



ISTIC 科学与管理译丛

# 信息专业人员 常用统计方法

用实用、轻松的方式来理解、使用和阐释统计学

---

## Statistical Methods for the Information Professional

A Practical, Painless Approach to Understanding,  
Using, and Interpreting Statistics

Liwen Vaughan(邱黎雯) 著

张爱霞 吴雯娜 等 译  
李 燕 潘晓蓉

张新民 等 校

 科学技术文献出版社

科学与管理译丛

# 信息专业人员常用 统计方法

用实用、轻松的方式来理解、使用和阐释统计学

Liwen Vaughan(邱黎雯) 著

张爱霞 吴雯娜 等 译  
李 燕 潘晓蓉  
张新民等 校

科学技术文献出版社

Scientific and Technical Documents Publishing House

北京

## 译者简介

**张爱霞** 北京大学信息管理系在读博士,中国科学技术信息研究所馆员

**吴雯娜** 硕士,中国科学技术信息研究所副研究馆员

**李 燕** 硕士,中国科学技术信息研究所馆员

**潘晓蓉** 硕士,中国科学技术信息研究所馆员

**尹盛鑫** 双学士,中国科学技术信息研究所工程师

**杨代庆** 硕士,中国科学技术信息研究所馆员

**刘春燕** 硕士,中国科学技术信息研究所馆员

**张新民** 博士,中国科学技术信息研究所副研究员,硕士研究生导师

## 译序

在这个“用数据说话”的年代，统计学基本知识与分析技能对于任何行当都非常重要。从正面说，掌握了统计学方法，就有利于收集适当的数据，开展深入的分析，得出有理有据毫不牵强的结论，从而说服别人；从反面说，不掌握统计学原理，就很容易受骗上当。曾任19世纪英国首相、同时又是小说家的本杰明·狄斯雷利(1804—1881)曾经幽默地说过，“世界上有三种类型的谎言：谎话、弥天大谎、还有统计数字”，就生动地说明了不懂统计学和误用统计学的危险。

在迈向知识社会的进程中，图书情报事业的重要性受到空前的重视，从而对图书情报人员的业务素质提出了更高的新要求，掌握基本统计学知识与技能便属于新要求之一。而我们面临的现状是，无论在中国还是在国外，图书情报学从业人员中能完善地掌握统计学方法的不多，一知半解的倒不算太少。英国诗人蒲伯(1688—1744)说，“一知半解可能是很危险的”，因为，拥有一知半解的人，往往自以为已经懂得不少了，容易率尔操觚，贻笑方家。

曾经在上海科技情报研究所工作、目前在加拿大西安大略大学任教的邱黎雯教授对于此类问题有深刻的观察与体会。为了帮助图书情报人员更有效地、更有针对性地学习和掌握统计学知识，她发奋

努力编写了这本教材。教材推出后，在国际图书情报学界获得好评，因为过去从来没有专门为图书情报人员编写的类似教材。我们组织力量将此书译出，也是为了填补中文图书市场上的这一空白，更好地满足中国图书情报专业的学生和从业人员的需求。在此，我们向本书作者邱黎雯教授表示深深的感谢和敬意。

在全球化的世界上，保持多样化显得更加重要。翻译工作是在多样性的文化之间架设桥梁的伟大工程。笔者曾写过一首打油词“西江月·译事咏”，上半阙说的是：“文化相隔如岛，翻译可作津梁。互通有无较短长，委实功德无量”，就反映了我们对翻译工作的认识。多年来，中国科学技术信息研究所一直注重科技翻译工作，推出了不少有影响的译著。遗憾的是，翻译劳动的价值迄今未受到应有的重视，翻译稿酬之低暂且不说，翻译成果在评职称、报奖时要么不被承认，要么大打折扣。但是，只要社会需要翻译作品，我们就会一如既往、义无反顾地做下去。

中国科学技术信息研究所总工程师

**武夷山**

2007年11月25日

## 译者的话

本书译自 Information Today 公司出版的邱黎雯教授所著的《Statistic Method for the Information Professional》。邱黎雯系加拿大西安大略大学教授,长期从事研究方法和统计学、网络搜索引擎、网络数据挖掘、信息检索以及情报计量学等领域的教学和科研工作。本书是作者教学和科研实践经验的总结,深入浅出地讲述统计学方法在信息科学领域中的应用,侧重统计方法的基本逻辑,而不是繁杂的数学推理;强调的是何时运用和怎样运用统计学方法以及如何理解和解释统计分析结果。对于我们而言,无论是做统计分析,还是理解在阅读文献中遇到的统计学问题,这本书都是非常有用的。

本书由张爱霞联合多人共同协作完成。各章的翻译分工如下:第1、2章和索引部分张爱霞;第3章刘春燕;第4、5章及前言部分吴雯娜;第6、7章李燕;第8、9章及第12章尹盛鑫;第10、11章杨代庆;第13章及附录部分潘晓蓉。为确保翻译质量,我们对译文进行了多次修改和校对。张新民、庞景安承担了译稿的审校工作,最后由张新民进行了终校终审。

在翻译中,为便于读者阅读查找,原著索引在译著中以中英对照的方式出现,此外,对于索引中出现的重要词汇,其在译著正文中首次或在重要位置出现时都标出了相应的英文,以便于理解。

在此,我们想对大家表示诚挚的谢意:感谢中国科学技术信息研究所总工武夷山研究员,是他推荐此书,并在百忙之余拨冗为本书作序,使得本书大为增色;感谢中国科学技术信息研究所信息资源中心主任沈玉兰研究馆员,她对本书的翻译和出版工作给予了极大的帮助和支持;感谢林明、陶冶、李建丽、蔡小英、刘智、蔡晓燕等人在本书校对中给予的支持。

译文虽然经过多次校对和修改,但由于译者水平有限,加上时间仓促,文中难免有不妥之处,我们真诚地希望同行和读者不吝赐教。

**译者**

2007年11月1日

## 关于作者

邱黎雯(Liwen Vaughan)博士于 1991 年在加拿大西安大略大学获得图书情报学博士学位。她对统计分析有着浓厚兴趣,同时,她又能以清晰和可理解的方式对复杂的概念进行解释,这使得她成为一个成功的教育家。她有 10 年以上面向不同专业和不同学位学生讲授统计学的教学经验,涉及专业包括图书情报学、新闻学和企业管理等,授课学生包括学士、硕士和博士。邱黎雯博士在情报学研究中有丰富的统计分析经验。她的研究成果发表在很多情报学期刊中,在这些成果中,她运用了许多统计方法。她的关于用 Markov 模型分析超文本信息的论文,发表在 JASIS 上,曾经吸引了全世界的博士生的关注。最近,她成功地利用 LISREL 模型(一种新的应用到情报学中的高级统计方法)定量地研究了信息对企业成功的影响。

## 致 谢

感谢我的丈夫 David Vaughan,他在本书编写过程中给予了我精神上和技术上的支持。他精心编辑了全书,不辞辛劳地把书中所有的数字和图表转换为要求的出版格式;更重要的是,他非常支持我的见解和想法。没有他的帮助,本书不可能完成。我的儿子 Ulysses,在我着手撰写此书时尚在腹中,到我完成这本书时他已经在蹒跚学步了。确切地说,他是伴随着我的写作一点点长大的。他一直都是我的灵感来源。写这本书常常使我顾不上照顾家庭,对于 David 和 Ulysses 我有着深深的负疚。谨将此书献给他们。

Liwen Vaughan(邱黎雯)

# 前　　言

## 不仅仅是又一本统计学书

统计学是一种最有用和最有力的数据分析工具，同时也是一个最令人生畏和最令人生厌的研究主题，这种现象很值得注意。之所以令人畏惧，是因为统计学著作通常充斥着大量公式、数学术语以及连篇累牍的推导和演算。本书与其他众多的同类书籍有所不同，不同之处在于其核心是理解统计学的基本逻辑，而不是错综复杂的数学难题；它强调的重点是统计学的含义、应用场合、应用方法和对结果的解释。本书反映了我从十年的统计学教学中提炼出来的非常成功的方法。这种方法具有以下三个特点：

第一，我采用逻辑推理而不是数学推导的方式去解释统计学概念和统计检验。统计学虽然是数学的一个分支，但很多统计检验的内在逻辑非常直接，不需要具备高等数学知识背景就可以理解。经常在统计学著作中见到的大量篇幅的公式和数学解释，往往掩盖了这种内在逻辑的简单性。所以，在本书中，你不会看到太多的公式。你看到的少数几个公式也会有详细的解释，而且，你也可以忽略这些公式，这并不会妨碍你理解正在讨论的基本原理。

第二，我采用信息科学研究中的实例来说明所涉及到的每个统计检验的过程。我先从构造假设开始，然后用计算机软件分析数据，整理得出的分析结果，最后得出结论。读者可以看到统计学方法应用的完整过程，而不会陷入到技术细节中去。这种强调统计学方法的实际应用的做法可以使读者明白统计学是一个多么有力和适用的工具，而不是与他们无关的晦涩的数学分支。

第三，我强调要利用计算机软件来进行演算和完成其他繁琐的数学工作。很多统计学著作花费大量时间去做公式推导和演算。事实上，现在几乎不会有人用手工方式来完成这些演算。现在已有各种各样的计算机软件包来完成繁

重的数学工作,因此我们就可以将重点放在理解结果的含义上。有鉴于此,我用了一章的篇幅来论述有关使用计算机软件进行统计分析的问题。

我采用这种方法来教我的学生学习统计学。他们来自各个学科,从本科生到博士,各种层次的都有。我的大部分学生几乎没有什么数学基础,许多学生刚接触这门课程时都感到畏惧和紧张。然而,令他们惊喜的是,统计学竟然如此使人轻松,弄懂一个对他们来说曾经是神秘和陌生的学科如此使人兴奋。事实上,认为统计学很复杂,非数学专业人员难以企及,是一种误解。我相信,一旦你读了这本书,就会赞同我的看法。

## 本书的读者对象

本书适用于信息专业人员,包括学术研究人员和实际从业人员。他们将会发现,这本书在理解统计学的概念和技巧上非常有用,这些概念和技巧在信息科学中的应用越来越广泛,而且在文献中也大量出现。作为信息专业人员,即使你不亲自去做统计分析,你至少也会使用别人做出的统计结果。本书的目的就是为你提供一种帮助你理解和应用统计学方法的策略。即使你不是信息专业人员,也可以从此书中受益。归根结底,统计学的基本理论和原理对于任何领域都是一样的。

## 如何使用本书

本书涵盖了信息科学研究中常用的基本统计学方法。这些方法包括:描述统计学、 $t$  检验、 $\chi^2$  检验、方差分析、相关、回归和基本非参数统计检验。书中也将介绍一些中级和高级的统计方法,比如多重回归分析和 LISREL。

如果你已经熟悉一般的统计分析方法,使用本书是为了重温某一主题,那你可直接跳到你感兴趣的章节。第 12 章中还特别给出了一个总图,你可根据实际情况,依据此图来决定要复习哪种检验方法,然后再去相应的章节阅读。

如果你是个统计学初学者,或者你过去曾经学过,但需要再复习一下以前学过的一些主要概念,我建议你依次读完前 6 章,然后再看其他章节。第 1 章讨论如何区分不同类型的数据,这是开始进行数据分析时最先遇到的问题。第 2 章涉及使用计算机软件进行统计分析的一般性问题,包括如何将数据组织到计算机文件中,以及怎样处理缺失数据。第 3 章和第 4 章分别讨论了图表和描

述统计学。第 5 章介绍了推断统计学的基本概念。因为本章是后面讨论的所有统计检验的基础,所以在转入后面的章节之前应该先阅读此章。第 6 章讨论了各种抽样方法。

前 6 章的性质决定了读者需要按章节顺序依次阅读,但从第 7 至 11 章就不必依次阅读了,因为每章介绍的是不同的统计检验方法。所以你可以直接转入你特别感兴趣的一章。第 12 章对第 7 至 11 章中所涉及到的推断统计学检验进行了总结,其中,有一个“路线图”总图,用以引导你根据特定情况选择正确的统计检验方法。第 13 章介绍了一些中级和高级的统计方法。建议你最后再阅读这一章,因为它以上述章节中介绍过的概念为基础。

Liwen Vaughan(邱黎雯)

(京)新登字 130 号

## 内 容 简 介

本书深入浅出地讲述了信息专业人员常用基本统计学方法——描述统计学、 $t$  检验、 $\chi^2$  检验、方差分析、相关、回归、基本非参数统计检验的理论及应用，也对中、高级统计方法——多重回归分析和线性关系模型(LISREL)作了介绍。本书的突出特点是，着眼于统计学方法应用的完整过程，着重于在对统计学基本逻辑的理解、统计学方法的选择和统计分析结果的解释上进行形象化描述，而不陷入技术细节和繁琐的数学公式之中，使读者感到轻松、实用。本书适合大专院校相关专业师生阅读和所有用到统计学方法的其他专业人员参考使用。

## 声 明

本书英文版原著由 Information Today, Inc. 出版。Information Today, Inc. 保留所有权利。

Third printing, July 2005

**Statistical Methods for the Information Professional: A Practical, Painless Approach to Understanding, Using, and Interpreting Statistics**

Copyright© 2001 by American Society for Information Science and Technology

All rights reserved. No part of this book may be reproduced in any form or by any electronic or mechanical means, including information storage and retrieval systems, without permission in writing from the publisher, except by a reviewer, who may quote brief passages in a review. Published by Information Today, Inc, 143 Old Marlton Pike, Medford, New Jersey 08055.

# 目 录

<b>第 1 章 起步——区分数据类型</b> .....	(1)
1.1 名义数据 .....	(2)
1.2 有序数据 .....	(2)
1.3 区间数据 .....	(3)
1.4 比例数据 .....	(4)
1.5 数据转换 .....	(5)
<b>第 2 章 避免手工计算和公式处理——使用软件</b> .....	(7)
2.1 软件类型 .....	(7)
2.2 软件选择 .....	(8)
2.3 如何将数据组织成计算机文档 .....	(9)
2.4 如何处理缺失数据 .....	(10)
<b>第 3 章 初步观察——利用图表观察数据特征</b> .....	(13)
3.1 图表的种类 .....	(13)
3.2 一种特殊的条图——直方图 .....	(19)
<b>第 4 章 将杂乱的数据汇总为整齐的数字——描述统计学</b> .....	(23)
4.1 集中趋势的度量 .....	(23)
4.1.1 均值——算术平均值 .....	(23)
4.1.2 中位数——中点 .....	(26)
4.1.3 众数——直方图峰值处横坐标值 .....	(27)
4.1.4 不同集中趋势度量指标的选用——小结 .....	(28)
4.2 变异性度量 .....	(31)

4.2.1 极差 .....	(31)
4.2.2 四分位数间距 .....	(32)
4.2.3 标准差 .....	(32)
4.2.4 方差 .....	(34)
4.2.5 变异性度量指标的选择——小结 .....	(34)
4.3 综合应用描述统计学度量指标——示例 .....	(35)
<b>第5章 什么是“统计显著性”——推断统计学的基本概念 .....</b>	<b>(40)</b>
5.1 描述统计学与推断统计学 .....	(40)
5.2 总体与样本 .....	(41)
5.3 参数与统计量 .....	(41)
5.4 概率与频率分布 .....	(42)
5.5 正态分布 .....	(44)
5.6 Z 分值 .....	(45)
5.7 标准正态分布 .....	(47)
5.8 置信区间 .....	(48)
5.9 假设检验——统计显著性与无统计显著性 .....	(52)
5.10 统计检验的错误——I类错误和II类错误 .....	(56)
<b>第6章 如何收集数据——抽样方法 .....</b>	<b>(59)</b>
6.1 简单随机抽样 .....	(59)
6.2 系统抽样 .....	(60)
6.3 分层抽样 .....	(62)
6.4 抽样偏倚 .....	(63)
<b>第7章 检验名义数据与有序数据的关系——<math>\chi^2</math> 检验 .....</b>	<b>(66)</b>
7.1 $\chi^2$ 检验的逻辑性 .....	(66)
7.2 预期频数计算 .....	(69)
7.3 $\chi^2$ 值 .....	(70)
7.4 $\chi^2$ 表 .....	(71)
7.5 检验关系模式 .....	(73)
7.6 使用软件进行 $\chi^2$ 检验的示例 .....	(74)

---

7.7 使用 $\chi^2$ 检验的前提条件 .....	(78)
<b>第 8 章 检验区间数据和比例数据的关系——相关和回归 .....</b>	(83)
8.1 相关类型 .....	(84)
8.2 用散点图观察关系模式 .....	(84)
8.3 度量关系强度——皮尔逊系数 $r$ .....	(87)
8.4 皮尔逊系数 $r_p$ 的显著性检验 .....	(88)
8.5 相关和因果 .....	(90)
8.6 回归方程和回归线 .....	(91)
8.7 预测 .....	(96)
8.8 进行相关分析和回归分析的前提条件 .....	(98)
<b>第 9 章 两个样本间存在显著差异吗? ——<math>t</math> 检验 .....</b>	(99)
9.1 独立 $t$ 检验与配对 $t$ 检验 .....	(100)
9.2 $t$ 检验逻辑 .....	(102)
9.3 $t$ 检验过程 .....	(104)
9.4 用软件进行 $t$ 检验示例 .....	(106)
9.5 使用 $t$ 检验的前提条件 .....	(109)
<b>第 10 章 三个或多个样本存在显著差异吗? ——方差分析 .....</b>	(111)
10.1 方差分析逻辑 .....	(112)
10.2 方差分析过程 .....	(115)
10.3 用软件进行方差分析示例 .....	(116)
10.4 检验差异模式 .....	(118)
10.5 使用方差分析的前提条件 .....	(121)
<b>第 11 章 数据不符合参数检验条件——使用非参数检验 .....</b>	(123)
11.1 Spearman 秩相关系数 .....	(124)
11.2 Mann-Whitney 检验 .....	(127)
11.3 Wilcoxon 符号秩检验 .....	(129)
11.4 Kruskal-Wallis 检验 .....	(132)
11.5 非参数检验的优缺点 .....	(136)
11.5.1 非参数检验的优点 .....	(136)

11.5.2 非参数检验的缺点 .....	(136)
11.5.3 何时使用非参数检验 .....	(137)
<b>第 12 章 如何选择检验法? ——路线图</b> .....	(139)
<b>第 13 章 多变量分析——使用高级统计学方法</b> .....	(144)
13.1 两因素方差分析 .....	(144)
13.2 多重回归 .....	(151)
13.2.1 我们为什么需要多重回归? .....	(151)
13.2.2 多重回归方程 .....	(152)
13.2.3 回归系数 .....	(153)
13.2.4 多重相关系数和多重决定系数 .....	(154)
13.2.5 偏相关系数 .....	(154)
13.3 LISREL .....	(155)
<b>附录</b> .....	(161)
附录 1 标准正态分布 .....	(161)
附录 2 随机数表 .....	(163)
附录 3 $\chi^2$ 临界值 .....	(164)
附录 4 皮尔逊系数 $r$ 的临界值 .....	(165)
附录 5 $t$ 的临界值 .....	(166)
附录 6 方差分析( $\alpha=0.05$ )中 $F$ 的临界值 .....	(167)
附录 7 Tukey's HSD( $\alpha=0.05$ )的临界值 .....	(168)
<b>参考书目</b> .....	(169)
<b>索引</b> .....	(171)