

高等 学 校 食 品 专 业 系 列 教 材



# 试验设计与数据处理

EXPERIMENTAL DESIGN & DATA PROCESSING

潘丽军 陈锦权 / 主编

東南大學 出版社  
SOUTHEAST UNIVERSITY PRESS

高等学校食品专业系列教材

# 试验设计与数据处理

主 编 潘丽军 陈锦权

副主编 刘建学

编写人员 (按姓氏笔画为序)

王 武 刘建学 邱春江

陈锦权 章银良 潘丽军

东南大学出版社  
·南京·

## 内 容 提 要

试验设计与数据处理是食品科学与工程及相关专业的专业基础课程。本教材从技术和应用的观点出发,重点阐述了试验设计和数据处理的基本原理与常用方法。主要内容包括试验设计基础,正交试验设计,均匀试验设计,回归正交试验设计,回归旋转试验设计;极差分析,方差分析,回归分析,以及SAS统计软件的基本知识与命令,SAS数据集、软件应用中的程序设计,SAS软件在简单计算、绘图、方差分析、多元回归设计等应用中的处理方法等。各章教学目标明确,配有的习题可帮助读者理解掌握。

本教材可作为食品、生物、化工、材料、制药等专业本科生和研究生的教学用书,也可供相关学科的科研、教学、实验、工程技术人员参考。

## 图书在版编目(CIP)数据

试验设计与数据处理/潘丽军,陈锦权主编. —南京:东南大学出版社,2008. 2  
ISBN 978 - 7 - 5641 - 1128 - 1

I . 试… II . ①潘… ②陈… III . ①试验设计(数学)  
②实验数据—数据处理 IV . 0212. 6 N33

中国版本图书馆 CIP 数据核字(2008)第 016337 号

## 试验设计与数据处理

---

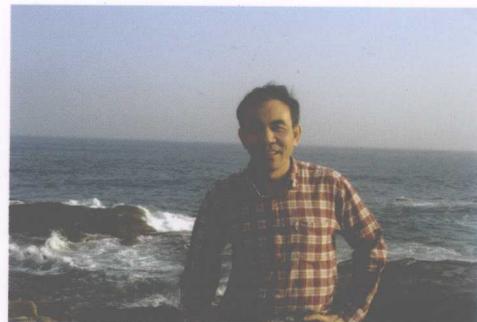
出版发行: 东南大学出版社  
社 址: 南京四牌楼 2 号 邮编: 210096  
出 版 人: 江汉  
网 址: <http://press. seu. edu. cn>  
电子邮箱: [press@seu.edu.cn](mailto:press@seu.edu.cn)  
经 销: 全国各地新华书店  
印 刷: 兴化印刷有限责任公司  
开 本: 787mm×1092mm 1/16  
印 张: 22. 75  
字 数: 554 千字  
版 次: 2008 年 2 月第 1 版  
印 次: 2008 年 2 月第 1 次印刷  
书 号: ISBN 978-7-5641-1128-1/TS · 28  
印 数: 1~4000 册  
定 价: 38. 00 元

---

本社图书若有印装质量问题,请直接与读者服务部联系。电话(传真):025 - 83792328



潘丽军，教授，安徽省优秀教师，享受国务院政府特殊津贴专家。现任合肥工业大学生物与食品工程学院副院长。长期从事本科生、研究生的教学工作，主持参加省部级以上科研及教改重点项目 10 余项，获省部级科技进步奖 3 项、教学成果奖 1 项，发表论文 30 余篇。



陈锦权，男，53岁，福建农林大学食品科学学院院长、教授、博导，致力于亚热带果蔬、茶叶等资源的深度加工与食品安全控制研究。目前研究的国家、省部级科研课题有5项，近年来发表学术论文12篇，其中国家一级刊物5篇，EI、SCI、ISTP收录3篇，获国家发明专利1项、省部级科技进步二等奖1项。

# 高等学校食品专业系列教材

## 编写委员会

(以姓氏笔画为序)

- 王晓曦 河南工业大学粮油食品学院副院长、教授  
王向东 山西师范大学工程学院院长、教授  
邓泽元 南昌大学生命科学学院副院长、教授,博士生导师  
毛多斌 郑州轻工业学院食品与生物工程学院院长、教授  
艾志录 河南农业大学食品科学技术学院副院长、副教授  
刘建学 河南科技大学食品与生物工程学院副院长、教授  
张 濞 江南大学食品学院院长、教授,博士生导师  
孟岳成 浙江工商大学食品科学与工程系主任、教授  
陆兆新 南京农业大学食品科技学院院长、教授,博士生导师  
陈正行 江南大学食品学院副院长、教授,博士生导师  
陈锦权 福建农林大学食品科学学院院长、教授,博士生导师  
杜云建 淮海工学院海洋学院副教授  
郑铁松 南京师范大学食品科学与营养系主任、副教授  
姜绍通 合肥工业大学生物与食品工程学院院长、教授,博士生导师  
赵丽芹 内蒙古农业大学食品科学与工程学院副院长、教授  
赵希荣 淮阴工学院食品系主任、副教授  
钱建亚 扬州大学食品科学与工程学院教授  
董 英 江苏大学食品与生物工程学院总支书记、教授,博士生导师  
蒋爱民 华南农业大学食品学院教授,博士生导师  
熊晓辉 南京工业大学食品科学与工程系主任、教授  
鞠兴荣 南京财经大学食品科学与工程学院院长、教授,博士生导师

# 总序

受编辑之托,为我等所著的高等学校食品专业系列教材作序,真是诚惶诚恐,迟迟难以下笔。苏轼《与孙子思》云:“……余空纸两幅,留与五百年后人跋尾也!”此一戏语道出了作序之尴尬。回想起当时来自各地高校食品院系的学者们共同讨论系列教材时认真而热烈的场景,我就勉为其难,介绍一下我们编写这套系列教材的来龙去脉和想法。

2005年11月18~20日,经东南大学出版社和江南大学食品学院的联合组织,在江苏无锡召开了“普通高等教育‘十一五’国家级教材规划·食品专业系列教材”编写和申报研讨会,来自江南大学、南昌大学、南京农业大学、合肥工业大学、江苏大学、内蒙古农业大学、福建农林大学、河南工业大学、郑州轻工业学院、河南农业大学、河南科技大学、浙江工商大学、扬州大学、华南农业大学、南京工业大学、南京财经大学、南京师范大学、淮阴工学院、淮海工学院等19所大学食品院系的30余名学者参加了会议。在两天的会议中,学者们探讨了近几年来食品专业教育的得失,研讨了新形势下为进一步推进食品学科创新型人才培养的系列教材的编写要求、体例和分工,明确了31部教材的编写任务。时间过去不到一年,硕果满园的金秋季节在望,这31部教材中已有5部列入普通高等教育“十一五”国家级教材规划,第一部教材《食品添加剂》将正式付梓,其他多部教材也将孕育而生,在近期内陆续出版,真是欣慰之极。

古人曰:教人以道者,师也。作为教师,不仅要教会学生如何掌握知识,更重要的是要教会学生如何运用知识和创造知识。这套系列教材的编者们,少则有十多年、多则有二十年左右从事相应课程教学和本专业领域科研的经历。我们一致的想法是希望把多年实践中的感悟和积累融入这套教材中,使本系列教材的阅读者在理解和掌握知识的同时,也能对知识的运用和创造有所领悟。

食品工业的GDP在我国国民经济中已连续几年居首位,现已接近2万亿元,食品科技进步与产业发展在国民经济发展中越来越发挥举足轻重的作用。目前全国约有200所高校办有食品专业,每年招收学生2万多人,食品专业的教育教学在一定程度上关系到我国食品工业的健康和可持续发展,编写一套反映当今科技发展现状、符合创新创业型人才培养要求的食品专业系列教材,是我们

所有编者的愿望，也是我们义不容辞的责任和义务。

愿我们的国家明天更美好，愿我们的食品工业发展更健康，愿我们在着力创建的和谐社会中享用的食品更安全。让我们所有编写和阅读本系列教材的同仁们共同为此尽绵薄之力！

张 瀚

2006年8月3日晚于无锡

# 前　　言

试验是开展科学的研究工作必不可少的手段。任何一种新产品、新工艺、新材料、新品种的产生,任何一项科研成果的获得,都需要进行多次反复试验,并通过对试验数据的分析处理与总结,来获取尽可能多的有用的信息,以达到解决实际问题的目的。因此,试验设计与数据处理作为一种通用现代技术,近年来在各领域的应用范围日益广泛、成效日益显著,是现今科研、工程技术、管理人员必备的一项技术。本教材的编写在保持系统性和科学性的前提下,注意引入相关学科发展的新知识、新成果;紧密联系生产、科研实际,以及统计分析与计算机科学的结合;用大量的实例,介绍一些常用的试验设计及其数据处理方法在科学试验和工农业生产中的应用;力求做到循序渐进,由浅入深,深入浅出,通俗易懂。对每一种设计或分析处理方法,在正确阐明基本原理的同时,注重交代其问题的提出,解决问题的思路与方法;各章的教学目标明确,配有的习题可帮助读者理解掌握。

全书共 10 章。第 1 章绪论(潘丽军编写),主要介绍试验设计与数据处理的作用和意义,以及试验设计与数据处理的发展历程和应用进展情况。第 2 章试验设计基础(潘丽军编写),主要介绍试验常用术语、试验研究的计划与方案、试验设计常用的优良性及应遵循的基本原则。第 3 章数据资料的特征数与误差分析(邱春江编写),介绍了不同类型资料的性质、整理方法,平均数和变异数的分类、性质及资料特征数的计算方法;试验误差的来源和分类,产生试验误差的原因及影响因素,试验误差的检测和判别方法;试验数据的列表法和图示法表达方式。第 4 章方差分析(刘建学编写),主要介绍单因素、双因素方差分析的基本原理和方法,无交互影响、有交互影响双因素方差分析及其简化算法与显著性检验。第 5 章试验数据的回归分析(王武编写),重点介绍线性回归方程的建立与回归效果显著性检验,最优线性回归方程的统计选择,利用偏回归系数的显著性程度或标准回归系数的大小判断试验因素重要程度。第 6 章正交试验设计(潘丽军编写),主要介绍正交试验设计的基本思想与正交表;单指标与多指标的正交试验设计、混合型正交试验设计、考虑交互作用的正交试验设计;正交试验设计直观分析、方差分析的基本原理和方法。第 7 章均匀试验设计(章银良编写),

介绍了均匀试验设计的概念与特点,均匀试验均匀性准则,均匀试验基本方法和应用。第8章回归正交试验设计(陈锦权编写),主要介绍一次回归正交试验设计和二次回归正交试验设计的原理、基本方法及统计分析。第9章回归旋转试验设计(陈锦权编写),主要介绍了回归旋转试验设计的基本原理、实现条件、组合设计和统计分析方法。第10章SAS统计软件在试验设计与数据处理中的应用(刘建学编写),重点介绍国际上最富知名度的三大统计软件之一的SAS统计软件,包括SAS软件的基本知识、基本命令、SAS数据集、软件应用中的程序设计;SAS软件在简单计算、绘图、方差分析、多元回归设计等应用中的处理方法等。书后附有常用统计数学用表。

本教材可作为食品、生物、化工、材料、制药等专业本科生和研究生的教学用书,也可供相关学科的科研、教学、实验、工程技术人员参考。

在本教材的编写过程中,参考了有关文献与资料,在此向这些文献资料的作者表示感谢。

限于编者的水平与经验,书中错误与不妥之处,恳望读者批评指正。

编 者

2007年8月

# 目 录

<b>1 绪论</b>	1
1.1 试验设计与数据处理的作用和意义	1
1.2 试验设计与数据处理的发展和应用	2
<b>2 试验设计基础</b>	8
2.1 常用术语	8
2.2 试验计划与方案	10
2.3 试验设计常用的优良性	15
2.4 试验设计应遵循的基本原则	16
<b>3 数据资料的特征数与误差分析</b>	19
3.1 数据资料的特征数	19
3.2 试验数据的误差	29
3.3 试验数据常用的表、图表达方式	50
<b>4 方差分析</b>	56
4.1 单因素方差分析	56
4.2 单因素试验方差分析的简化计算	62
4.3 双因素试验的方差分析	63
<b>5 试验数据的回归分析</b>	77
5.1 基本概念	77
5.2 直线回归方程的建立与回归效果显著性检验	78
5.3 多元线性回归分析	86
5.4 试验因素重要程度(主次顺序)的判别方法	96
5.5 能直线化的曲线回归分析	97
<b>6 正交试验设计</b>	100
6.1 正交试验设计的基本思想	100
6.2 正交表	104
6.3 正交试验设计的基本步骤	108
6.4 正交试验设计的直观分析	110
6.5 正交试验设计的方差分析	140
<b>7 均匀试验设计</b>	165
7.1 均匀试验设计的概念与特点	165
7.2 均匀设计的思想	166
7.3 均匀设计表	167

7.4 均匀性准则 .....	170
7.5 均匀试验设计的基本方法 .....	175
7.6 均匀试验设计的应用 .....	177
7.7 含有定性因素的均匀设计 .....	184
7.8 均匀试验设计特别注意的几个问题 .....	190
<b>8 回归正交试验设计 .....</b>	<b>192</b>
8.1 回归正交试验设计简介 .....	192
8.2 一次回归正交试验设计原理 .....	193
8.3 二次回归正交组合设计 .....	209
<b>9 回归旋转试验设计 .....</b>	<b>221</b>
9.1 回归旋转试验设计的基本原理 .....	221
9.2 二次回归正交旋转组合设计及统计分析 .....	227
9.3 通用旋转组合设计及统计分析 .....	232
<b>10 SAS 软件在试验设计与数据处理中的应用 .....</b>	<b>239</b>
10.1 SAS 软件简介 .....	239
10.2 SAS 软件的应用基础 .....	243
10.3 SAS 数据集的使用 .....	250
10.4 SAS 程序的创建和运行 .....	252
10.5 SAS 常用语句 .....	257
10.6 SAS 服务过程 .....	259
10.7 SAS 软件的应用 .....	261
<b>附录 .....</b>	<b>319</b>
一、F 分布表 .....	319
二、t 分布表 .....	329
三、相关系数 $\gamma$ 与 $R$ 的临界值 ( $\gamma_c$ 与 $R_c$ ) 表 .....	331
四、常用正交表 .....	334
五、均匀试验设计表 .....	342
<b>参考文献 .....</b>	<b>353</b>

# 1

# 绪论

试验设计与数据处理是数理统计学中的一个重要分支。它是以概率论、数理统计及线性代数为理论基础,结合一定的专业知识和实践经验,研究如何经济、合理地安排试验方案以及如何系统、科学地分析处理试验结果的一项科学技术,从而解决了长期以来在试验领域中,传统的试验方法对于多因素试验往往只能被动地处理试验数据,而对试验方案的设计及试验过程的控制显得无能为力这一问题。近代创立和发展起来的试验设计方法,将试验方案的最优化设计与数据处理方法的最优化选择进行有机地结合,并将其思想和要求贯穿于试验的全过程,使试验领域发生了深刻的变化,有力地推动了科学研究和生产实践的进程。

## 1.1 试验设计与数据处理的作用和意义

试验是开展科学的研究工作必不可少的手段。任何一种新产品、新工艺、新材料、新品种的产生以及任何一项科研成果的获得,往往需要做多次反复试验,并通过对试验数据的分析与总结,来获取尽可能多的有用的信息,以达到解决实际问题的目的。因此,只要做试验,就存在着如何科学合理地安排试验方案和分析处理试验结果的问题,即进行试验设计与数据处理的问题。所以说试验设计与数据处理在科学的研究和工农业生产中的作用与意义是极其重要和深远的。

### 1.1.1 试验设计

试验研究可分为试验设计、试验的实施、收集整理和分析实验数据等步骤。而试验设计是影响研究成功与否最关键的一个环节,是提高试验质量的重要基础。试验设计是在试验开始之前,根据某项研究的目的和要求,制定试验研究进程计划和具体的试验实施方案。其主要内容是研究如何合理地安排试验、取得数据,然后进行综合的科学分析,从而达到尽快获得最优方案的目的。

如果试验安排得合理,就能用较少的试验次数,在较短的时间内达到预期的试验目标;反之,试验次数既多,其结果还往往不能令人满意。试验次数过多,不仅浪费大量的人力和物力,有时还会由于时间拖得很长,使试验条件发生变化而导致试验失败。因此,如何合理地安排试验方案是值得研究的一个重要课题。一项科学合理的试验设计应能做到如下几个方面:①尽可能采用具有一定优良性的试验方案,以最少的人力、物力和试验次数,实现预期目的;②能运用试验设计的基本原则,有效控制试验干扰,提高试验精度;③通过简便的计算和分析,可直接获得整个试验区域内较多的、有价值的信息;④试验研究结果具有较好的重演性和推广性。

目前,已建立起许多试验设计方法。如我们大家比较熟悉的,常用单因素试验设计方法的有黄金分割法(0.618法)、分数法、交替法、等差法、等比法、对分法和随机法等,这些方法为多因素试验水平范围的选取提供了重要的依据,并在生产中取得了显著成效。而多因素试验设计方法有正交试验设计、均匀试验设计、稳健试验设计、完全随机化设计、随机区组试验设计、回归正交试验设计、回归正交旋转试验设计、回归通用旋转试验设计、混料回归试验设计、 $D$ -最优回归设计等,其中最基础的、在各领域应用最广泛的多因素试验设计方法是正交试验设计、均匀试验设计、回归正交试验设计以及回归正交旋转试验设计。

### 1.1.2 数据处理

合理的试验方案只是试验成功的充分条件,如果结合系统科学的数据处理与分析,就能对所研究的问题有一个明确的认识,即能从大量的、带有偶然性误差的试验观测值中找出可靠规律性的结论。数据处理主要是研究试验测量或观察值的分析计算处理方法,并依据所得到的规律和结果对工农业生产等进行预报和控制,从而掌握和主宰客观事物的发展规律,使之更好地服从和服务于人类。

目前常用的数据处理方法中,参数估计主要是对某些重要参数进行点估计和区间估计;假设检验是判断各种数据处理结果的可靠性程度;极差分析是利用一组数据中最大值与最小值之间差值的大小,判断因素对指标的影响程度;方差分析是分析各影响因素对考察指标的显著性程度;回归分析则是描述如何获得反映事物客观规律性的数学表达式等。

因此,对试验获取的大量数据,采用不同的数据处理方法,不仅能够方便、有效地揭示产品质量和性能指标与众多影响因素之间的内在关系,还可以简捷地求得回归方程、确定最佳工艺条件(或最优参数组合),如预报气象和病虫害、建立自动控制中的数学模型、制定产品的生产工艺参数等。

### 1.1.3 试验设计与数据处理的作用

- (1) 有助于研究者掌握试验因素对试验考察指标影响的规律性,即各因素的水平改变时指标的变化情况。
- (2) 有利于分清试验因素对试验考察指标影响的大小顺序,找出主要因素。
- (3) 有助于反映试验因素之间的相互影响情况,即因素间是否存在交互作用。
- (4) 能正确估计和有效控制试验误差,提高试验的精度。
- (5) 能较为迅速地优选出最佳工艺条件(或称最优方案),并能预估或控制一定条件下的试验指标值及其波动范围。
- (6) 根据试验因素对试验考察指标影响规律的分析,可以深入揭示事物内在规律,明确进一步试验研究的方向。

## 1.2 试验设计与数据处理的发展和应用

### 1.2.1 试验设计与数据处理的发展历程

试验设计与数据处理是在概率论和数理统计的基础上不断完善和发展起来的。试验统

计的方法最早起源于对农业及生物遗传研究的应用统计方法,故一般称为生物统计学,它是应用数理统计学原理来研究生物界数量现象的科学方法,是一门数理统计学与生物科学相结合的交叉学科。

在 20 世纪 20 年代,英国生物统计学家费歇(R. A. Fisher, 1890—1962)运用均衡排列的拉丁方,解决了长期未能解决的试验条件不均衡的问题,首创了“试验设计”方法。开始时该法主要应用于农业、生物学和遗传学等方面,取得了丰硕成果,后用于田间试验,使农业大幅度增产。费歇于 1935 年出版了他的著作 *Design of Experiments*,从此开创了“试验设计”这门新的应用技术科学。以后随着农业和生物学研究的发展,生物统计、试验设计和抽样理论也得到了快速的同步发展,伴随着工业试验研究设计和数理科学理论研究的不断深入,进而推动了应用数理统计学的发展,反过来又促进了试验统计学研究水平的不断提升。其发展历程大致可分为三个阶段。

### 1) 古典记录统计学阶段

古典记录统计学形成期间大致在 17 世纪中叶至 19 世纪中叶。统计学在这个兴起阶段,只还是一门意义和范围不太明确的学问,在它用文字或数字如实记录与分析国家社会经济状况的过程中,初步建立了统计研究的方法和规则。在这一期间的研究成果有 17 世纪帕斯和费马的概率论,18 世纪德莫弗、拉普拉斯和高斯的正态分布理论。其中最卓有成效地把古典概率论引进统计学的是法国天文学家、数学家、统计学家拉普拉斯(P. S. Laplace, 1749—1827),因此,后来比利时大统计学家凯特勒指出,统计学应从拉普拉斯开始。

#### (1) 拉普拉斯

拉普拉斯的主要贡献:①发展了对概率论的研究。拉普拉斯第一篇关于概率论的表述发表于 1774 年,1812 年发表的《概率分析理论》(先后出过 4 版)是他的代表作。拉普拉斯最早系统地把数学分析方法运用到概率论研究中,建立了严密的概率数学理论。②推广了概率论在统计中的应用,主要表现在人口统计、观察误差理论和概率论对于天文问题的应用。1809 年~1813 年,拉普拉斯结合概率分布模型和中心极限思想来研究最小二乘法,首次为这项后来最常用的统计手段奠定了理论基础。③明确了统计学的大数法则。拉普拉斯发现在观察天体运动现象时,当次数足够多时,能使个体的特征趋于消失,而呈现出某种同一现象,认为这其中一定存在着某些原因,而绝非出于偶然。④进行了大样本推断的尝试。在统计发展史上,人口的推算问题多少年来一直是统计学家耿耿于怀的难题,直到 19 世纪初,拉普拉斯才用概率论的原理迈出了关键一步。1781 年~1786 年他提出了“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。在统计发展史上,他利用样本来推断总体的思想方法,开创了一条抽样调查的新思路。

#### (2) 高斯

另一位在概率论与统计学结合的研究上作出贡献的是德国数学家高斯(C. F. Gauss, 1777—1855)。他的主要贡献:①建立最小二乘法。1795 年,高斯设想以残差平方和  $\sum(Y_i - a - bx_i)^2$  为最小的情况下,求得的  $a$  与  $b$  来估计  $\alpha$  与  $\beta$ ;1798 年他完成了最小二乘法的整个构思与结构,并正式发表于 1809 年。②发现高斯分布。高斯以他丰富的数学实践经验,发现观察值  $x$  与真正值  $\mu$  的误差变异服从现代人们最熟悉的正态分布,他运用极大似然法及其他数学方法推导出测量误差的概率分布公式。“误差概率分布曲线”这个术语就是高斯提出来的,后人为了纪念他,称该分布曲线为高斯曲线,也就是今天的正态分布曲

线。高斯所发现的一般误差概率分布曲线以及据此来测定误差的方法,不仅在理论上,而且在应用上都有极为重要的意义。

## 2) 近代描述统计学阶段

近代描述统计学形成期间大致在 19 世纪中叶至 20 世纪上半叶。由于生物学家们为了解决达尔文进化论中的复杂问题,经常需要借助统计学手段,而在这个过程中,原有的统计学方法的不足与局限性逐步地暴露出来。因此,许多学者在改善手段方面做了许多工作。19 世纪达尔文应用统计方法研究生物界地连续性变异;孟德尔应用统计方法发现显性、分离、独立分配等遗传定律。由于这种“描述”特色由一批研究生物进化的学者们提炼而成,因此历史上称他们为生物统计学派。生物统计学派的创始人是英国的高登(F. Galton, 1822—1911),而主要发展是由高登的得意门生泊松(K. Poisson, 1857—1936)完成的。

### (1) 高登

高登的主要贡献:①初创生物统计学。高登自 1882 年起开设“人体测量实验室”,在连续 6 年中,共测量了 9 337 人的“身高、体质量、阔度、呼吸力、效力和压力、手击的速率、听力、视力、色觉及个人的其他资料”,他深入钻研这些资料中隐藏着的内在联系,最终得出“祖先遗传法则”;在极其广泛地收集资料的同时,为了能使他的遗传理论建立在比较精确的基础上,他出色地引入了中位数、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和方法;他在著作 *Natural Inheritance* 中首先提出了“生物统计学”一词,指出“所谓生物统计学,是应用于生物学中的现代统计方法”。②对统计学的贡献。变异是进化论中的重要概念,高登首先以统计方法加以处理,最终在英国创立了生物统计学派。1889 年,高登把总体的定量测定法引入遗传研究中。通过总体测量发现,对动物或植物的每一个种别都可以决定一个平均类型,在一个种别中,所有个体都围绕着这个平均类型,并把它当作轴心向多方面变异。这就是他提出的“平均离差法则”。关于“相关”,统计相关法是由高登创造的。关于相关研究的起因,最早是他因度量甜豌豆的大小,觉察到子代在遗传后有“返于中亲”的现象。1877 年,他搜集了大量人体身高数据后,计算分析高个子父母以及一高一矮父母的后代各有多少个高个子和矮个子子女,从而把“父母高的后代高个子比较多,父母矮的其后代高个子比较少”这一认识具体化为父母与子女之间在身高方面的定量关系。高登在研究人类身高的遗传时发现,高个子父母的子女,其身高有低于他们父母身高的趋势;相反,矮个子父母的子女,其身高却往往有高于他们父母身高的趋势。这就是统计学上“回归”的最初含义。1886 年,高登在论文“在遗传的身高中向中等身高的回归”中,正式提出了“回归”概念。

### (2) 泊松

对生物统计学倾注心血,并把它上升到通用方法论高度的是泊松。他对统计学的主要贡献:①变异数据的处理。生物统计中所取得的数据常常是零乱的,很难看出其规律,泊松首创的频数分布图成为统计方法中最基本的手段之一。②分布曲线的选配。19 世纪以前,人们认为以频数分布描述变异值,最终都表现为正态分布曲线。但是,泊松从生物统计资料的经验分布中,注意到许多生物上的度量不具有正态分布,而常常呈偏态分布,甚至倾斜度很大,而且也不一定都是单峰,也有非单峰的。1894 年,他在“关于不对称频率曲线的分解”一文中首先把非对称的观察曲线分解为几个正态曲线,利用所谓“相对斜率”的方法得到 12 种分布函数型,其中包括正态分布、矩形分布、J 形分布、U 形分布或铃形分布等。③卡方检

验的提出。1900 年泊松发现了  $\chi^2$  分布，并提出了有名的“卡方检验法”。在自然现象的范围内， $\chi^2$  检验法运用得很广泛，以后经费歇补充，成了小样本推断统计的早期方法之一。<sup>④</sup> 回归与相关的发展。回归与相关，经泊松进一步发展后，这两个出自于生物统计学领域的概念，便被推广为一般统计方法论的重要概念。此外，泊松还提出复相关、总相关、相关比等概念，不仅发展了高登的相关理论，还为之建立了数学基础。

### 3) 现代推断统计学阶段

现代推断统计学形成期间大致是 20 世纪初叶至 20 世纪中叶。人类历史进入 20 世纪后，无论社会领域还是自然领域都向统计学提出了更多的要求。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，相继产生了“推断”的方法来掌握事物总体的真正联系以及预测未来的发展。从描述统计学到推断统计学，这是统计发展过程中的一个大飞跃。统计学发展中的这场深刻变革是在农业田间试验领域中完成的，因此，历史上称之为农业试验学派。对现代推断统计的建立贡献最大的是英国统计学家哥赛特(Willian Seely Gosset, 1876—1937)和费歇。

#### (1) 哥赛特的 t 检验与小样本思想

1908 年，哥赛特首次以 Student 为笔名，在 Biometrika 杂志上发表了“平均数的概率误差”一文。由于这篇文章提供了“学生 t 检验”的基础，为此，许多统计学家把 1908 年看做是统计推断理论发展史上的里程碑。此后，哥赛特又连续发表了“相关系数的概率误差”、“非随机抽样的样本平均数分布”、“从无限总体随机抽样平均数的概率估算表”等论文。他在这些论文中，比较了平均误差与标准误差的两种计算方法，研究了泊松分布应用中的样本误差问题，建立了相关系数的抽样分布，导入了“学生氏分布”，即 t 分布。这些论文的完成，为“小样本理论”奠定了基础。由于哥赛特开创的理论使统计学开始由大样本向小样本、由描述向推断发展，因此有人把哥赛特推崇为推断统计学的先驱者。人们认为哥赛特研究成果的战略意义远比其战术意义大，它打开了人们的思路，启发后人发展出许多统计方法。

#### (2) 费歇的统计理论与方法

费歇一生先后共写作论文 329 篇，他在统计学方面的贡献是多方面的：① 通用方法论。费歇非常强调统计学是一门通用方法论，他认为无论对各种自然现象还是社会生活现象的研究，统计方法及其计算公式“正如同其他数学课目一样，这里同一公式适用于一切问题的研究”。② 假设无限总体。费歇认为，在研究各种事物现象包括社会经济现象时，必须把具体物质内容的信息舍弃掉，使统计处理的只是“统计总体”，例如我们已有关于 1 万名新兵身长的资料，那么，统计研究的对象不是新兵的整体，而是各种身长尺寸的总体。他在 1922 年所写的“关于理论统计学的数学基础”一文中，提出了一个重要的概念：假设无限总体，所谓假设无限总体，即现有的资料就是它的随机样本。③ 抽样分布。费歇跨进统计学界就是从研究概率分布开始的。1915 年，他在 Biometrika 杂志上发表了“无限总体样本相关系数值的概率分布”一文。由于这篇论文对相关系数的一般公式做了论证，对以后整个推断统计的发展也有一定贡献，因此，有人把这篇论文称为现代推断统计学的第一篇论文。1922 年，费歇导出相关系数  $r$  的 Z 分布，后来还编制了“Z 曲线末端面积为 0.05, 0.01 和 0.001 的 Z 数值分布表”；1924 年，费歇对 t 分布、 $\chi^2$  分布和 Z 分布加以综合研究，使哥赛特的 t 检验也能适用于大样本，使泊松的  $\chi^2$  检验也能适用于小样本；1938 年，费歇与耶茨合编了“F 分布显著性水平表”，为该分布的研究与应用提供了方便。④ 方差分析。方差和方差分析两词，由