

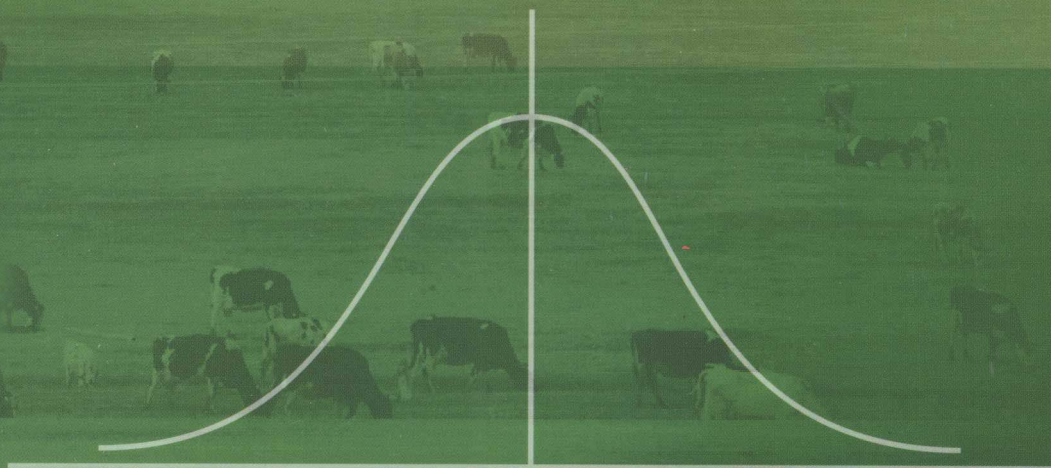


“十二五”普通高等教育本科国家级规划教材
普通高等教育“十一五”国家级规划教材

生物统计学

(第五版)

李春喜 姜丽娜 邵云 张黛静 编著



科学出版社

Q-332-43

4054

“十二五”普通高等教育本科国家级规划教材
普通高等教育“十一五”国家级规划教材

生物统计学

(第五版)

李春喜 姜丽娜 邵云 张黛静 编著

科学出版社

北京

内 容 简 介

本书较为系统地介绍了生物统计学的基本原理和方法,在简要叙述了生物统计学的产生、发展及其研究对象与作用以及生物学研究中抽样方法、试验资料的整理、特征数的计算、概率和概率分布、抽样分布的基础上,着重介绍平均数的统计推断、 χ^2 检验、方差分析、直线回归与相关分析、可直线化的非线性回归分析、协方差分析、多元线性回归与多元相关分析、逐步回归与通径分析和多项式回归分析,同时对试验设计原理及对比设计、随机区组设计、拉丁方设计、裂区设计、正交设计等常用试验设计及其统计分析也进行了详细叙述。

本书可供综合性大学、师范院校生物类及其相关专业的本科生作为教材使用,也可作为从事生命科学、生物工程、农业科学、林业科学、医学、畜牧兽医、水产科学等专业的科研工作者、教师和研究生的参考书。

图书在版编目(CIP)数据

生物统计学/李春喜等编著.—5版.—北京:科学出版社,2013
“十二五”普通高等教育本科国家级规划教材·普通高等教育“十一五”国家级规划教材

ISBN 978-7-03-037502-5

I.①生… II.①李… III.①生物统计-高等学校-教材 IV.①Q-332

中国版本图书馆CIP数据核字(2013)第103373号

责任编辑:丛楠 贺密青/责任校对:桂伟利

责任印制:阎磊/封面设计:迷底书装

科学出版社出版

北京东黄城根北街16号

邮政编码:100717

<http://www.sciencep.com>

化学工业出版社印刷厂印刷

科学出版社发行 各地新华书店经销

*

1997年8月第 一 版 开本:787×1092 1/16

2013年6月第 五 版 印张:19,1/4

2013年6月第一次印刷 字数:442 000

定价:38.00元

(如有印装质量问题,我社负责调换)

第五版前言

作为研究生命科学最基础的工具性课程之一,生物统计学越来越被从事生物学基础教学、生命科学研究的教师和科技工作者所高度重视。生物学研究中会产生大量试验与调查数据资料,这些资料在计算机出现之前是很难进行系统整理和深入分析的。计算机的普及应用,不仅能够将生命科学问题推向深化和精确化,而且还能对其理论及其假说进行论证和说明。尽管数据计算手段已不是限制因素,但不少学生在学习生物统计学时仍然存在很大焦虑,面对大量数学公式和统计数据无所适从。多年的生物统计学教学实践经验表明,上述现象是由于学生没有真正掌握生物统计学基本知识和框架内容所致。针对上述问题,本书按照“循序渐进、强化基础、先易后难、注重实践”的基本思路从生物统计学的概念、发展历史、基本作用出发,在介绍抽样方法、数据资料整理、特征数计算、概率和概率分布、统计推断的基础上,重点对方差分析、回归与相关分析及常用试验设计与统计分析进行详细叙述,以加强高校生物学类专业生物统计学的教学工作,为促进新型生物学高级人才的培养提供基础性知识和工具性方法。

本书在前四版的基础上,对全书内容进行了调整和精编,并获教育部“十二五”普通高等教育本科国家级规划教材立项。本书在第四版的基础上进行了系统修编,主要修编内容有以下几个方面:一是将第四版的14章缩编为13章,取消了第四版第九章抽样原理和方法,将该章部分概念移至第一章,将抽样方法放在第二章中进行介绍;二是在常用试验设计与统计分析中增加了拉丁方设计,兼顾了诸如动物实验等来自双向系统误差试验设计的需要;三是更换了部分例题,使其对相关方法实际运用的介绍更具有代表性;四是对各章部分内容进行了修改,补充了一些难点分析,统一了各章相同术语的表述,订正了个别错误内容;五是配合本书的出版,编写了《生物统计学》(第五版)立体化教材《生物统计学学习指导》一书,为进一步帮助读者理解和学习生物统计学课程内容提供了学习资料。

本书在编写和出版过程中,得到了河南省教育厅、河南师范大学和科学出版社的大力支持。科学出版社的王国栋先生、甄文全先生为书稿的内容编排提出了建设性建议,丛楠女士在本书编审方面作了大量工作。河南省教育厅在规划教材申报立项方面给予了指导和支持。河南师范大学教务处、生命科学学院在适应新形势生物统计学教学改革和教材建设方面给予了指导性意见,对本书再版修订给予了大力支持。在此一并感谢。

本书的再版,离不开广大读者的支持和帮助。恳请各位读者继续对本书给予关心和帮助,对本书存在的不足之处及时提出意见,以便进一步完善。

李春喜

2013年3月于河南师范大学

第一版前言

生物统计学是运用数理统计的原理和方法来分析和解释生物界各种现象和试验调查资料的一门科学。随着生物学的不断发展,对生物体的研究和观察已不再局限于定性的描述,而是需要从大量调查和测定数据中,应用统计学方法,分析和解释其数量上的变化,以正确制订试验计划,科学地对试验结果进行分析,从而作出符合科学实际的推断。目前,生物统计学在农学、林学、畜牧、医药、卫生、生态、环保等领域已有广泛应用,但在纯生物学研究方面的应用,不管是在深度上还是广度上都不及上述领域。有鉴于此,在生物学研究中,迫切需要加强生物统计学的应用,对高校生物类专业,它也是一门应被十分重视的工具课程。本书正是为了满足这些需要而编写的。

本书的写作是在作者多年从事生物统计学教学和应用研究的基础上完成的。书中的内容主要侧重于各种统计方法的应用,在统计原理方面,一般只作概念上的介绍和公式的简单推导,对有些较复杂的统计公式则只给出公式,其目的主要是为了让读者不但对统计学原理有较全面的了解,更重要的是结合实例了解和掌握各种常用统计方法。在本书的安排上,全书共分12章,概括起来主要有五个方面:第一章至第三章介绍统计和概率的基础知识,包括生物统计学的概念和内容、数据的搜集与整理、平均数和变异数的计算、概率和概率分布等;第四章、第五章介绍统计推断,包括样本平均数的检验、样本频率的检验、方差同质性检验、非参数检验和 χ^2 检验;第六章至第九章介绍统计分析方法,主要内容有方差分析、直线回归与相关分析、可直线化的曲线回归分析、多元回归与相关分析、逐步回归分析、多项式回归、协方差分析;第十章、第十一章介绍抽样与试验设计,主要包括抽样误差估计、抽样方法、抽样方案制订及常见的试验设计,如对比设计、随机区组设计、正交设计及其相应的统计分析方法;第十二章对近年来应用越来越多的多元统计分析进行了简单介绍。每章都附有一定数量的思考练习题,供读者参考。

本书中的例子主要有两个来源,一个是近年来有关生物学、农学、林学、医学、畜牧、水产、环保等领域或学科的实际研究资料,另一个是有关著作中的一些例题。崔党群教授在百忙中通审了全书,并提出了富有建设性的建议。贾玉书同志承担了本书的大部分绘图工作。姜丽娜同志在本书的录排中做了大量工作。在本书的出版过程中,得到了科学出版社的大力支持,特别是张晓春同志在书稿的编审和发行方面作了大量工作,在此一并表示谢意。

本书通俗易懂,具有一定的深度和广度,适合生物学、农学、医学、畜牧、水产、环保等领域或学科的科学工作者阅读,也可供本、专科院校生物类专业作为教材使用。

由于作者水平的限制和资料占有的局限性,本书难免会有错误和不妥之处,敬请读者批评指正,以便日后修订完善。

李春喜 王文林

1997年3月

目 录

第五版前言	
第一版前言	
第一章 概论	1
第一节 生物统计学的概念	1
第二节 统计学发展概况	2
一、古典记录统计学	2
二、近代描述统计学	2
三、现代推断统计学	3
第三节 常用统计学术语	4
一、总体与样本	4
二、参数与统计数	4
三、变量与资料	5
四、因素与水平	5
五、处理与重复	6
六、效应与互作	7
七、准确性与精确性	7
八、误差与错误	7
第四节 生物统计学的内容与作用	8
思考练习题	8
第二章 试验资料整理与特征数计算	9
第一节 试验资料的搜集与整理	9
一、试验资料的类型	9
二、试验资料的搜集	10
三、试验资料的整理	15
第二节 试验资料特征数的计算	20
一、平均数	21
二、变异数	24
思考练习题	27
第三章 概率与概率分布	29
第一节 概率基础知识	29
一、概率的概念	29
二、概率的计算	31
三、概率分布	32

四、大数定律	34
第二节 几种常见的理论分布	35
一、二项分布	35
二、泊松分布	39
三、正态分布	41
第三节 统计数的分布	46
一、抽样试验与无偏估计	46
二、样本平均数的分布	47
三、样本平均数差数的分布	49
四、 t 分布	50
五、 χ^2 分布	51
六、 F 分布	52
思考练习题	53
第四章 统计推断	54
第一节 假设检验的原理与方法	54
一、假设检验的概念	54
二、假设检验的步骤	55
三、双尾检验与单尾检验	57
四、假设检验中的两类错误	58
第二节 样本平均数的假设检验	59
一、一个样本平均数的假设检验	59
二、两个样本平均数的假设检验	62
第三节 样本频率的假设检验	69
一、一个样本频率的假设检验	69
二、两个样本频率的假设检验	71
第四节 参数的区间估计与点估计	73
一、参数区间估计与点估计的原理	73
二、一个总体平均数 μ 的区间估计与点估计	74
三、两个总体平均数差数 $\mu_1 - \mu_2$ 的区间估计与点估计	75
四、一个总体频率 p 的区间估计与点估计	76
五、两总体频率差数 $p_1 - p_2$ 的区间估计与点估计	77
第五节 样本方差的同质性检验	78
一、一个样本方差的同质性检验	78
二、两个样本方差的同质性检验	79
三、多个样本方差的同质性检验	80
思考练习题	81
第五章 χ^2 检验	83
第一节 χ^2 检验的原理与方法	83

第二节 适合性检验	85
第三节 独立性检验	88
一、 2×2 列联表的独立性检验	88
二、 $2 \times c$ 列联表的独立性检验	90
三、 $r \times c$ 列联表的独立性检验	91
思考练习题	92
第六章 方差分析	94
第一节 方差分析的基本方法	95
一、方差分析的基本原理	95
二、数学模型	95
三、平方和与自由度的分解	97
四、统计假设的显著性检验—— F 检验	99
五、多重比较	100
第二节 单因素方差分析	105
一、组内观测次数相等的方差分析	105
二、组内观测次数不相等的方差分析	107
第三节 二因素方差分析	109
一、无重复观测值的二因素方差分析	109
二、具有重复观测值的二因素方差分析	113
第四节 多因素方差分析	119
第五节 方差分析缺失数据的估计	124
一、缺失一个数据的估计方法	124
二、缺失两个数据的估计方法	125
第六节 方差分析的基本假定和数据转换	126
一、方差分析的基本假定	126
二、数据转换	126
思考练习题	130
第七章 直线回归与相关分析	132
第一节 回归和相关的概念	133
第二节 直线回归分析	134
一、直线回归方程的建立	134
二、直线回归的数学模型和基本假定	137
三、直线回归的假设检验	138
四、直线回归的区间估计	140
五、直线回归的应用及注意事项	144
第三节 直线相关	145
一、相关系数和决定系数	145
二、相关系数的假设检验	146

82	三、相关系数的区间估计	147
83	四、应用直线相关的注意事项	148
82	思考练习题	149
82	第八章 可直线化的非线性回归分析	150
82	第一节 非线性回归的直线化	151
88	一、曲线类型的确定	151
87	二、数据变换的方法	152
82	第二节 倒数函数曲线	153
82	第三节 指数函数曲线	156
82	第四节 对数函数曲线	158
82	第五节 幂函数曲线	160
82	第六节 Logistic 生长曲线	163
89	一、Logistic 生长曲线的由来和基本特征	163
89	二、Logistic 生长曲线方程的配合	164
89	思考练习题	166
89	第九章 试验设计及其统计分析	167
89	第一节 试验设计的基本原理	167
89	一、试验设计的意义	168
89	二、生物学试验的基本要求	168
89	三、试验设计的基本要素	169
89	四、试验误差及其控制途径	169
89	五、试验设计的基本原则	171
89	第二节 对比设计及其统计分析	172
89	一、对比设计	172
89	二、对比设计试验结果的统计分析	173
89	第三节 随机区组设计及其统计分析	175
89	一、随机区组设计	175
89	二、随机区组设计试验结果的统计分析	176
89	第四节 拉丁方设计及其统计分析	182
89	一、拉丁方设计	182
89	二、拉丁方设计试验结果的统计分析	184
89	第五节 裂区设计及其统计分析	187
89	一、裂区设计	187
89	二、裂区设计试验结果的统计分析	187
89	第六节 正交设计及其统计分析	194
89	一、正交表及其特点	194
89	二、正交试验的基本方法	196
89	三、正交设计试验结果的统计分析	198

思考练习题·····	201
第十章 协方差分析 ·····	204
第一节 协方差分析的作用 ·····	205
一、降低试验误差,实现统计控制·····	205
二、分析不同变异来源的相关关系·····	205
三、估计缺失数据·····	206
第二节 单因素试验资料的协方差分析 ·····	206
一、计算变量各变异来源的平方和、乘积和与自由度·····	208
二、检验 x 和 y 是否存在直线回归关系·····	209
三、检验矫正平均数 $\bar{y}_{i(x=\bar{x})}$ 间的差异显著性·····	210
四、矫正平均数 $\bar{y}_{i(x=\bar{x})}$ 间的多重比较·····	211
第三节 二因素试验资料的协方差分析 ·····	213
一、乘积和与自由度的分解·····	215
二、检验 x 和 y 是否存在直线回归关系·····	216
三、检验矫正平均数 $\bar{y}_{i(x=\bar{x})}$ 间的差异显著性·····	216
第四节 协方差分析的数学模型和基本假定 ·····	217
一、协方差分析的数学模型·····	217
二、协方差分析的基本假定·····	217
思考练习题·····	218
第十一章 多元线性回归与多元相关分析 ·····	219
第一节 多元线性回归分析 ·····	219
一、多元线性回归模型·····	220
二、多元线性回归方程的建立·····	220
三、多元线性回归的假设检验和置信区间·····	226
第二节 多元相关分析 ·····	230
一、多元相关分析·····	230
二、偏相关分析·····	231
思考练习题·····	235
第十二章 逐步回归与通径分析 ·····	237
第一节 逐步回归分析 ·····	237
一、逐个淘汰不显著自变量的回归方法·····	238
二、逐个选入显著自变量的回归方法·····	243
第二节 通径分析 ·····	247
一、通径与通径系数的概念·····	247
二、通径系数的求解方法·····	248
三、通径分析的假设检验·····	251
思考练习题·····	253

第十三章 多项式回归分析	255
第一节 多项式回归的数学模型	255
第二节 多项式回归方程的建立	256
一、多项式回归方程的建立与求解	256
二、多项式回归方程的图示	259
第三节 多项式回归方程的假设检验	259
第四节 相关指数	261
第五节 正交多项式回归分析	261
一、正交多项式回归分析原理	261
二、正交多项式回归分析示例	263
思考练习题	265
主要参考文献	266
附表	268
索引	293

第一章

概 论

本章提要

生物统计学是把数学的语言引入具体的生命科学领域,运用数理统计的原理和方法对生物有机体开展调查和试验,目的是以样本的统计数估计总体的参数,对所研究的总体进行合理的推论。生物统计学主要包括试验设计和统计分析两部分内容,其作用主要有4个方面:提供整理、描述数据资料的科学方法并确定其数量特征,判断试验结果的可靠性,提供由样本推断总体的方法,提供试验设计的原则。统计学的发展经历了古典记录统计学、近代描述统计学和现代推断统计学3个阶段。本章还介绍了统计学中几组常用的术语。

第一节 生物统计学的概念

统计学(statistics)是把数学的语言引入具体的科学研究领域,将所研究的问题抽象为数学问题的过程,是搜集、分析和解释数据的一门科学。生物统计学(biostatistics)是数理统计(mathematical statistics)在生物学研究中的应用,它是用数理统计的原理和方法来分析和解释生物界各种现象和试验调查资料的一门学科,属于应用统计学的一个分支。随着生物学研究的不断发展,生物统计学的应用也越来越广泛。

生物学研究的对象是生物有机体,与非生物相比,它具有特殊的变异性、随机性和复杂性。生物有机体的生长发育、生理活动、生化变化及有机体受外界各种随机因素的影响等,都使生物学研究的试验结果有较大的差异性,这种差异性往往会掩盖生物体本身的特殊规律。在生物学研究中,大量试验资料内在的规律性也容易被杂乱无章的数据所迷惑,从而被人们忽视。因此,在生物学研究中,应用生物统计学就显得特别重要。生物学研究的实践证明,只有正确地应用生物统计学的原理和分析方法对生物学试验进行合理设计,对数据进行客观分析,才能得出科学的结论。

生物统计学是在生物学研究过程中,逐渐与数学的发展相结合而形成的,它是应用数学的一个分支,属于生物数学的范畴。生物统计学是把数学的方法引入具体的生命科学

领域,把生命科学领域中具体的研究问题抽象为数学问题,从大量试验数据中探寻其规律的过程,以数学的概率论和数理统计为基础,涉及数列、排列、组合、矩阵、微积分等知识。作为一门工具课,生物统计学一般不讨论数学原理,而主要偏重于统计原理的介绍和具体分析方法的应用。

第二节 统计学发展概况

人类的统计实践是随着记数活动而产生的。因此,对统计发展的历史可追溯到远古的原始社会。但是,使人类的统计实践上升到理论予以总结和概括成一门系统的统计学,起源于17世纪英国,其代表人物 W. Petty (1623~1687)是政治算术学派的奠基人,代表作是《政治算术》。政治算术学派主张用大量观察和数量分析等方法对社会经济现象进行研究,为统计学的发展开辟了广阔的前景。由于 W. Petty 对统计学的形成有着巨大的贡献,马克思称他为“统计学的创始人”。统计学的发展经历了古典记录统计学、近代描述统计学和现代推断统计学3个阶段。

一、古典记录统计学

古典记录统计学(record statistics)形成于17世纪中叶至19世纪中叶。在最初兴起时,通过用文字或数字如实记录与分析国家社会经济状况,初步建立了统计研究的方法和规则。概率论被引进之后,逐渐成为一种较为成熟的方法。

瑞士数学家 J. Bernoulli(1654~1705)系统论证了大数定律。后来, J. Bernoulli 的后代 D. Bernoulli(1700~1782)将概率论的理论应用到医学和人类保险。

法国天文学家、数学家、统计学家 P. S. Laplace(1749~1827)发展了概率论的研究,建立了严密的概率数学理论,并在天文学、物理学的研究中进行了推广应用。他研究了最小二乘法,提出了“拉普拉斯定理”(中心极限定理的一部分),初步建立了大样本推断的理论基础,为后人开创了抽样调查的方法。

正态分布理论对研究生物统计学的理论十分重要,它最早是由法国数学家 De Moivre 于1733年发现的。德国天文学家和数学家 G. F. Gauss(1777~1855)在研究观察误差理论时,也独立推导出测量误差的概率分布方程,并提出了“误差分布曲线”。这条分布曲线称为 Gauss 分布曲线,也就是正态分布曲线。

二、近代描述统计学

近代描述统计学(description statistics)形成于19世纪中叶至20世纪上半叶,这个时期也是统计学应用于生物学研究的开始和发展时期,其“描述”特色是由一批原来研究生物进化的学者们提炼而成的。英国遗传学家 F. Galton(1822~1911)自1882年起开设“人体测量实验室”,分析父母与子女的变异,探寻其遗传规律,应用统计方法研究人种特性和遗传,探索了能把大量数据加以描述与比较的方法和途径,引入了中位数、百分位数、四分位数,以及分布、相关、回归等重要的统计学概念与方法,开辟了生物学研究的新领域。尽管他的研究当时并未成功,但由于他开创性地将统计方法应用于生物学研究,后人推崇他为生物统计学的创始人。

F. Galton 和他的继承人 K. Pearson (1857~1936) 经过共同努力于 1895 年成立了伦敦大学生物统计实验室, 1889 年发表了《自然界的遗传》一文, 并于 1901 年创办了 *Biometrika* (《生物统计学报》或《生物计量学报》) 这一权威杂志。在该杂志的创刊词中, F. Galton 和 K. Pearson 首次为他们所运用的统计方法明确提出了“生物统计”(biometry) 一词, F. Galton 解释为: 所谓生物统计学, 就是应用于生物学科中的统计方法。在《自然界的遗传》一文中, K. Pearson 提出了相关与回归分析问题, 并给出了简单相关系数和复相关系数的计算公式。1900 年, K. Pearson 在研究样本误差效应时, 提出了 χ^2 检验, 它在属性资料的统计分析中有着广泛的应用。

三、现代推断统计学

现代推断统计学(inference statistics)形成于 20 世纪初至 20 世纪中叶。随着社会科学和自然科学领域研究的不断深入, 各种事物与现象之间繁杂的数量关系以及一系列未知的数量变化, 单靠记录或描述的统计方法已难以奏效。因此, 要求采用推断的方法来掌握事物之间的真正联系并对事物进行预测。从描述统计学到推断统计学, 这是统计学发展过程中的一个巨大飞跃。

K. Pearson 的学生 W. S. Gosset (1876~1937) 对样本标准差进行了大量研究, 于 1908 年以笔名“Student”在 *Biometrika* 杂志上发表了论文《平均数的概率误差》, 创立了小样本检验的理论和方法, 即 t 分布和 t 检验法。 t 检验已成为当代生物统计工作的基本工具之一, 它也为多元分析的理论形成和应用奠定了基础。因此, 许多统计学家把 1908 年看成是统计推断理论发展史上的里程碑, 也有人推崇 W. S. Gosset 为推断统计学(尤其是小样本研究理论)的先驱者。

英国统计学家 R. A. Fisher (1890~1962) 于 1923 年发展了显著性检验及估计理论, 提出了 F 分布和 F 检验, 创立了方差和方差分析。在从事农业试验及数据分析研究时, 他提出了随机区组法、拉丁方法和正交试验的方法。1915 年, R. A. Fisher 在 *Biometrika* 上发表论文《无限总体样本相关系数值的频率分布》, 被称为现代推断统计学的第一篇论文。1925 年, R. A. Fisher 发表了《试验研究工作中的统计方法》, 对方差分析及协方差分析进一步作了完整的解释, 从而推动和促进了农业科学、生物学及遗传学的研究与发展。自方差分析问世以来, 各种数理统计方法不但在实验室中成为研究人员的析因工具, 而且在田间试验、饲养试验、临床试验等农学、医学和生物学领域也得到了广泛应用。

J. Newman (1894~1981) 和 E. S. Pearson 进行了统计理论的研究工作, 分别于 1936 年和 1938 年提出了一种统计假设检验学说。假设检验和区间估计作为数学上的最优化问题, 对促进统计理论研究和对试验作出正确结论具有非常实用的价值。

另外, P. C. Mabeilinobis 对作物抽样调查、A. Waecl 对序贯抽样、K. Mather 对群体遗传学、F. Yates 对田间试验设计等都作出了杰出的贡献。

我国对生物统计学的应用始于 1913 年顾澄教授翻译的英国统计学家 G. U. Yule 在 1911 年出版的关于描述统计学的名著《统计学之理论》, 这标志着英国、美国数理统计学传入中国的开始。之后, 许多生物学研究工作者积极从事统计学理论和实践的应用研究, 使生物统计学在农业科学、医学科学、生物学、遗传学、生态学等学科领域发挥了重要作

用。应用试验设计方法和统计分析理论,进行农作物品种产量比较试验、病虫害的预测预报、动物饲养试验、饲料配方、毒理试验、动植物资源的调查与分析、动植物育种中遗传资源及亲代和子代遗传的分析等都取得了较好成果。

近年来,生物统计学发展迅速,从中又分支出群体遗传学、生态统计学、生物分类统计学、毒理统计学等。由于数学与生物学、医学和农学的应用,使生物数学成为一门新的学科,生物统计学只是它的一个分支学科。1974年,联合国教育、科学及文化组织在编制学科分类目录时,第一次把生物数学作为一门独立的学科列入生命科学类中。随着计算机的普及和网络技术的发展,SAS(statistical analysis system)、SPSS(statistical package for the social science)等国际通用统计软件的开发和应用,以及生命科学研究领域的不断深入,生物统计学的研究和应用必将越来越广泛,越来越深入。

第三节 常用统计学术语

一、总体与样本

具有相同性质的个体所组成的集合称为总体(population),它是指研究对象的全体,而组成总体的基本单元称为个体(individual)。

总体按所含个体的数目可分为有限总体和无限总体。个体极多或无限多的总体称为无限总体(infinite population)。例如,某一棉田棉铃虫的头数,可以认为是无限总体。另外,也可从抽象意义上来理解无限总体。例如,通过临床试验来推断某种药品比另一种药品治愈剪率高,这里无限总体是指一个理论性总体。个体有限的总体称为有限总体(finite population)。例如,对某一班学生身高进行调查,这时总体是指这一班中每位学生的身高。

要研究总体的性质,一般情况下我们无法对总体中的个体全部取出进行调查或研究。因为在实际研究过程中,常会遇到两种难以克服的困难:一是总体的个体数目较多,甚至无限多;二是总体的数目虽然不多,但试验具有破坏性,或者试验费用很高,不允许做更多的试验。在这种情况下,只能采取抽样的方法,从总体中抽取一部分个体进行研究。

从总体中抽出的若干个体所构成的集合称为样本(sample),构成样本的每个个体称为样本单位(sample unit),样本中个体的数目称为样本容量(sample size),记为 n 。样本的作用在于估计总体。例如,可以调查某一地区棉田100株棉花上的棉铃虫头数,来推断该地区棉铃虫的发生状况,以采取相应的对策。一般在生物学研究中, $n < 30$ 的样本称为小样本, $n \geq 30$ 的样本称为大样本。在一些计算和分析检验方法上,大、小样本是不同的。

在对事物的研究过程中,人们常通过某事物的一部分(样本)来估计事物全部(总体)的特征,目的是为了以样本的特征对未知总体进行推断,从特殊推导一般,对所研究的总体作出合乎逻辑的推论,得到对客观事物的本质和规律性的认识。在生物学研究中,我们所期望的是总体,而不是样本。但是在具体的试验过程中,我们所得到的却是样本而不是总体。因此,从某种意义上讲,生物统计学是研究生命过程中以样本来推断总体的一门学科。

二、参数与统计数

参数(parameter)也称为参量,是对一个总体特征的度量,常用希腊字母表示,如总体

平均数 μ 、总体标准差 σ 等均为参数。统计数(statistic)也称为统计量,是由样本计算所得的数值,它是描述样本特征的数量,常用英文字母表示,如样本平均数 \bar{x} 、样本标准差 s 等。由于总体一般都很大,有的甚至不可能取得,所以总体参数通常是未知的。正因为如此,我们才进行抽样,由于样本是已经抽出来的,所以统计数是可以计算出来的,我们可以根据样本统计数来估计总体的参数。此外,还有一些统计量是为了进行统计分析而构造出来的,如后续章节中的 u 统计量、 t 统计量及 F 统计量等。

三、变量与资料

相同性质的事物间表现差异性的某项特征或性状称为变量或变数(variable),是研究者在确定了研究目的之后,所观测的试验指标。由于试验目的不同,所选择的变量也不相同,如植物叶片叶绿素的含量,人体身高、体重、血糖含量、血型等。变量通常记为 x ,如 10 个人的身高为 155~180cm,共有 158,167,173,155,180,165,175,178,170,162(cm) 10 个变量值,记作 $x_i (i=1,2,\dots,10)$,表示 x_1 到 x_{10} 之间任一数值。变量的观察结果可以是定量的,也可以是定性的,其结果称为变量值(value of variable)或观测值(observed value),也称为数据、资料(data)。

根据获取观测值的方式及测量方法所提供的数值信息的差异,变量可以分为定量变量和定性变量。通过测量所获得的、用具体数值与特定计量单位表达的数据称为定量变量(quantitative variable),也称为数值变量(numerical variable)。其变量值是定量的,表现为数值大小,一般有度量衡单位,如人的身高(cm)、体重(kg)、脉搏计数(次/min)等。定量变量根据取值的不同,可以分为连续变量和非连续变量。连续变量(continuous variable)表示在变量范围内可抽出某一范围的所有值,变量之间是连续的、无限的。例如,小麦的株高为 80~90cm,在此范围内可以取得无数个变量。非连续变量(discontinuous variable)也称为离散型变量(discrete variable),表示在变量数列中仅能取得固定数值,并且通常是整数,如菌落中的菌数、单位面积水稻的茎数、小白鼠每胎产仔数等。

定性变量(qualitative variable)也称为分类变量(categorical variable)、名义变量(nominative variable),其变量值是定性的,表示某个体属于几种互不相容的类型中的一种。例如,果蝇的翅有长翅与残翅,人的血型有 A、B、AB 和 O 型,豌豆花的颜色有白色、红色和紫色,等等。

变量的类型是根据研究目的而确定的。根据需要,各类变量可以互相转化。例如,以人作为研究对象,观察某人群成年男子的血红蛋白含量(mg/L),属于定量变量;若按血红蛋白含量正常与偏低分为两类,则属于定性变量。

对应于变量,常量(constant)是不能给予不同数值的变量,它是代表事物特征和性质的数值,通常由变量计算而来,在一定过程中是不变的,如总体平均数、标准差、变异系数等。只有在事物的总体发生变动时,常量才随之变化。

四、因素与水平

试验中所研究的影响试验指标的原因或原因组合称为试验因素(experimental factor)或处理因素(treatment factor),简称为因素或因子(factor)。试验因素常用大写字

母,如 A、B、C 等来表示。

每个试验因素的不同状态(处理的某种特定状态或数量上的差别)称为因素水平(level of factor),简称为水平(level)。例如,研究温度对某种酶活性的影响,所设置的 15℃、20℃、25℃、30℃ 分别称为温度因素的一个水平。可见,因素是一个抽象的概念,而水平则是一个较为具体的概念。水平常用代表该因素的字母添加下标(如 1、2、3 等)来表示,如 A_1 、 A_2 、 B_1 、 B_2 等。

按照性质不同,因素可以分为可控因素和非控因素。在试验中可以人为调控的因素称为可控因素(controllable factor)或固定因素(fixed factor)。该因素的水平可准确控制,且水平固定后,其效应也固定,同时在试验进行重复时可以得到相同的结果。例如,研究 3 种温度对胰蛋白酶水解产物的影响,因为温度是可以严格控制的,所以在重复该试验时对于相同的温度其水解产物的量也是固定的。温度在此例中即为固定因素。

在试验中不能人为调控的因素称为非控因素(uncontrollable factor)或随机因素(random factor)。该因素的水平不能严格控制,或虽水平能控制,但其效应仍为随机变量,同时在试验进行重复时不易得到相同的结果。例如,研究农家肥不同施用量对作物产量的影响,由于农家肥有效成分较为复杂,不能像控制温度那样,将农家肥有效成分严格地控制在某一固定值上,在重复试验时即使施用相同数量的农家肥,也得不到一个固定的效应值。农家肥在此例中即为随机因素。

五、处理与重复

试验处理(experimental treatment)通常也称为处理(treatment),是指对受试对象给予的某种外部干预(或措施)。其中受试对象(tested subject)又称为试验单位或试验单元(experimental unit),是指在试验中能接受不同试验处理的独立的试验载体。植物个体、动物个体,以及不同的组织、器官等都可以作为试验单位。

处理根据所涉及的因素数可以分为单因素处理和多因素处理。当试验中涉及的因素只有一个时,称为单因素处理(single factor treatment)。在单因素处理中,实施在试验单位上的具体项目就是试验因素的某一水平。例如,饲料的比较试验,实施在试验单位(如某种畜禽)上的具体项目就是饲喂某一种饲料。进行单因素试验时,试验因素的一个水平就是一个处理。

如果试验中涉及两个或两个以上的因素,则称为多因素处理(multiple factors treatment)。可依处理因素数进行具体命名,如二因素试验处理、三因素试验处理等。在多因素试验处理中,实施在试验单位上的具体项目是各因素的某一水平组合。例如,3 个播种密度对 4 个小麦品种的产量影响试验,就是一个二因素试验处理,试验共有 $3 \times 4 = 12$ 个水平组合,实施在试验单位(小麦)上的具体项目就是某个种植密度与某个小麦品种的组合。进行多因素试验时,试验因素的一个水平组合就是一个处理。相对于单因素试验,多因素试验不但可以研究因素的主效,同时也可研究因素之间的交互作用。

重复(repetition)是指在试验中,将一个处理实施在两个或两个以上的试验单位上。处理实施的试验单位数即为处理的重复数。例如,研究某种饲料对猪的增重效果,将该种饲料饲喂 5 头猪,则表明这个处理(饲料)有 5 次重复。