

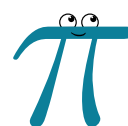
# 数理统计 - 茆诗松等 - 笔记

作者：若水

邮箱：ethanmxzhou@163.com

主页：helloethanzhou.github.io

时间：July 18, 2024



## 致谢

感谢 勇敢的 自己

# 目录

<b>第一章 统计量及其分布</b>	<b>1</b>
1.1 总体与概率	1
1.1.1 总体与个体	1
1.1.2 样本	1
1.2 样本数据的整理与表示	1
1.2.1 经验分布函数	1
1.2.2 频数频率分布表	2
1.3 统计量及其分布	2
1.3.1 统计量与抽样分布	2
1.3.2 样本均值及其抽样分布	2
1.3.3 样本方差	3
1.3.4 样本矩及其函数	3
1.3.5 次序统计量及其分布	4
1.3.6 样本分位数与分位数	5
1.3.7 五数概括与箱线图	6
1.4 三大抽样分布	6
1.4.1 $\chi^2$ 分布	7
1.4.2 $F$ 分布	7
1.4.3 $T$ 分布	8
1.5 充分统计量	8
1.5.1 充分性	8
1.5.2 因子分解定理	8
<b>第二章 参数估计</b>	<b>9</b>
2.1 点估计的概念	9
2.1.1 点估计及无偏性	9
2.1.2 有效性	9
2.2 矩估计及相关性	9
2.2.1 替换原理和矩法估计	9
2.2.2 概率函数已知时未知参数的矩估计	10
2.2.3 相合性	10
2.3 最大似然估计	11
2.3.1 最大似然估计	11
2.3.2 渐进正态性	11
2.4 最小方差无偏估计	12
2.4.1 均方误差	12
2.4.2 一致最小方差无偏估计	12
2.4.3 充分性原则	13
2.4.4 Cramer-Rao 不等式	13
2.5 Bayes 估计	14
2.6 统计判断的基础	14

2.6.1 Bayes 公式的密度函数形式 . . . . .	14
2.6.2 Bayes 估计 . . . . .	14
2.6.3 共轭先验分布 . . . . .	15
2.7 区间估计 . . . . .	15
2.7.1 区间估计的概念 . . . . .	15
2.7.2 枢轴量法 . . . . .	16
2.7.3 单个正态总体参数的置信区间 . . . . .	16
2.7.4 大样本置信区间 . . . . .	16
2.7.5 样本量的确定 . . . . .	17
2.7.6 两个正态总体下的置信区间 . . . . .	17
<b>第三章 假设检验</b>	<b>18</b>
3.1 假设检验的基本思想与概念 . . . . .	18
3.1.1 假设检验问题 . . . . .	18
3.1.2 假设检验的基本步骤 . . . . .	18
3.1.3 检验的 $p$ 值 . . . . .	19
3.2 正态总体参数假设检验 . . . . .	20
3.2.1 单个正态总体均值的检验 . . . . .	20
3.2.2 两个正态总体均值差的检验 . . . . .	21
3.2.3 成对数据检验 . . . . .	22
3.2.4 正态总体方差的检验 . . . . .	22
3.3 其他分布参数的假设检验 . . . . .	23
3.4 似然比检验与分布拟合检验 . . . . .	23
3.4.1 似然比检验的思想 . . . . .	23
3.4.2 分布数据的 $\chi^2$ 拟合优度检验 . . . . .	24
3.4.3 分布的 $\chi^2$ 拟合优度检验 . . . . .	24
3.4.4 列联表的独立性检验 . . . . .	25
3.5 正态性检验 . . . . .	25
3.5.1 正态概率纸 . . . . .	25
3.5.2 W 检验 . . . . .	25
3.5.3 EP 检验 . . . . .	26
3.6 非参数检验 . . . . .	26
3.6.1 游程检验 . . . . .	26
3.6.2 符号检验 . . . . .	26
3.6.3 秩和检验 . . . . .	27
<b>第四章 方差分析与回归分析</b>	<b>28</b>
4.1 方差分析 . . . . .	28
4.1.1 问题的提出 . . . . .	28
4.1.2 单因子方差分析的统计模型 . . . . .	28
4.1.3 平方和分解 . . . . .	29
4.1.4 检验方法 . . . . .	30
4.1.5 参数估计 . . . . .	31
4.1.6 重复数不等情形 . . . . .	31
4.2 多重比较 . . . . .	32

---

4.2.1 水平均值差的置信区间 . . . . .	32
4.2.2 多重比较问题 . . . . .	32
4.2.3 重复数相等的 T 法 . . . . .	33
4.2.4 重复数不等场合的 S 法 . . . . .	33
4.3 方差齐性检验 . . . . .	34
4.3.1 Hartley 检验 . . . . .	34
4.3.2 Bartlett 检验 . . . . .	34
4.3.3 修正的 Bartlett 检验 . . . . .	35
4.4 一元线性回归 . . . . .	35
4.4.1 变量间的两类关系 . . . . .	35
4.4.2 一元线性回归模型 . . . . .	35
4.4.3 回归系数的最小二乘估计 . . . . .	35
4.4.4 回归模型的显著性检验 . . . . .	36
4.4.5 估计与预测 . . . . .	37
4.4.6 曲线回归方程的比较 . . . . .	38
<b>附录 A 概率模型</b>	<b>39</b>

# 第一章 统计量及其分布

## 1.1 总体与概率

### 1.1.1 总体与个体

总体：研究对象的全体

个体：构成总体的每个成员

### 1.1.2 样本

样本：从总体中随机的抽取  $n$  个个体，记其指标值为

$$x_1, \cdots, x_n$$

那么此称为总体的一个样本， $n$  称为**样本容量**，或简称样本量，样本中的个体称为**样品**。

简单随机抽样原则：随机性，独立性

#### 定义 1.1.1 (联合分布函数)

总体  $X$  具有分布函数  $F(x)$ ， $x_1, \cdots, x_n$  为取自该总体的容量为  $n$  的样本，那么样本联合分布函数为

$$F(x_1, \cdots, x_n) = \prod_{k=1}^n F(x_k)$$



## 1.2 样本数据的整理与表示

### 1.2.1 经验分布函数

#### 定义 1.2.1 (经验分布函数)

对于取自总体分布函数为  $F(x)$  的样本  $x_1, \cdots, x_n$ ，记其对应的次序统计量为  $x_{(1)}, \cdots, x_{(n)}$ ，定义该样本的经验分布函数为

$$F_n(x) = \begin{cases} 0, & x < x_{(1)} \\ k/n, & x_{(k)} \leq x < x_{(k+1)}, k = 1, \cdots, n-1 \\ 1, & x \geq x_{(n)} \end{cases}$$



#### 命题 1.2.1 (经验分布函数的性质)

- $F_n(x)$  非减且右连续。
- $F_n(-\infty) = 0, \quad F_n(+\infty) = 1$



#### 定理 1.2.1 (Glivenko 定理)

对于取自总体分布函数为  $F(x)$  的样本  $x_1, \cdots, x_n$ ，记其经验分布函数为  $F_n(x)$ ，那么

$$P\left(\lim_{n \rightarrow \infty} \sup_{x \in \mathbb{R}} |F_n(x) - F(x)| = 0\right) = 1$$



## 1.2.2 频数频率分布表

### 定义 1.2.2 (频数频率分布表)

1. 对样本进行分组：通常为 5 ~ 20 个。
2. 确定每组组距：

$$d = \frac{\text{样本最大观测值} - \text{样本最小观测值}}{\text{组数}}$$

3. 确定每组组限：

$$(a_0, a_1], \dots, (a_{n-1}, a_n]$$

4. 统计样本数据落入每个区间的个数——频数



表 1.1: 频数频率分布表

分组区间	频数	频率
$(a_0, a_1]$	$f_1$	$\frac{f_1}{n}$
$\vdots$	$\vdots$	$\sum_{k=1} f_k$
$(a_{n-1}, a_n]$	$f_n$	$\frac{f_n}{n}$
		$\sum_{k=1} f_k$

## 1.3 统计量及其分布

### 1.3.1 统计量与抽样分布

#### 定义 1.3.1 (统计量)

对于取自总体的样本  $x_1, \dots, x_n$ ，若  $T = T(x_1, \dots, x_n)$  中不含有任何位置参数，那么称  $T$  为统计量。



#### 定义 1.3.2 (抽样分布)

统计量的分布称为抽样分布。



### 1.3.2 样本均值及其抽样分布

#### 定义 1.3.3 (样本均值)

对于取自总体的样本  $x_1, \dots, x_n$ ，其算术平均值称为样本均值，记为  $\bar{x}$ ，即

$$\bar{x} = \frac{1}{n} \sum_{k=1}^n x_k$$

特别的，在分组样本中，样本均值的近似公式为

$$\bar{x} = \frac{1}{n} \sum_{k=1}^m x_k f_k$$

其中  $m$  为组数,  $x_k$  为第  $k$  组的组中值,  $f_k$  为第  $k$  组的频数, 同时

$$n = \sum_{k=1}^n f_k$$



### 命题 1.3.1 (样本均值的性质)

1. 样本的所有偏差之和为 0, 即

$$\sum_{k=1}^n (x_k - \bar{x}) = 0$$

2. 对于任意  $c \in \mathbb{R}$ , 成立

$$\sum_{k=1}^n (x_k - \bar{x})^2 \leq \sum_{k=1}^n (x_k - c)^2$$

当且仅当  $c = \bar{x}$  时等号成立。



### 定理 1.3.1

对于取自总体的样本  $x_1, \dots, x_n$ , 记其样本均值为  $\bar{x}$ 。

1. 如果总体分布为  $N(\mu, \sigma^2)$ , 那么  $\bar{x}$  满足分布  $N(\mu, \frac{\sigma^2}{n})$ 。

2. 对于一般的总体的分布, 记  $E(x) = \mu, \text{Var}(x) = \sigma^2$ , 那么当  $n \rightarrow \infty$  时,  $\bar{x}$  满足近似分布  $N(\mu, \frac{\sigma^2}{n})$ , 记作

$$\bar{x} \sim N(\mu, \frac{\sigma^2}{n})$$



## 1.3.3 样本方差

### 定义 1.3.4 (样本方差)

对于取自总体的样本  $x_1, \dots, x_n$ , 定义其样本方差为

$$s^2 = \frac{1}{n-1} \sum_{k=1}^n (x_k - \bar{x})^2 = \frac{1}{n-1} \left( \sum_{k=1}^n x_k^2 - n\bar{x}^2 \right)$$



### 定理 1.3.2

对于具有二阶矩的总体  $X$ , 即  $E(X) = \mu, \text{Var}(X) = \sigma^2 < +\infty$ , 取自总体的样本  $x_1, \dots, x_n$ , 记  $\bar{x}$  和  $s^2$  分别为样本均值和样本方差, 那么

$$\begin{aligned} E(\bar{x}) &= \mu, & \text{Var}(\bar{x}) &= \frac{\sigma^2}{n} \\ E(s^2) &= \sigma^2 \end{aligned}$$



## 1.3.4 样本矩及其函数

### 定义 1.3.5 (样本原点矩)

对于取自总体的样本  $x_1, \dots, x_n$ , 定义其样本  $k$  阶原点矩为

$$a_k = \frac{1}{n} \sum_{i=1}^n x_i^k$$





**定义 1.3.6 (样本中心矩)**

对于取自总体的样本  $x_1, \dots, x_n$ , 定义其样本  $k$  阶中心矩为

$$b_k = \frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})^k$$

**定义 1.3.7 (样本偏差)**

对于取自总体的样本  $x_1, \dots, x_n$ , 定义其样本偏差为

$$\hat{\beta}_s = \frac{b_3}{b_2^{3/2}}$$

**命题 1.3.2 (样本偏差的性质)**

样本偏差反应总体分布密度曲线的对称性。

1.  $\hat{\beta}_s = 0$ : 完全对称
2.  $\hat{\beta}_s > 0$ : 存在右长尾
3.  $\hat{\beta}_s < 0$ : 存在左长尾

**定义 1.3.8 (样本峰度)**

对于取自总体的样本  $x_1, \dots, x_n$ , 定义其样本峰度为

$$\hat{\beta}_k = \frac{b_4}{b_2^2} - 3$$

**命题 1.3.3 (样本峰度的性质)**

样本峰度反应总体分布密度曲线在其峰值附近的陡峭程度。

1.  $\hat{\beta}_k > 0$ : 比正态分布陡峭, 称为尖顶型
2.  $\hat{\beta}_k < 0$ : 比正态分布平缓, 称为平顶型



### 1.3.5 次序统计量及其分布

**定义 1.3.9 (次序统计量)**

对于取自总体的样本  $x_1, \dots, x_n$ , 称其次序统计量为

$$x_{(1)}, \dots, x_{(n)}$$

其中  $x_{(1)} \leq \dots \leq x_{(n)}$ 。

**定理 1.3.3 (单个次序统计量的分布)**

对于取自总体的样本  $x_1, \dots, x_n$ , 如果  $X$  的密度函数为  $p(x)$ , 分布函数为  $F(x)$ , 那么第  $k$  个次序统计量  $x_{(k)}$  的密度函数为

$$p_k(x) = \frac{n!}{(k-1)!(n-k)!} (F(x))^{k-1} (1-F(x))^{n-k} p(x)$$



**推论 1.3.1**

1.  $x_{(1)}$  的密度函数为

$$p_1(x) = np(x)(1 - F(x))^{n-1}$$

分布函数为

$$F_1(x) = 1 - (1 - F(x))^n$$

2.  $x_{(n)}$  的密度函数为

$$p_n(x) = np(x)(F(x))^{n-1}$$

分布函数为

$$F_n(x) = (F(x))^n$$

**定理 1.3.4 (两个次序统计量的联合分布)**

对于取自总体的样本  $x_1, \dots, x_n$ , 如果  $X$  的密度函数为  $p(x)$ , 分布函数为  $F(x)$ , 那么第  $i$  个次序统计量  $x_{(i)}$  和第  $j$  个次序统计量  $x_{(j)}$  的联合分布密度函数为

$$p_{ij}(x, y) = \frac{n!}{(i-1)!(j-i-1)!(n-j)!} (F(x))^{i-1} (F(x) - F(y))^{j-i-1} (1 - F(y))^{n-j} p(x)p(y)$$

其中  $i < j$ 。

**1.3.6 样本分位数与分位数****定义 1.3.10 (样本  $p$  分位数)**

对于取自总体的样本  $x_1, \dots, x_n$ , 定义其样本  $p$  分位数为

$$m_p = \begin{cases} x_{([np+1])}, & np \notin \mathbb{Z} \\ \frac{1}{2} (x_{(np)} + x_{(np+1)}), & np \in \mathbb{Z} \end{cases}$$

其中  $p \in (0, 1)$ 。

**定义 1.3.11 ( $\alpha$  分位数)**

对于随机变量  $X$ , 称  $x_\alpha$  为其  $\alpha$  分位数, 如果

$$P(X \leq x_\alpha) = \alpha$$

**定理 1.3.5**

如果总体密度函数为  $p(x)$ ,  $x_p$  为其  $p$  分位数,  $p(x)$  在  $x_p$  处连续且  $p(x_p) > 0$ , 那么当  $n \rightarrow +\infty$  时, 样本  $p$  分位数  $m_p$  的渐进分布为

$$m_p \sim N\left(x_p, \frac{p(1-p)}{np^2(x_p)}\right)$$



## 1.3.7 五数概括与箱线图

## 定义 1.3.12 (五数概括)

$$x_{\min}, \quad Q_1 = m_{0.25}, \quad m_{0.5}, \quad Q_3 = m_{0.75}, \quad x_{\max}$$



## 定义 1.3.13 (箱线图)

1. 画一个箱子，其两侧恰为第一 4 分位数和第三 4 分位数，在中位数位置上画一条竖线，其在箱子内，这个箱子包含了样本中 50% 的数据。
2. 在箱子左右两侧各引出一条水平线，分别至最小值和最大值为止。每条线段中包含了样本 25% 的数据。



## 1.4 三大抽样分布

表 1.2: 三大抽样分布

分布名称	表示	统计量的构造	抽样分布密度函数	期望	方差	特征函数
正态分布	$N(\mu, \sigma^2)$		$p(x) = \frac{1}{\sqrt{2\pi}\sigma} e^{-\frac{(x-\mu)^2}{2\sigma^2}}, x \in \mathbb{R}$	$\mu$	$\sigma^2$	$e^{i\mu t - \frac{1}{2}\sigma^2 t^2}$
$\Gamma$ 分布	$\Gamma(\alpha, \lambda)$		$p(x) = \frac{\lambda^\alpha}{\Gamma(\alpha)} x^{\alpha-1} e^{-\lambda x}, x > 0$	$\frac{\alpha}{\lambda}$	$\frac{\alpha}{\lambda^2}$	$(1 - \frac{it}{\lambda})^{-\alpha}$
$\chi^2$ 分布	$\chi^2(n)$	$\chi^2(n) = \sum_{k=1}^n (N(0, 1))^2$	$p(x) = \frac{1}{\Gamma(\frac{n}{2})2^{\frac{n}{2}}} x^{\frac{n}{2}-1} e^{-\frac{x}{2}}, x > 0$	$n$	$2n$	$(1 - 2it)^{-\frac{n}{2}}$
$F$ 分布	$F(m, n)$	$F(m, n) = \frac{\chi^2(m)/m}{\chi^2(n)/n}$	$p(x) = \frac{\Gamma(\frac{m+n}{2})\Gamma(\frac{n}{2})^{\frac{m}{2}}}{\Gamma(\frac{m}{2})\Gamma(\frac{n}{2})} x^{\frac{m}{2}-1} (1 + \frac{m}{n}x)^{-\frac{m+n}{2}}, x > 0$	$\frac{n}{n-2}$	$\frac{2n^2(m+n-2)}{m(n-2)^2(n-4)}$	
$t$ 分布	$t(n)$	$t(n) = \frac{N(0, 1)}{\sqrt{\chi^2(n)/n}}$	$p(x) = \frac{\Gamma(\frac{n+1}{2})}{\sqrt{n\pi}\Gamma(\frac{n}{2})} \left(1 + \frac{x^2}{n}\right)^{-\frac{n+1}{2}}, x \in \mathbb{R}$	0	$\frac{n}{n-2}$	

定义 1.4.1 ( $\Gamma$  函数)

$$\Gamma(x) = \int_0^\infty t^{x-1} e^{-t} dt$$



## 命题 1.4.1 (分布间的联系)

1. 若  $X \sim N(0, 1)$ ，那么

$$X^2 \sim \Gamma\left(\frac{1}{2}, \frac{1}{2}\right)$$

因此

$$\chi^2(n) = \Gamma\left(\frac{n}{2}, \frac{1}{2}\right)$$

2. 如果  $X \sim N(\mu, \sigma^2)$ ，那么

$$aX + b \sim N(a\mu + b, a^2\sigma^2)$$

3. 如果独立分布  $X \sim N(\mu_1, \sigma_1^2)$  和  $Y \sim N(\mu_2, \sigma_2^2)$ ，那么

$$X + Y \sim N(\mu_1 + \mu_2, \sigma_1^2 + \sigma_2^2)$$

4. 如果  $X \sim \Gamma(\alpha, \lambda)$ , 那么

$$aX \sim \Gamma(\alpha, \frac{\lambda}{a})$$

5. 如果独立分布  $X \sim \Gamma(\alpha_1, \lambda)$  和  $Y \sim \Gamma(\alpha_2, \lambda)$ , 那么

$$X + Y \sim \Gamma(\alpha_1 + \alpha_2, \lambda)$$



### 1.4.1 $\chi^2$ 分布

#### 定义 1.4.2 ( $\chi^2$ 分布)

对于独立同分布于标准正态分布的  $N(0, 1)$  的随机变量  $X_1, \dots, X_n$ , 称随机变量  $X = X_1^2 + \dots + X_n^2$  的分布为自由度为  $n$  的  $\chi^2$  分布, 记作  $X \sim \chi^2(n)$ , 其密度函数为

$$p(x) = \frac{1}{\Gamma\left(\frac{n}{2}\right) 2^{\frac{n}{2}}} x^{\frac{n}{2}-1} e^{-\frac{x}{2}}, \quad x > 0$$



#### 命题 1.4.2 ( $\chi^2$ 分布的性质)

对于来自正态总体  $N(\mu, \sigma^2)$  的样本  $x_1, \dots, x_n$ , 记其样本均值和样本方差分别为  $\bar{x}$  和  $s^2$ , 那么

1.  $\bar{x}$  和  $s^2$  相互独立。

2.

$$\bar{x} \sim N\left(\mu, \frac{\sigma^2}{n}\right)$$

3.

$$\frac{(n-1)s^2}{\sigma^2} \sim \chi^2(n-1) = \Gamma\left(\frac{n-1}{2}, \frac{1}{2}\right)$$

即

$$s^2 \sim \Gamma\left(\frac{n-1}{2}, \frac{n-1}{2\sigma^2}\right)$$



### 1.4.2 $F$ 分布

#### 定义 1.4.3 ( $F$ 分布)

对于独立的随机变量  $X \sim \chi^2(m)$  和  $Y \sim \chi^2(n)$ , 称随机变量  $F = \frac{X/m}{Y/n}$  的分布为自由度为  $m$  和  $n$  的  $F$  分布, 记作  $F \sim F(m, n)$ , 其密度函数为

$$p(x) = \frac{\Gamma\left(\frac{m+n}{2}\right) \left(\frac{m}{n}\right)^{\frac{m}{2}}}{\Gamma\left(\frac{m}{2}\right) \Gamma\left(\frac{n}{2}\right)} x^{\frac{m}{2}-1} \left(1 + \frac{m}{n}x\right)^{-\frac{m+n}{2}}, \quad x > 0$$



#### 推论 1.4.1

对于独立的分别来自  $N(\mu_1, \sigma_1^2)$  和  $N(\mu_2, \sigma_2^2)$  的样本  $x_1, \dots, x_m$  和  $y_1, \dots, y_n$ , 那么记其样本方差分别为  $s_x^2$  和  $s_y^2$ , 那么

$$\frac{s_x^2/\sigma_1^2}{s_y^2/\sigma_2^2} \sim F(m-1, n-1)$$



### 1.4.3 $T$ 分布

#### 定义 1.4.4 ( $T$ 分布)

对于独立的随机变量  $X \sim N(0, 1)$  和  $Y \sim \chi^2(n)$ , 称随机变量  $t = \frac{X}{\sqrt{\frac{Y}{n}}}$  的分布为自由度为  $n$  的  $t$  分布, 记作  $t \sim t(n)$ , 其密度函数为

$$p(x) = \frac{\Gamma\left(\frac{n+1}{2}\right)}{\sqrt{n\pi}\Gamma\left(\frac{n}{2}\right)} \left(1 + \frac{x^2}{n}\right)^{-\frac{n+1}{2}}, x \in \mathbb{R}$$



#### 推论 1.4.2

对于来自正态分布  $N(\mu, \sigma^2)$  的样本  $x_1, \dots, x_n$ , 记其样本均值和样本方差分别为  $\bar{x}$  和  $s^2$ , 那么

$$\frac{\sqrt{n}(\bar{x} - \mu)}{s} \sim T(n-1)$$



#### 推论 1.4.3

对于独立的分别来自  $N(\mu_1, \sigma^2)$  和  $N(\mu_2, \sigma^2)$  的样本  $x_1, \dots, x_m$  和  $y_1, \dots, y_n$ , 那么记其样本方差分别为  $s_x^2$  和  $s_y^2$ , 且

$$s_w^2 = \frac{(m-1)s_x^2 + (n-1)s_y^2}{m+n-2}$$

那么

$$\frac{(\bar{x} - \bar{y}) - (\mu_1 - \mu_2)}{s_w \sqrt{\frac{1}{m} + \frac{1}{n}}} \sim T(m+n-2)$$



## 1.5 充分统计量

### 1.5.1 充分性

#### 定义 1.5.1 (充分统计量)

对于来自总体分布函数为  $F(x; \theta)$  的样本  $x_1, \dots, x_n$ , 称统计量  $T = T(x_1, \dots, x_n)$  为  $\theta$  的充分统计量, 如果给定  $T$  的取值后, 样本  $x_1, \dots, x_n$  的条件分布与  $\theta$  无关。



### 1.5.2 因子分解定理

#### 定理 1.5.1 (Fischer-Neyman 因子分解定理)

对于来自总体概率函数为  $f(x; \theta)$  的样本  $x_1, \dots, x_n$ , 那么  $T = T(x_1, \dots, x_n)$  为充分统计量的充分必要条件为, 存在函数  $g(t, \theta)$  和  $h(x_1, \dots, x_n)$ , 使得对于任意的  $\theta$  和  $x_1, \dots, x_n$ , 成立

$$f(x_1, \dots, x_n; \theta) = g(T(x_1, \dots, x_n); \theta)h(x_1, \dots, x_n)$$



#### 定理 1.5.2

对于充分统计量  $T$ , 如果存在函数  $h$ , 使得  $T = h(S)$ , 那么统计量  $S$  也为充分统计量。



## 第二章 参数估计

### 2.1 点估计的概念

#### 2.1.1 点估计及无偏性

##### 定义 2.1.1 (点估计)

对于来自总体的样本  $x_1, \dots, x_n$ , 用于估计未知参数  $\theta$  的统计量  $\hat{\theta} = \hat{\theta}(x_1, \dots, x_n)$  称为  $\theta$  的点估计。

##### 定义 2.1.2 (无偏估计)

对于  $\theta$  的点估计  $\hat{\theta} = \hat{\theta}(x_1, \dots, x_n)$ ,  $\theta$  的参数空间为  $\Theta$ , 称  $\hat{\theta}$  为  $\theta$  的无偏估计, 如果对于任意  $\theta \in \Theta$ , 成立

$$E_{\theta}(\hat{\theta}) = \theta$$

##### 定义 2.1.3 (可估参数)

称参数  $\theta$  为可估参数, 如果存在无偏估计  $\hat{\theta} = \hat{\theta}(x_1, \dots, x_n)$ 。

#### 2.1.2 有效性

##### 定义 2.1.4 (有效性)

对于  $\theta$  的两个无偏估计  $\hat{\theta}_1$  和  $\hat{\theta}_2$ , 称  $\hat{\theta}_1$  比  $\hat{\theta}_2$  有效, 如果对于任意  $\theta \in \Theta$ , 成立

$$\text{Var}(\hat{\theta}_1) \leq \text{Var}(\hat{\theta}_2)$$

且存在  $\theta_0 \in \Theta$ , 使得成立

$$\text{Var}(\hat{\theta}_1) < \text{Var}(\hat{\theta}_2)$$

### 2.2 矩估计及相关性

#### 2.2.1 替换原理和矩法估计

##### 定义 2.2.1 (替换原理)

1. 用样本矩替换总体矩。
2. 用样本矩的函数替换总体矩的函数。

根据替换原理, 在总体分布形式未知场合对参数作出估计:

1. 用样本均值  $\bar{x}$  估计总体均值  $E(X)$ 。
2. 用样本方差  $s^2$  估计总体方差  $\text{Var}(X)$ 。
3. 用事件  $A$  出现的频率估计事件  $A$  发生的概率。
4. 用样本  $p$  分位数估计总体的  $p$  分位数。

**定理 2.2.1 (Khinchin 大数定律)**

对于独立同分布的随机变量序列  $X_1, \dots, X_n$ , 如果对于任意  $i = 1, \dots, n$ , 总体  $X$  的  $k$  阶原点矩  $E(X^k)$  存在, 那么对于任意  $\varepsilon > 0$ , 成立

$$\lim_{n \rightarrow \infty} P \left\{ \left| \frac{1}{n} \sum_{i=1}^n X_i^k - E(X^k) \right| \geq \varepsilon \right\} = 0$$

**2.2.2 概率函数已知时未知参数的矩估计****定义 2.2.2 (矩估计)**

对于具有概率函数  $p(x; \theta_1, \dots, \theta_k)$  的总体, 以及样本  $x_1, \dots, x_n$ , 其中  $(\theta_1, \dots, \theta_k) \in \Theta$  是未知参数或参数向量, 如果总体的  $i$  阶原点矩  $\mu_i$  存在, 而且  $\theta_i = \theta_i(\mu_1, \dots, \mu_k)$ , 其中  $1 \leq i \leq k$ , 那么  $\theta_i$  的矩估计为

$$\hat{\theta}_i = \theta_i(a_1, \dots, a_k), \quad i = 1, \dots, k$$

其中  $a_i$  为样本  $i$  阶原点矩

$$a_i = \frac{1}{n} \sum_{j=1}^n x_j^i, \quad i = 1, \dots, k$$

进一步, 对于  $\theta_1, \dots, \theta_k$  的函数  $\eta = g(\theta_1, \dots, \theta_k)$  的矩估计为

$$\hat{\eta} = g(\hat{\theta}_1, \dots, \hat{\theta}_k)$$

**2.2.3 相合性****定义 2.2.3 (相合性)**

对于未知参数  $\theta$ , 以及  $\theta$  的一个估计量  $\hat{\theta}_n = \hat{\theta}_n(x_1, \dots, x_n)$ , 称  $\hat{\theta}_n$  为参数  $\theta$  的相合估计, 如果对于任意  $\varepsilon > 0$ , 成立

$$\lim_{n \rightarrow \infty} P \left( \left| \hat{\theta}_n - \theta \right| \geq \varepsilon \right) = 0$$

**定理 2.2.2 (相合估计的充分条件)**

对于  $\theta$  的一个估计量  $\hat{\theta}_n = \hat{\theta}_n(x_1, \dots, x_n)$ , 如果

$$\lim_{n \rightarrow \infty} E(\hat{\theta}_n) = \theta, \quad \lim_{n \rightarrow \infty} \text{Var}(\hat{\theta}_n) = 0$$

那么  $\hat{\theta}_n$  为参数  $\theta$  的相合估计。

**定理 2.2.3 (相合估计在连续函数下的像为相合估计)**

如果  $\hat{\theta}_{n_1}, \dots, \hat{\theta}_{n_k}$  分别是  $\theta_1, \dots, \theta_k$  的相合估计,  $\eta = g(\theta_1, \dots, \theta_k)$  是连续函数, 那么  $\hat{\eta} = g(\hat{\theta}_{n_1}, \dots, \hat{\theta}_{n_k})$  是  $\eta$  的相合估计。



## 2.3 最大似然估计

### 2.3.1 最大似然估计

#### 定义 2.3.1 (似然函数)

对于概率函数为  $p(x; \theta)$  的总体，其中  $\theta \in \Theta$  为一个或多个未知参数组成的参数向量， $\Theta$  为参数空间， $x_1, \dots, x_n$  是来自该总体的样本，称样本的联合概率函数

$$L(\theta) = L(\theta; x_1, \dots, x_n) = \prod_{k=1}^n p(x_k; \theta)$$

为样本的似然函数。



#### 定义 2.3.2 (最大似然估计 MLE)

对于概率函数为  $p(x; \theta)$  的总体，其中  $\theta \in \Theta$  为一个或多个未知参数组成的参数向量， $\Theta$  为参数空间， $x_1, \dots, x_n$  是来自该总体的样本，统计量  $\hat{\theta} = \hat{\theta}(x_1, \dots, x_n)$  为  $\theta$  的最大似然估计，如果对于任意  $\theta \in \Theta$ ，成立

$$L(\hat{\theta}) \geq L(\theta)$$



#### 定理 2.3.1 (最大似然估计的不变性)

如果  $\hat{\theta}$  为  $\theta$  的最大似然估计，那么对于任意函数  $g$ ， $g(\hat{\theta})$  是  $g(\theta)$  的最大似然估计。



#### 定理 2.3.2 (正态分布参数的最大似然估计)

对于来自正态分布  $N(\mu, \sigma^2)$  的样本  $x_1, \dots, x_n$ ，记样本均值为  $\bar{x}$ ，样本方差为  $s^2$ ，那么  $\mu$  和  $\sigma^2$  的最大似然估计分别为

$$\hat{\mu} = \bar{x}, \quad \hat{\sigma}^2 = \frac{n-1}{n} s^2$$



### 2.3.2 渐进正态性

#### 定义 2.3.3 (渐进正态分布)

参数  $\theta$  的相合估计  $\hat{\theta}_n$  称为渐进正态的，如果存在趋于 0 的非负常数序列  $\sigma_n(\theta)$ ，使得成立  $\frac{\hat{\theta}_n - \theta}{\sigma_n(\theta)}$  依分布收敛于标准正态分布。此时也称  $\hat{\theta}_n$  服从渐进正态分布  $N(\theta, \sigma_n^2(\theta))$ ，记为  $\hat{\theta}_n \sim AN(\theta, \sigma_n^2(\theta))$ 。 $\sigma_n^2(\theta)$  称为  $\hat{\theta}_n$  的渐近方差。



#### 定理 2.3.3

对于密度函数为  $p(x; \theta)$  的总体  $X$ ，其中  $\theta \in \Theta$ ，如果

1. 对于任意  $x$ ，以及任意  $\theta \in \Theta$ ，偏导数  $\frac{\partial \ln p}{\partial \theta}$ ， $\frac{\partial^2 \ln p}{\partial \theta^2}$  和  $\frac{\partial^3 \ln p}{\partial \theta^3}$  都存在。
2. 对于任意  $\theta \in \Theta$ ，成立

$$\left| \frac{\partial p}{\partial \theta} \right| < F_1(x), \quad \left| \frac{\partial^2 p}{\partial \theta^2} \right| < F_2(x), \quad \left| \frac{\partial^3 p}{\partial \theta^3} \right| < F_3(x)$$



其中函数  $F_1(x), F_2(x), F_3(x)$  满足

$$\int_{-\infty}^{\infty} F_1(x) dx < \infty, \quad \int_{-\infty}^{\infty} F_2(x) dx < \infty$$

$$\sup_{\theta \in \Theta} \int_{-\infty}^{\infty} F_3(x) p(x; \theta) dx < \infty$$

3. 对于任意  $\theta \in \Theta$ , 成立

$$0 < I(\theta) = \int_{-\infty}^{\infty} \left( \frac{\partial \ln p}{\partial \theta} \right)^2 p(x; \theta) dx < \infty$$

那么对于来自该总体的样本  $x_1, \dots, x_n$ , 存在未知参数  $\theta$  的最大似然估计  $\hat{\theta}_n = \hat{\theta}_n(x_1, \dots, x_n)$ , 且  $\hat{\theta}_n$  具有相合性和渐近正态性, 同时

$$\hat{\theta}_n \sim AN \left( \theta, \frac{1}{nI(\theta)} \right)$$



## 2.4 最小方差无偏估计

### 2.4.1 均方误差

#### 定义 2.4.1 (均方误差)

对于  $\theta$  的点估计  $\hat{\theta} = \hat{\theta}(x_1, \dots, x_n)$ , 称  $E(\hat{\theta} - \theta)^2$  为  $\hat{\theta}$  关于  $\theta$  的均方误差, 记为  $MSE(\hat{\theta}, \theta)$ , 或  $M_\theta(\hat{\theta})$ 。



#### 命题 2.4.1 (均方误差的性质)

1. - 对于  $\theta$  的任意估计  $\hat{\theta}$  而言, 成立

$$MSE(\hat{\theta}, \theta) = \text{Var}(\hat{\theta}) + (E(\hat{\theta}) - \theta)^2$$

2. 对于  $\theta$  的无偏估计  $\hat{\theta}$  而言, 成立

$$MSE(\hat{\theta}, \theta) = \text{Var}(\hat{\theta})$$



#### 定义 2.4.2 (一致最小均方误差估计)

对于样本  $x_1, \dots, x_n$ , 以及待估参数  $\theta$  的一个估计类, 称  $\hat{\theta}(x_1, \dots, x_n)$  是该估计类中  $\theta$  中的一致最小均方误差估计, 如果对于该估计类中另外任意一个  $\theta$  的估计  $\tilde{\theta}$ , 在参数空间  $\Theta$  上均成立

$$MSE_\theta(\hat{\theta}) \leq MSE_\theta(\tilde{\theta})$$



### 2.4.2 一致最小方差无偏估计

#### 定义 2.4.3 (一致最小方差无偏估计 UMVUE)

对于  $\theta$  的一个无偏估计  $\hat{\theta}$ , 称  $\hat{\theta}$  是  $\theta$  的一致最小方差无偏估计, 如果对于  $\theta$  的任意无偏估计  $\tilde{\theta}$ , 在参数空间  $\Theta$  上均成立

$$\text{Var}_\theta(\hat{\theta}) \leq \text{Var}_\theta(\tilde{\theta})$$



**定理 2.4.1**

对于来自某总体的样本  $X = (x_1, \dots, x_n)$ , 如果  $\hat{\theta} = \hat{\theta}(X)$  是  $\theta$  的一个无偏估计,  $\text{Var}(\hat{\theta}) < \infty$ , 那么  $\hat{\theta}$  是  $\theta$  的一致最小方差无偏估计的充分必要条件是, 对于任意满足  $E(\varphi(X)) = 0$  和  $\text{Var}(\varphi(X)) < \infty$  的  $\varphi(X)$ , 以及任意  $\theta \in \Theta$ , 成立

$$\text{Cov}_{\theta}(\hat{\theta}, \varphi) = 0$$

即

$$E(\hat{\theta}\varphi) = 0$$

**2.4.3 充分性原则****定理 2.4.2**

对于来自总体概率密度函数为  $p(x; \theta)$  的样本  $x_1, \dots, x_n$ , 如果  $T = T(x_1, \dots, x_n)$  是  $\theta$  的充分统计量, 那么对于  $\theta$  的任意无偏估计  $\hat{\theta} = \hat{\theta}(x_1, \dots, x_n)$ , 成立  $\tilde{\theta} = E(\hat{\theta}|T)$  是  $\theta$  的无偏估计, 且

$$\text{Var}(\tilde{\theta}) \leq \text{Var}(\hat{\theta})$$

**2.4.4 Cramer-Rao 不等式****定义 2.4.4 (Fisher 信息量)**

对于满足如下条件的概率函数为  $p(x; \theta), \theta \in \Theta$  的总体

1. 参数空间  $\Theta$  是直线上的一个开区间。
2. 支撑  $S = \{x : p(x; \theta) > 0\}$  与  $\theta$  无关。
3. 导数  $\frac{\partial}{\partial \theta} p(x; \theta)$  对任意  $\theta \in \Theta$  均存在。
4. 对于  $p(x; \theta)$ , 积分与微分运算可交换次序, 即

$$\int_{-\infty}^{\infty} \frac{\partial}{\partial \theta} p(x; \theta) dx = \frac{\partial}{\partial \theta} \int_{-\infty}^{\infty} p(x; \theta) dx = 0$$

5. 期望  $E\left(\frac{\partial}{\partial \theta} \ln p(x; \theta)\right)^2$  存在。

称

$$I(\theta) = E\left(\frac{\partial}{\partial \theta} \ln p(x; \theta)\right)^2$$

为总体分布的 Fisher 信息量。如果二阶导数  $\frac{\partial^2}{\partial \theta^2} p(x; \theta)$  对于任意  $\theta \in \Theta$  存在, 那么

$$I(\theta) = -E\left(\frac{\partial^2}{\partial \theta^2} \ln p(x; \theta)\right)$$

**定理 2.4.3 (Cramer-Rao 不等式)**

对于满足 Fisher 信息量定义的总体分布  $p(x; \theta)$ ,  $X = (x_1, \dots, x_n)$  是来自该总体的样本, 如果  $T = T(X)$  是  $g(\theta)$  的任意无偏估计, 即

$$g(\theta) = \int_{\mathbb{R}^n} T(X) L(X; \theta) dX$$

其中  $L(x_1, \dots, x_n; \theta)$  为  $X = (x_1, \dots, x_n)$  的总体概率密度函数

$$L(X; \theta) = \prod_{k=1}^n p(x_k; \theta)$$

并且  $g'(\theta) = \frac{\partial g(\theta)}{\partial \theta}$  存在, 同时对于任意  $\theta \in \Theta$ ,  $g(\theta)$  的微商可在积分号下进行, 即

$$g'(\theta) = \int_{\mathbb{R}^n} T(X) \frac{\partial}{\partial \theta} L(X; \theta) dX$$

(对于离散总体, 将上述积分号改为求和符号) 那么

$$\text{Var}(T) \geq \frac{(g'(\theta))^2}{nI(\theta)}$$

其中  $I(\theta)$  为总体分布的 Fisher 信息量,  $\frac{(g'(\theta))^2}{nI(\theta)}$  称为  $g(\theta)$  的无偏估计的方差的 C-R 下界。当等号成立时, 称  $T = T(X)$  为  $g(\theta)$  的有效估计, 有效估计一定是一致最小方差无偏估计。



## 2.5 Bayes 估计

## 2.6 统计判断的基础

**Bayes 学派基本观点:** 任意未知量都可看作随机变量, 可用一个概率分布去描述, 这个分布称为先验分布。

### 2.6.1 Bayes 公式的密度函数形式

#### 定理 2.6.1 (Bayes 公式的密度函数形式)

1.  $p(x | \theta)$  表示随机变量  $\theta$  取给定值时总体的条件概率函数。
2. 根据参数  $\theta$  的先验信息确定先验分布  $\pi(\theta)$ 。
3. 样本  $X = (x_1, \dots, x_n)$  的产生分两步进行, 首先设想从先验分布  $\pi(\theta)$  产生一个个体  $\theta_0$ , 其次从  $p(X | \theta)$  中产生一组样本, 此时样本  $X$  的联合条件概率函数为

$$P(X | \theta_0) = \prod_{k=1}^n p(x_k | \theta)$$

4. 由于  $\theta_0$  是设想出来的, 因此需要考虑  $\pi(\theta)$ , 那么样本  $X$  和参数  $\theta$  的联合分布为

$$h(X, \theta) = P(X | \theta) \pi(\theta)$$

5. 将  $h(X, \theta)$  分解为

$$h(X, \theta) = \pi(\theta | X) m(X)$$

其中  $m(X)$  为  $X$  的边缘概率函数

$$m(X) = \int_{\Theta} h(X, \theta) d\theta = \int_{\Theta} P(X | \theta) \pi(\theta) d\theta$$

进而  $\theta$  的后验分布为

$$\pi(\theta | X) = \frac{h(X, \theta)}{m(X)} = \frac{P(X | \theta) \pi(\theta)}{\int_{\Theta} P(X | \theta) \pi(\theta) d\theta}$$



### 2.6.2 Bayes 估计

#### 定义 2.6.1 (Bayes 估计)

由后验分布  $\pi(\theta | X)$  估计  $\theta$  有三种常用的方法:

1. 最大后验估计: 后验分布的密度函数的最大值点。
2. 后验中位数估计: 后验分布的中位数。

3. 后验期望估计：后验分布的均值。  
称后验期望估计为 Bayes 估计，记为  $\hat{\theta}$ 。



### 2.6.3 共轭先验分布

#### 定义 2.6.2 (共轭先验分布)

对于总体分布  $p(x; \theta)$  中的参数  $\theta$ ， $\pi(\theta)$  是其先验分布，如果对于任意来自该总体的样本观测值得到的后验分布  $\pi(\theta | X)$  与  $\pi(\theta)$  属于同一个分布族，那么称该分布族为  $\theta$  的共轭先验分布（族）。



## 2.7 区间估计

### 2.7.1 区间估计的概念

#### 定义 2.7.1 (置信区间)

对于总体的参数  $\theta \in \Theta$ ，以及来自该总体的样本  $x_1, \dots, x_n$ ，给定  $\alpha \in (0, 1)$ ，如果两个统计量  $\hat{\theta}_L = \hat{\theta}_L(x_1, \dots, x_n)$  和  $\hat{\theta}_U = \hat{\theta}_U(x_1, \dots, x_n)$ ，满足对于任意  $\theta \in \Theta$ ，成立

$$P_{\theta}(\hat{\theta}_L \leq \theta \leq \hat{\theta}_U) \geq 1 - \alpha$$

那么称随机区间  $[\hat{\theta}_L, \hat{\theta}_U]$  为  $\theta$  的置信水平为  $1 - \alpha$  的置信区间，或简称  $[\hat{\theta}_L, \hat{\theta}_U]$  为  $\theta$  的  $1 - \alpha$  置信区间。其中  $\hat{\theta}_L$  和  $\hat{\theta}_U$  分别称为  $\theta$  的（双侧）置信下限和置信上限。



#### 定义 2.7.2 (同等置信区间)

对于总体的参数  $\theta \in \Theta$ ，以及来自该总体的样本  $x_1, \dots, x_n$ ，给定  $\alpha \in (0, 1)$ ，如果两个统计量  $\hat{\theta}_L = \hat{\theta}_L(x_1, \dots, x_n)$  和  $\hat{\theta}_U = \hat{\theta}_U(x_1, \dots, x_n)$ ，满足对于任意  $\theta \in \Theta$ ，成立

$$P_{\theta}(\hat{\theta}_L \leq \theta \leq \hat{\theta}_U) = 1 - \alpha$$

那么称随机区间  $[\hat{\theta}_L, \hat{\theta}_U]$  为  $\theta$  的置信水平为  $1 - \alpha$  的同等置信区间。



#### 定义 2.7.3 (单侧置信下限)

对于总体的参数  $\theta \in \Theta$ ，以及来自该总体的样本  $x_1, \dots, x_n$ ，给定  $\alpha \in (0, 1)$ ，如果统计量  $\hat{\theta}_L = \hat{\theta}_L(x_1, \dots, x_n)$  满足对于任意  $\theta \in \Theta$ ，成立

$$P_{\theta}(\hat{\theta}_L \leq \theta) \geq 1 - \alpha$$

那么称  $\hat{\theta}_L$  为  $\theta$  的（单侧）置信下限。



#### 定义 2.7.4 (单侧置信上限)

对于总体的参数  $\theta \in \Theta$ ，以及来自该总体的样本  $x_1, \dots, x_n$ ，给定  $\alpha \in (0, 1)$ ，如果统计量  $\hat{\theta}_U = \hat{\theta}_U(x_1, \dots, x_n)$  满足对于任意  $\theta \in \Theta$ ，成立

$$P_{\theta}(\hat{\theta}_U \geq \theta) \geq 1 - \alpha$$

那么称  $\hat{\theta}_U$  为  $\theta$  的（单侧）置信上限。



## 2.7.2 枢轴量法

## 定理 2.7.1 (构造枢轴量的方法)

1. 构造函数  $G = G(x_1, \dots, x_n, \theta)$ , 使得  $G$  的分布不依赖于  $\theta$ , 此函数  $G$  称为枢轴量。
2. 选择常数  $a, b$ , 使得对于给定  $\alpha \in (0, 1)$ , 使得成立

$$P(a \leq G \leq b) = 1 - \alpha$$

3. 将不等式  $a \leq G \leq b$  等价变形为  $\hat{\theta}_L \leq \theta \leq \hat{\theta}_U$ , 即

$$P_\theta(\hat{\theta}_L \leq \theta \leq \hat{\theta}_U) = 1 - \alpha$$

那么区间  $[\hat{\theta}_L, \hat{\theta}_U]$  为  $\theta$  的置信水平为  $1 - \alpha$  同等置信区间。

4. 其中常数  $a, b$  的选择应该使得区间  $[\hat{\theta}_L, \hat{\theta}_U]$  的长度最短, 否则使得成立

$$P(G < a) = P(G > b) = \frac{\alpha}{2}$$

称这样得到的置信区间  $[\hat{\theta}_L, \hat{\theta}_U]$  为等尾置信区间。



## 2.7.3 单个正态总体参数的置信区间

表 2.1: 单个正态总体参数的置信区间

目标	条件	枢轴量	分布	置信区间
$\mu$	$\sigma$ 已知	$\frac{\sqrt{n}(\bar{x}-\mu)}{\sigma}$	$N(0, 1)$	$\left[ \bar{x} - \frac{\sigma}{\sqrt{n}} n_{1-\frac{\alpha}{2}}, \quad \bar{x} + \frac{\sigma}{\sqrt{n}} n_{1-\frac{\alpha}{2}} \right]$
$\mu$	$\sigma$ 未知	$\frac{\sqrt{n}(\bar{x}-\mu)}{s}$	$T(n-1)$	$\left[ \bar{x} - \frac{s}{\sqrt{n}} t_{1-\frac{\alpha}{2}}, \quad \bar{x} + \frac{s}{\sqrt{n}} t_{1-\frac{\alpha}{2}} \right]$
$\sigma^2$	$\mu$ 未知	$\frac{(n-1)s^2}{\sigma^2}$	$\chi^2(n-1)$	$\left[ \frac{(n-1)s^2}{\chi^2_{1-\frac{\alpha}{2}}}, \quad \frac{(n-1)s^2}{\chi^2_{\frac{\alpha}{2}}} \right]$

## 2.7.4 大样本置信区间

在有些场合, 寻找枢轴量及其分布比较困难。在样本量充分大时, 可用渐进分布来构造近似的置信区间。以下为二点分布关于比例  $p$  的置信区间。

对于来自二点分布  $b(1, p)$  的样本  $x_1, \dots, x_n$ , 由中心极限定理

$$N = \frac{\sqrt{n}(\bar{x} - p)}{\sqrt{p(1-p)}} \sim N(0, 1)$$

因此置信水平为  $1 - \alpha$  的同等置信区间为

$$\left[ \frac{1}{1 + \frac{n_{1-\frac{\alpha}{2}}^2}{n}} \left( \bar{x} + \frac{n_{1-\frac{\alpha}{2}}^2}{2n} - \sqrt{\frac{\bar{x}(1-\bar{x})}{n} n_{1-\frac{\alpha}{2}}^2 + \left( \frac{n_{1-\frac{\alpha}{2}}^2}{2n} \right)^2} \right), \quad \frac{1}{1 + \frac{n_{1-\frac{\alpha}{2}}^2}{n}} \left( \bar{x} + \frac{n_{1-\frac{\alpha}{2}}^2}{2n} + \sqrt{\frac{\bar{x}(1-\bar{x})}{n} n_{1-\frac{\alpha}{2}}^2 + \left( \frac{n_{1-\frac{\alpha}{2}}^2}{2n} \right)^2} \right) \right]$$

其中  $n_{1-\frac{\alpha}{2}}$  为  $N(0, 1)$  的  $1 - \frac{\alpha}{2}$  分位数。由于  $n$  充分大, 略去  $\frac{n_{1-\frac{\alpha}{2}}^2}{n}$  项, 因此置信水平为  $1 - \alpha$  的同等置信区间近似为

$$\left[ \bar{x} - n_{1-\frac{\alpha}{2}} \sqrt{\frac{\bar{x}(1-\bar{x})}{n}}, \quad \bar{x} + n_{1-\frac{\alpha}{2}} \sqrt{\frac{\bar{x}(1-\bar{x})}{n}} \right]$$

## 2.7.5 样本量的确定

## 定义 2.7.5 (保证概率)

称置信水平  $1 - \alpha$  为保证概率。



## 定义 2.7.6 (绝对误差)

称置信区间的半径 (即长度的一半) 为绝对误差。



## 2.7.6 两个正态总体下的置信区间

$x_1, \dots, x_m$  是取自  $N(\mu_1, \sigma_1^2)$  的样本,  $y_1, \dots, y_n$  是取自  $N(\mu_2, \sigma_2^2)$  的样本, 两个样本相互独立, 记  $\bar{x}$  和  $\bar{y}$  分别记为两者的样本均值,  $s_x^2$  和  $s_y^2$  分别记为两者的样本方差。

表 2.2: 两个正态总体下的置信区间

目标	条件	枢轴量	分布	置信区间
$\mu_1 - \mu_2$	$\sigma_1^2$ 和 $\sigma_2^2$ 已知	$\frac{(\bar{x} - \bar{y}) - (\mu_1 - \mu_2)}{\sqrt{\frac{\sigma_1^2}{m} + \frac{\sigma_2^2}{n}}}$	$N(0, 1)$	$\left[ (\bar{x} - \bar{y}) - n_{1-\frac{\alpha}{2}} \sqrt{\frac{\sigma_1^2}{m} + \frac{\sigma_2^2}{n}}, (\bar{x} - \bar{y}) + n_{1-\frac{\alpha}{2}} \sqrt{\frac{\sigma_1^2}{m} + \frac{\sigma_2^2}{n}} \right]$
$\mu_1 - \mu_2$	$\sigma_1^2 = \sigma_2^2$ 未知	$\frac{(\bar{x} - \bar{y}) - (\mu_1 - \mu_2)}{s_w \sqrt{\frac{1}{m} + \frac{1}{n}}}$	$T(m + n - 2)$	$\left[ (\bar{x} - \bar{y}) - s_w t_{1-\frac{\alpha}{2}} \sqrt{\frac{1}{m} + \frac{1}{n}}, (\bar{x} - \bar{y}) + s_w t_{1-\frac{\alpha}{2}} \sqrt{\frac{1}{m} + \frac{1}{n}} \right]$
$\mu_1 - \mu_2$	$\frac{\sigma_2^2}{\sigma_1^2} = c$ 已知	$\frac{(\bar{x} - \bar{y}) - (\mu_1 - \mu_2)}{s_{wc} \sqrt{\frac{1}{m} + \frac{c}{n}}}$	$T(m + n - 2)$	$\left[ (\bar{x} - \bar{y}) - s_{wc} t_{1-\frac{\alpha}{2}} \sqrt{\frac{1}{m} + \frac{c}{n}}, (\bar{x} - \bar{y}) + s_{wc} t_{1-\frac{\alpha}{2}} \sqrt{\frac{1}{m} + \frac{c}{n}} \right]$
$\mu_1 - \mu_2$	$n_1$ 和 $n_2$ 充分大	$\frac{(\bar{x} - \bar{y}) - (\mu_1 - \mu_2)}{\sqrt{\frac{s_x^2}{m} + \frac{s_y^2}{n}}}$	$N(0, 1)$	$\left[ (\bar{x} - \bar{y}) - n_{1-\frac{\alpha}{2}} \sqrt{\frac{s_x^2}{m} + \frac{s_y^2}{n}}, (\bar{x} - \bar{y}) + n_{1-\frac{\alpha}{2}} \sqrt{\frac{s_x^2}{m} + \frac{s_y^2}{n}} \right]$
$\mu_1 - \mu_2$	一般情况	$\frac{(\bar{x} - \bar{y}) - (\mu_1 - \mu_2)}{s_0}$	$T(l)$	$\left[ (\bar{x} - \bar{y}) - s_0 t_{1-\frac{\alpha}{2}}, (\bar{x} - \bar{y}) + s_0 t_{1-\frac{\alpha}{2}} \right]$
$\sigma_2^2 / \sigma_1^2$	一般情况	$\frac{s_x^2 / \sigma_1^2}{s_y^2 / \sigma_2^2}$	$F(m - 1, n - 1)$	$\left[ \frac{s_x^2}{s_y^2} \frac{1}{f_{1-\frac{\alpha}{2}}}, \frac{s_x^2}{s_y^2} \frac{1}{f_{\frac{\alpha}{2}}} \right]$

其中

$$s_w^2 = \frac{(m-1)s_x^2 + (n-1)s_y^2}{m+n-2}$$

$$s_{wc}^2 = \frac{(m-1)s_x^2 + (n-1)\frac{s_y^2}{c}}{m+n-2}$$

$$s_0^2 = \frac{s_x^2}{m} + \frac{s_y^2}{n}$$

$$l = \left[ \frac{s_0^4}{\frac{s_x^4}{m^2(m-1)} + \frac{s_y^4}{n^2(n-1)}} \right]$$

## 第三章 假设检验

### 3.1 假设检验的基本思想与概念

#### 3.1.1 假设检验问题

**基本思想：**如果试验结果与假设  $H$  发生矛盾，那么拒绝原假设  $H$ ，否则接受原假设  $H$ 。

**假设检验问题：**

1. **假设：**两个非空不交参数集合。
2. **检验：**通过样本对一个假设作出“对”或“不对”的具体判断规则。
3. **参数假设检验问题：**假设可用一个参数的集合表示的检验问题。

#### 3.1.2 假设检验的基本步骤

##### 一、建立假设

对于来自参数分布族  $\{F(x, \theta) : \theta \in \Theta\}$  的样本  $x_1, \dots, x_n$ ，其中  $\Theta$  为参数空间，如果非空集合  $\Theta_0 \subset \Theta$ ，那么命题  $H_0 : \theta \in \Theta_0$  称为原假设或零假设，命题  $H_a : \theta \in \Theta \setminus \Theta_0$  称为对立假设或备择假设，那么  $H_0$  对  $H_a$  的假设检验问题记为

$$H_0 : \theta \in \Theta_0 \quad \text{vs} \quad H_a : \theta \in \Theta \setminus \Theta_0$$

如果  $\Theta_0$  仅含有一个点，那么称  $H_0$  为简单原假设，否则称为复杂原假设或复合原假设。当  $H_0$  为简单假设时，其形式可写为  $H_0 : \theta = \theta_0$ ，此时备择假设通常有如下三种可能：

$$H_1 : \theta \neq \theta_0, \quad H_2 : \theta < \theta_0, \quad H_3 : \theta > \theta_0$$

称  $H_0$  vs  $H_1$  为双侧假设或双边假设， $H_0$  vs  $H_2$  以及  $H_0$  vs  $H_3$  为单侧假设或单边假设。

在假设检验中，通常将不宜轻易否定的假设作为原假设。

##### 二、选择检验统计量，给出拒绝域形式

当有了具体的样本后，将样本空间划分为两个互不相交的部分  $W$  和  $\bar{W}$ ，当样本属于  $W$  时，拒绝  $H_0$ ，否则接受  $H_0$ 。称  $W$  为该检验的拒绝域， $\bar{W}$  为该检验的接受域。事实上，在拒绝域和接受域外，还有保留域，但通常将保留域合并于接受域内。

选择分布已知的检验统计量  $T(X)$ ，确定拒绝域  $W$  的形式。

##### 三、选择显著性水平

当  $\theta \in \Theta_0$  时，样本由于随机性却落入了拒绝域  $W$ ，于是采取了拒绝  $H_0$  的错误决策，称之为第一类错误或拒真错误，记第一类错误概率为

$$\alpha(\theta) = P\{X \in W \mid H_0\}, \quad \theta \in \Theta_0$$

当  $\theta \in \Theta \setminus \Theta_0$  时，样本由于随机性却落入了接受域  $\bar{W}$ ，于是采取了接受  $H_0$  的错误决策，称之为第二类错误或取伪错误，记第二类错误概率为

$$\beta(\theta) = P\{X \in \bar{W} \mid H_a\}, \quad \theta \in \Theta \setminus \Theta_0$$

#### 定义 3.1.1 (势函数)

对于检验问题

$$H_0 : \theta \in \Theta_0 \quad \text{vs} \quad H_a : \theta \in \Theta \setminus \Theta_0$$

其拒绝域为  $W$ ，那么定义势函数为

$$\rho(\theta) = P_\theta(X \in W), \quad \theta \in \Theta$$

即

$$\rho(\theta) = \begin{cases} \alpha(\theta), & \theta \in \Theta_0 \\ 1 - \beta(\theta), & \theta \in \Theta \setminus \Theta_0 \end{cases}$$



### 定义 3.1.2 (显著性检验)

对于检验问题

$$H_0 : \theta \in \Theta_0 \quad \text{vs} \quad H_a : \theta \in \Theta \setminus \Theta_0$$

其势函数为  $\rho(\theta)$ ，如果一个检验满足对于任意  $\theta \in \Theta_0$ ，成立

$$\rho(\theta) \leq \alpha$$

那么称该检验为显著性水平为  $\alpha$  的显著性检验，简称水平为  $\alpha$  的检验。



### 四、给出拒绝域

依据显著性水平  $\alpha$  以及拒绝域  $W$  的形式，确定具体的拒绝域。

### 五、做出判断

由拒绝域  $W$  唯一相互确定的判断准则为

1. 如果  $(x_1, \dots, x_n) \in W$ ，那么拒绝  $H_0$ 。
2. 如果  $(x_1, \dots, x_n) \in \bar{W}$ ，那么接受  $H_0$ 。

## 3.1.3 检验的 $p$ 值

### 定义 3.1.3 (检验的 $p$ 值)

在假设检验问题中，利用样本观测值能够作出拒绝原假设的最小显著性水平称为检验的  $p$  值。

1. 如果  $p \leq \alpha$ ，那么在显著性水平  $\alpha$  下拒绝  $H_0$ 。
2. 如果  $p > \alpha$ ，那么在显著性水平  $\alpha$  下接受  $H_0$ 。





## 3.2 正态总体参数假设检验

### 3.2.1 单个正态总体均值的检验

表 3.1: 单个正态总体均值的检验

检验	条件	$H_0$	$H_a$	统计检验量	分布	拒绝域	$p$ 值
		$\mu \leq \mu_0$	$\mu > \mu_0$			$\{u \geq u_{1-\alpha}\}$	$1 - \Phi(u_0)$
$u$ 检验	$\sigma$ 已知	$\mu \geq \mu_0$	$\mu < \mu_0$	$u = \frac{\sqrt{n}(\bar{x} - \mu_0)}{\sigma}$	$N(0, 1)$	$\{u \leq u_\alpha\}$	$\Phi(u_0)$
		$\mu = \mu_0$	$\mu \neq \mu_0$			$\{ u  \geq u_{1-\frac{\alpha}{2}}\}$	$2(1 - \Phi( u_0 ))$
		$\mu \leq \mu_0$	$\mu > \mu_0$			$\{t \geq t_{1-\alpha}\}$	$P(T \geq t_0)$
$t$ 检验	$\sigma$ 未知	$\mu \geq \mu_0$	$\mu < \mu_0$	$t = \frac{\sqrt{n}(\bar{x} - \mu_0)}{s}$	$T(n - 1)$	$\{t \leq t_\alpha\}$	$P(T \leq t_0)$
		$\mu = \mu_0$	$\mu \neq \mu_0$			$\{ t  \geq t_{1-\frac{\alpha}{2}}\}$	$P( T  \geq  t_0 )$

## 3.2.2 两个正态总体均值差的检验

表 3.2: 两个正态总体均值差的检验

检验	条件	$H_0$	$H_a$	检验统计量	分布	拒绝域	$p$ 值
		$\mu_1 \leq \mu_2$	$\mu_1 > \mu_2$			$\{u \geq u_{1-\alpha}\}$	$1 - \Phi(u_0)$
$u$ 检验	$\sigma_1, \sigma_2$ 已知	$\mu_1 \geq \mu_2$	$\mu_1 < \mu_2$	$u = \frac{\bar{x} - \bar{y}}{\sqrt{\frac{\sigma_1^2}{m} + \frac{\sigma_2^2}{n}}}$	$N(0, 1)$	$\{u \leq u_\alpha\}$	$\Phi(u_0)$
		$\mu_1 = \mu_2$	$\mu_1 \neq \mu_2$			$\{ u  \geq u_{1-\frac{\alpha}{2}}\}$	$2(1 - \Phi( u_0 ))$
		$\mu_1 \leq \mu_2$	$\mu_1 > \mu_2$			$\{t \geq t_{1-\alpha}\}$	$P(T \geq t_0)$
$t$ 检验	$\sigma_1 = \sigma_2$ 未知	$\mu_1 \geq \mu_2$	$\mu_1 < \mu_2$	$t = \frac{\bar{x} - \bar{y}}{s_w \sqrt{\frac{1}{m} + \frac{1}{n}}}$	$T(m + n - 2)$	$\{t \leq t_\alpha\}$	$P(T \leq t_0)$
		$\mu_1 = \mu_2$	$\mu_1 \neq \mu_2$			$\{ t  \geq t_{1-\frac{\alpha}{2}}\}$	$P( T  \geq  t_0 )$
		$\mu_1 \leq \mu_2$	$\mu_1 > \mu_2$			$\{u \geq u_{1-\alpha}\}$	$1 - \Phi(u_0)$
$u$ 检验	$m, n$ 充分大	$\mu_1 \geq \mu_2$	$\mu_1 < \mu_2$	$u = \frac{\bar{x} - \bar{y}}{\sqrt{\frac{s_x^2}{m} + \frac{s_y^2}{n}}}$	$N(0, 1)$	$\{u \leq u_\alpha\}$	$\Phi(u_0)$
		$\mu_1 = \mu_2$	$\mu_1 \neq \mu_2$			$\{ u  \geq u_{1-\frac{\alpha}{2}}\}$	$2(1 - \Phi( u_0 ))$
		$\mu_1 \leq \mu_2$	$\mu_1 > \mu_2$			$\{t \geq t_{1-\alpha}\}$	$P(T \geq t_0)$
$t$ 检验	一般情况	$\mu_1 \geq \mu_2$	$\mu_1 < \mu_2$	$t = \frac{\bar{x} - \bar{y}}{\sqrt{\frac{s_x^2}{m} + \frac{s_y^2}{n}}}$	$T(l)$	$\{t \leq t_\alpha\}$	$P(T \leq t_0)$
		$\mu_1 = \mu_2$	$\mu_1 \neq \mu_2$			$\{ t  \geq t_{1-\frac{\alpha}{2}}\}$	$P( T  \geq  t_0 )$

其中

$$s_w^2 = \frac{(m-1)s_x^2 + (n-1)s_y^2}{m+n-2}$$

$$l = \left[ \frac{\left( \frac{s_x^2}{m} + \frac{s_y^2}{n} \right)^2}{\frac{s_x^4}{m^2(m-1)} + \frac{s_y^4}{n^2(n-1)}} \right]$$

## 3.2.3 成对数据检验

表 3.3: 成对数据检验

$H_0$	$H_a$	统计检验量	分布	拒绝域	$p$ 值
$\mu \leq 0$	$\mu > 0$			$\{t \geq t_{1-\alpha}\}$	$P(T \geq t_0)$
$\mu \geq 0$	$\mu < 0$	$t = \frac{\sqrt{n}\bar{d}}{s_d}$	$T(n-1)$	$\{t \leq t_\alpha\}$	$P(T \leq t_0)$
$\mu = 0$	$\mu \neq 0$			$\{ t  \geq t_{1-\frac{\alpha}{2}}\}$	$P( T  \geq  t_0 )$

## 3.2.4 正态总体方差的检验

表 3.4: 正态总体方差的检验

检验	条件	$H_0$	$H_a$	统计检验量	分布	拒绝域	$p$ 值
		$\sigma^2 \leq \sigma_0^2$	$\sigma^2 > \sigma_0^2$			$\{\chi^2 \geq \chi_{1-\alpha}^2\}$	$P(\chi^2 \geq \chi_0^2)$
$\chi^2$ 检验	一个	$\sigma^2 \geq \sigma_0^2$	$\sigma^2 < \sigma_0^2$	$\chi^2 = \frac{(n-1)s^2}{\sigma_0^2}$	$\chi^2(n-1)$	$\{\chi^2 \leq \chi_\alpha^2\}$	$P(\chi^2 \leq \chi_0^2)$
		$\sigma^2 = \sigma_0^2$	$\sigma^2 \neq \sigma_0^2$			$\{\chi^2 \leq \chi_{\frac{\alpha}{2}}^2\} \cup \{\chi^2 \geq \chi_{1-\frac{\alpha}{2}}^2\}$	$2 \min\{P(\chi^2 \leq \chi_0^2), P(\chi^2 \geq \chi_0^2)\}$
		$\sigma_1^2 \leq \sigma_2^2$	$\sigma_1^2 > \sigma_2^2$			$\{F \geq F_{1-\alpha}\}$	$P(F \geq F_0)$
$F$ 检验	两个	$\sigma_1^2 \geq \sigma_2^2$	$\sigma_1^2 < \sigma_2^2$	$F = \frac{s_x^2}{s_y^2}$	$F(m-1, n-1)$	$\{F \leq F_\alpha\}$	$P(F \leq F_0)$
		$\sigma_1^2 = \sigma_2^2$	$\sigma_1^2 \neq \sigma_2^2$			$\{F \leq F_{\frac{\alpha}{2}}\} \cup \{F \geq F_{1-\frac{\alpha}{2}}\}$	$2 \min\{P(F \leq F_0), P(F \geq F_0)\}$

### 3.3 其他分布参数的假设检验

表 3.5: 其他分布参数的假设检验

检验	条件	$H_0$	$H_a$	统计检验量	分布	拒绝域	$p$ 值
		$\lambda \leq \lambda_0$	$\lambda > \lambda_0$			$\{\chi^2 \geq \chi_{1-\alpha}^2\}$	$P(\chi^2 \geq \chi_0^2)$
$\chi^2$ 分布	$\text{Exp}(\frac{1}{\lambda})$	$\lambda \geq \lambda_0$	$\lambda < \lambda_0$	$\frac{2n\bar{x}}{\lambda_0}$	$\chi^2(2n)$	$\{\chi^2 \leq \chi_\alpha^2\}$	$P(\chi^2 \leq \chi_0^2)$
		$\lambda = \lambda_0$	$\lambda \neq \lambda_0$			$\{\chi^2 \leq \chi_{\frac{\alpha}{2}}^2\} \cup \{\chi^2 \geq \chi_{1-\frac{\alpha}{2}}^2\}$	$2 \min\{P(\chi^2 \leq \chi_0^2), P(\chi^2 \geq \chi_0^2)\}$
		$p \leq p_0$	$p > p_0$				$P(x \geq x_0)$
$B$ 检验	$B(1, p)$	$p \geq p_0$	$p < p_0$	$x$	$B(n, p)$		$P(x \leq x_0)$
		$p = p_0$	$p \neq p_0$				$2 \min\{P(x \leq x_0), P(x \geq x_0)\}$
		$\theta \leq \theta_0$	$\theta > \theta_0$			$\{u \geq u_{1-\alpha}\}$	$1 - \Phi(u_0)$
$u$ 检验	大样本分布 $F(x; \theta)$	$\theta \geq \theta_0$	$\theta < \theta_0$	$\frac{\sqrt{n}(\bar{x} - \theta_0)}{\sqrt{\sigma^2(\hat{\theta})}}$	$N(0, 1)$	$\{u \leq u_\alpha\}$	$\Phi(u_0)$
		$\theta = \theta_0$	$\theta \neq \theta_0$			$\{ u  \geq u_{1-\frac{\alpha}{2}}\}$	$2(1 - \Phi( u_0 ))$

其中分布  $F(x; \theta)$  的均值为  $\theta$ , 方差为  $\sigma(\theta)$ ,  $\hat{\theta}$  为  $\theta$  的最大似然估计。

### 3.4 似然比检验与分布拟合检验

#### 3.4.1 似然比检验的思想

##### 定义 3.4.1 (似然比)

对于来自密度函数为  $p(x; \theta), \theta \in \Theta$  的总体的样本  $x_1, \dots, x_n$ , 对于如下检验问题

$$H_0 : \theta \in \Theta_0 \quad \text{vs} \quad H_a : \theta \in \Theta \setminus \Theta_0$$

定义改假设检验问题的似然比统计量为

$$\Lambda(x_1, \dots, x_n) = \frac{\sup_{\theta \in \Theta} p(x_1, \dots, x_n; \theta)}{\sup_{\theta \in \Theta_0} p(x_1, \dots, x_n; \theta)}$$

即

$$\Lambda(x_1, \dots, x_n) = \frac{p(x_1, \dots, x_n; \hat{\theta})}{p(x_1, \dots, x_n; \hat{\theta}_0)}$$

其中  $\hat{\theta}$  和  $\hat{\theta}_0$  分别为参数空间  $\Theta$  和  $\Theta_0$  上的最大似然估计。



##### 定义 3.4.2 (似然比检验 LRT)

对于来自密度函数为  $p(x; \theta), \theta \in \Theta$  的总体的样本  $x_1, \dots, x_n$ , 对于如下检验问题

$$H_0 : \theta \in \Theta_0 \quad \text{vs} \quad H_a : \theta \in \Theta \setminus \Theta_0$$

其似然比统计量

$$\Lambda(x_1, \dots, x_n) = \frac{p(x_1, \dots, x_n; \hat{\theta})}{p(x_1, \dots, x_n; \hat{\theta}_0)}$$

作为检验问题的检验统计量，且取拒绝域为  $W = \{\Lambda(x_1, \dots, x_n) \geq \lambda_0\}$ ，其中临界值  $\lambda_0$  满足对于任意  $\theta \in \Theta_0$ ，成立

$$P_\theta(\Lambda(x_1, \dots, x_n) \geq \lambda_0) \leq \alpha$$

那么称此检验为显著性水平为  $\alpha$  的似然比检验。



### 3.4.2 分布数据的 $\chi^2$ 拟合优度检验

#### 定理 3.4.1

总体被分为  $r$  类  $A_1, \dots, A_r$ ，考虑假设检验

$$H_0: A_k \text{ 所占的比率为 } p_k, \quad k = 1, \dots, r$$

其中  $p_k$  已知且  $\sum_{k=1}^r p_k = 1$ 。从该总体抽出  $n$  个样本， $n_k$  为样本中属于  $A_k$  的样本个数，记检验统计量为

$$\chi^2 = \sum_{k=1}^r \frac{(n_k - np_k)^2}{np_k}$$

那么当  $H_0$  成立时，成立

$$\chi^2 \xrightarrow{L} \chi^2(r-1)$$

因此对于显著性水平  $\alpha$ ，拒绝域为  $W = \{\chi^2 \geq \chi_{1-\alpha}^2\}$ ，检验的  $p$  值为  $p = P(\chi^2 \geq \chi_0^2)$ 。

如果  $A_k$  出现的概率含有  $s$  个参数，那么可用最大似然估计方法估计出该  $s$  个参数，然后再算出  $p_k$  的估计值  $\hat{p}_k$ ，于是统计检验量

$$\chi^2 = \sum_{k=1}^r \frac{(n_k - n\hat{p}_k)^2}{n\hat{p}_k} \xrightarrow{L} \chi^2(r-s-1)$$



### 3.4.3 分布的 $\chi^2$ 拟合优度检验

对于来自分布函数为  $F(x)$  的总体的样本  $x_1, \dots, x_n$ ，考虑假设检验问题

$$H_0: F(x) = F_0(x)$$

其中  $F_0(x)$  为可含参的理论分布。

#### 一、总体 $X$ 为离散分布

如果总体  $X$  为至多可数个值  $a_1, a_2, \dots$ ，将其分为  $r$  类  $A_1, \dots, A_r$ ，使得每一个  $A_k$  中的样本个数  $n_k$  不小于 5，记  $P(X \in A_k) = p_k$ ，那么原假设检验转化为

$$H_0: A_k p_k, \quad k = 1, \dots, r$$

#### 二、总体 $X$ 为连续分布

如果总体  $X$  的分布为  $F_0$ ，选取  $-\infty = a_0 < a_1 < \dots < a_{r-1} < a_r = \infty$ ，记  $A_k = (a_{k-1}, a_k]$ ，那么

$$p_k = P(X \in A_k) = F_0(a_k) - F_0(a_{k-1}), \quad k = 1, \dots, r$$

于是原假设转化为

$$H_0: A_k p_k, \quad k = 1, \dots, r$$

### 3.4.4 列联表的独立性检验

将总体分为两个属性  $A$  和  $B$ , 其中  $A$  有  $r$  个类  $A_1, \dots, A_r$ ,  $B$  有  $s$  个类  $B_1, \dots, B_s$ , 从总体中抽取  $n$  个样本, 设其中有  $n_{ij}$  个个体属于  $A_i$  和  $B_j$ , 构造列联表  $\{n_{ij}\}_{r \times s}$ 。

记总体中的个体仅属于  $A_i$  和仅属于  $B_j$  的概率分别为  $p_{i\cdot}$  和  $p_{\cdot j}$ , 总体中的个体同时属于  $A_i$  和  $B_j$  的概率为  $p_{ij}$ , 那么得到二维离散分布表  $\{p_{ij}\}_{r \times s}$ ,  $A$  和  $B$  两属性度量的假设可表述为

$$H_0: p_{ij} = p_{i\cdot} p_{\cdot j}, \quad i = 1, \dots, r; j = 1, \dots, s$$

$H_0$  成立时  $p_{ij}$  的最大似然估计为

$$\hat{p}_{ij} = \frac{1}{n^2} \sum_{k=1}^r n_{kj} \sum_{k=1}^s n_{ik}$$

那么检验统计量为

$$\chi^2 = \sum_{i,j} \frac{(n_{ij} - n\hat{p}_{ij})^2}{n\hat{p}_{ij}} \xrightarrow{L} \chi^2((r-1)(s-1))$$

因此对于显著性水平  $\alpha$ , 拒绝域为  $W = \{\chi^2 \geq \chi_{1-\alpha}^2\}$ , 检验的  $p$  值为  $p = P(\chi^2 \geq \chi_0^2)$ 。

## 3.5 正态性检验

### 3.5.1 正态概率纸

对于给定的样本观测值  $x_1, \dots, x_n$ , 做点

$$\left( x_{(k)}, \frac{k - 0.375}{n + 0.25} \right), \quad k = 1, \dots, n$$

如果诸点在一条直线附近, 那么认为该批数据来自正态总体; 否则不认为该批数据来自正态总体。

### 3.5.2 W 检验

对于来自正态分布总体  $N(\mu, \sigma^2)$  的样本  $x_1, \dots, x_n$ , 其中  $8 \leq n \leq 50$ , 定义  $W$  统计量为

$$W = \frac{\sum_{k=1}^n (w_k - \bar{w})^2 (x_{(k)} - \bar{x})^2}{\sum_{k=1}^n (w_k - \bar{w})^2 \sum_{k=1}^n (x_{(k)} - \bar{x})^2} = \frac{\sum_{k=1}^{[\frac{n}{2}]} w_k^2 (x_{(k)} - x_{(n+1-k)})^2}{\sum_{k=1}^n (x_{(k)} - \bar{x})^2}$$

其中

$$\mathbf{e} = \begin{pmatrix} E\left(\frac{x_{(1)} - \mu}{\sigma}\right) \\ \vdots \\ E\left(\frac{x_{(n)} - \mu}{\sigma}\right) \end{pmatrix}, \quad \mathbf{C} = \begin{pmatrix} \text{Cov}\left(\frac{x_{(1)} - \mu}{\sigma}, \frac{x_{(1)} - \mu}{\sigma}\right) & \cdots & \text{Cov}\left(\frac{x_{(1)} - \mu}{\sigma}, \frac{x_{(n)} - \mu}{\sigma}\right) \\ \vdots & \ddots & \vdots \\ \text{Cov}\left(\frac{x_{(n)} - \mu}{\sigma}, \frac{x_{(1)} - \mu}{\sigma}\right) & \cdots & \text{Cov}\left(\frac{x_{(n)} - \mu}{\sigma}, \frac{x_{(n)} - \mu}{\sigma}\right) \end{pmatrix}$$

$$\begin{pmatrix} w_1 \\ \vdots \\ w_n \end{pmatrix} = \frac{\mathbf{C}^{-1} \mathbf{e}}{\sqrt{\mathbf{e}^T (\mathbf{C}^{-1})^2 \mathbf{e}}}$$

拒绝域为  $\{W \leq W_\alpha\}$ , 其中  $W_\alpha$  为  $\alpha$  分位数。

### 3.5.3 EP 检验

对于来自正态分布总体  $N(\mu, \sigma^2)$  的样本  $x_1, \dots, x_n$ , 其中  $n \geq 8$ , 定义 EP 检验统计量为

$$T_{EP} = 1 + \frac{n}{\sqrt{3}} + \frac{2}{n} \sum_{i=2}^n \sum_{j=1}^{i-1} \exp\left(-\frac{(x_j - x_i)^2}{2^{\frac{n-1}{n}} s^2}\right) - \sqrt{2} \sum_{k=1}^n \exp\left(-\frac{(x_k - \bar{x})^2}{4^{\frac{n-1}{n}} s^2}\right)$$

拒绝域为  $\{T_{EP} \geq T_{EP_{1-\alpha}}\}$ , 其中  $T_{EP_{1-\alpha}}$  为  $1 - \alpha$  分位数。

## 3.6 非参数检验

### 3.6.1 游程检验

对于依时间顺序连续得到的样本观测值  $x_1, \dots, x_n$ , 记样本中位数为  $m_e$ , 对于  $k = 1, \dots, n$ , 记

$$y_k = \begin{cases} 1, & x_k \geq m_e \\ 0, & x_k < m_e \end{cases}$$

$y_1, \dots, y_n$  构成 0-1 序列。

记 0-1 序列中 0 和 1 的个数分别为  $n_1$  和  $n_2$ , 游程总数为  $R$ , 那么  $1 < n_1, n_2 < n$  且  $2 \leq R \leq n$ 。同时  $|n_1 - n_2|$  为 0 或 1。原假设为

$H_0$ : 样本序列符合随机抽取的原则

$R$  的分布如下

$$P(R = 2k) = \frac{2 \binom{n_1-1}{k-1} \binom{n_2-1}{k-1}}{\binom{n_1+n_2}{n_1}}, \quad k = 1, \dots, \left\lfloor \frac{n_1 + n_2}{2} \right\rfloor$$

$$P(R = 2k + 1) = \frac{\binom{n_1-1}{k-1} \binom{n_2-1}{k} + \binom{n_1-1}{k} \binom{n_2-1}{k-1}}{\binom{n_1+n_2}{n_1}}, \quad k = 1, \dots, \left\lfloor \frac{n_1 + n_2 - 1}{2} \right\rfloor$$

拒绝域为  $\{R \leq R_{\frac{\alpha}{2}}\} \cup \{R \geq R_{1-\frac{\alpha}{2}}\}$ , 检验的  $p$  值为  $2 \min\{P(R \leq R_0), P(R \geq R_0)\}$ 。

### 3.6.2 符号检验

表 3.6: 符号检验

$H_0$	$H_a$	拒绝域	检验的 $p$ 值
$x_p \leq x_0$	$x_p > x_0$	$\{S^+ \geq c\}$	$\sum_{k=S_0^+}^n \binom{n}{k} (1-p)^k p^{n-k}$
$x_p \geq x_0$	$x_p < x_0$	$\{S^+ \leq c\}$	$\sum_{k=0}^{S_0^+} \binom{n}{k} (1-p)^k p^{n-k}$
$x_p = x_0$	$x_p \neq x_0$	$\{S^+ \leq c_1\} \cup \{S^+ \geq c_2\}$	$2 \min \left\{ \sum_{k=0}^{S_0^+} \binom{n}{k} (1-p)^k p^{n-k}, \sum_{k=S_0^+}^n \binom{n}{k} (1-p)^k p^{n-k} \right\}$

其中  $S^+$  为  $x_1 - x_0, \dots, x_n - x_0$  中正数的个数, 即

$$S^+ = \sum_{k=1}^n I_{x_k > x_0}$$

## 3.6.3 秩和检验

## 定义 3.6.1 (秩)

对于来自连续分布  $F(x)$  的简单随机样本  $x_1, \dots, x_n$ , 次序样本为  $x_{(1)}, \dots, x_{(n)}$ , 称  $x_k$  秩为  $r_k$ , 如果  $x_k = x_{(r_k)}$ , 记作  $R_k = r_k$ 。



## 定义 3.6.2 (秩统计量)

对于来自连续分布  $F(x)$  的简单随机样本  $x_1, \dots, x_n$ ,  $R_k$  为  $x_k$  的秩, 那么称  $R = (R_1, \dots, R_n)$  为  $x_1, \dots, x_n$  的秩统计量。



对于来自连续分布  $F(x-\theta)$  的简单随机样本  $x_1, \dots, x_n$ , 其中  $\theta$  为总体的中位数, 记  $R_k$  为  $|x_k|$  在  $|x_1|, \dots, |x_n|$  中的秩, 定义符号秩和统计量为

$$W^+ = \sum_{k=1}^n R_k I_{x_k > 0} \sim W^+(n)$$

$H_0$	$H_a$	拒绝域
$\theta \leq 0$	$\theta > 0$	$\{W^+ \leq W_\alpha^+\}$
$\theta \geq 0$	$\theta < 0$	$\{W^+ \geq W_\alpha^+\}$
$\theta = 0$	$\theta \neq 0$	$\{W^+ \leq W_{\frac{\alpha}{2}}^+\} \cup \{W^+ \geq W_{1-\frac{\alpha}{2}}^+\}$

其中  $W_\alpha^+ + W_{1-\alpha}^+ = \frac{1}{2}n(n-1)$ 。

对于来自连续分布  $F(x - \theta_1)$  的简单随机样本  $x_1, \dots, x_m$  和对于来自连续分布  $F(x - \theta_2)$  的简单随机样本  $y_1, \dots, y_n$ , 产生的秩为

$$R = (Q_1, \dots, Q_m, R_1, \dots, R_n)$$

那么秩和统计量为

$$W = \sum_{k=1}^n R_k \sim W(m, n)$$

$H_0$	$H_a$	拒绝域
$\theta_1 \leq \theta_2$	$\theta_1 > \theta_2$	$\{W \leq W_\alpha\}$
$\theta_1 \geq \theta_2$	$\theta_1 < \theta_2$	$\{W \geq W_\alpha\}$
$\theta_1 = \theta_2$	$\theta_1 \neq \theta_2$	$\{W \leq W_{\frac{\alpha}{2}}\} \cup \{W \geq W_{1-\frac{\alpha}{2}}\}$

其中  $W_\alpha + W_{1-\alpha} = n(m+n-1)$ 。



## 第四章 方差分析与回归分析

### 4.1 方差分析

#### 4.1.1 问题的提出

因子:  $A$

水平:  $A_1, \dots, A_r$

结果:  $y_{ij}$ , 其中  $i = 1, \dots, r$

#### 4.1.2 单因子方差分析的统计模型

在单因子试验中, 记因子为  $A$ , 设其由  $r$  个水平, 记为  $A_1, \dots, A_r$ , 在每一个水平下考察的指标可以看成是一个总体, 现有  $r$  个水平, 故有  $r$  个总体, 假定:

1. 每一个总体均为正态分布, 记为  $N(\mu_k, \sigma_k^2)$ , 其中  $k = 1, \dots, r$ 。
2. 各总体的方差相同, 记为  $\sigma_1^2 = \dots = \sigma_r^2 = \sigma^2$ 。
3. 从每一总体中抽取的样本是互相独立的, 即所有的试验结果  $y_{ij}$  都相互独立。

作假设检验:

$$H_0: \mu_1 = \dots = \mu_r \quad \text{vs} \quad H_a: \mu_1, \dots, \mu_r$$

如果  $H_0$  成立, 称因子  $A$  的  $r$  个水平没有显著差异, 简称因此  $A$  不显著。

对  $r$  个总体每个作  $m$  次重复实现, 得到试验结果  $\{y_{ij}\}_{r \times m}$ , 定义随机误差为

$$\varepsilon_{ij} = y_{ij} - \mu_i$$

那么试验结果  $y_{ij}$  的数据结构式为

$$y_{ij} = \mu_i + \varepsilon_{ij}$$

单因子方差分析的统计模型为

$$\begin{cases} y_{ij} = \mu_i + \varepsilon_{ij} \\ \varepsilon_{ij} \text{相互独立} \\ \varepsilon_{ij} \sim N(0, \sigma^2) \end{cases}$$

总均值

$$\mu = \frac{1}{r} \sum_{k=1}^r \mu_k$$

因子  $A$  的第  $k$  个水平的主效应

$$a_i = \mu_i - \mu$$

容易知道

$$\sum_{k=1}^r a_k = 0$$

$$\mu_k = \mu + a_k$$

于是统计模型改写为

$$\begin{cases} y_{ij} = \mu + a_i + \varepsilon_{ij} \\ \sum_{i=1}^r a_i = 0 \\ \varepsilon_{ij} \text{相互独立} \\ \varepsilon_{ij} \sim N(0, \sigma^2) \end{cases}$$

统计假设改写为

$$H_0: a_1 = \cdots = a_r = 0 \quad \text{vs} \quad H_a: a_1, \cdots, a_r \text{不全为0}$$

### 4.1.3 平方和分解

#### 一、实验数据

表 4.1: 符号检验

因子水平	试验数据	和	均值
$A_1$	$y_{11}, \cdots, y_{1m}$	$T_1 = \sum_{j=1}^m y_{1j}$	$\bar{y}_1 = \frac{T_1}{m}$
$\vdots$	$\vdots$	$\vdots$	$\vdots$
$A_r$	$y_{r1}, \cdots, y_{rm}$	$T_r = \sum_{j=1}^m y_{rj}$	$\bar{y}_r = \frac{T_r}{m}$
		$T = \sum_{i=1}^r T_i$	$\bar{y} = \frac{T}{rm} = \frac{T}{n}$

#### 二、组内偏差与组间方差

记

$$y_{ij} - \bar{y} = (y_{ij} - \bar{y}_i) + (\bar{y}_i - \bar{y})$$

$$\bar{\varepsilon}_i = \frac{1}{m} \sum_{j=1}^m \varepsilon_{ij}$$

$$\bar{\varepsilon} = \frac{1}{r} \sum_{j=1}^m \bar{\varepsilon}_i = \frac{1}{n} \sum_{i,j} \varepsilon_{ij}$$

组内偏差为

$$y_{ij} - \bar{y}_i = \varepsilon_{ij} - \bar{\varepsilon}_i$$

组间偏差为

$$\bar{y}_i - \bar{y} = a_i + \bar{\varepsilon}_i - \bar{\varepsilon}$$

#### 三、偏差平方和及其自由度

偏差平方和

$$Q = \sum_{k=1}^n (y_k - \bar{y})^2$$

自由度

$$f_Q = n - 1$$

#### 四、总平方和分解公式

总偏差平方和

$$S_T = \sum_{i,j} (y_{ij} - \bar{y})^2 = \sum_{i,j} y_{ij}^2 - \frac{T^2}{n}, \quad f_T = n - 1$$

组内偏差平方和（因子 A 的偏差平方和）

$$S_A = m \sum_{i=1}^r (\bar{y}_i - \bar{y})^2 = \frac{1}{m} \sum_{i=1}^r T_i^2 - \frac{T^2}{n}, \quad f_A = r - 1$$

组内偏差平方和（误差偏差平方和）

$$S_e = \sum_{i,j} (y_{ij} - \bar{y}_i)^2 = \sum_{i,j} y_{ij}^2 - \frac{1}{m} \sum_{i=1}^r T_i^2, \quad f_e = n - r$$

总平方和分解式

$$S_T = S_A + S_e$$

#### 4.1.4 检验方法

均方

$$MS = \frac{Q}{f_Q}$$

因子均方和误差均方

$$MS_A = \frac{S_A}{f_A}, \quad MS_e = \frac{S_e}{f_e}$$

##### 定理 4.1.1

1.

$$\frac{S_e}{\sigma^2} \sim \chi^2(n - r), \quad E(S_e) = (n - r)\sigma^2$$

2.

$$E(S_A) = (r - 1)\sigma^2 + m \sum_{i=1}^r a_i^2$$

3. 若  $H_0$  成立, 那么

$$\frac{S_A}{\sigma^2} \sim \chi^2(r - 1), \quad E(S_e) = (r - 1)\sigma^2$$

4.  $S_A$  与  $S_e$  相互独立。



检验统计量

$$F = \frac{MS_A}{MS_e} \sim F(r - 1, n - r)$$

拒绝域

$$W = \{F \geq F_{1-\alpha}\}$$

1.  $F \geq F_{1-\alpha}$ : 拒绝原假设, 认为因子 A 显著。

2.  $F \leq F_{1-\alpha}$ : 接受原假设, 认为因子 A 不显著。

检验的  $p$  值为

$$p = P(F \geq F_0)$$

表 4.2: 单因子方差分析表

来源	平方和	自由度	均方	$F$ 比	$p$ 值
因子 $A$	$S_A$	$f_A = r - 1$	$MS_A = \frac{S_A}{f_A}$	$F = \frac{MS_A}{MS_e}$	$p = P(F \geq F_0)$
误差 $e$	$S_e$	$f_e = n - r$	$MS_e = \frac{S_e}{f_e}$		
总和 $T$	$S_T$	$f_T = n - 1$			

### 4.1.5 参数估计

#### 一、点估计

1.

$$y_{ij} \sim N(\mu + a_i, \sigma^2)$$

2.  $\mu$  的最大似然估计为

$$\hat{\mu} = \bar{y}$$

3.  $a_i$  的最大似然估计为

$$\hat{a}_i = \bar{y}_i - \bar{y}$$

4.  $\sigma^2$  的最大似然估计为

$$\hat{\sigma}^2 = MS_e$$

#### 二、置信区间

由于

$$\bar{y}_i \sim N(\mu_i, \frac{\sigma^2}{m}), \quad \frac{S_e}{\sigma^2} \sim \chi^2(n-r)$$

因此

$$\frac{\sqrt{m}(\bar{y}_i - \mu_i)}{\sqrt{\hat{\sigma}^2}} \sim T(f_e)$$

进而  $\mu_i$  的  $1 - \alpha$  的置信区间为

$$\left[ \bar{y}_i - t_{1-\frac{\alpha}{2}} \frac{\hat{\sigma}}{\sqrt{m}}, \quad \bar{y}_i + t_{1-\frac{\alpha}{2}} \frac{\hat{\sigma}}{\sqrt{m}} \right]$$

### 4.1.6 重复数不等情形

#### 一、数据

记从第  $i$  个水平下的总体获得  $m_i$  个试验结果, 记为  $y_{i1}, \dots, y_{im_i}$ , 其中  $i = 1, \dots, r$ , 实验总次数为  $n = m_1 + \dots + m_r$ , 统计模型为

$$\begin{cases} y_{ij} = \mu_i + \varepsilon_{ij} \\ \varepsilon_{ij} \text{ 相互独立} \\ \varepsilon_{ij} \sim N(0, \sigma^2) \end{cases}$$

#### 二、总均值

加权均值

$$\mu = \frac{1}{n} \sum_{i=1}^r m_i \mu_i$$

水平效应:

$$a_i = \mu_i - \mu$$

统计模型为

$$\begin{cases} y_{ij} = \mu + a_i + \varepsilon_{ij} \\ \sum_{i=1}^r m_i a_i = 0 \\ \varepsilon_{ij} \\ \varepsilon_{ij} \sim N(0, \sigma^2) \end{cases}$$

#### 四、各平方和的计算

$$T_i = \sum_{j=1}^{m_i} y_{ij}, \quad \bar{y}_i = \frac{T_i}{m_i}$$

$$T = \sum_{i,j} y_{ij} = \sum_{i=1}^r T_i, \quad \bar{y} = \frac{T}{n}$$

$$S_T = \sum_{i,j} (y_{ij} - \bar{y})^2 = \sum_{i,j} y_{ij}^2 - \frac{T^2}{n}, \quad f_T = n - 1$$

$$S_A = \sum_{i=1}^r m_i (\bar{y}_i - \bar{y})^2 = \sum_{i=1}^r \frac{T_i^2}{m_i} - \frac{T^2}{n}, \quad f_A = r - 1$$

$$S_e = \sum_{i,j} (y_{ij} - \bar{y}_i)^2 = \sum_{i,j} y_{ij}^2 - \sum_{i=1}^r \frac{T_i^2}{m_i}, \quad f_e = n - r$$

## 4.2 多重比较

### 4.2.1 水平均值差的置信区间

检验问题

$$H_0: \mu_i - \mu_j = 0 \quad \text{vs} \quad H_a: \mu_i - \mu_j \neq 0$$

由于

$$\bar{y}_i - \bar{y}_j \sim N\left(\mu_i - \mu_j, \left(\frac{1}{m_i} + \frac{1}{m_j}\right) \sigma^2\right)$$

而  $\frac{S_e}{\sigma^2} \sim \chi^2(n-r)$ , 因此

$$\frac{(\bar{y}_i - \bar{y}_j) - (\mu_i - \mu_j)}{\sqrt{\left(\frac{1}{m_i} + \frac{1}{m_j}\right) \hat{\sigma}^2}} \sim T(n-r)$$

那么置信水平为  $1 - \alpha$  的置信区间为

$$\left[ \bar{y}_i - \bar{y}_j - t_{1-\frac{\alpha}{2}} \sqrt{\left(\frac{1}{m_i} + \frac{1}{m_j}\right) \hat{\sigma}^2}, \quad \bar{y}_i - \bar{y}_j + t_{1-\frac{\alpha}{2}} \sqrt{\left(\frac{1}{m_i} + \frac{1}{m_j}\right) \hat{\sigma}^2} \right]$$

这也是检验问题的接受域  $\bar{W}$ 。如果包含 0, 那么接受原假设, 认为  $\mu_i$  和  $\mu_j$  无显著差异; 反之拒绝原假设, 认为  $\mu_i$  和  $\mu_j$  存在显著差异。

### 4.2.2 多重比较问题

首先经过方差检验, 表明因子 A 是显著的, 即  $r$  个水平均值不全相等, 那么考虑如下多重比较问题检验

$$H_0^{ij}: \mu_i = \mu_j, \quad 1 \leq i < j \leq r$$

拒绝域

$$W = \bigcup_{1 \leq i < j \leq r} \{|\bar{y}_i - \bar{y}_j| \geq c_{ij}\}$$

### 4.2.3 重复数相等的 T 法

当  $m_1 = \cdots = m_r = m$  时, 记  $c_{ij} = c$ , 于是检验统计量为

$$\frac{\sqrt{m}(\bar{y}_i - \mu_i)}{\hat{\sigma}} \sim T(n - r)$$

当原假设成立时,  $\mu_1 = \cdots = \mu_r = \mu$ , 此时

$$P(W) = P\left(q(r, n - r) \geq \frac{c\sqrt{m}}{\hat{\sigma}}\right)$$

其中  $t$  化极差统计量为

$$q(r, n - r) = \max_{1 \leq i \leq r} \frac{\sqrt{m}(\bar{y}_i - \mu_i)}{\hat{\sigma}} - \min_{1 \leq j \leq r} \frac{\sqrt{m}(\bar{y}_j - \mu_j)}{\hat{\sigma}}$$

仅与  $n$  和  $r$  有关。由  $P(W) = \alpha$ , 可知

$$c = q_{1-\alpha} \frac{\hat{\sigma}}{\sqrt{m}}$$

因此, 如果

$$|\bar{y}_i - \bar{y}_j| \geq q_{1-\alpha} \frac{\hat{\sigma}}{\sqrt{m}}$$

那么认为水平  $A_i$  和  $A_j$  存在显著差异; 反之认为水平  $A_i$  和  $A_j$  无显著差异。

### 4.2.4 重复数不等场合的 S 法

由于

$$\frac{(\bar{y}_i - \bar{y}_j) - (\mu_i - \mu_j)}{\sqrt{\left(\frac{1}{m_i} + \frac{1}{m_j}\right) \hat{\sigma}^2}} \sim T(n - r)$$

当原假设成立时,  $\mu_1 = \cdots = \mu_r = \mu$ , 此时

$$\frac{(\bar{y}_i - \bar{y}_j)^2}{\left(\frac{1}{m_i} + \frac{1}{m_j}\right) \hat{\sigma}^2} \sim F(1, n - r)$$

令  $c_{ij} = c\sqrt{\frac{1}{m_i} + \frac{1}{m_j}}$ , 那么

$$P(W) = P\left(\max_{1 \leq i < j \leq r} \frac{(\bar{y}_i - \bar{y}_j)^2}{\left(\frac{1}{m_i} + \frac{1}{m_j}\right) \hat{\sigma}^2} \geq \frac{c^2}{\hat{\sigma}^2}\right)$$

其中

$$\frac{\max_{1 \leq i < j \leq r} \frac{(\bar{y}_i - \bar{y}_j)^2}{\left(\frac{1}{m_i} + \frac{1}{m_j}\right) \hat{\sigma}^2}}{r - 1} \sim F(r - 1, n - r)$$

由  $P(W) = \alpha$ , 可知

$$\frac{c^2}{\hat{\sigma}^2} = (r - 1)f_{1-\alpha}$$

即

$$c_{ij} = \sqrt{(r - 1)f_{1-\alpha}\hat{\sigma}^2 \left(\frac{1}{m_i} + \frac{1}{m_j}\right)}$$

其中  $f_{1-\alpha}$  为  $F(r-1, n-r)$  的  $1-\alpha$  分位数。因此, 如果

$$|\bar{y}_i - \bar{y}_j| \geq \sqrt{(r-1)f_{1-\alpha}\hat{\sigma}^2 \left( \frac{1}{m_i} + \frac{1}{m_j} \right)}$$

那么认为水平  $A_i$  和  $A_j$  存在显著差异; 反之认为水平  $A_i$  和  $A_j$  无显著差异。

### 4.3 方差齐性检验

方差齐性检验

$$H_0: \sigma_1^2 = \cdots = \sigma_r^2$$

#### 4.3.1 Hartley 检验

对于单因子方差分析中含有  $r$  个样本, 当  $m_1 = \cdots = m_r = m$  时, 设第  $i$  个样本方差为

$$s_i^2 = \frac{1}{m-1} \sum_{k=1}^m (y_{ik} - \bar{y}_i)^2$$

检验统计量为

$$H = \frac{\max\{s_i^2\}}{\min\{s_i^2\}} \sim H(r, m-1)$$

拒绝域为

$$W = \{H \geq H_{1-\alpha}\}$$

#### 4.3.2 Bartlett 检验

对于单因子方差分析中含有  $r$  个样本, 设第  $i$  个样本方差为

$$s_i^2 = \frac{1}{m_i-1} \sum_{k=1}^{m_i} (y_{ik} - \bar{y}_i)^2 = \frac{Q_i}{f_i}$$

其中  $m_i$  为第  $i$  个样本的容量且  $m_i \geq 5$ ,  $Q_i = \sum_{k=1}^{m_i} (y_{ik} - \bar{y}_i)^2$  和  $f_i = m_i - 1$  为该样本的偏差平方和自由度。 $s_i^2$  的算术加权平均即为均方误差

$$MS_e = \frac{1}{f_e} \sum_{i=1}^r Q_i = \sum_{i=1}^r \frac{f_i}{f_e} s_i^2$$

其加权几何平均为

$$GMS_e = \left( \prod_{i=1}^r (s_i^2)^{f_i} \right)^{\frac{1}{f_e}}$$

其中  $f_e = \sum_{i=1}^r f_i = n - r$ 。由算术-几何平均不等式

$$MS_e \geq GMS_e$$

当且仅当  $s_1^2 = \cdots = s_r^2$  时等号成立。而

$$B = \frac{f_e}{C} \ln \frac{MS_e}{GMS_e} \sim \chi^2(r-1)$$

其中

$$C = 1 + \frac{1}{3(r-1)} \left( \sum_{i=1}^r \frac{1}{f_i} - \frac{1}{f_e} \right)$$

因此拒绝域为

$$W = \{B \geq \chi_{1-\alpha}^2\}$$

### 4.3.3 修正的 Bartlett 检验

修正的检验统计量

$$B' = \frac{f_2 BC}{f_1(A - BC)} \sim F(f_1, f_2)$$

其中

$$f_1 = r - 1, \quad f_2 = \frac{r + 1}{(C - 1)^2}, \quad A = \frac{f_2}{2 - C + \frac{2}{f_2^0}}$$

拒绝域为

$$W = \{B' \geq F_{1-\alpha}\}$$

## 4.4 一元线性回归

### 4.4.1 变量间的两类关系

确定性关系，相关关系

### 4.4.2 一元线性回归模型

第一类回归问题

$$f(x) = E(Y | x) = \int_{-\infty}^{\infty} yp(y | x)dx$$

第二类回归问题

$$y = f(x) + \varepsilon$$

其中  $\varepsilon \sim N(0, \sigma^2)$ 。

一元回归模型：

$$\begin{cases} y_i = \beta_0 + \beta x_i + \varepsilon_i \\ \varepsilon_i \text{相互独立} \\ \varepsilon_i \sim N(0, \sigma^2) \end{cases}$$

由数据  $(x_i, y_i)$  得到的  $\beta_0$  和  $\beta$  的估计  $\hat{\beta}_0$  和  $\hat{\beta}$ ，称

$$\hat{y} = \hat{\beta}_0 + \hat{\beta}x$$

为  $y$  关于  $x$  的回归函数。给定  $x = x_0$ ，称  $\hat{y}_0 = \hat{\beta}_0 + \hat{\beta}x_0$  为回归值。

### 4.4.3 回归系数的最小二乘估计

$$\begin{aligned} l_{xy} &= \sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y}) = \sum_{i=1}^n x_i y_i - \frac{1}{n} \sum_{i=1}^n x_i \sum_{i=1}^n y_i = \sum_{i=1}^n x_i y_i - n\bar{x}\bar{y} \\ l_{xx} &= \sum_{i=1}^n (x_i - \bar{x})^2 = \sum_{i=1}^n x_i^2 - \frac{1}{n} \left( \sum_{i=1}^n x_i \right)^2 = \sum_{i=1}^n x_i^2 - n\bar{x}^2 \\ l_{yy} &= \sum_{i=1}^n (y_i - \bar{y})^2 = \sum_{i=1}^n y_i^2 - \frac{1}{n} \left( \sum_{i=1}^n y_i \right)^2 = \sum_{i=1}^n y_i^2 - n\bar{y}^2 \end{aligned}$$



$\beta_0$  和  $\beta$  的最小二乘估计 (LSE)  $\hat{\beta}_0$  和  $\hat{\beta}$  为

$$\hat{\beta} = \frac{l_{xy}}{l_{xx}}$$

$$\hat{\beta}_0 = \bar{y} - \hat{\beta}\bar{x}$$

#### 定理 4.4.1

在如下模型下, 成立

$$\begin{cases} y_i = \beta_0 + \beta x_i + \varepsilon_i \\ \varepsilon_i \\ \varepsilon_i \sim N(0, \sigma^2) \end{cases}$$

1.

$$\hat{\beta}_0 \sim N\left(\beta_0, \left(\frac{1}{n} + \frac{\bar{x}^2}{l_{xx}}\right)\sigma^2\right), \quad \hat{\beta} \sim N\left(\beta, \frac{\sigma^2}{l_{xx}}\right)$$

2.

$$\text{Cov}(\hat{\beta}_0, \hat{\beta}) = -\frac{\bar{x}}{l_{xx}}\sigma^2$$

3. 给定  $x_0$

$$\hat{y}_0 = \hat{\beta}_0 + \hat{\beta}x_0 \sim N\left(\beta_0 + \beta x_0, \left(\frac{1}{n} + \frac{(\bar{x} - x_0)^2}{l_{xx}}\right)\sigma^2\right)$$



### 4.4.4 回归模型的显著性检验

**显著性:**  $\beta \neq 0$  称为显著, 否则称为不显著。

显著性假设检验:

$$H_0: \beta = 0 \quad \text{vs} \quad H_a: \beta \neq 0$$

表 4.3: 方差分析表

来源	平方和	自由度	均方	F 比
回归	$S_R$	$f_R = 1$	$MS_R = \frac{S_R}{f_R}$	$F = \frac{MS_A}{MS_e}$
残差	$S_e$	$f_e = n - 2$	$MS_e = \frac{S_e}{f_e}$	
总和	$S_T$	$f_T = n - 1$		

#### 一、F 检验

总偏差平方和:

$$S_T = \sum_{i=1}^n (y_i - \bar{y})^2 = l_{yy}$$

回归平方和:

$$S_R = \sum_{i=1}^n (\hat{y}_i - \bar{y})^2 = \frac{l_{xy}^2}{l_{xx}}$$

残差平方和:

$$S_e = \sum_{i=1}^n (y_i - \hat{y}_i)^2$$

平方和分解:

$$S_T = S_R + S_e$$

**定理 4.4.2**

在如下模型下, 成立

$$\begin{cases} y_i = \beta_0 + \beta x_i + \varepsilon_i \\ \varepsilon_i \text{ 相互独立} \\ \varepsilon_i \sim N(0, \sigma^2) \end{cases}$$

1.

$$E(S_R) = \sigma^2 + \hat{\beta} l_{xx}, \quad E(S_e) = (n-2)\sigma^2$$

2.

$$\frac{S_e}{\sigma^2} \sim \chi^2(n-2)$$

3. 如果  $H_0$  成立, 那么

$$\frac{S_R}{\sigma^2} \sim \chi^2(1)$$

4.  $S_R$  与  $S_e$ 、 $\bar{y}$  独立。



统计检验量

$$F = \frac{(n-2)S_R}{S_e} \sim F(1, n-2)$$

拒绝域为  $W = \{F \geq F_{1-\alpha}\}$ 。

二、 $T$  检验

检验统计量

$$T = \frac{\sqrt{(n-2)l_{xx}}\hat{\beta}}{\sqrt{S_e}} \sim T(n-2)$$

拒绝域为  $W = \{|t| \geq t_{1-\frac{\alpha}{2}}\}$ 。

三、相关系数检验

相关系数假设检验:

$$H_0: \rho = 0 \quad \text{vs} \quad H_a: \rho \neq 0$$

检验统计量: 相关系数

$$r = \frac{l_{xy}}{\sqrt{l_{xx}l_{yy}}} = \sqrt{\frac{F}{F + (n-2)}} \sim r(n-2) = \sqrt{\frac{F(1, n-2)}{F(1, n-2) + (n-2)}}$$

1.  $|r| = 1$ :  $(x_i, y_i)$  共线。

2.  $r > 0$ :  $(x_i, y_i)$  正相关。

3.  $r < 0$ :  $(x_i, y_i)$  负相关。

4.  $r = 0$ :  $(x_i, y_i)$  不相关。

拒绝域为  $W = \{|r| \geq r_{1-\alpha}\}$ , 其中

$$r_{1-\alpha} = \sqrt{\frac{F_{1-\alpha}}{F_{1-\alpha} + (n-2)}}$$

#### 4.4.5 估计与预测

一、 $E(y_0)$  的置信区间

枢轴量为

$$\frac{\hat{y}_0 - E(y_0)}{\sqrt{\frac{S_e}{n-2} \sqrt{\frac{1}{n} + \frac{(x_0 - \bar{x})^2}{l_{xx}}}}} \sim T(n-2)$$

$1 - \alpha$  的置信区间为

$$\left[ \hat{y}_0 - t_{1-\frac{\alpha}{2}} \sqrt{\frac{S_e}{n-2}} \sqrt{\frac{1}{n} + \frac{(x_0 - \bar{x})^2}{l_{xx}}}, \quad \hat{y}_0 + t_{1-\frac{\alpha}{2}} \sqrt{\frac{S_e}{n-2}} \sqrt{\frac{1}{n} + \frac{(x_0 - \bar{x})^2}{l_{xx}}} \right]$$

二、 $y_0$  的预测区间

枢轴量为

$$\frac{y_0 - \hat{y}_0}{\sqrt{\frac{S_e}{n-2}} \sqrt{1 + \frac{1}{n} + \frac{(x_0 - \bar{x})^2}{l_{xx}}}} \sim T(n-2)$$

预测区间为

$$\left[ \hat{y}_0 - t_{1-\frac{\alpha}{2}} \sqrt{\frac{S_e}{n-2}} \sqrt{1 + \frac{1}{n} + \frac{(x_0 - \bar{x})^2}{l_{xx}}}, \quad \hat{y}_0 + t_{1-\frac{\alpha}{2}} \sqrt{\frac{S_e}{n-2}} \sqrt{1 + \frac{1}{n} + \frac{(x_0 - \bar{x})^2}{l_{xx}}} \right]$$

#### 4.4.6 曲线回归方程的比较

决定系数：越大说明残差越小，回归曲线拟合越好。

$$R^2 = 1 - \frac{\sum (y_i - \hat{y}_i)^2}{\sum (y_i - \bar{y})^2}$$

剩余标准差：越小，回归曲线拟合越好。

$$s = \sqrt{\frac{\sum (y_i - \hat{y}_i)^2}{n-2}}$$

## 附录 A 概率模型

概率模型	密度函数 $p(x)$	参数范围	数学期望 $E\xi$	方差 $D\xi$	特征函数 $f(t)$
退化分布 $I_c(x)$	$p(x) = \begin{cases} 1, & x = c \\ 0, & x \neq c \end{cases}$		$c$	0	$e^{ict}$
Bernoulli 分布	$p(x) = \begin{cases} 1-p, & x = 0 \\ p, & x = 1 \end{cases}$	$0 < p < 1$	$p$	$p(1-p)$	$pe^{it} + 1 - p$
二项分布 $B(n, p)$	$b(k; n, p) = \binom{n}{k} p^k (1-p)^{n-k}$	$0 \leq k \leq n; 0 < p < 1$	$np$	$np(1-p)$	$(pe^{it} + 1 - p)^n$
Poisson 分布 $P(\lambda)$	$p(k; \lambda) = \frac{\lambda^k}{k!} e^{-\lambda}$	$k \in \mathbb{N}; \lambda > 0$	$\lambda$	$\lambda$	$e^{\lambda(e^{it} - 1)}$
几何分布	$g(k; p) = p(1-p)^{k-1}$	$k \in \mathbb{N}^*, 0 < p < 1$	$\frac{1}{p}$	$\frac{q}{p^2}$	$\frac{pe^{it}}{1 - (1-p)e^{it}}$
超几何分布	$p_k = \frac{\binom{M}{k} \binom{N-M}{n-k}}{\binom{N}{n}}$	$M, n \leq N; 0 \leq k \leq \min\{M, n\}$	$\frac{nM}{N}$	$\frac{nM(N-M)(N-n)}{N^2(N-1)}$	$\sum_{k=0}^n \frac{\binom{M}{k} \binom{N-M}{n-k}}{\binom{N}{n}} e^{ikt}$
Pascal 分布	$p_k = \binom{k-1}{r-1} p^r (1-p)^{k-r}$	$k \geq r, 0 < p < 1$	$\frac{r}{p}$	$\frac{r(1-p)}{p^2}$	$(\frac{(1-p)e^{it}}{1 - (1-p)e^{it}})^r$
负二项分布	$p_k = \binom{k-1}{r-1} p^r (1-p)^{k-r}$	$k \in \mathbb{N}, 0 < p < 1, r > 0$	$\frac{r(1-p)}{p}$	$\frac{r(1-p)}{p^2}$	$(\frac{p}{1 - (1-p)e^{it}})^r$
正态分布 $N(\mu, \sigma^2)$	$p(x) = \frac{1}{\sqrt{2\pi}\sigma} e^{-\frac{(x-\mu)^2}{2\sigma^2}}$		$\mu$	$\sigma$	$e^{i\mu t - \frac{1}{2}\sigma^2 t^2}$
均匀分布 $U[a, b]$	$p(x) = \begin{cases} \frac{1}{b-a}, & a \leq x \leq b \\ 0, & \text{其他} \end{cases}$	$a < b$	$\frac{a+b}{2}$	$\frac{(b-a)^2}{12}$	$\frac{e^{ibt} - e^{iat}}{i(b-a)t}$
指数分布 $\text{Exp}(\lambda)$	$p(x) = \begin{cases} \lambda e^{-\lambda x}, & x \geq 0 \\ 0, & x < 0 \end{cases}$	$\lambda > 0$	$\lambda^{-1}$	$\lambda^{-2}$	$(1 - \frac{it}{\lambda})^{-1}$
$\chi^2$ 分布	$p(x) = \begin{cases} \frac{1}{2^{\frac{n}{2}} \Gamma(\frac{n}{2})} x^{\frac{n}{2}-1} e^{-\frac{x}{2}}, & x \geq 0 \\ 0, & x < 0 \end{cases}$	$n \in \mathbb{N}^*$	$n$	$2n$	$(1 - 2it)^{-\frac{n}{2}}$
$\Gamma$ 分布 $\Gamma(r, \lambda)$	$p(x) = \begin{cases} \frac{\lambda^r}{\Gamma(r)} x^{r-1} e^{-\lambda x}, & x \geq 0 \\ 0, & x < 0 \end{cases}$	$r, \lambda > 0$	$\frac{r}{\lambda}$	$\frac{r}{\lambda^2}$	$(1 - \frac{it}{\lambda})^{-r}$
Cauchy 分布	$p(x) = \frac{1}{\pi} \frac{\lambda}{\lambda^2 + (x-\mu)^2}$	$\mu \in \mathbb{R}, \lambda > 0$	不存在	不存在	$e^{i\mu t - \lambda t }$
$t$ 分布	$p(x) = \frac{\Gamma(\frac{n+1}{2})}{\sqrt{n\pi} \Gamma(\frac{n}{2})} (1 + \frac{x^2}{n})^{-\frac{n+1}{2}}$	$n \in \mathbb{N}^*$	$0 (n > 1)$	$\frac{n}{n-2} (n > 2)$	
Pareto 分布	$p(x) = \begin{cases} rA^r \frac{1}{x^{r+1}}, & x \geq A \\ 0, & x < A \end{cases}$	$r, A > 0$	$(r > 1 \text{ 时存在})$	$(r > 2 \text{ 时存在})$	
$F$ 分布	$p(x) = \begin{cases} \frac{\Gamma(\frac{m+n}{2})}{\Gamma(\frac{m}{2})\Gamma(\frac{n}{2})} m^{\frac{m}{2}} n^{\frac{n}{2}} \frac{x^{\frac{m}{2}-1}}{(n+mx)^{\frac{m+n}{2}}}, & x \geq 0 \\ 0, & x < 0 \end{cases}$	$m, n \in \mathbb{N}^*$	$\frac{n}{n-2} (n > 2)$	$\frac{2n^2(m+n-2)}{m(n-2)^2(n-4)} (n > 4)$	
$\beta$ 分布	$p(x) = \begin{cases} \frac{\Gamma(p+q)}{\Gamma(p)\Gamma(q)} x^{p-1} (1-x)^{q-1}, & 0 < x < 1 \\ 0, & \text{其他} \end{cases}$	$p, q > 0$	$\frac{p}{p+q}$	$\frac{pq}{(p+q)^2(p+q+1)}$	$\frac{\Gamma(p+q)}{\Gamma(p)} \sum_{k=0}^{\infty} \frac{\Gamma(p+k)(it)^k}{\Gamma(p+q+k)\Gamma(k+1)}$
对数正态分布	$p(x) = \begin{cases} \frac{1}{\sqrt{2\pi}\sigma x} e^{-\frac{(\ln \frac{x-\mu}{\sigma})^2}{2\sigma^2}}, & x > 0 \\ 0, & x \leq 0 \end{cases}$	$\alpha, \sigma > 0$	$e^{\alpha + \frac{\sigma^2}{2}}$	$e^{2\alpha + \sigma^2} (e^{\sigma^2} - 1)$	
Weibull 分布	$p(x) = \begin{cases} \alpha \lambda x^{\alpha-1} e^{-\lambda x^\alpha}, & x > 0 \\ 0, & x \leq 0 \end{cases}$	$\lambda, \alpha > 0$	$\Gamma(\frac{1}{\alpha} + 1) \lambda^{-\frac{1}{\alpha}}$	$\lambda^{-\frac{2}{\alpha}} (\Gamma(\frac{2}{\alpha} + 1) - (\Gamma(\frac{1}{\alpha} + 1))^2)$	
Rayleigh 分布	$p(x) = \begin{cases} xe^{-\frac{x^2}{2}}, & x \geq 0 \\ 0, & x < 0 \end{cases}$		$\sqrt{\frac{\pi}{2}}$	$2 - \frac{\pi}{2}$	