

当代国外语言学与应用语言学文库

Statistics in Language Studies

语言研究中的统计学

Anthony Woods

Paul Fletcher

Arthur Hughes

外语教学与研究出版社

Foreign Language Teaching and Research Press

剑桥大学出版社

Cambridge University Press



当代国外语言学与应用语言学文库

Statistics in Language Studies

语言研究中的统计学

Anthony Woods, Paul Fletcher and Arthur Hughes 著

林连书 导读

外语教学与研究出版社

剑桥大学出版社

(京)新登字 155 号

京权图字: 01-2000-0157

图书在版编目(CIP)数据

语言研究中的统计学/(英)安东尼(Anthony, W. S.)等著;林连书导读. - 北京:外语教学与研究出版社, 2000. 8

ISBN 7-5600-1926-9

I. 语… II. ①安… ②林… III. 统计学-应用-应用-语言-研究-英文 IV. H0-059

中国版本图书馆 CIP 数据核字(2000)第 34967 号

English edition © Cmbridge University Press 1986

All rights reserved. No part of this publication may be reproduced, stored or transmitted by any means without the prior permission of the publishers.

This edition of Statistics in Language Studies by Anthony Woods, Paul Fletcher and Arthur Hughes is published by arrangement with the Syndicate of the Press of the University of Cambridge, Cambridge, England. It is for sale in the People's Republic of China only. Not for export elsewhere.

本书由剑桥大学出版社授权外语教学与研究出版社出版

版权所有 翻印必究

语言研究中的统计学

Anthony Woods, Paul Fletch and Arthur Hughes 著

林连书 导读

* * *

责任编辑: 李景峰

出版发行: 外语教学与研究出版社

社 址: 北京市西三环北路 19 号 (100089)

网 址: <http://www.fltrp.com.cn>

印 刷: 北京市鑫鑫印刷厂

开 本: 650×980 1/16

印 张: 23

版 次: 2000 年 8 月第 1 版 2000 年 8 月第 1 次印刷

印 数: 1—5000 册

书 号: ISBN 7-5600-1926-9/G·829

定 价: 29.90 元

* * *

如有印刷、装订质量问题出版社负责调换

当代国外语言学与 应用语言学文库



专家委员会

主任 王宗炎

副主任 (以姓氏笔画为序)

刘润清	吴一安	李朋义	沈家煊	陆俭明
陈国华	胡文仲	胡壮麟	徐烈炯	桂诗春
顾曰国	戴炜栋			

委员 (以姓氏笔画为序)

文秋芳	方立	王才仁	王立弟	王克非
王初明	王逢鑫	王嘉龄	史宝辉	宁春岩
田贵森	申丹	刘世生	朱永生	何兆熊
何自然	张绍杰	张柏然	张德禄	李宇明
李延福	李行德	李筱菊	杨永林	杨信彰
杨惠中	杜学增	汪榕培	邵永真	陈治安
周流溪	林连书	罗选民	姚小平	祝畹瑾
徐盛桓	秦秀白	贾玉新	顾阳	高远
高一虹	黄国文	惠宇	董燕萍	蒋祖康
韩宝成	蓝纯	潘永樑		

策划 霍庆文

Preface by Halliday

Foreign Language Teaching & Research Press is to be congratulated on its initiative in making these publications in linguistics available to foreign language teachers and postgraduate students of linguistics in China.

The books are a representative selection of up-to-date writings on the most important branches of linguistic studies, by scholars who are recognized as leading authorities in their fields.

The availability of such a broad range of materials in linguistics will greatly help individual teachers and students to build up their own knowledge and understanding of the subject. At the same time, it will also contribute to the development of linguistics as a discipline in Chinese universities and colleges, helping to overcome the divisions into “English linguistics”, “Chinese linguistics” and so on which hinder the progress of linguistics as a unified science.

The series is to be highly commended for what it offers to all those wanting to gain insight into the nature of language, whether from a theoretical point of view or in application to their professional activities as language teachers. It is being launched at a time when there are increasing opportunities in China for pursuing linguistic studies, and I am confident that it will succeed in meeting these new requirements.

M. A. K. Halliday
Emeritus Professor
University of Sydney

王宗炎序

近年来，国际交往日益频繁，国际贸易急速发展，出现了一种前所未有的现象：学外语、教外语、用外语的人多了；研究语言学和应用语言学的人多了；开设这方面专业的高校也多了，语言学硕士生和博士生也多了。就是不以此为专业，学习语言学和应用语言学的也不乏其人。为了给从事这个专业的师生提供便利，同时又帮助一般外语教师、涉外工作者以及汉语研究者开阔思路，扩大视野，提高效率，我们献上这套内容崭新而丰富的丛书——英文版《当代国外语言学与应用语言学文库》。

文库首批推出 54 部外国英文原著，它覆盖了语言学与应用语言学 28 个分支学科。这批书是我们与各地有关专家教授反复研究之后精选出来的。出版这样大规模的语言学与应用语言学丛书，这在我国语言学界和外语教学界是破天荒第一次。

我们这样做，抱着什么希望呢？总的说来，是遵循教育部关于加强一级学科教育的指示，在世纪之交，推出一套书来给中国的外语教育领航，同时也给一般外语工作者和汉语研究者提供信息，拓宽思路。

我们希望这个文库能成为进一步带动外语教学改革和科研的发动机；我们希望它能成为运载当代外国语言学理论、语言研究方法和语言教学方法来到中国的特快列车；我们希望，有了这套书，语言学与应用语言学专业师生就能顺利地进行工作；我们希望，通过读这套书，青年外语教师和外语、汉语研究者能迅速把能力提高，把队伍不断扩大。

以上是我们的愿望，可是从广大读者看来，这个文库

是否真的有出台的有必要呢？我们想，只要大家看一下今天的客观情况，就知道这套书有填空补缺的作用，是让大家更上一层楼的扶梯。

我们跟许多人一样，认为国内的外语教学和语言学与应用语言学研究是成绩斐然的，但是某些不足之处也无庸讳言。

在语言研究方面，有大量工作还等着大家去做。汉语语法研究，过去由于结构主义的启示，已经成绩卓著，可是现在虽则引进了功能主义，还看不出什么出色的成果。语料语言学是新兴学科，在我国刚刚起步，机器翻译从50年代就有人搞，然而其进展至今不能令人满意。

在语言理论方面，我们不时听到一些片面的、所见不全的论调。有人说，1957年前西方根本没有什么理论语言学，其创始者是 Chomsky；也有人说，语言纯属社会文化范畴；还有人说，搞语言研究只有量化方法才是科学方法，定性方法不值得一提。

谈到外语教学，某些看法做法是分明不值得赞许的。有人以为交际教学只管听说，不管读写，也有人以为教精读课就是教阅读，不管口语。在分析课文时老师满堂灌，学生开口不得，是常见的；教听力课时老师只管放录音，对学生不给半点提示点拨，也并非罕有现象。

上述这些缺点，我们早有所知，现在我们更加明白，必须力图改进，再也不能安于现状了。为了改进，我们就得参考国外的先进理论，借鉴国外的有效措施。眼前这个文库，就是我们上下求索的结果。

在编辑这个文库时，我们在两方面下了功夫。

一方面，在选书时，我们求全，求新，求有代表性和前瞻性。我们不偏爱一家之言，也不只收一家外国出版社之书。语言学与应用语言学的主干学科固然受到了应有的重视，分支学科可也不忽视。语料语言学、语言统计学是新兴学科，我们收入了专著；句法学、语义学久已有人研

究，我们也找到了有关的最新著作。

另一方面，我们邀请了国内知名的博士生导师、硕士生导师为各书撰文导读，为读者铺平道路。语言学和应用语言学专著包罗宏富，初学者读起来可能觉得茫无头绪。为了助他们一臂之力，本文库中每一种书我们都请专家写了一万字左右的导读材料。哪怕书中内容比较陌生，谁只要在读书前看一下导读材料，读书后把材料再看一遍，一定能弄清脉络，掌握要点。

在结束本文时，我们想向爱好泛读的人们提个建议。语言和社会生活息息相关；我们靠语言与他人协作；通过语言继承传统文化，接受外国先进思想和科学知识；利用语言来教育下一代，帮助他们创造美好的未来；语言又反过来表达着我们的个性和我们充当的各种角色。学一点语言学和应用语言学，有助于增强我们的语言意识，对我们的工作和生活都是有利的。我们不妨把此事作为一个项目，列入自己的日程。持之以恒，必有所获。

王宗炎

中山大学教授
博士生导师

导 读

《语言研究中的统计学》是一本全面讨论统计学原理、方法及其在语言研究中应用的专著，作者是统计学家 Anthony Woods 和两位语言学家 Paul Fletcher 和 Arthur Hughes。该书于 1986 年在剑桥大学出版社初版之后，受到语言界的重视和欢迎，并分别于 1989、1991、1993 和 1996 年再版。在我国，不少外语院校和其他高校外语专业将其定为语言学研究生的教材和参考书，一些关于科研方法和统计学应用的书或论文也以它为参考书，其影响是不言而喻的。

读者也许会问：统计学向来是自然科学和社会科学的研究工具，为什么作为人文学科之一的语言学也需要统计学呢？我认为有 3 个原因：第一是语言学研究范围的扩大。现代语言学研究已扩大到包括语言习得、功能语言学、社会语言学、心理语言学、语言教学法、语言测试、数理语言学、语病学、神经语言学、话语分析、语料库研究、甚至研究方法等范围，所涉及的往往已不限于单纯的语言问题，而是扩展到社会学、心理学、人类学、信息科学、概率论、计算机应用等其他学科。现代语言学研究正在成为一门综合性学科，部分已成为交叉学科和边缘学科，所以，它应该从自然科学和社会科学研究方法中吸收营养，以应付日益增加的研究问题的需要。第二是语言学研究方法的不断更新。以往的语言学研究，往往偏重于文献资料的研究，包括对文献资料的综合、总结、分析、分类、比较、对比等。但现代语言学研究中，人们已不满足于对二手材料的研究。人们要建立语言理论体系和探索语言习得与运用的模式，要对实际发生的语言现象进行描写，要通过语言学实验对某些理论和假设进行验证，找出规律，所以，他们更多地用到第一手材料，并广泛运用科学实验、社会调查、自然观察等方法。在这些方法中，概率论和统计学是归纳分析数据、进行统计推理的重要工具。第三是语言学研究的科学性的要求。现代语言学研究比较强调信度和效度，信度主要是指研究的稳定性和可重复性，效度主要是指研究的可解释性和可推广性，缺乏这些特性的研究是不可信或无效的。而这些特性的检验，也是以统计学为基础的。

《语言研究中的统计学》一书共分 15 章。第一章是一个引子，通过几个例子说明统计学在归纳研究结果和进行统计推理中的作用，实际上是在介绍统计学的两个组成部分：描述性统计学和推论性统计学。本章最后部分对实验研究的过程也有所提及。

第二、三章谈的是描述性统计学的内容。第二章介绍如何以表格和图形的方式将研究结果形象地展示给读者，其中包括建立频率分布表（频数、频率、累计频率等）、条形图、频率分布直方图和累计频率曲线等。第三章介绍分析数据的集中趋势和离散趋势的统计量。前者包括平均数、中数和众数；后者包括全距、标准差和方差。掌握这些统计量的基本意义和计算方法是学习其余各章的基础。

第四章提出统计推理的概念及其重要性。实验研究的目的往往是通过分析样本来推知总体的情况，找出某些规律，所以，抽样的代表性和推理的合理性具有重大的意义。而要进行统计推理，就必须掌握有关概率的基本知识，第五章就讨论了概率、随机变量和随机抽样等问题。

第六章讨论总体的数学模型，即概率分布。总体的概率分布有若干种模型，其中最常用的是正态分布。作者介绍了正态分布曲线的形状、性质和如何使用标准正态分布表，最后介绍了抽样分布的概念。抽样分布是学习推论性统计学的基础。第七章主要介绍如何利用抽样分布对总体的参数进行估计，即从观察到的样本平均数的大小来推知不能观察到的总体平均数落在什么范围内。

从第八章至第十二章，每章都涉及推论性统计学的内容。第八章举例说明如何用 Z 检验法和 t 检验法进行单样本假设检验，即研究某个样本是否属于某个总体。第九章讨论 χ^2 检验法在假设检验中的两个作用：1. 检验样本的频数分布是否与某个数学模式（如正态分布）相似合；2. 检验两个或多个变量间的相互独立性，即以频数表示的某个变量的分布情况是否与另一个变量有关。第十章介绍相关系数的意义、作用和检验其显著性的方法。第十一章讨论如何用 t 检验法来检验两个样本平均数之间差异的显著性，即差异是由随机因素引起的，还是由某个研究中的因素引起的，这包括独立样本 t 检验法和配对样本（或称相关样本）t 检验法两种情况。第十二章介绍如何用方差分析对两个或多个样本平均数的差异的显著性进行检验，方差分析可能是单向的，也可能是多向的。学习了上述各章，对推论性统计学就会有比较清楚的了解。

从第十三章至第十五章，每章都涉及多变量研究问题。第十三章所谈的线性回归和第十章所谈的相关系数，都是进行多变量分析的有效方法。第十四章讨论聚类分析和判别分析。第十五章介绍主成分分析和因素分析。这些章为进行多变量分析提供了非常有用的工具。

本书的特点是：1. 条理分明，层次清楚，令读者了解到统计学的全貌；2. 深入浅出，解释详细，每章都有大量的例子来解释一些基本概念、统计技术及其应用，避开复杂的数学推导，这样，在有教师指导的情况下，读者不难掌握全书的主要内容；3. 每章配有练习，使读者可对每章的内容进行回忆，通过做练习巩固学到的知识，并避免只会空谈别人的研究而不会自己动手做研究的现象。书后附录附有练习答案和可供查阅的统计表，统计表是进行统计推理时必不可少的工具。针对本书的特点，读者在阅读本书的时候，应该做到：1. 按步就班，循序渐进，在掌握前一章的内容之后，再开始阅读下一章的内容；2. 先弄清一些基本概念，然后了解某些公式的由来、意义和应用，最后学会能对得到的结果进行解释；3. 在可能的情况下，参加以本书为教材或主要参考书的统计学或研究方法课，以增进对所学内容的理解和提高学习效率；四、多做练习，多进行实验研究的实践，这是巩固所学内容、学以致用最好方法。

我们感谢 Woods 等人为语言学研究者写了一本谈统计学在语言研究中应用的书。如果我国的语言研究者能从阅读这本书中获益，我想无论读者或作者都会高兴的。

第一章 为什么语言学家需要统计学？

本章讨论语言学家需要学习统计学的原因。作者指出，在不少语言学研究领域，很多数据需要用统计学来进行处理。统计学包含两个主要方面：一个是描述性统计学（descriptive statistics），用以归纳复杂的数据，即把数据表达为某些具有代表性的统计量（如平均值、标准差等）；另一个是推论性统计学（inferential statistics），有了推论性统计学，从样本（sample）中得到的统计量可用来推知总体（population）的情况，这样就便于语言学家从研究中得到一些规律性的东西。

作者用两个例子来说明描述性统计学的应用。一是爆破音从开

始发音到后面元音出现之间所花的时间 (VOT, 即 voice onset time)。语音学家找了 20 个人, 让他们发出 10 个以 /p/ 开头的单词, 并重复 10 次。可以想象, 人与人之间, 或同一个人在重复读同一个音时, 发音长度都是有差异 (variation) 的, 所以, 必须进行如下统计: 不同的人发音长度平均是多少, 有多大的差异; 同一个人对不同单词的发音的平均长度是多少, 有多大的差异。另一个例子是心理语言学家研究能力倾向测验 (aptitude test) 与成绩测验 (achievement test) 之间的关系。其关系可用相关系数来表示, 计算方法详见第 10 章。

推论性统计学常常用来比较两组或多组数据间有没有显著差异, 即它们平均数间的差异是由随机因素引起的, 还是由实验条件因素引起的, 也就是说, 它们之间是否有真正的差异 (real difference)。例如, 在 /p/ 开头的单词和 /b/ 开头的单词中, /p/ 的 VOT 和 /b/ 的 VOT 是否有显著的差异。根据有关文献, 清爆破音的 VOT 比浊爆破音的长。儿童在语音习得过程中, 常常混淆清浊爆破音的发音。所以, 研究这个语音问题还是有意义的。

作者提出了进行诸如 VOT 之类的研究需要注意的几个问题: 1. 研究对象的选择; 2. 样本量 (sample size) 的大小; 3. 收集数据和进行测量之后, 所得的结果以何种形式表示; 4. 描述性统计数字带来的结论; 5. 如何将样本中得到的结果推广到总体中。

基于上述讨论, 作者指出了语言学家学习统计学的必要性: 第一, 能够对使用统计技术的文献进行评价, 不但能看懂这些文章, 而且能对方法的科学性进行评论; 第二, 能够设计自己的研究方案, 本书所介绍的统计分析技术可帮助读者实现这个目的。

第二章 表格和图形

首先, 作者将数据分为分类数据 (categorical data) 和数量数据 (numerical data) 两种。不同类型数据的列表和制图的方法是不一样的。

其次作者举例说明分类数据列表和制图的方法。例如, 研究美国 560 个有语言障碍 (language-impaired) 的病人, 其中男性 364 人, 女性 196 人。语言障碍可分成 4 类: 口吃 (stuttering)、发音不良 (phonological disability)、具体语言表达障碍 (specific language

disorder)、听觉受损 (impaired hearing)。可以根据这 4 类列出男女病人每一类的频数 (frequency) 和频率 (relative frequencies), 进而根据这些频数和频率列出男女病人语言障碍对比的条形图 (bar chart)。

与分类数据不同, 数量数据是连续型的 (continuous), 所以, 其列表和制图方式也不同。例如书中第 16 页表 2.5 (a) 列出 1980 年 6 月剑桥英语水平考试中 108 个考生的成绩, 这些成绩从 117 分至 255 分不等。最高分与最低分之差称为全距 (range)。在列表时, 一般将全距分成 8—15 个分数段 (class interval), 数出在每个分数段里分数的个数, 这些个数便是每个分数段的频数。将每个频数除以考生总人数, 所得便是频率。将某个分数段前面的频数全部加起来, 就是累积频数 (cumulative frequency)。将其前面所有的频率加起来, 便是累积频率。

根据频数或频率, 可以作出如第 17 页所示的频率分布直方图 (histogram)。频率分布直方图将各分数段的频数或频率形象地以条形的高度表示出来。用频率表示的直方图是最常用的, 它广泛地用来研究样本的频率分布。同时, 当样本数量无限增加时, 频率分布会形成某种规律, 例如正态分布, 或其他类型的分布。也可根据累积频率作出累积频率曲线 (cumulative frequency curve), 累积频率曲线的作用之一是用来找出某个分数的百分位 (percentile)。例如在第 18 页, 找到横坐标轴上 150 分的位置, 从该点作垂直线与曲线相交于一点, 再从这一点向纵坐标轴作一与横坐标轴相平行的直线, 与纵坐标轴相交于一点, 这一点是 0.1, 它就是 150 分的百分位, 即得 150 分的考生成绩比 10% 的考生高, 比 90% 的考生低。

最后作者介绍了如何列出更复杂的表格的方法和一些特例。只要有前面的基础, 理解、使用表格和图形并不困难。

第三章 归纳性测量

描述性统计学对样本数据常常作两方面的分析: 一是数据的集中趋势 (central tendency), 二是数据的离散趋势 (dispersion 或 variability)。本章讨论了这两个问题。

对数据的集中趋势, 作者首先介绍了中数 (median) 的概念。中数就是位于一组数据中间的数, 它高于一半数据, 但又低于另一

半数据，也就是说，它处于第 50 个百分位。算术平均数 (arithmetic mean) 是读者都熟悉的统计量，其计算公式是

$$\bar{X} = \frac{1}{n} \sum X$$

式中的 \bar{X} 是平均数， X 代表各个数据， \sum 是连加符号，读作 sigma，表示将所有数据连加。 n 是数据的个数。这样，整个式子表示所有数据的和除以数据的个数，便是平均数。

作者接着比较了平均数和中数的特点，并指出在什么情况下，哪一个统计量是更合适的量度。如果数据的分布是中间大、两头小，或者更精确地说，近似于正态分布，那么平均数和中数很接近，都是数据集中趋势的良好量度。但是，尽管平均数是计算最容易、使用最方便的统计量，它却容易受到极端数值 (extreme values) 的影响。例如计算 30 个家庭的平均收入，如果其中有一个家庭是百万富翁，得到的平均数就不是典型的、能够代表这 30 个家庭的统计量。这时用中数比较好，因为中数比较固定，不会受到极端数值的影响。

数据的分布有时是向一边倾斜的 (skewed)，例如，第 33 页中图 3.4 (a) 的频率分布直方图向右倾斜，图 3.4 (b) 的频率分布直方图向左倾斜。这表示：在英语考试中，第一种情况的考生考低分的特别多，第二种情况的考生考高分的特别多。这时平均数和中数都不是很好的量度。所以，作者引入了第 3 个统计量——众数 (mode)。众数表示频数最高的数据。

对数据的离散趋势，作者首先介绍了中四分间距 (interquartile distance or range) 的概念。把数据的个数分成相等的 4 组，即每组的数据占数据总量的 25%，每组分数段的距离就叫中四分间距。其次，作者介绍了测量离散趋势最重要的一个统计量——标准差 (standard deviation)，其计算公式是

$$s = \sqrt{\frac{\sum (X - \bar{X})^2}{n - 1}}$$

其中， X 是一个数据， \bar{X} 是平均数， $\sum (X - \bar{X})^2$ 表示将每个数据减去平均数，然后平方，再将其连加起来 (即求平方和)，再将平方和拿去平均，平均时用 $(n-1)$ 作除数而不用 n 是为了得到无偏统计量。最后求平方根，便可得到标准差。

上述公式如果不开方，便是方差 (variance)，记作 V 或 s^2 。即

$$V = s^2 = \frac{\sum (X - \bar{X})^2}{n - 1}$$

最后，作者介绍了分数标准化的方法，但他只介绍了 Z 分数。Z 分数是一种平均数为零，标准差为 1 的标准分数。计算公式很简单：

$$Z = \frac{X - \bar{X}}{s}$$

即把每个分数减去平均分数（这实际上是将每个分数向左平移 \bar{X} 个单位，使平均数变成零），再除以标准差（实际上是将得到的分数每个缩小 s 个单位，从而使标准差变为 1）。

得到 Z 分数后，可根据需要用下列公式将 Z 分数化成平均数为 500、标准差为 100 的标准分数：

$$\text{标准分数} = 500Z + 100$$

我国广东、海南等省每年高考就是根据此原理将原始分数（raw score）变成这种标准分数的。美国 TOEFL 分数的平均数为 500，标准差为 50。2000 年以后，TOEFL CBT（computer-based test）的标准分会有较大的变化。这些在本书中没有谈及，我们略为介绍，以便读者知道分数标准化的过程。

第四章 统计推理

本章简单地介绍了统计推理（statistical inference）的意义。进行统计推理，必须了解总体和样本这两个概念。总体即研究对象的全体，也就是实验结果所能推广到的最大一个实验对象群体，样本是从总体中抽出的一部分实验对象。在实际研究中，由于时间、人力、物力的限制，研究总体是很困难的。所以，往往从总体中抽出有限的一些样本来研究。这就存在下列的问题：1. 抽样（sampling）是否有代表性？如果抽样有代表性，从样本得到的某些数据就有可能推广到总体中去。例如，如果知道样本的平均数和标准差，就可以用统计学的方法推知总体的平均数和标准差落在什么范围，准确性达到 95% 或 99% 以上。2. 如何做到抽样有代表性？书中介绍了简单随机抽样（simple random sampling）的方法。所谓简单随机抽样就是总体中的每个元素（element）都有同等的机会或概率（probability）被抽到，这样做可避免抽样过程中人为的偏差（bias）。(3) 在随机抽样难以做到的时候怎么办？作者介绍了分阶段（stages）抽样的方法。

在语言学研究中，很多实际问题并不一定涉及到以某种形式进

行抽样的问题，但这并不意味着统计技术对这些研究不适用。统计技术为实验人员提供籍以测量和比较证据（evidence）的共同基础和标尺。如果通过研究得到某些有趣的结果，例如暗示某个新的假设与以前的某个假设相矛盾，那么研究人员可进而了解其样本是如何得到的和它对所得结论的效度（validity）有何影响等问题。

若仔细阅读本章特别是作者举出的例子，读者就会对抽样和统计推理的关系有比较清楚的了解。

第五章 概率

概率（probability），即事件发生的可能性的的大小，是统计学的基础。本章一开始，作者就举例解释这个术语的含义。例如，一个箱子装着编号从1至10大小形状都一样的10个小盘，有3个是红的（编号为1，2，3），有7个是白的（编号为7—10），那么，随便从箱子里抽出一个盘子，其编号为4的概率为 $\frac{1}{10} = 0.1$ ，抽到偶数号码的盘子的概率是 $\frac{5}{10} = 0.5$ ，抽到白盘子的概率是 $\frac{7}{10} = 0.7$ 。接着作者介绍了事件独立性（independence）和条件概率（conditional probability）的概念。条件概率是在事件Y发生的条件下事件X发生的概率，记作 $P(X/Y)$ 。在两个事件X和Y互相独立的情况下，他们同时发生的概率

$$P(X \text{ 和 } Y \text{ 同时发生}) = P(X) \times P(Y)$$

但是，在一般情况下，

$$P(X \text{ 和 } Y \text{ 同时发生}) = P(X/Y) \times P(Y) = P(Y/X) \times P(X)$$

这些关系在介绍 χ^2 检验法时会用到。

其次，作者又介绍了随机变量（random variable）的概念。随机变量是其取值在实验前不能确定、而要靠实验结果来确定的变量。随机变量有两种，一种是离散型的（discrete），一种是连续型的（continuous）。在作抽样分析时，离散型随机变量的图形一般用条形图，连续型随机变量的图形一般用频率分布直方图。

最后，作者介绍了简单随机抽样（simple random sampling）的原理和使用随机数字（random numbers）表的方法。例如，要从7832个实验对象中随机抽出10个，首先，将他们从1至7832编号，然后从随机数字表任何一个地方开始连续记下10个4位数（重复或

超过 7832 的舍去), 这 10 个 4 位数所代表的实验对象就是随机抽出的实验对象。这个过程相当于从摇珠中摇出的 10 个数字。

读者学习本章要记住上述几个基本概念, 这样可为下几章的学习打下基础。

第六章 建立统计总体的模型

这一章主要讨论总体的数学模型。大家知道, 总体有平均数和标准差, 这两个参数 (parameters) 一般分别用希腊字母 μ 和 σ 表示。在进行实验研究时, 从总体中抽出来的样本也有平均数和标准差, 这两个统计量 (statistics) 一般分别用罗马字母 \bar{X} 和 s 表示。一般来说, 总体的平均数和标准差是不知道的。例如, 让同一个人无限多次地读 /tu:/ 这个音, 用精密的仪器把每次清爆破音 /t/ 的 VOT 记下来, 可以发现每次读音的 VOT 都会有细微的差别, 而我们无法通过测量无限多次来获得 VOT 的平均数 μ 。我们只能让他读有限多次 (例如 10 次, 20 次等等), 这有限多次测得的 VOT 就是一个样本。从这个样本中可以算出平均数 \bar{X} 和标准差 s 。在统计学中, 我们常常用样本的统计量来推知总体的参数。

所以, 在阅读本章的时候, 必须分清样本的频率分布 (frequency distribution) 和总体的概率分布 (probability distribution)。总体的概率分布有若干种模型, 其中最重要而且最常用的是正态分布 (normal distribution), 用来表示这种分布的曲线叫做正态分布曲线 (normal distribution curve)。正态分布曲线看起来象一个倒转的钟形, 正态分布的数值有 68% 左右落在正负 1 个标准差之内, 有 95% 左右落在正负 2 个标准差之内, 而且在平均数两边的数值是对称的。为研究所有正态分布的分布规律, 可以用下列公式把所有 X 数值化成 Z 数值:

$$Z = \frac{X - \mu}{\sigma}$$

这样所有 Z 数值就会形成平均数是零、标准差是 1 的正态分布, 这种正态分布称为标准正态分布 (standard normal distribution)。标准正态分布曲线如 88 页图 6.6 所示。通过查 298 页的标准正态分布表 A2, 可以查出从 -4 至 +4 的每个 Z 数值左侧的概率, 即小于这个数值的所有数值在总体中占多大的比例。299 页的表 A3 列出 Z 数值是