

APPLIED STATISTICS

MBA CLASSICS

MBA 精品系列

MBA

应用统计学

(第三版)

主编 / 贾俊平 谭英平

 中国人民大学出版社

MBA CLASSICS

MBA 精品系列

MBA

应用统计学

(第三版)

主编 / 贾俊平 谭英平

中国人民大学出版社
· 北京 ·

图书在版编目 (CIP) 数据

应用统计学/贾俊平, 谭英平主编. —3 版. —北京: 中国人民大学出版社, 2017. 3
MBA 精品系列
ISBN 978-7-300-23934-7

I. ①应… II. ①贾…②谭… III. ①应用统计学-教材 IV. ①C8

中国版本图书馆 CIP 数据核字 (2017) 第 016811 号

MBA 精品系列

应用统计学 (第三版)

主编 贾俊平 谭英平

Yingyong Tongjixue

出版发行 中国人民大学出版社

社 址 北京中关村大街 31 号

邮政编码 100080

电 话 010-62511242 (总编室)

010-62511770 (质管部)

010-82501766 (邮购部)

010-62514148 (门市部)

010-62515195 (发行公司)

010-62515275 (盗版举报)

网 址 <http://www.crup.com.cn>

<http://www.ttrnet.com> (人大教研网)

经 销 新华书店

印 刷 北京鑫丰华彩印有限公司

版 次 2005 年 5 月第 1 版

规 格 185 mm×260 mm 16 开本

2017 年 3 月第 3 版

印 张 14.75 插页 1

印 次 2017 年 3 月第 1 次印刷

字 数 323 000

定 价 36.00 元

版权所有 侵权必究 印装差错 负责调换

前 言

一本什么样的教材能让学生更好地理解统计呢？根据笔者对统计的理解及多年的教学经验，尽可能少使用那些专业的统计术语、少去纠缠那些复杂的公式、少去用晦涩的词汇表述统计问题和结果，或许是个不错的选择。本书在写法上做了一些新的尝试：力图把统计方法的思想用书中标题的形式表达出来，尽管这种表达不一定确切；在书中内容的表述上，每种方法都尽力用实际问题引出，而不是从概念开始，尽量不使用更专业的统计术语；书中例题的解答直接使用计算机的输出结果，尽可能抛弃手工计算过程，书中例题的计算使用 SPSS 和 Excel 两种软件，但以 SPSS 为主，对软件操作的一些说明放在每章后的附录里。

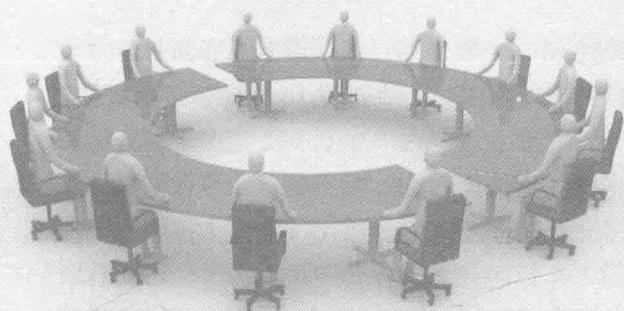
作为一门应用性很强的学科，多数人学习的目的也主要是应用。但初学者学习统计时面临的主要困惑是学完不会用。问题在于学习过程中多把注意力集中在公式和计算上，而忽视对统计思想的理解。学习统计关键在于理解。记住公式，不等于学会统计；学会计算，不等于会用统计。统计的真谛在于它所体现的思想，在于它所提供的思维方式。学好统计的关键是掌握如何运用统计思维来思考问题，而不是简单地记住那些死的统计知识。有些初学者对统计课程往往感到畏惧，被书中的公式吓倒。实际上，抛开公式照样可以学会统计。特别是在计算机应用已经普及的今天，所有的计算都可以由计算机来完成。只要清楚统计方法使用的前提，理解统计方法的实质，要应用统计并不困难。

本书的初衷是作为统计学专业学生的入门课程教材，以替代过去的描述统计内容。作为本专业的学生，在最初接触统计时，应该让他们对统计有一个较全面的认识，了解一些统计思想，为后续的专业课学习奠定基础。当然，本书也可以作为非统计专业学生通开课的教材使用。由于书中的有些提法只是笔者的个人看法，不一定恰当，希望读者多提意见和建议，以便进一步修改和完善。

贾俊平 谭英平

第 1 章	统计能为你做些什么	1
	1.1 统计无处不在	1
	1.2 统计学研究数据	4
	1.3 怎样获得数据	6
第 2 章	用图表和统计量看数据	10
	2.1 用图表描述数据	10
	2.2 用统计量描述数据	20
第 3 章	用概率分布描述随机变量	33
	3.1 度量事件发生的可能性	33
	3.2 随机变量的概率分布	34
	3.3 由正态分布导出的几个重要分布	40
	3.4 样本统计量的抽样分布	42
第 4 章	用样本推断总体	49
	4.1 怎样进行推断	49
	4.2 估计总体参数	52
	4.3 检验总体假设	56
第 5 章	类别变量分析	75
	5.1 某个类别变量的观测频数与期望频数是否一致	75
	5.2 两个类别变量是否独立	79
	5.3 度量两个类别变量的关系强度	81
第 6 章	类别变量对数值变量的影响	86
	6.1 方差分析解决什么问题	86
	6.2 考虑一个类别变量的影响	89
	6.3 考虑两个类别变量的影响	90
第 7 章	利用变量间的关系进行预测	98
	7.1 变量之间有什么样的关系	98
	7.2 建立变量之间的数学表达式	102

7.3	拟合效果的度量和回归检验	104
7.4	所有自变量都有必要放进模型中吗	106
7.5	用自变量预测因变量	110
7.6	含有定性自变量的回归	110
第8章	根据过去的模式预测未来	120
8.1	时间序列的组成要素	120
8.2	时间序列预测的程序	122
8.3	平滑法预测	126
8.4	趋势模型预测	127
8.5	多成分序列的预测	132
第9章	用少数变量代表多个变量	143
9.1	主成分分析	143
9.2	因子分析	149
第10章	把对象分成不同的类别	156
10.1	聚类分析	156
10.2	判别分析	162
第11章	不依赖于分布的检验	172
11.1	关于非参数检验	172
11.2	单样本的非参数检验	174
11.3	两样本的非参数检验	178
附录	各章习题答案	189
	参考书目	228



第 1 章

统计能为你做些什么

统计思维总有一天会像读与写一样成为一个有效率公民的必备能力。

——H. G. 韦尔斯 (H. G. Wells)

1.1 统计无处不在

1.1.1 每个人都离不开统计

了解一些统计知识对每个人都是必要的。比如，在外出旅游时，你需要关心一段时间内的详细天气预报；在投资股票时，你需要了解股票市场价格的信息，了解某只特定股票的有关财务信息；在观看足球比赛时，除了关心进球数的多少外，你还要知道各支球队的技术统计，等等。要正确阅读并理解下面的一些统计研究结论，就更需要具备一些统计知识。

- 吸烟对健康是有害的。
- 不结婚的男性会早逝 10 年。
- 身材高的父亲，其子女的身高也较高。
- 第二个出生的子女没有第一个聪明，第三个出生的子女没有第二个聪明，依此类推。
- 两天服一片阿司匹林会减少心脏病第二次发作的机会。
- 身体超重 30% 会使寿命减少 1 300 天。

- 每天摄取 500 毫升维生素 C, 生命可延长 6 年。
- 学生们在听了莫扎特钢琴曲 10 分钟后的推理测试会比他们听 10 分钟娱乐磁带或其他曲目做得更好。
- 上课坐在前面的学生平均考试分数比坐在后面的学生高。

看懂这些结论并不困难,但这些结论是怎样得出来的?你相信这些结论吗?学点统计知识你就会正确理解它们。

了解一些统计知识,对政策的制定者或企业的管理者来说同样重要,这有助于他们作出正确决策,否则会因为不懂统计而弄出笑话来。一个统计办公室的主管与一些统计学者开会,统计学者抱怨从其他部门收到的一些估计值没有给出标准误差(估计时的误差大小,表示估计的精度),这个主管马上问道:“对误差也有标准吗?”健康部门的一位官员看到一个统计学者提供的报告,报告中提到去年由于某种疾病,平均 1 000 人中死亡人数为 3.2 人,这位官员对这个数字产生了兴趣。他问他的私人秘书,3.2 个人是如何死法?他的秘书说:“先生,当一个统计学家说死了 3.2 个人时,意味着 3 个人已经死了,两个人正要死。”

有点统计知识就不会闹出这样的笑话来。一位学者写道:“假定你是市场部的新任经理,一次广告活动的统计结果摆到了你面前,声称某个结果是‘统计显著’的。你如何解释这份报告而又不暴露你对该术语的无知?赶快学点统计,这对你和你的事业都非常有用。”^①

1.1.2 几乎所有的领域都要用统计

统计是适用于所有学科领域的通用数据分析方法,是一种通用的数据分析语言。只要有数据的地方就会用到统计方法。这里,我们不想列举统计的应用领域,只想通过几个简单的例子说明统计的应用。

【例 1.1】用统计识别作者。1787—1788 年,三位作者亚历山大·汉密尔顿(Alexander Hamilton)、约翰·杰伊(John Jay)、詹姆斯·麦迪逊(James Madison)为了说服纽约人认可宪法,匿名发表了著名的 85 篇论文。这些论文的作者大多数已经得到了确认,但是,其中 12 篇论文的作者身份引起了争议。通过对不同单词的频数进行统计分析,得出的结论是詹姆斯·麦迪逊最有可能是这 12 篇论文的作者。现在,对于这些存在争议的论文,认为詹姆斯·麦迪逊是原创作者的说法占主导地位,而且几乎可以肯定这种说法是正确的。

【例 1.2】用统计进行产品质量管理。统计在企业产品质量管理中的应用是统计应用的一个重要方面。在统计中, σ 表示一个总体的标准差,它表示的是数据之间的差异程度。比如,在企业生产的产品中,同一种产品没有两个是完全一样的,因为在生产过程中,由于各种因素的影响而使产品质量产生波动。产品的这种差异称为质量的波动性,也正是由

^① Gudmund R. Iversen, Mary Gergen. 统计学:基本概念和方法. 北京:高等教育出版社,2000.

于波动性的存在才需要进行质量管理。 6σ 是质量管理中使用的一个术语，它的含义是指偏离正态分布的中心6个标准差。它表示在生产过程中缺陷率不超过百万分之三点四，通俗地说，如果生产100万个产品，不合格产品平均来说不超过3.4个。这样的不合格率非常低，以至于可以忽略不计。 6σ 质量管理已成为最新的质量管理理念，近年来，它已成为一些著名国际大企业的质量管理方法，并且使企业受益匪浅。例如，实行 6σ 质量管理，使摩托罗拉公司在3年中节省的资金超过9.4亿美元。实行 6σ 管理的大公司还有美国通用电气公司（GE）、宝利来（Polaroid）和得州仪器（Texas Instruments）等。GE的前CEO杰克·韦尔奇1999年4月曾说过这样一段话：“ 6σ 培训计划是GE下一个世纪领导层得以产生繁衍的园地， 6σ 是我们曾经尝试过的最重要的管理培训方法，它胜过到哈佛商学院就读，也胜过到克顿维尔（克顿维尔是GE公司内部的质量培训部）进修，它教会你一种完全与众不同的思维方式。”在推广 6σ 质量管理不到10年的时间内，GE的总市值从世界排名第十位跃升到第二位。

【例 1.3】用简单的描述量得到一个重要发现。费希尔（R. A. Fisher）在1952年的一篇文章中举了一个例子，说明如何由基本的描述统计量知识引出一个重要的发现。20世纪早期，哥本哈根卡尔堡实验室的施米特（J. Schmidt）发现不同地区所捕获的同种鱼类的脊椎骨和鳃腺的数量有很大不同，甚至在同一海湾内不同地点所捕获的同种鱼类也有这样的倾向。然而，鳗鱼的脊椎骨数量变化不大。施米特在从欧洲冰岛、亚速尔群岛等地以及尼罗河等几乎分离的海域里所捕获的鳗鱼样本中，计算发现了几乎一样的均值和标准偏差值。由此，施米特推断所有各个不同海域内的鳗鱼是由海洋中某公共场所繁殖的，后来名为“戴纳”（Dana）的科学考察船在一次远征中发现了这个场所。

这些例子表明统计不仅在许多领域都有广泛应用，而且在生产、生活、科学研究等各个领域都发挥着日趋重要的作用。

1.1.3 统计的误用与滥用

马克·吐温（Mark Twain）有一句名言：“有三种谎言：谎言、该死的谎言和统计数据。”历史学家安德鲁·兰（Andrew Lang）说，一些人使用统计“就像喝醉酒的人使用街灯柱——支撑的功能多于照明”。统计常常被人们有意或无意地滥用，比如错误的统计定义、错误的图表展示、不合理的样本、数据的遗漏或逻辑错误，等等。这些误用有些是常识性的，有些是技术性的，有些则是故意的。作为从数据中寻找事实的统计，却被有些人变成了歪曲事实的工具。你也许常常看到这样的产品质检报告：某某产品的抽样合格率是80%。乍看上去没什么问题，但实际上只抽查了5件产品，有4件合格。这样的合格率能说明什么问题呢？在马路上随便采访几个人，他们的看法能代表大多数人的观点吗？“调查结果表明……”，调查了多少个人？是随机调查的吗？样本是怎样选取的？这看上去是在用事实说话，实际上成了统计陷阱。

此外，统计也往往被作为两个极端使用：一个极端是不懂或不太懂统计的人认为统计没什么用，他们因为不懂统计而瞧不起统计，他们不用或几乎不用统计方法分析数据，即

使做些统计分析, 往往也是表面上的。走入这一极端的人, 他们的决策依据就是自己的大脑: 一些杂乱无章的信息组合出的某种直觉。如果他们的决策是正确的, 更增加了他们的自信, 更加感到不用统计也挺好; 如果他们的决策出了毛病, 他们会找出一大堆开脱理由: 市场难测, 环境突变, 竞争激烈, 价格下跌, 需求疲软, 管理不善, 成本上升, 出口下降……另一个极端是把简单问题复杂化。特别是在管理领域, 一些管理者把本来可以用简单方法解决的问题故意复杂化, 他们不用简单的分析方法, 而是用复杂的分析方法; 他们为证明管理的科学性, 建立一个别人看不懂模型, 编一大堆程序, 输出了一大堆数字和符号; 他们得出用统计语言陈述的结论, 提出一些似是而非的建议……这样的分析往往脱离了管理问题, 对实际决策也未必有用。在工商管理中, 这两个极端都是不可取的。管理决策中不用统计几乎不可想象; 把简单问题复杂化对管理决策也未必有用。从统计的实际应用来看, 简单的方法不一定没用, 复杂的方法也不一定有用。统计应该恰当地应用到它能起作用的地方。不能把统计神秘化, 不能歪曲统计, 更不能把统计作为掩盖事实的陷阱。

曲解统计是一种常见的现象。在有些人的心目中, 使用统计就是寻找支持: 他们的心目中可能有了某种“结论”性的东西, 或者说他们希望看到符合他们需要的某种结论, 而后去找些数据来支持他们的结论。如果数据分析的结果与他们预期的结论一致, 他们就会声称自己是用科学方法得到的结论; 如果与预期的不一致, 他们要么篡改数据, 要么对统计弃而不用。这恰恰歪曲了数据分析的本质。数据分析的真正目的是从数据中找出结论、从数据中寻找启发, 而不是寻找支持。真正的数据分析事先是没有结论的, 通过对数据的分析才得出结论。

1.2 统计学研究数据

你问身边的人, GDP (国内生产总值) 是什么? CIP (消费者价格指数) 是什么, 似乎都能说上几句。但要是仔细追问它们究竟代表了什么, 就不是每个人都能够说清楚的。统计也是一样。你要问一个人统计是什么, 似乎没有人不知道, 但多数人会将其与统计工作相联系。要问统计学是什么, 就不是每个人都能够说明白的, 要搞清楚统计学研究什么就更困难了。

1.2.1 有数据的地方就需要统计学

物理学研究的是像热、光、电这类自然现象的运动规律。化学家测定物质的组成及化学元素之间的交互作用。生物学家研究植物和动物的生活。数学家则在给出的假定之下推演各种命题。这些学科都有它们自己特定的问题, 而且各自有解决这些问题的方法, 各学科因此而成为一门单独的学科。

统计学是一门独立的学科, 这似乎没人怀疑。但统计学究竟研究什么, 可能就有不同

的看法。有人认为，统计学是一门独特的学问，没有任何固定的对象。乍听起来似乎难以理解，但仔细想想，也许有道理。统计学研究的是来自各领域的的数据，靠解决其他领域的问题而存在和发展。按萨维奇 (L. J. Savage) 的说法：“统计学基本上是寄生的。靠研究其他领域内的工作而生存。这不是对统计学的轻视，这是因为对很多寄主来说，如果没有寄生虫就会死。对有的动物来说，如果没有寄生虫就不能消化它们的食物。因此，人类奋斗的很多领域，如果没有统计学，虽然不会死亡，但一定会变得很弱。”^① 看上去统计似乎被边缘化了，实际上这也正说明了统计在各学科领域的独特地位和作用，同时表明统计作为一门独立存在的学科而具有的特点。

统计学研究的是数据，只要有数据的地方，就需要用统计方法进行分析。一堆数字不去分析，它也仅仅是数字而已，没什么价值。要分析这些数据，就一定要用到统计方法。没有数据，统计学就没有存在的必要了。

1.2.2 统计学提供研究数据的方法

按统计学家 C. R. 劳的说法：“今天，统计学已发展成为一门媒介科学，它研究的对象是其他学科的逻辑和方法论——做出决策的逻辑和试验这些决策的逻辑。统计学的未来依赖于向其他学习领域内的研究者正确传授统计学的观点；依赖于如何能够在其他知识领域内将其主要问题模式化。”^② 因此，在他看来，统计学是一门科学、一种工艺和一门艺术。

统计学是一门科学。它提供一套方法和技术，这些方法和技术并不是一成不变的，使用者在给定的情况下必须根据所掌握的专门知识选择使用这些方法，而且，如果需要，还要进行必要的修正。统计方法是通用的数据分析方法，这些方法不是为某个特定的问题领域而构造的。

统计学是一种工艺。如同工业生产过程中的质量控制程序一样，统计方法是在为保证产品达到所希望的质量和保持其稳定性的管理系统中建立起来的。统计方法也能用于控制、减少和考察不确定性。

统计学是一门艺术。它提供一种归纳推理的方法，推理就是一种艺术。既然是归纳推理，就不能保证结论百分之百正确，就不能没有争议。怎样让别人看懂并理解统计结论，就要看统计表达这些结论的技巧和艺术性了。

统计学提供的是一套通用于所有学科领域的的数据方法。它是为自然科学、社会科学的多个领域而发展起来的，它为多个学科提供了一种通用的数据分析方法。从某种意义上说，统计仅仅是一种数据分析的方法。与数学一样，统计学是一种工具，是一种数据分析的工具。

统计研究数据所使用的方法通常分为描述统计 (descriptive statistics) 和推断统计 (inferential statistics) 两大类。描述统计研究的是数据收集、处理、汇总、图表描述、概括与分析等统计方法。推断统计研究的是利用样本数据来推断总体特征的统计方法。如何

^{①②} C. R. 劳. 统计与真理：怎样运用偶然性. 北京：科学出版社，2004.

划分其实并不重要,重要的是当你面对所研究的数据时,如何选择适当的统计方法进行分析,并对结果作出合理解释。

1.2.3 统计方法不是万能的

无论是作为一个工商管理人员,还是一个研究人员,你都会面对大量的数据,也都需要分析这些数据。通过分析找到隐藏在数据里面的有用信息。比如,你知道一个地区每个家庭的收入数据,难以给出这个地区收入状况的一个概括性认识;你知道20只灯泡的使用寿命、知道50件产品的合格率,这显然不够,因为你要知道的是这批灯泡的使用寿命、这批产品的合格率。要得到这样的结果,就需要用统计方法去分析。

但是,统计并不是万能的,它不能解决你面临的所有问题。吸烟能引起肺癌,这是一个统计结论,但吸烟为什么能引起肺癌,就不是统计所能回答的问题。统计能帮助你进行数据分析,并从分析中得出某种结论,但对统计结论的进一步解释,则需要你的专业知识。对于你所面对的数据,统计并不能告诉你应该用什么方法去分析,而是你自己选择你认为适合的方法。这就需要你除了具备统计知识外,还应具备你所研究问题领域的专业知识。用灵活的头脑分析数据,是统计方法所一贯强调的。

大多数统计方法都有一定的前提。有些人在使用统计方法时,往往忽视这些前提,特别是在社会科学领域,多数数据都是非实验性的,不一定能满足统计的假定。教条地使用统计方法,往往不能得到预期的结果。这不是统计方法的错,而是你的数据不满足统计方法使用的前提。不好的数据,再好的统计方法也无济于事;再好的数据,错误地选择所用的方法,统计不能得到你要的结论;如果你希望证实早已存在于你心中的某种结论,再好的数据,再好的方法,统计对此也无能为力。

统计不能提供想要的一切技巧和方法。当你把它用到自己的研究领域时,统计能做的和你所需要的之间或许还有差距。针对你所研究问题的特殊性,你需要灵活使用统计方法,而不是教条地照搬。必要时你需要对统计方法做出修正,以适应你所研究的问题。

统计是一种分析数据的工具。当你不需要这种工具时,它对你就是没有用的。当你需要它时,它也只能帮你做它能做的事情。你不能指望统计成为你解决问题的灵丹妙药。

1.3 怎样获得数据

统计学研究数据,就要先有数据。数据是什么?到哪儿去找数据?

1.3.1 变量与数据

观测一个企业的销售额,你会发现这个月和上个月有所不同。观测股票市场上涨股票的家数,今天与昨天数量不一样。观测一个班学生的生活费支出,一个人和另一个人不一样。投掷一枚骰子观测其出现的点数,这次投掷的结果和下一次也不一样。这里的“企业

销售额”“上涨股票的家数”“生活费支出”“投掷一枚骰子出现的点数”等就是变量 (variable)，它们的特点是从一次观测到下一次观测会出现不同结果。把观测到的结果记录下来就是数据 (data)。

“企业销售额”“上涨股票的家数”“生活费支出”“投掷一枚骰子出现的点数”这些变量可以用阿拉伯数字来记录其观测结果，这样的变量称为数值变量 (metric variable) 或定量变量 (quantitative variable)。定量变量的观测结果称为数值型数据 (metric data) 或定量数据。但你要观测人的性别、企业所属的行业、学生所在的学院等，这些变量的观测结果就不是数字，而是表现为不同的类别。比如“性别”表现为“男”或“女”；“企业所属的行业”表现为“制造业”“零售业”“旅游业”，等等；“学生所在的学院”则可能是“商学院”“法学院”，等等，这些表现为不同类别的变量称为类别变量 (categorical variable)。类别变量的观测结果就是类别数据 (categorical data)。由于类别数据在坐标轴上的位置是任意的，因此也称为无序类别数据。如果类别具有一定的顺序，这样的类别变量也称为顺序变量 (rank variable)，相应的观测结果就是顺序数据 (rank data)，也称为有序类别数据。比如考试成绩按等级分为优、良、中、及格、不及格，一个人对事物的态度分为赞成、中立、反对。这里的“考试成绩等级”“态度”等就是顺序变量。类别变量和顺序变量也统称为定性变量 (qualitative variable)。

1.3.2 怎样得到一个样本

从哪里取得所需的数据呢？对大多数人来说，研究社会科学问题，可以使用已有的数据。比如公开出版或公开报道的数据，统计部门公开出版的各种统计年鉴，分布在各种报纸、杂志、图书、广播、电视传媒中的数据，其他管理部门已有的数据，等等。也可以在网络上获取所需的数据，比如各种金融产品的交易数据、官方统计网站的各种宏观经济数据，等等。

当已有的数据不能满足需要时，可以亲自去调查。比如，你了解全校学生的生活费支出状况，可以从中抽出一个样本获得样本数据。这里“全校所有学生”是你所关心的总体 (population)，它是包含所研究的全部个体 (数据) 的集合。从全校学生中抽取 200 人进行调查，这就是一个样本 (sample)，它是从总体中抽取的一部分元素的集合。构成样本的元素的数目称为样本量 (sample size)。

怎样获得一个样本呢？要在全校学生中抽取 200 人组成一个样本，如果全校学生中每一个学生被抽中与否完全是随机的，而且每个学生被抽中的概率是已知的，这样的抽样方法称为概率抽样。概率抽样方法有简单随机抽样、分层抽样、系统抽样、整群抽样等。

简单随机抽样 (simple random sampling) 是从含有 N 个元素的总体中，抽取 n 个元素组成一个样本，使得总体中的每一个元素都有相同的机会 (概率) 被抽中。采用简单随机抽样时，如果抽取一个个体记录下数据后，再把这个个体放回到原来的总体中参加下一次抽选，叫做重复抽样 (sampling with replacement)；如果抽中的个体不再放回，再从剩下的个体中抽取第二个元素，直到抽取 n 个个体为止，这样的抽样方法叫做不重复抽样



(sampling without replacement)。由简单随机抽样得到的样本称为简单随机样本 (simple random sample)。

分层抽样 (stratified sampling) 也称分类抽样，它是在抽样之前先将总体的元素划分为若干层 (类)，然后从各个层中抽取一定数量的元素组成一个样本。比如，要研究学生的生活费支出，可先将学生按地区进行分类，然后从各类中抽取一定数量的学生组成一个样本。分层抽样的优点是可以使样本分布在各个层内，从而使样本在总体中的分布比较均匀。

系统抽样 (systematic sampling) 也称等距抽样，它是先将总体各元素按某种顺序排列，并按某种规则确定一个随机起点，然后，每隔一定的间隔抽取一个元素，直至抽取 n 个元素组成一个样本。比如，要从全校学生中抽取一个样本，可以找到全校学生的花名册，按花名册中的学生顺序，用随机数找到一个随机起点，然后依次抽取就得到一个样本。

整群抽样 (cluster sampling) 是先将总体划分成若干群，然后以群作为抽样单元从中抽取部分群组成一个样本，再对抽中的每个群中包含的所有元素进行观测。比如，可以把每一个学生宿舍看作一个群，在全校学生宿舍中抽取一定数量的宿舍，然后对抽中的宿舍中的每一个学生进行调查。整群抽样的误差相对要大一些。

软件应用

1. Excel【数据分析】工具的安装 (Office 2013)

第 1 步：在 Excel 工作表界面中点击【文件】，点击【选项】，在弹出的对话框中选择找到【分析工具库】选项并单击。

第 2 步：点击【转到】，在弹出的对话框中选中【分析工具库】，然后单击【确定】，即可完成安装。

2. 用 Excel 中的【数据分析】工具抽取随机样本

如果用于抽取样本的元素是类别数据，比如学生名单，需要先将类别数据用数字代码来表示（数值型数据不用指定代码）。然后按下列步骤操作：

第 1 步：点击【工具】，并点击【数据分析】，在【数据分析】选项中选择【抽样】。

第 2 步：在【抽样】对话框中的【输入区域】中输入代码区域（数值型数据直接输入数据区域）；在【抽样方法】中单击【随机】；在【样本数】中输入需要抽取的样本量；在【输出区域】中选择抽样结果放置的区域；单击【确定】后即得到要抽取的一个随机样本。

3. 用 Excel 中的【RANDBETWEEN】函数生成位于两个指定数之间的一个随机数

第 1 步：在 Excel 表格界面中点击插入函数【fx】。

第 2 步：在“函数分类”中点击【全部】选项，并在“函数名”中点击【RANDBETWEEN】函数，然后单击【确定】。

第 3 步：在【Bottom】输入指定的最小整数；在【Top】输入指定的最大整数；单击

【确定】即可得到一个随机数（要得到多个随机数，向下复制即可）。

4. 用 Excel 中的【RAND】函数生成位于 0~1 之间的均匀分布随机数

第 1 步：在 Excel 表格界面中点击插入函数【fx】。

第 2 步：在“函数分类”中点击【全部】选项，并在“函数名”中点击【RAND】函数，然后单击【确定】即可得到一个随机数（要得到多个随机数，向下复制即可）。

习 题

1.1 请举出统计的一些应用领域。

1.2 举例说明定量变量和定性变量。

1.3 你怎样理解统计的研究对象？

1.4 获得数据的概率抽样方法有哪些？

1.5 指出下列变量的类型：

(1) 年龄。

(2) 性别。

(3) 汽车产量。

(4) 员工对企业某项改革措施的态度（赞成、中立、反对）。

(5) 购买商品时的支付方式（现金、信用卡、支票）。

1.6 一家研究机构从 IT 从业者中随机抽取 1 000 人作为样本进行调查，其中 60% 的人回答他们的月收入在 5 000 元以上，50% 的人回答他们的消费支付方式是使用信用卡。

(1) 这一研究的总体是什么？

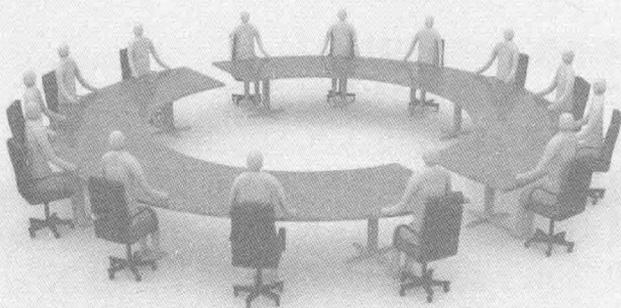
(2) 月收入是类别变量、顺序变量还是数值变量？

(3) 消费支付方式是类别变量、顺序变量还是数值变量？

1.7 一项调查表明，消费者每月在网上购物的平均花费是 200 元，他们选择在网上购物的主要原因是价格便宜。

(1) 这一研究的总体是什么？

(2) “消费者在网上购物的原因”是类别变量、顺序变量还是数值变量？



第 2 章

用图表和统计量看数据

图并没有说谎，是说谎者在画图。

——本杰明·迪斯雷利 (Benjamin Disraeli)

当你获得了一个地区各年的 GDP（国内生产总值）数据，如何观测经济的走势？当你有一个班级学生考试分数的数据，如何知道全班学生的学习状况？如果你是企业的薪酬设计人员，如何根据已有的职工工资数据，进行合理的薪酬设计？要回答这些问题，首先需要弄清楚你要用这些数据做什么？你关心这些数据的哪些特征？一堆杂乱的数据看不出它有什么特征，也很难用它说明问题。当你把这些数据用图形展示出来或计算出它们的平均数，就可能有令人惊喜的发现。本章介绍的用图表和统计量看数据会告诉你怎样对数据做初步的描述性分析。

2.1 用图表描述数据

把数据汇总在一张表格里面，用它来看一组数据的分布状况，这就是频数分布（frequency distribution）表。当然，你也可以把数据画成图，通过图形来看数据的分布。

2.1.1 用图表展示定性数据

如果你想知道不同品牌饮料的市场占有率、一所大学不同学院或专业的学生人数、一个社会中不同收入阶层的人数、一个地区不同类型的企业的数量，等等，这很容易，只要把

各个类别列出来并给出各类别的数据个数就可以了，这就是一张频数分布表。通过频数分布表可以看出不同类型数据的分布状况。一组数据的分布包含了很多有用的信息。下面通过一个例子说明怎样用 Excel 生成一张频数分布表。

【例 2.1】一家市场调查公司为研究不同类型饮料的市场占有率，对随机抽取的一家超市进行调查。调查员在某天对 50 名顾客购买饮料的类型进行了记录，如果一个顾客购买某一类型的饮料，就将这一类型的饮料记录一次。下面的表 2—1 就是记录的原始数据。

表 2—1 顾客购买的饮料类型

绿茶	碳酸饮料	绿茶	果汁	矿泉水
矿泉水	绿茶	碳酸饮料	矿泉水	碳酸饮料
绿茶	碳酸饮料	碳酸饮料	其他	绿茶
碳酸饮料	其他	绿茶	碳酸饮料	其他
矿泉水	矿泉水	矿泉水	其他	矿泉水
碳酸饮料	绿茶	绿茶	果汁	果汁
果汁	绿茶	碳酸饮料	碳酸饮料	碳酸饮料
碳酸饮料	其他	矿泉水	果汁	其他
矿泉水	碳酸饮料	其他	碳酸饮料	矿泉水
碳酸饮料	绿茶	其他	果汁	绿茶

表 2—2 是根据表 2—1 用 Excel 建立的频数分布表。

表 2—2 不同类型饮料的频数分布

	A	B	C
1	饮料类型	频数	百分比%
2	果汁	6	12
3	矿泉水	10	20
4	碳酸饮料	15	30
5	绿茶	11	22
6	其他	8	16
7	总计	50	100

利用图形来看数据的分布更直观和形象。定性数据的常用图示方法主要有条形图 (bar chart) (或称柱形图 (column chart))、饼图 (pie chart)、环形图 (doughnut chart) 等。图 2—1 是不同类型饮料的条形图。

饼图主要用于表示一个样本 (或总体) 中各组成部分的数据占全部数据的比例, 对于研究结构性问题十分有用。例如, 根据表 2—2 中的数据绘制的饼图如图 2—2 所示。

把一个地区的人口按高收入、中等收入和低收入划分成三部分。如果要比较五个地区的人口构成, 你需要绘制五个饼图, 这种做法既不经济也不便于比较, 能否用一个图形比较出五个地区的人口构成呢? 把饼图叠在一起, 挖去中间的部分就可以了, 这就是环形图。环形图与饼图类似, 但又有区别。环形图中间有一个“空洞”, 样本中的每一部分数据用环中的一段表示。饼图只能显示一个样本各部分所占的比例, 而环形图则可以同时绘制多个样本的数据系列, 每一个样本的数据系列为一个环。因此环形图可显示多个样本各部分所占的相应比例, 从而有利于进行比较研究。