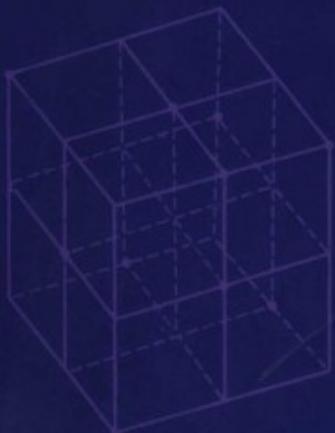


研究生应用数学丛书

# 应用统计

● 刘达民 程岩 编



化学工业出版社

PDG

# 研究生应用数学丛书

- 高等应用数学  
——非线性分析
- 矩阵论及应用
- 应用统计
- 张量分析及应用

ISBN 7-5025-5696-6



9 787502 556969 >

ISBN 7-5025-5696-6/O · 54 定价：30.00元

PDG

研究生应用数学丛书

# 应    用    统    计

刘达民 程岩 编



(京)新登字039号

**图书在版编目(CIP)数据**

应用统计/刘达民, 程岩编. —北京: 化学工业出版社,  
2004. 6

(研究生应用数学丛书)

ISBN 7-5025-5696-6

I. 应… II. ①刘… ②程… III. 应用统计学-研究  
生-教材 IV. C8

中国版本图书馆CIP数据核字(2004)第070747号

---

研究生应用数学丛书

**应 用 统 计**

刘达民 程岩 编

责任编辑:任文斗

文字编辑:宋 薇

责任校对:李 林

封面设计:蒋艳君

\*

化学工业出版社出版发行

(北京市朝阳区惠新里3号 邮政编码100029)

发行电话:(010) 64982530

<http://www.cip.com.cn>

\*

新华书店北京发行所经销

北京市彩桥印刷厂印刷

三河市前程装订厂装订

开本 787mm×960mm 1/16 印张 17 1/4 字数 315 千字

2004年8月第1版 2004年8月北京第1次印刷

ISBN 7-5025-5696-6/O·54

定 价: 30.00 元

---

版权所有 违者必究

该书如有缺页、倒页、脱页者, 本社发行部负责退换

PDG

## 前　　言

应用统计对于工科研究生是一门十分重要的课程，它既是许多专业的数学基础，又直接提供实用数学方法，特别是随着计算机的广泛应用，应用统计的理论和方法为处理科学的研究和工程技术中众多受随机因素影响的实际问题，提供了有力的工具，日益显示其重要性。

本书编写过程中针对工科研究生的特点，在教材的材料安排、习题的选用上，突出实用性，体现了本课程应用性强的特色。为适应解决实际问题时，处理大量数据的需要，专辟一章介绍统计软件 SPSS 的使用，使读者能应用它去解决诸多的实际问题。全书行文力求通俗易懂，便于读者自学。

本书由刘达民作整体安排并编写了第 1 章至第 6 章，程岩编写了第 7 章。

本书得到北京化工大学“新世纪教学改革项目”及“研究生教育创新基金”的资助；在编写过程中还得到刘慧、王云的大力支持和帮助，在此表示感谢。

限于编者学识水平，书中难免存在不当之处，欢迎读者批评指正。

编者

2004 年 4 月于北京化工大学



## 内 容 提 要

本书是根据教委颁布的“工学硕士研究生应用统计课程基本要求”编写的，内容除了包括数理统计的基本概念、参数估计、假设检验、回归分析、方差分析、正交设计外，还补充了统计软件 SPSS 的应用简介。本书注重阐明统计思想和介绍各种统计方法，强调统计实际应用，并通过材料的选择安排、问题的引入、内容的阐述、例题和习题的配置等环节体现上述特色。由于是以工学硕士研究生为主要对象，全书论述严谨，不仅追求数学上的严密性和完整性，而且对未给出推导的结论指明出处。本书行文深入浅出，注意启发性，书末附有概率论等基础知识，便于读者自学。

本书可作为工学硕士研究生应用统计课程的教材，也可作为高等院校高年级学生、教师和科研人员、工程技术人员的参考书。



# 目 录

<b>第1章 数理统计的基本概念</b> .....	1
1.1 数理统计的基本问题 .....	1
1.2 总体和样本 .....	2
1.2.1 总体 .....	2
1.2.2 样本 .....	3
1.3 经验分布函数和直方图 .....	4
1.3.1 经验分布函数 .....	4
1.3.2 直方图 .....	6
1.4 统计量 .....	8
1.4.1 统计量 .....	8
1.4.2 顺序统计量 .....	9
1.5 由正态分布导出的抽样分布 .....	10
1.5.1 $\chi^2$ 分布 .....	11
1.5.2 t 分布和 F 分布 .....	13
1.5.3 分位数 .....	14
1.5.4 正态总体的抽样分布 .....	15
习题 1 .....	19
<b>第2章 参数估计</b> .....	22
2.1 参数估计问题 .....	22
2.2 点估计与求估计量的方法 .....	24
2.2.1 矩法估计 .....	24
2.2.2 极大似然估计 .....	27
2.3 估计量的评价准则 .....	30
2.3.1 无偏性 .....	30
2.3.2 有效性 .....	31
2.3.3 其他准则 .....	35
2.4 贝叶斯估计 .....	36
2.5 区间估计 .....	38
2.5.1 区间估计的概念 .....	38
2.5.2 单个正态总体参数的区间估计 .....	39

2.5.3 两个正态总体参数的区间估计	42
2.5.4 非正态总体参数的区间估计	44
习题 2	45
<b>第 3 章 假设检验</b>	<b>48</b>
3.1 假设检验思想	48
3.1.1 假设检验的提出	48
3.1.2 假设检验的基本原理	49
3.1.3 假设检验的步骤	50
3.1.4 尾概率	51
3.2 总体参数的假设检验	51
3.2.1 单个正态总体的参数检验	51
3.2.2 两个正态总体参数的检验	55
3.3 非正态总体参数的假设检验	59
3.3.1 非正态总体的均值检验	59
3.3.2 指数分布总体的参数检验	61
3.4 检验的实际意义及两类错误	62
3.4.1 检验结果的实际意义	62
3.4.2 两类错误和犯两类错误的概率	63
3.4.3 样本容量的确定	66
3.5 非参数假设检验	67
3.5.1 皮尔逊 $\chi^2$ 拟合检验	68
3.5.2 柯尔莫哥洛夫检验法	75
3.5.3 斯米尔诺夫检验法	77
3.5.4 秩和检验法	79
习题 3	81
<b>第 4 章 回归分析</b>	<b>84</b>
4.1 相关关系与回归	84
4.2 一元线性回归	85
4.2.1 基本概念	85
4.2.2 最小二乘估计及其性质	86
4.2.3 相关系数与回归显著性检验	93
4.2.4 预测与控制	101
4.3 多元线性回归	104
4.3.1 多元线性模型	104
4.3.2 $\beta$ 的估计及其主要性质	105

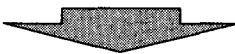
4.3.3 回归方程和回归系数的显著性检验 .....	110
4.3.4 最优回归方程的选择 .....	117
4.4 可化为线性的非线性回归 .....	119
4.4.1 多项式回归 .....	119
4.4.2 可线性化的曲线问题 .....	121
习题 4 .....	127
<b>第 5 章 方差分析.....</b>	<b>130</b>
5.1 单因子试验方差分析 .....	131
5.1.1 基本概念 .....	131
5.1.2 数学模型与分布假设 .....	132
5.1.3 统计分析 .....	133
5.2 双因子试验方差分析 .....	140
5.2.1 基本概念与数学模型 .....	140
5.2.2 平方和分解式 .....	142
5.2.3 检验统计量的分布 .....	145
5.2.4 无交互作用的方差分析 .....	148
习题 5 .....	151
<b>第 6 章 正交试验设计.....</b>	<b>153</b>
6.1 多因素试验 .....	153
6.1.1 全面试验 .....	154
6.1.2 单因素试验轮换法 .....	154
6.1.3 部分因子试验 .....	155
6.2 正交表和正交试验方案 .....	156
6.2.1 正交表 .....	156
6.2.2 正交试验方案 .....	157
6.2.3 正交试验方案的合理性 .....	158
6.3 正交试验的数据分析和统计模型 .....	159
6.3.1 试验的极差分析 .....	160
6.3.2 统计模型和参数估计 .....	161
6.3.3 方差分析和假设检验 .....	164
6.4 有交互作用的正交试验 .....	167
6.4.1 有交互作用的正交试验设计 .....	167
6.4.2 实例分析 .....	169
习题 6 .....	174
<b>第 7 章 统计软件 SPSS 应用简介 .....</b>	<b>176</b>

7.1 SPSS 的窗口、菜单、命令和对话框 .....	176
7.1.1 SPSS 的窗口 .....	176
7.1.2 SPSS 中的菜单 .....	179
7.1.3 SPSS 中的命令和对话框 .....	181
7.2 数据的建立、编辑和基本统计分析 .....	184
7.2.1 数据文件的建立 .....	185
7.2.2 数据的删除或移动 .....	185
7.2.3 插入一个新变量 .....	185
7.2.4 插入一个新观测个体 .....	185
7.2.5 查找指定的观测个体 .....	186
7.2.6 给观测个体排序 .....	186
7.2.7 数据的转置 .....	187
7.2.8 已存在的变量生成新变量 .....	187
7.2.9 生成秩变量 .....	189
7.2.10 变量重新编码与自动重新编码 .....	190
7.2.11 数据的频数统计 .....	193
7.2.12 描述统计量 .....	196
7.2.13 数据探索与统计图 .....	197
7.3 参数的假设检验与区间估计 .....	202
7.3.1 点估计量 .....	202
7.3.2 单样本均值的 $t$ 检验 .....	202
7.3.3 相互独立两样本的 $t$ 检验 .....	205
7.3.4 成对数据均值的 $t$ 检验 .....	207
7.4 线性回归模型 .....	209
7.4.1 一元线性回归 .....	209
7.4.2 多元线性回归 .....	212
7.4.3 曲线估计 .....	216
7.5 方差分析模型 .....	219
7.5.1 单因子方差分析 .....	219
7.5.2 双因子方差分析 .....	220
部分习题答案 .....	224
附录 I 概率论基本知识 .....	228
附录 II $\Gamma$ 函数和 B 函数 .....	242
附录 III 附表 .....	243
附表 1 泊松分布表 .....	243

附表 2 标准正态分布表 .....	244
附表 3 t 分布分位数表 .....	246
附表 4 $\chi^2$ 分布分位数表 .....	247
附表 5 F 分布分位数表 .....	250
附表 6 柯尔莫哥洛夫检验的临界值 ( $D_{n,\alpha}$ ) 表 .....	261
附表 7 $\hat{D}_n$ 的临界值 ( $\hat{D}_{n,\alpha}$ ) 表 .....	263
附表 8 秩和检验表 .....	264
附表 9 常用正交表 .....	265
<b>参考文献</b> .....	<b>272</b>



## 第 1 章



### 数理统计的基本概念

随机性在自然现象和社会现象的普遍存在，促进了数理统计的发展。近几十年来，特别是计算机技术的长足进步，为数理统计的广泛应用提供了广阔的前景。数理统计首先就是因为生物学、遗传学和农业科学的研究的需要而兴起的。在近一个世纪的发展中，数理统计几乎不同程度地渗到所有人类活动的领域。在农业方面，方差分析已经是农业试验的常规手段；在工业生产中，正交试验方法在新产品、新工艺、新材料的开发研究过程中得到广泛应用；在医学中，显著性检验是说明一些药物和治疗手段疗效的典型方法；在国防尖端武器的研制中，精度分析主要也是用数理统计的方法；特别是随着电子计算机的发展，诸如回归分析、多元统计分析等数理统计方法更是在测量、通信、质量控制、气象、水文、地震预报、地质探矿、市场预测、保险等各方面得到越来越多的应用。因此，可以说数理统计是一门应用性很强的数学学科。

#### 1.1 数理统计的基本问题

数理统计的研究对象是带随机性影响的数据，任务是如何有效地、合理地收集、整理、分析这些数据，并利用所得数据对所观察的现象作出推断或预测，直到为采取决策提供依据为止。

数理统计不同于工农业生产中各项经济指标完成情况的统计，也不同于政府部门中的各项资料的汇总。它侧重于应用随机现象本身的统计规律性来解决问题，研究如何合理、有效地获得数据资料，并利用有限资料对所关心的问题，做出尽可能精确、可靠的结论。

由于随机现象的统计规律性必然在大量重复试验中呈现出来，因而从理论上讲，只要对随机现象进行足够多次的观察便能揭示出规律性。然而，实际上所允许的观察或试验往往是有限的，甚至是少量的。如在质量管理中，若对产品的检验是破坏性的（如产品是弹药，必须通过击发来检验是否合格），这时要确定该批弹药是否合格，只能从中随机抽取若干件，通过击发

检验其是否是合格的，得出合格品数，并由此对该批弹药的合格率做出推测（即由局部来推断整体）；有时这种检验虽不是破坏性的，但是限于财力、人力和时间，也无法对该批产品中的所有产品逐个检验以确定该批产品的合格率，不得不采用随机抽取一些产品进行检验。

在产品检验中的上述两类情况，均由局部来推断整体，显而易见，这样做的代价是无法得到完全正确的结论，只能得到某种意义上正确的推断。

在实际问题中，数据的随机性通常是无法避免的。它的来源大致有两个方面：一方面是由“局部性”带来的，这类问题往往是研究的对象数量很大，不可能也没有必要对它们全部加以考察，只能抽取一部分来加以研究，尽管从抽取方式来说应力求能较全面地反映对象的全部信息，但由于只是抽取其中的一部分，所以难免有偶然性；另一方面则是由“不确定性”带来的，如在产品生产中，即使使用同样的材料、设备、工艺流程，所生产的产品质量仍然有差异，这是因为实际上总有一些因素无法控制（如生产环境中温度、湿度的微小变化等）或不便控制，这就使质量指标数据具有不确定性。

数理统计的基本内容大致包括数据采集和统计推断两个方面，两者在应用中都很重要，而且关系密切。数据采集主要包含抽样方法和试验设计等内容，统计推断则包括估计和检验两类大问题。本书主要介绍统计推断的基本原理和方法。

## 1.2 总体和样本

### 1.2.1 总体

通常把研究对象的全体称为**总体**（Population），把总体中的每个元素称为**个体**。如一批产品、一个学校中的全体学生都构成各自的总体，其中每一件产品、每一个学生则是该总体中的个体。

在实际中，人们关心的往往是总体的某个特征指标及其分布状况。如研究显像管的质量，关心的是其寿命（反映质量好坏的一个重要指标）而不是显像管本身。由于各显像管（即使属同一批、同一类型）的寿命不全相同，不可能也没有必要逐个地测出每个显像管的寿命，而只需了解全体显像管的寿命分布状况。由于任一个显像管的寿命测试前是不能确定的，但每一个显像管都确实对应着一个寿命值，所以可以认为显像管的寿命是个随机变量，而人们关心的正是这个随机变量的概率分布。一般说来，都可以认为所考察的总体是用一个随机变量来代表的。由这个观点总体可描述为：**总体就是一**

一个具有确定概率分布的随机变量. 以后, 可以说总体  $F(x)$  或总体  $X$ , 其含义是总体是一个以  $F(x)$  为分布函数的随机变量  $X$ . 当然这个随机变量也可能是二维的, 这时就可以说二维总体  $F(x, y)$ , 其含义就是总体是一个以  $F(x, y)$  为分布函数的二维随机变量  $(X, Y)$ .

### 1.2.2 样本

在数理统计中, 总体的分布总是未知的. 为了获得总体的分布, 就必须对总体进行抽样观察. 从总体中抽取一个个体, 就是做一次试验, 如果进行了  $n$  次抽样观察, 并记录其试验结果, 就得到总体的一组数值, 记作  $(x_1, x_2, \dots, x_n)$ , 称之为样本容量为  $n$  的样本观测值, 也简称为样本值. 人们自然希望这一组样本值能很好地反映总体的情况, 这就要求对抽取方法加上一定的限制. 容易想到, 如果总体中每个个体被抽到的机会均等, 并且在抽取一个个体后总体的成分不改变, 那么, 抽得的个体就能很好地反映总体的情况. 基于这种想法的抽取个体的方法称为简单随机抽样, 抽得的一组个体称为一组简单样本观测值. 换句话说, 简单随机抽样就是独立地、重复地做随机试验, 从而能使试验结果既具有独立性, 又有代表性, 便于理论分析处理. 简单随机抽样相当于概率论中的有放回抽样, 当从总体中抽取容量为  $n$  的样本 (Sample) 时, 每抽取一个个体作观测后立即放回并搅匀, 然后再抽取下一个个体; 但实际情形中往往采取的是无放回抽样, 特别是当总体所含的个体个数  $N$  很大, 而样本容量  $n$  相对较小 (一般是  $n/N$  不大于 0.1 时), 可以近似地把无放回抽样看做是有放回抽样.

设  $(x_1, x_2, \dots, x_n)$  为总体  $X$  的一组样本值, 由于抽样的随机性与独立性, 每个  $x_i$  都可以看做某一个随机变量  $X_i$  ( $i=1, 2, \dots, n$ ) 所取的观测值, 这里  $X_1, X_2, \dots, X_n$  相互独立, 且皆与总体  $X$  具有相同分布. 因此,  $(x_1, x_2, \dots, x_n)$  又可以看做  $n$  维随机变量  $(X_1, X_2, \dots, X_n)$  的观测值, 就一次观察结果而言,  $(x_1, x_2, \dots, x_n)$  是完全确定的一组值, 但它又随每次抽样观察而改变, 具有随机性, 这就是说抽样观察结果具有两重性, 在不会混淆的情况下, 今后, 总是将来自总体  $X$  的容量为  $n$  的样本记为  $(X_1, X_2, \dots, X_n)$ .

以上讨论的概念以定义和定理的形式表述如下.

**定义 1.2.1** 设  $(X_1, X_2, \dots, X_n)$  为来自总体  $X$  的容量为  $n$  的样本, 如果  $X_1, X_2, \dots, X_n$  相互独立且每个都是与总体  $X$  具有相同分布的随机变量, 则称  $(X_1, X_2, \dots, X_n)$  为总体  $X$  的容量为  $n$  的简单随机样本, 简称为样本. 它们的观测值  $(x_1, x_2, \dots, x_n)$  称为总体  $X$  的  $n$  个独立的观测值.

随机向量  $(X_1, X_2, \dots, X_n)$  所有可能取值的全体称为样本空间, 一个

样本值  $(x_1, x_2, \dots, x_n)$  就是样本空间中的一个点.

**定理 1.2.1** 若  $(X_1, X_2, \dots, X_n)$  是来自具有分布函数  $F(x)$  的总体  $X$  的样本, 则  $(X_1, X_2, \dots, X_n)$  的联合分布函数为  $F_n(x_1, x_2, \dots, x_n) = \prod_{i=1}^n F(x_i)$ .

又若总体  $X$  具有概率密度函数  $f(x)$ , 则  $(X_1, X_2, \dots, X_n)$  的联合概率密度函数为  $f_n(x_1, x_2, \dots, x_n) = \prod_{i=1}^n f(x_i)$ .

**例 1.2.1** 总体  $X \sim b(1, p)$ , 它的概率函数为  $f(x) = p^x (1-p)^{1-x}$ ,  $x=0, 1$ .

因此, 样本  $(X_1, X_2, \dots, X_n)$  的联合概率函数为

$$\begin{aligned} f_n(x_1, x_2, \dots, x_n) &= \prod_{i=1}^n [p^{x_i} (1-p)^{1-x_i}] \\ &= p^{\sum_{i=1}^n x_i} (1-p)^{n-\sum_{i=1}^n x_i}, \quad x_i = 0, 1, i = 1, 2, \dots, n \end{aligned}$$

**例 1.2.2** 总体  $X \sim N(\mu, \sigma^2)$ , 它的概率密度函数为  $f(x) = \frac{1}{\sqrt{2\pi}\sigma} e^{-\frac{(x-\mu)^2}{2\sigma^2}}$ ,  $-\infty < x < +\infty$ .

因此, 样本  $(X_1, X_2, \dots, X_n)$  的联合概率密度函数为

$$\begin{aligned} f_n(x_1, x_2, \dots, x_n) &= \prod_{i=1}^n \frac{1}{\sqrt{2\pi}\sigma} e^{-\frac{(x_i-\mu)^2}{2\sigma^2}} \\ &= \frac{1}{(2\pi\sigma^2)^{\frac{n}{2}}} e^{-\frac{\sum_{i=1}^n (x_i-\mu)^2}{2\sigma^2}}, \quad -\infty < x_i < +\infty, \quad i = 1, 2, \dots, n \end{aligned}$$

### 1.3 经验分布函数和直方图

总体分布未知时, 要用样本对总体进行非参数推断, 常用的方法是经验分布函数和直方图.

#### 1.3.1 经验分布函数

总体  $X$  的分布函数, 称为理论分布或总体分布, 取自总体的一个容量为  $n$  的样本的分布函数称为经验分布函数.

**定义 1.3.1** 从总体  $X$  中抽取容量为  $n$  的样本  $(X_1, X_2, \dots, X_n)$ , 将其样本值从小到大排列后为  $x_{(1)} \leq x_{(2)} \leq \dots \leq x_{(n)}$ , 对任何实数  $x$ , 定义函数

$$F_n^*(x) = \begin{cases} 0, & x < x_{(1)} \\ \frac{k}{n}, & x_{(k)} \leq x < x_{(k+1)}, k=1, 2, \dots, n-1 \\ 1, & x \geq x_{(n)} \end{cases} \quad (1.3.1)$$

称  $F_n^*(x)$  为总体  $X$  的经验分布函数.

$F_n^*(x)$  的图形如图 1.1 所示, 是呈跳跃上升的一条阶梯形折线, 若样本值均不同, 则在每一观测值  $x_{(k)}$  处有一跳跃为  $\frac{1}{n}$ , 若有相同的观测值, 则  $F_n^*(x)$  在此处按  $\frac{1}{n}$  的倍数 (即相同观测值的个数) 跳跃上升.

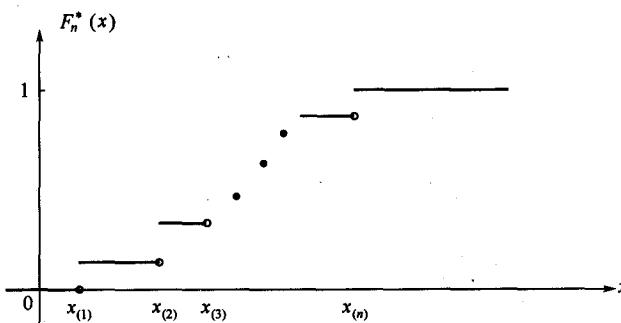


图 1.1 经验分布函数

显然,  $F_n^*(x)$  具有分布函数的一切性质, 即有  $0 \leq F_n^*(x) \leq 1$ ,  $F_n^*(x)$  非降且右连续.

**例 1.3.1** 对某厂生产的电子仪器作寿命测试, 得样本值 (单位: 100h) 为 5, 3, 7, 5, 4, 试求经验分布函数.

解 样本值 (按从小到大顺序排列) 的频率分布为

样本值	3	4	5	7
频率	$\frac{1}{5}$	$\frac{1}{5}$	$\frac{2}{5}$	$\frac{1}{5}$

将上表视为某个离散型随机变量的概率分布列, 与其相应的分布函数便是经验分布函数

$$F_5^*(x) = \begin{cases} 0, & x < 3 \\ \frac{1}{5}, & 3 \leq x < 4 \\ \frac{2}{5}, & 4 \leq x < 5 \\ \frac{4}{5}, & 5 \leq x < 7 \\ 1, & x \geq 7 \end{cases}$$

自然会想到未知的总体分布函数是否可以用经验分布函数去近似？显然取决于 $|F_n^*(x) - F_n(x)|$ 这个量的大小。下面的理论结果表明，这种做法是合理的。

**定理 1.3.1** 设 $(X_1, X_2, \dots, X_n)$ 为来自总体 $X$ 的样本，总体分布函数为 $F(x)$ ，对任意的 $x \in (-\infty, +\infty)$ 与任意给定的 $\epsilon > 0$ ，有

$$\lim_{n \rightarrow \infty} P\{|F_n^*(x) - F(x)| \geq \epsilon\} = 0$$

**证明** 对 $\forall x \in (-\infty, +\infty)$ ，定义随机变量

$$Y_i = \begin{cases} 1 & x_i \leq x \\ 0 & x_i > x \end{cases} \quad i=1, 2, \dots, n$$

显然 $Y_1, Y_2, \dots, Y_n$ 是独立同分布的随机变量，每个 $Y_i \sim b(1, p)$ ，其中 $p = P\{Y_i = 1\} = P\{X_i \leq x\} = F(x)$ 。

另外由 $F_n^*(x)$ 的定义得到 $F_n^*(x) = \frac{1}{n} \sum_{i=1}^n Y_i$ 。

由于 $E(Y_i) = P\{X_i \leq x\} = F(x)$ ，易得 $E[F_n^*(x)] = F(x)$ 。

于是，由伯努利（Bernoulli）大数定律即得定理成立。

定理表明，当 $n$ 足够大时，用经验分布函数 $F_n^*(x)$ 来近似未知的总体分布函数 $F(x)$ 的效果是比较理想的。

实际上， $F_n^*(x)$ 还一致地收敛于 $F(x)$ 。格列汶科在 1933 证明了这一更深刻的结果。

**定理 1.3.2** (W. Glivenko 定理)

$$P\{\lim_{n \rightarrow \infty} D_n = 0\} = 1, \text{ 其中 } D_n = \sup_{-\infty < x < +\infty} |F_n^*(x) - F(x)|$$

证明请参阅文献 [1]。

由定理可知，当 $n$ 足够大时， $F_n^*(x)$ 与 $F(x)$ 可以足够地接近，它们的差别的最大值（定理中的 $\sup$ 表示最大值），也会随 $n$ 增大而趋于零，这样的结局是以概率 1 发生的事件。因而，当 $n$ 足够大时，就可以用 $F_n^*(x)$ 来近似代替 $F(x)$ ，这就是以后可以用样本来推断总体的最基本的理论依据。

### 1.3.2 直方图

当总体 $X$ 为连续型随机变量时，总体分布可以用总体密度函数 $f(x)$ 来刻画，当然 $f(x)$ 是未知的，需要用样本来对它进行推断，直方图是一种简便易行的方法，而且比较直观，特别适合于统计现场使用。下面给出绘制直方图的一般步骤。设 $(x_1, x_2, \dots, x_n)$ 为取自总体 $X$ 的样本值。

步骤 1 求出 $x_{(1)} = \min_{1 \leq i \leq n} x_i$ ,  $x_{(n)} = \max_{1 \leq i \leq n} x_i$ .

步骤 2 取 $a$ 略小于 $x_{(1)}$ ,  $b$ 略大于 $x_{(n)}$ ，并把区间 $[a, b]$ 分为 $m$ 个