向下数, 此后再沿第二组的四列数字向下数, 余下类推。凡是学生编号范围(0001 – 1917)之外的数字, 以及虽在此范围之内但不与在校学生相联的数字都被忽略。表中

前四个数字(6728, 8586, 4010, 9455)都未选出样本学生, 因此第一个被选出的学生是 1163 号(假定这个学生仍在校)。这样继续数到底, 从表的这一部分仅能选出的其它样本学生是 0588 号和 0385。至此已经可以清楚地看出, 以这种方式抽出 250 个学生是一件令人乏味的工作, 需要大量的随机数字, 而其中大部分是无用的。

将每个学生与几个随机号码相联系, 而不是仅与一个相联系, 可以避免使这么多的随机数字被浪费; 只要所有学生都与同样多的随机数字相联系, 样本就仍然是一个简单随机样本。在此例中, 每个学生可以与五个四位的随机数字相联系。一种简单的方法是使 0001 号学生也与 2001、4001、6001 和 8001 相联系; 0002 号学生也与 2002、4002、6002、8002 相联系; 依此类推, 直至 1917 号学生也与 3917、5917、7917 和 9917 相联系。然后再从表的左上角开始, 被抽选的学生为 6728 = 0728 号学生; 8586 = 0586 号学生; 4010 = 0010 号学生; 9455 = 1455 号学生; 1163 = 1163 号学生, 等等。

在使用随机数字表抽选样本的过程中, 一个元素有可能被抽中一次以上。而在使用前面介绍的抽彩方法时不存在这种可能性, 因为当一个元素的碟片被从坛中抽取出来后就不再放回, 不存在再次中选的机会。然而, 在使用抽彩方法时, 若在下一次抽选进行之前将被抽中的圆碟片放回坛中, 则也存在被抽中一次以上的可能性。若在进行抽样时不放回, 样本就一定包含 n 个不同的元素, 而在放回抽样时, 大小为 n 的样本则可能包含少于 n 个不同的元素。若此处描述的抽样程序是有放回的, 抽样方法即被称为无限制随机抽样, 或有放回简单随机抽样。若这些抽样程序是不放回的, 则这种方法被称为无放回简单随机抽样, 或者就称为简单随机抽样。要从随机数字表中抽选一个无放回的简单随机样本, 只要将已被选入样本的元素的重复中选忽略即可。由于无放回的抽样能够得出比放回的抽样更为精