

数理统计

致谢

感谢邢小玉老师对于本课程学习的帮助。

目录

数理统计

致谢

目录

第五章：统计量及其分布

- 5.1 总体与概率
 - 5.1.1 总体与个体
 - 5.1.2 样本
- 5.2 样本数据的整理与表示
 - 5.2.1 经验分布函数
 - 5.2.2 频数频率分布表
 - 5.2.3 样本数据的图形显示
- 5.3 统计量及其分布
 - 5.3.1 统计量与抽样分布
 - 5.3.2 样本均值及其抽样分布
 - 5.3.3 样本方差与样本标准差
 - 5.3.4 样本矩及其函数
 - 5.3.5 次序统计量及其分布
 - 5.3.6 样本分位数与分位数
 - 5.3.7 五数概括与箱线图
- 5.4 三大抽样分布
 - 5.4.1 χ^2 分布
 - 5.4.2 F 分布
 - 5.4.3 T 分布
- 5.5 充分统计量
 - 5.5.1 充分性的概念
 - 5.5.2 因子分解定理

第六章：参数估计

- 6.1 点估计的概念
 - 6.1.1 点估计及无偏性
 - 6.1.2 有效性
- 6.2 矩估计及相关性
 - 6.2.1 替换原理和矩法估计
 - 6.2.2 概率函数已知时未知参数的矩估计
 - 6.2.3 相合性
- 6.3 最大似然估计与EM算法
 - 6.3.1 最大似然估计
 - 6.3.2 EM算法
 - 6.3.3 渐进正态性
- 6.4 最小方差无偏估计
 - 6.4.1 均方误差
 - 6.4.2 一致最小方差无偏估计
 - 6.4.3 充分性原则
 - 6.4.4 Cramer-Rao不等式
- 6.5 Bayes估计
 - 6.5.1 统计判断的基础
 - 6.5.2 Bayes公式的密度函数形式
 - 6.5.3 Bayes估计
 - 6.5.4 共轭先验分布
- 6.6 区间估计
 - 6.6.1 区间估计的概念
 - 6.6.2 枢轴量法
 - 6.6.3 单个正态总体参数的置信区间
 - 6.6.4 大样本置信区间

- 6.6.5 样本量的确定
- 6.6.6 两个正态总体下的置信区间

第七章：假设检验

- 7.1 假设检验的基本思想与概念
 - 7.1.1 假设检验问题
 - 7.1.2 假设检验的基本步骤
 - 7.1.3 检验的 p 值
- 7.2 正态总体参数假设检验
 - 7.2.1 单个正态总体均值的检验
 - 7.2.2 假设检验与置信区间的关系
 - 7.2.3 两个正态总体均值差的检验
 - 7.2.4 成对数据检验
 - 7.2.5 正态总体方差的检验
- 7.3 其他分布参数的假设检验
- 7.4 似然比检验与分布拟合检验
 - 7.4.1 似然比检验的思想
- 7.4.2 分布数据的 χ^2 拟合优度检验
- 7.4.3 分布的 χ^2 拟合优度检验
- 7.4.4 列联表的独立性检验
- 7.5 正态性检验
 - 7.5.1 正态概率纸
 - 7.5.2 W检验
 - 7.5.3 EP检验
- 7.6 非参数检验
 - 7.6.1 游程检验
 - 7.6.2 符号检验
 - 7.6.3 秩和检验

第八章：方差分析与回归分析

- 8.1 方差分析
 - 8.1.1 问题的提出
 - 8.1.2 单因子方差分析的统计模型
 - 8.1.3 平方和分解
 - 8.1.4 检验方法
 - 8.1.5 参数估计
 - 8.1.6 重复数不等情形
- 8.2 多重比较
 - 8.2.1 水平均值差的置信区间
 - 8.2.2 多重比较问题
 - 8.2.3 重复数相等的T法
 - 8.2.4 重复数不等场合的S法
- 8.3 方差齐性检验
 - 8.3.1 Hartley检验
 - 8.3.2 Bartlett检验
 - 8.3.3 修正的Bartlett检验
- 8.4 一元线性回归
 - 8.4.1 变量间的两类关系
 - 8.4.2 一元线性回归模型
 - 8.4.3 回归系数的最小二乘估计
 - 8.4.4 回归模型的显著性检验
 - 8.4.5 估计与预测
- 8.5 一元非线性回归
 - 8.5.1 确定可能的函数形式
 - 8.5.2 参数估计
 - 8.5.3 曲线回归方程的比较

附录：概率模型

第五章：统计量及其分布

5.1 总体与概率

5.1.1 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

5.1.2 样本

样本：从总体中随机的抽取 n 个个体，记其指标值为

$$x_1, \cdots, x_n \quad (1)$$

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

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

$$F(x_1, \cdots, x_n) = \prod_{k=1}^n F(x_k) \quad (2)$$

5.2 样本数据的整理与表示

5.2.1 经验分布函数

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

$$F_n(x) = \begin{cases} 0, & x < x_{(1)} \\ \frac{k}{n}, & x_{(k)} \leq x < x_{(k+1)}, k = 1, \cdots, n-1 \\ 1, & x \geq x_{(n)} \end{cases} \quad (3)$$

- $F_n(x)$ 非减且右连续。

$$\bullet \quad F_n(-\infty) = 0 \quad F_n(+\infty) = 1 \quad (4)$$

定理5.2.1 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 $F_n(x)$ ，那么当 $n \rightarrow +\infty$ 时，成立

$$P\left(\sup_{x \in \mathbb{R}} |F_n(x) - F(x)| \rightarrow 0\right) = 1 \quad (5)$$

5.2.2 频数频率分布表

- 对样本进行分组：通常为5 ~ 20个。
- 确定每组组距：

$$\text{组距 } d = \frac{\text{样本最大观测值} - \text{样本最小观测值}}{\text{组数}} \quad (6)$$

- 确定每组组限：

$$(a_0, a_1], \cdots, (a_{n-1}, a_n] \quad (7)$$

- 统计样本数据落入每个区间的个数——频数

频数频率分布表

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

5.2.3 样本数据的图形显示

- 直方图
- 茎叶图

5.3 统计量及其分布

5.3.1 统计量与抽样分布

定义5.3.1 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。统计量的分布称为**抽样分布**。

5.3.2 样本均值及其抽样分布

定义5.3.2 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

$$\bar{x} = \frac{1}{n} \sum_{k=1}^n x_k \quad (8)$$

特别的，在分组样本中，样本均值的近似公式为

$$\bar{x} = \frac{1}{n} \sum_{k=1}^m x_k f_k \quad (9)$$

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

$$n = \sum_{k=1}^m f_k \quad (10)$$

性质5.3.1：样本的所有偏差之和为0，即

$$\sum_{k=1}^n (x_k - \bar{x}) = 0 \quad (11)$$

性质5.3.2：对于任意 $c \in \mathbb{R}$ ，成立

$$\sum_{k=1}^n (x_k - \bar{x})^2 \leq \sum_{k=1}^n (x_k - c)^2 \quad (12)$$

当且仅当 $c = \bar{x}$ 时等号成立。

定理5.3.1：对于取自总体的样本 x_1, \dots, x_n ，记其样本均值为 \bar{x} 。

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

$$\bar{x} \sim N(\mu, \frac{\sigma^2}{n}) \quad (13)$$

5.3.3 样本方差与样本标准差

定义5.3.3 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

$$s^2 = \frac{1}{n-1} \sum_{k=1}^n (x_k - \bar{x})^2 = \frac{1}{n-1} \left(\sum_{k=1}^n x_k^2 - n\bar{x}^2 \right) \quad (14)$$

定理5.3.2：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

$$E(\bar{x}) = \mu, \quad \text{Var}(\bar{x}) = \frac{\sigma^2}{n} \quad (15)$$

$$E(s^2) = \sigma^2 \quad (16)$$

5.3.4 样本矩及其函数

定义5.3.4 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

$$a_k = \frac{1}{n} \sum_{i=1}^n x_i^k \quad (17)$$

定义5.3.5 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

$$b_k = \frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})^k \quad (18)$$

定义5.3.6 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

$$\hat{\beta}_s = \frac{b_3}{b_2^{\frac{3}{2}}} \quad (19)$$

样本偏差反应总体分布密度曲线的对称性。

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

定义5.3.7 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

$$\hat{\beta}_k = \frac{b_4}{b_2^2} - 3 \quad (20)$$

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。

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

5.3.5 次序统计量及其分布

定义5.3.8 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

$$x_{(1)}, \dots, x_{(n)} \quad (21)$$

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

定理5.3.3 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

$$p_k(x) = \frac{n!}{(k-1)!(n-k)!} (F(x))^{k-1} (1-F(x))^{n-k} p(x) \quad (22)$$

- $x_{(1)}$ 的密度函数为

$$p_1(x) = np(x)(1-F(x))^{n-1} \quad (23)$$

分布函数为

$$F_1(x) = 1 - (1-F(x))^n \quad (24)$$

- $x_{(n)}$ 的密度函数为

$$p_n(x) = np(x)(F(x))^{n-1} \quad (25)$$

分布函数为

$$F_n(x) = (F(x))^n \quad (26)$$

定理5.3.4 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 $i < j$ 。

5.3.6 样本分位数与分位数

定义5.3.9 样本 p 分位数: 对于取自总体的样本 x_1, \cdots, 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5.3.10 α 分位数: 对于随机变量 X , 称 x_α 为其 α 分位数, 如果

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

定理5.3.5: 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

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

5.3.7 五数概括与箱线图

定义5.3.111 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5.3.12 箱线图:

- 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。
- 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。

5.4 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	$N(\mu,\sigma^2)$		$p(x)=\frac{1}{\sqrt{2\pi}\sigma}e^{-\frac{(x-\mu)^2}{2\sigma^2}},x\in\mathbb{R}$	μ	σ^2	$e^{it\mu-\frac{1}{2}\sigma^2t^2}$
Γ 分布	$\Gamma(\alpha,\lambda)$		$p(x)=\frac{\lambda^\alpha}{\Gamma(\alpha)}x^{\alpha-1}e^{-\lambda x},x>0$	$\frac{\alpha}{\lambda}$	$\frac{\alpha}{\lambda^2}$	$(1-\frac{it}{\lambda})^{-\alpha}$
χ^2 分布	$\chi^2(n)$	$\chi^2(n)=\sum_{i=1}^n(N(0,1))^2$	$p(x)=\frac{1}{\Gamma(\frac{n}{2})2^{\frac{n}{2}}}x^{\frac{n}{2}-1}e^{-\frac{x}{2}},x>0$	n	$2n$	$(1-2it)^{-\frac{n}{2}}$
F 分布	$F(m,n)$	$F(m,n)=\frac{\chi^2(m)/m}{\chi^2(n)/n}$	$p(x)=\frac{\Gamma(\frac{m+n}{2})\Gamma(\frac{n}{2})}{\Gamma(\frac{m}{2})\Gamma(\frac{n}{2})}x^{\frac{m}{2}-1}(1+\frac{m}{n}x)^{-\frac{m+n}{2}},x>0$	$\frac{n}{m-2}$	$\frac{2n^2(m+n-2)}{m(m-2)(m-4)}$	
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	$\frac{n}{n-2}$	

Γ 函数

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

分布间的联系:

- 若 $X\sim N(0,1)$, 那么

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

因此

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

- 如果 $X\sim N(\mu,\sigma^2)$, 那么

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

- 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

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

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

$$aX \sim \Gamma(\alpha, \frac{\lambda}{a}) \quad (37)$$

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

$$X + Y \sim \Gamma(\alpha_1 + \alpha_2, \lambda) \quad (38)$$

5.4.1 χ^2 分布

定义5.4.1 χ^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 χ^2 分布, 记作 $X \sim \chi^2(n)$, 其密度函数为

$$p(x) = \frac{1}{\Gamma\left(\frac{n}{2}\right)2^{\frac{n}{2}}} x^{\frac{n}{2}-1} e^{-\frac{x}{2}}, \quad x > 0 \quad (39)$$

定理5.4.1: 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

- \bar{x} 和 s^2 相互独立。

$$\bar{x} \sim N\left(\mu, \frac{\sigma^2}{n}\right) \quad (40)$$

-

$$\frac{(n-1)s^2}{\sigma^2} \sim \chi^2(n-1) = \Gamma\left(\frac{n-1}{2}, \frac{1}{2}\right) \quad (41)$$

-

即

$$s^2 \sim \Gamma\left(\frac{n-1}{2}, \frac{n-1}{2\sigma^2}\right) \quad (42)$$

5.4.2 F 分布

定义5.4.2 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

$$p(x) = \frac{\Gamma\left(\frac{m+n}{2}\right)\left(\frac{m}{n}\right)^{\frac{m}{2}}}{\Gamma\left(\frac{m}{2}\right)\Gamma\left(\frac{n}{2}\right)} x^{\frac{m}{2}-1} \left(1 + \frac{m}{n}x\right)^{-\frac{m+n}{2}}, \quad x > 0 \quad (43)$$

推论5.4.1: 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

$$\frac{s_x^2/\sigma_1^2}{s_y^2/\sigma_2^2} \sim F(m-1, n-1) \quad (44)$$

5.4.3 T 分布

定义5.4.3 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

$$p(x) = \frac{\Gamma\left(\frac{n+1}{2}\right)}{\sqrt{n\pi}\Gamma\left(\frac{n}{2}\right)} \left(1 + \frac{x^2}{n}\right)^{-\frac{n+1}{2}}, \quad x \in \mathbb{R} \quad (45)$$

推论5.4.2: 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

$$\frac{\sqrt{n}(\bar{x} - \mu)}{s} \sim T(n-1) \quad (46)$$

推论5.4.3: 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

$$s_w^2 = \frac{(m-1)s_x^2 + (n-1)s_y^2}{m+n-2} \quad (47)$$

那么

$$\frac{(\bar{x} - \bar{y}) - (\mu_1 - \mu_2)}{s_w \sqrt{\frac{1}{m} + \frac{1}{n}}} \sim T(m+n-2) \quad (48)$$

5.5 充分统计量

5.5.1 充分性的概念

定义5.5.1 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 θ 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 θ 无关。

5.5.2 因子分解定理

定理5.5.1 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 θ 和 x_1, \dots, x_n ，成立

$$f(x_1, \dots, x_n; \theta) = g(T(x_1, \dots, x_n); \theta)h(x_1, \dots, x_n) \quad (49)$$

定理5.5.2：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。

第六章：参数估计

6.1 点估计的概念

6.1.1 点估计及无偏性

定义6.1.1 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 θ 的统计量 $\hat{\theta} = \hat{\theta}(x_1, \dots, x_n)$ 称为 θ 的点估计。

定义6.1.2 无偏估计：对于 θ 的点估计 $\hat{\theta} = \hat{\theta}(x_1, \dots, x_n)$ ， θ 的参数空间为 Θ ，称 $\hat{\theta}$ 为 θ 的无偏估计，如果对于任意 $\theta \in \Theta$ ，成立

$$E_{\theta}(\hat{\theta}) = \theta \quad (50)$$

当参数存在无偏估计时，称其为**可估的**，否则称为**不可估的**。

6.1.2 有效性

定义6.1.3 有效性：对于 θ 的两个无偏估计 $\hat{\theta}_1$ 和 $\hat{\theta}_2$ ，如果对于任意 $\theta \in \Theta$ ，成立

$$\text{Var}(\hat{\theta}_1) \leq \text{Var}(\hat{\theta}_2) \quad (51)$$

且存在 $\theta_0 \in \Theta$ ，使得成立

$$\text{Var}(\hat{\theta}_1) < \text{Var}(\hat{\theta}_2) \quad (52)$$

那么称 $\hat{\theta}_1$ 比 $\hat{\theta}_2$ 有效。

6.2 矩估计及相关性

6.2.1 替换原理和矩法估计

替换原理 矩法：

- 用样本矩替换总体矩。
- 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：

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

Хинчин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

$$\lim_{n \rightarrow \infty} P \left\{ \left| \frac{1}{n} \sum_{i=1}^n X_i^k - E(X^k) \right| \geq \varepsilon \right\} = 0 \quad (53)$$

即

$$\frac{1}{n} \sum_{i=1}^n X_i^k \xrightarrow{P} E(X^k) \quad (54)$$

6.2.2 概率函数已知时未知参数的矩估计

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 μ_i 存在，而且 $\theta_i = \theta_i(\mu_1, \dots, \mu_k)$ ，其中 $1 \leq i \leq k$ ，那么 θ_i 的矩估计为

$$\hat{\theta}_i = \theta_i(a_1, \dots, a_k), \quad i = 1, \dots, k \quad (55)$$

其中 a_i 为样本 i 阶原点矩

$$a_i = \frac{1}{n} \sum_{j=1}^n x_j^i, \quad i = 1, \dots, k \quad (56)$$

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

$$\hat{\eta} = g(\hat{\theta}_1, \dots, \hat{\theta}_k) \quad (57)$$

6.2.3 相合性

定义6.2.1 相合性: 对于未知参数 θ , 以及 θ 的一个估计量 $\hat{\theta}_n = \hat{\theta}_n(x_1, \dots, x_n)$, 称 $\hat{\theta}_n$ 为参数 θ 的相合估计, 如果对于任意 $\varepsilon > 0$, 成立

$$\lim_{n \rightarrow \infty} P(|\hat{\theta}_n - \theta| \geq \varepsilon) = 0 \quad (58)$$

定理6.2.1: 对于 θ 的一个估计量 $\hat{\theta}_n = \hat{\theta}_n(x_1, \dots, x_n)$, 如果

$$\lim_{n \rightarrow \infty} E(\hat{\theta}_n) = \theta, \quad \lim_{n \rightarrow \infty} \text{Var}(\hat{\theta}_n) = 0 \quad (59)$$

那么 $\hat{\theta}_n$ 为参数 θ 的相合估计。

定理6.2.2: 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 η 的相合估计。

6.3 最大似然估计与EM算法

6.3.1 最大似然估计

定义6.3.1 最大似然估计: 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, Θ 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

$$L(\theta) = L(\theta; x_1, \dots, x_n) = \prod_{k=1}^n p(x_k; \theta) \quad (60)$$

为样本的**似然函数**。如果某统计量 $\hat{\theta} = \hat{\theta}(x_1, \dots, x_n)$ 满足, 对于任意 $\theta \in \Theta$, 成立

$$L(\hat{\theta}) \geq L(\theta) \quad (61)$$

那么称 $\hat{\theta}$ 为 θ 的最大似然估计, 简记为MLE。

定理6.3.1 最大似然估计的不变性: 如果 $\hat{\theta}$ 为 θ 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。

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 μ 和 σ^2 的最大似然估计分别为

$$\hat{\mu} = \bar{x}, \quad \hat{\sigma}^2 = \frac{n-1}{n} s^2 \quad (62)$$

6.3.2 EM算法

6.3.3 渐进正态性

定义6.3.2 渐进正态分布: 参数 θ 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6.3.2: 对于密度函数为 $p(x; \theta)$ 的总体 X , 其中 $\theta \in \Theta$, 如果

- 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。
- 对于任意 $\theta \in \Theta$, 成立

$$\left| \frac{\partial p}{\partial \theta} \right| < F_1(x), \quad \left| \frac{\partial^2 p}{\partial \theta^2} \right| < F_2(x), \quad \left| \frac{\partial^3 p}{\partial \theta^3} \right| < F_3(x) \quad (63)$$

其中函数 $F_1(x), F_2(x), F_3(x)$ 满足

$$\int_{-\infty}^{\infty} F_1(x) dx < \infty, \quad \int_{-\infty}^{\infty} F_2(x) dx < \infty \quad (64)$$

$$\sup_{\theta \in \Theta} \int_{-\infty}^{\infty} F_3(x) p(x; \theta) dx < \infty \quad (65)$$

- 对于任意 $\theta \in \Theta$, 成立

$$0 < I(\theta) = \int_{-\infty}^{\infty} \left(\frac{\partial \ln p}{\partial \theta} \right)^2 p(x; \theta) dx < \infty \quad (66)$$

那么对于来自该总体的样本 x_1, \dots, x_n , 存在未知参数 θ 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

$$\hat{\theta}_n \sim AN \left(\theta, \frac{1}{nI(\theta)} \right) \quad (67)$$

6.4 最小方差无偏估计

6.4.1 均方误差

定义6.4.1 均方误差: 对于 θ 的点估计 $\hat{\theta} = \hat{\theta}(x_1, \dots, x_n)$, 称 $E(\hat{\theta} - \theta)^2$ 为 $\hat{\theta}$ 关于 θ 的均方误差, 记为 $MSE(\hat{\theta}, \theta)$, 或 $M_\theta(\hat{\theta})$.

- 对于 θ 的任意估计 $\hat{\theta}$ 而言, 成立

$$MSE(\hat{\theta}, \theta) = \text{Var}(\hat{\theta}) + (E(\hat{\theta}) - \theta)^2 \quad (68)$$

- 对于 θ 的无偏估计 $\hat{\theta}$ 而言, 成立

$$MSE(\hat{\theta}, \theta) = \text{Var}(\hat{\theta}) \quad (69)$$

定义6.4.2 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 θ 的一个估计类, 称 $\hat{\theta}(x_1, \dots, x_n)$ 是该估计类中 θ 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 θ 的估计 $\tilde{\theta}$, 在参数空间 Θ 上均成立

$$MSE_\theta(\hat{\theta}) \leq MSE_\theta(\tilde{\theta}) \quad (70)$$

6.4.2 一致最小方差无偏估计

定义6.4.3 一致最小方差无偏估计: 对于 θ 的一个无偏估计 $\hat{\theta}$, 称 $\hat{\theta}$ 是 θ 的一致最小方差无偏估计(简记为UMVUE), 如果对于 θ 的任意无偏估计 $\tilde{\theta}$, 在参数空间 Θ 上均成立

$$\text{Var}_\theta(\hat{\theta}) \leq \text{Var}_\theta(\tilde{\theta}) \quad (71)$$

定理6.4.1: 对于来自某总体的样本 $X = (x_1, \dots, x_n)$, 如果 $\hat{\theta} = \hat{\theta}(X)$ 是 θ 的一个无偏估计, $\text{Var}(\hat{\theta}) < \infty$, 那么 $\hat{\theta}$ 是 θ 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

$$\text{Cov}_\theta(\hat{\theta}, \varphi) = 0 \quad (72)$$

即

$$E(\hat{\theta}\varphi) = 0 \quad (73)$$

6.4.3 充分性原则

定理6.4.2: 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 θ 的充分统计量, 那么对于 θ 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 θ 的无偏估计, 且

$$\text{Var}(\tilde{\theta}) \leq \text{Var}(\hat{\theta}) \quad (74)$$

6.4.4 Cramer-Rao不等式

定义6.4.4 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

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

$$\int_{-\infty}^{\infty} \frac{\partial}{\partial \theta} p(x; \theta) dx = \frac{\partial}{\partial \theta} \int_{-\infty}^{\infty} p(x; \theta) dx = 0 \quad (75)$$

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

称

$$I(\theta) = E \left(\frac{\partial}{\partial \theta} \ln p(x; \theta) \right)^2 \quad (76)$$

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

$$I(\theta) = -E \left(\frac{\partial^2}{\partial \theta^2} \ln p(x; \theta) \right) \quad (77)$$

定理6.4.3 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

$$g(\theta) = \int_{\mathbb{R}^n} T(X) L(X; \theta) dX \quad (78)$$

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

$$L(X; \theta) = \prod_{k=1}^n p(x_k; \theta) \quad (79)$$

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

$$g'(\theta) = \int_{\mathbb{R}^n} T(X) \frac{\partial}{\partial \theta} L(X; \theta) dX \quad (80)$$

(对于离散总体, 将上述积分号改为求和符号) 那么

$$\text{Var}(T) \geq \frac{(g'(\theta))^2}{nI(\theta)} \quad (81)$$

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。

6.5 Bayes估计

6.5.1 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。

6.5.2 Bayes公式的密度函数形式

- $p(x | \theta)$ 表示随机变量 θ 取给定值时总体的条件概率函数。
- 根据参数 θ 的先验信息确定先验分布 $\pi(\theta)$ 。
- 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 θ_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

$$P(X | \theta_0) = \prod_{k=1}^n p(x_k | \theta) \quad (82)$$

- 由于 θ_0 是设想出来的, 因此需要考虑 $\pi(\theta)$, 那么样本 X 和参数 θ 的联合分布为

$$h(X, \theta) = P(X | \theta) \pi(\theta) \quad (83)$$

- 将 $h(X, \theta)$ 分解为

$$h(X, \theta) = \pi(\theta | X) m(X) \quad (84)$$

其中 $m(X)$ 为 X 的边际概率函数

$$m(X) = \int_{\Theta} h(X, \theta) d\theta = \int_{\Theta} P(X | \theta) \pi(\theta) d\theta \quad (85)$$

进而 θ 的后验分布为

$$\pi(\theta | X) = \frac{h(X, \theta)}{m(X)} = \frac{P(X | \theta) \pi(\theta)}{\int_{\Theta} P(X | \theta) \pi(\theta) d\theta} \quad (86)$$

6.5.3 Bayes估计

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

- **最大后验估计**：后验分布的密度函数的最大值点。
- **后验中位数估计**：后验分布的中位数。
- **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 $\hat{\theta}$ 。

6.5.4 共轭先验分布

定义6.5.1 共轭先验分布：对于总体分布 $p(x; \theta)$ 中的参数 θ ， $\pi(\theta)$ 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 θ 的共轭先验分布（族）。

6.6 区间估计

6.6.1 区间估计的概念

定义6.6.1 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

$$P_{\theta}(\hat{\theta}_L \leq \theta \leq \hat{\theta}_U) \geq 1 - \alpha \quad (87)$$

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

定义6.6.2 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

$$P_{\theta}(\hat{\theta}_L \leq \theta \leq \hat{\theta}_U) = 1 - \alpha \quad (88)$$

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

定义6.6.3 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

$$P_{\theta}(\hat{\theta}_L \leq \theta) \geq 1 - \alpha \quad (89)$$

那么称 $\hat{\theta}_L$ 为 θ 的（单侧）置信下限。

定义6.6.4 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

$$P_{\theta}(\hat{\theta}_U \geq \theta) \geq 1 - \alpha \quad (90)$$

那么称 $\hat{\theta}_U$ 为 θ 的（单侧）置信上限。

6.6.2 枢轴量法

构造枢轴量的方法：

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

$$P(a \leq G \leq b) = 1 - \alpha \quad (91)$$

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

$$P_{\theta}(\hat{\theta}_L \leq \theta \leq \hat{\theta}_U) = 1 - \alpha \quad (92)$$

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

- 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

$$P(G < a) = P(G > b) = \frac{\alpha}{2} \quad (93)$$

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

6.6.3 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
μ	σ 已知	$\frac{\sqrt{n}(\bar{x}-\mu)}{\sigma}$	$N(0,1)$	$\left[\bar{x}-\frac{\sigma}{\sqrt{n}}n_{1-\frac{\alpha}{2}}, \quad \bar{x}+\frac{\sigma}{\sqrt{n}}n_{1-\frac{\alpha}{2}}\right]$
μ	σ 未知	$\frac{\sqrt{n}(\bar{x}-\mu)}{s}$	$T(n-1)$	$\left[\bar{x}-\frac{s}{\sqrt{n}}t_{1-\frac{\alpha}{2}}, \quad \bar{x}+\frac{s}{\sqrt{n}}t_{1-\frac{\alpha}{2}}\right]$
σ^2	μ 未知	$\frac{(n-1)s^2}{\sigma^2}$	$\chi^2(n-1)$	$\left[\frac{(n-1)s^2}{\chi_{1-\frac{\alpha}{2}}^2}, \quad \frac{(n-1)s^2}{\chi_{\frac{\alpha}{2}}^2}\right]$

6.6.4 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, \cdots, x_n ，由中心极限定理

$$N = \frac{\sqrt{n}(\bar{x} - p)}{\sqrt{p(1 - p)}} \sim N(0, 1) \tag{94}$$

因此置信水平为 $1 - \alpha$ 的同等置信区间为

$$\left[\frac{1}{1 + \frac{n_{1-\frac{\alpha}{2}}^2}{n}} \left(\bar{x} + \frac{n_{1-\frac{\alpha}{2}}^2}{2n} - \sqrt{\frac{\bar{x}(1 - \bar{x})}{n} n_{1-\frac{\alpha}{2}}^2 + \left(\frac{n_{1-\frac{\alpha}{2}}^2}{2n} \right)^2} \right), \quad \frac{1}{1 + \frac{n_{1-\frac{\alpha}{2}}^2}{n}} \left(\bar{x} + \frac{n_{1-\frac{\alpha}{2}}^2}{2n} + \sqrt{\frac{\bar{x}(1 - \bar{x})}{n} n_{1-\frac{\alpha}{2}}^2 + \left(\frac{n_{1-\frac{\alpha}{2}}^2}{2n} \right)^2} \right) \right] \tag{95}$$

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

$$\left[\bar{x} - n_{1-\frac{\alpha}{2}} \sqrt{\frac{\bar{x}(1 - \bar{x})}{n}}, \quad \bar{x} + n_{1-\frac{\alpha}{2}} \sqrt{\frac{\bar{x}(1 - \bar{x})}{n}} \right] \tag{96}$$

6.6.5 样本量的确定

置信水平 $1 - \alpha$ 称为**保证概率**。

置信区间的半径（即长度的一半）称为**绝对误差**。

6.6.6 两个正态总体下的置信区间

x_1, \cdots, x_m 是取自 $N(\mu_1, \sigma_1^2)$ 的样本， y_1, \cdots, 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。

目标	条件	枢轴量	分布	置信区间
$\mu_1 - \mu_2$	σ_1^2 和 σ_2^2 已知	$\frac{(\bar{x}-\bar{y})-(\mu_1-\mu_2)}{\sqrt{\frac{\sigma_1^2}{m}+\frac{\sigma_2^2}{n}}}$	$N(0,1)$	$\left[(\bar{x}-\bar{y})-n_{1-\frac{\alpha}{2}}\sqrt{\frac{\sigma_1^2}{m}+\frac{\sigma_2^2}{n}}, \quad (\bar{x}-\bar{y})+n_{1-\frac{\alpha}{2}}\sqrt{\frac{\sigma_1^2}{m}+\frac{\sigma_2^2}{n}}\right]$
$\mu_1 - \mu_2$	$\sigma_1^2 = \sigma_2^2$ 未知	$\frac{(\bar{x}-\bar{y})-(\mu_1-\mu_2)}{s_w\sqrt{\frac{1}{m}+\frac{1}{n}}}$	$T(m+n-2)$	$\left[(\bar{x}-\bar{y})-s_w t_{1-\frac{\alpha}{2}}\sqrt{\frac{1}{m}+\frac{1}{n}}, \quad (\bar{x}-\bar{y})+s_w t_{1-\frac{\alpha}{2}}\sqrt{\frac{1}{m}+\frac{1}{n}}\right]$
$\mu_1 - \mu_2$	$\frac{\sigma_1^2}{\sigma_2^2} = c$ 已知	$\frac{(\bar{x}-\bar{y})-(\mu_1-\mu_2)}{s_w\sqrt{\frac{1}{m}+\frac{c}{n}}}$	$T(m+n-2)$	$\left[(\bar{x}-\bar{y})-s_w t_{1-\frac{\alpha}{2}}\sqrt{\frac{1}{m}+\frac{c}{n}}, \quad (\bar{x}-\bar{y})+s_w t_{1-\frac{\alpha}{2}}\sqrt{\frac{1}{m}+\frac{c}{n}}\right]$
$\mu_1 - \mu_2$	n_1 和 n_2 充分大	$\frac{(\bar{x}-\bar{y})-(\mu_1-\mu_2)}{\sqrt{\frac{s_x^2}{m}+\frac{s_y^2}{n}}}$	$N(0,1)$	$\left[(\bar{x}-\bar{y})-n_{1-\frac{\alpha}{2}}\sqrt{\frac{s_x^2}{m}+\frac{s_y^2}{n}}, \quad (\bar{x}-\bar{y})+n_{1-\frac{\alpha}{2}}\sqrt{\frac{s_x^2}{m}+\frac{s_y^2}{n}}\right]$
$\mu_1 - \mu_2$	一般情况	$\frac{(\bar{x}-\bar{y})-(\mu_1-\mu_2)}{s_0}$	$T(l)$	$\left[(\bar{x}-\bar{y})-s_0 t_{1-\frac{\alpha}{2}}, \quad (\bar{x}-\bar{y})+s_0 t_{1-\frac{\alpha}{2}}\right]$
$\frac{\sigma_1^2}{\sigma_2^2}$	一般情况	$\frac{s_x^2/\sigma_1^2}{s_y^2/\sigma_2^2}$	$F(m-1, n-1)$	$\left[\frac{s_x^2}{s_y^2} \frac{1}{f_{1-\frac{\alpha}{2}}}, \quad \frac{s_x^2}{s_y^2} \frac{1}{f_{\frac{\alpha}{2}}}\right]$

其中

$$s_w^2 = \frac{(m-1)s_x^2 + (n-1)s_y^2}{m+n-2} \tag{97}$$

$$s_{w_c}^2 = \frac{(m-1)s_x^2 + (n-1)\frac{s_y^2}{c}}{m+n-2} \tag{98}$$

$$s_0^2 = \frac{s_x^2}{m} + \frac{s_y^2}{n} \tag{99}$$

$$l = \left\lceil \frac{s_0^4}{\frac{s_x^4}{m^2(m-1)} + \frac{s_y^4}{n^2(n-1)}} \right\rceil \tag{100}$$

$\lceil \cdot \rceil$ 表示最近整数。

第七章：假设检验

7.1 假设检验的基本思想与概念

7.1.1 假设检验问题

假设检验的**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：

- **假设**：两个非空不交参数集合。
- **检验**：通过样本对一个假设作出“对”或“不对”的具体判断规则。
- **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。

7.1.2 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 Θ 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 $H_a : \theta \in \Theta - \Theta_0$ 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

$$H_0 : \theta \in \Theta_0 \quad \text{vs} \quad H_a : \theta \in \Theta - \Theta_0 \quad (101)$$

如果 Θ_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：

$$H_1 : \theta \neq \theta_0, \quad H_2 : \theta < \theta_0, \quad H_3 : \theta > \theta_0 \quad (102)$$

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 \overline{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**， \overline{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

$$\alpha(\theta) = P\{X \in W \mid H_0\}, \quad \theta \in \Theta_0 \quad (103)$$

当 $\theta \in \Theta - \Theta_0$ 时，样本由于随机性却落入了接受域 \overline{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

$$\beta(\theta) = P\{X \in \overline{W} \mid H_a\}, \quad \theta \in \Theta - \Theta_0 \quad (104)$$

定义7.1.1 势函数：对于检验问题

$$H_0 : \theta \in \Theta_0 \quad \text{vs} \quad H_a : \theta \in \Theta - \Theta_0 \quad (105)$$

其拒绝域为 W ，那么定义势函数为

$$\rho(\theta) = P_\theta(X \in W), \quad \theta \in \Theta \quad (106)$$

即

$$\rho(\theta) = \begin{cases} \alpha(\theta), & \theta \in \Theta_0 \\ 1 - \beta(\theta), & \theta \in \Theta - \Theta_0 \end{cases} \quad (107)$$

定义7.1.2 显著性检验：对于检验问题

$$H_0 : \theta \in \Theta_0 \quad \text{vs} \quad H_a : \theta \in \Theta - \Theta_0 \quad (108)$$

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

$$\rho(\theta) \leq \alpha \quad (109)$$

那么称该检验为显著性水平为 α 的显著性检验，简称水平为 α 的检验。

四、给出拒绝域

依据显著性水平 α 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

- 如果 $(x_1, \cdots, x_n) \in W$ ，那么拒绝 H_0 。
- 如果 $(x_1, \cdots, x_n) \in \overline{W}$ ，那么接受 H_0 。

7.1.3 检验的 p 值

定义7.1.3 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。

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

7.2 正态总体参数假设检验

7.2.1 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	σ 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	σ 未知	$\mu \geq \mu_0$	$\mu < \mu_0$	$t = \frac{\sqrt{n}(\bar{x} - \mu_0)}{s}$	$T(n - 1)$	$\{t \leq t_\alpha\}$	$P(T \leq t_0)$
		$\mu = \mu_0$	$\mu \neq \mu_0$			$\{ t \geq t_{1-\frac{\alpha}{2}}\}$	$P(T \geq t_0)$

7.2.2 假设检验与置信区间的关系

7.2.3 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	σ_1, σ_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	$\mu_1 \geq \mu_2$	$\mu_1 < \mu_2$	$t = \frac{\bar{x} - \bar{y}}{s_p \sqrt{\frac{1}{m} + \frac{1}{n}}}$	$T(m + n - 2)$	$\{t \leq t_\alpha\}$	$P(T \leq t_0)$
		$\mu_1 = \mu_2$	$\mu_1 \neq \mu_2$			$\{ t \geq t_{1-\frac{\alpha}{2}}\}$	$P(T \geq t_0)$
		$\mu_1 \leq \mu_2$	$\mu_1 > \mu_2$			$\{u \geq u_{1-\alpha}\}$	$1 - \Phi(u_0)$
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

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

(110)

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

(111)

[.]表示最近整数。

7.2.4 成对数据检验

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

7.2.5 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
		$\sigma^2 \leq \sigma_0^2$	$\sigma^2 > \sigma_0^2$			$\{\chi^2 \geq \chi_{1-\alpha}^2\}$	$P(\chi^2 \geq \chi_0^2)$
χ^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	$\sigma_1^2 \geq \sigma_2^2$	$\sigma_1^2 < \sigma_2^2$	$F = \frac{s_1^2}{s_2^2}$	$F(m-1, n-1)$	$\{F \leq F_\alpha\}$	$P(F \leq F_0)$
		$\sigma_1^2 = \sigma_2^2$	$\sigma_1^2 \neq \sigma_2^2$			$\{F \leq F_{\frac{\alpha}{2}}\} \cup \{F \geq F_{1-\frac{\alpha}{2}}\}$	$2 \min\{P(F \leq F_0), P(F \geq F_0)\}$

7.3 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
		$\lambda \leq \lambda_0$	$\lambda > \lambda_0$			$\{\chi^2 \geq \chi_{1-\alpha}^2\}$	$P(\chi^2 \geq \chi_0^2)$
χ^2 分布	$\text{Exp}(\frac{1}{\lambda})$	$\lambda \geq \lambda_0$	$\lambda < \lambda_0$	$-\frac{2n}{\lambda_0}$	$\chi^2(2n)$	$\{\chi^2 \leq \chi_\alpha^2\}$	$P(\chi^2 \leq \chi_0^2)$
		$\lambda = \lambda_0$	$\lambda \neq \lambda_0$			$\{\chi^2 \leq \chi_{\frac{\alpha}{2}}^2\} \cup \{\chi^2 \geq \chi_{1-\frac{\alpha}{2}}^2\}$	$2 \min\{P(\chi^2 \leq \chi_0^2), P(\chi^2 \geq \chi_0^2)\}$
		$p \leq p_0$	$p > p_0$				$P(x \geq x_0)$
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 $F(x; \theta)$	$\theta \geq \theta_0$	$\theta < \theta_0$	$\frac{\sqrt{n}(\hat{\theta} - \theta_0)}{\sqrt{v'(\theta)}}$	$N(0, 1)$	$\{u \leq u_\alpha\}$	$\Phi(u_0)$
		$\theta = \theta_0$	$\theta \neq \theta_0$			$\{ u \geq u_{1-\frac{\alpha}{2}}\}$	$2(1 - \Phi(u_0))$

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

7.4 似然比检验与分布拟合检验

7.4.1 似然比检验的思想

定义7.4.1 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

$$H_0 : \theta \in \Theta_0 \quad \text{vs} \quad H_a : \theta \in \Theta - \Theta_0 \tag{112}$$

定义改假设检验问题的似然比统计量为

$$\Lambda(x_1, \dots, x_n) = \frac{\sup_{\theta \in \Theta} p(x_1, \dots, x_n; \theta)}{\sup_{\theta \in \Theta_0} p(x_1, \dots, x_n; \theta)} \tag{113}$$

即

$$\Lambda(x_1, \dots, x_n) = \frac{p(x_1, \dots, x_n; \hat{\theta})}{p(x_1, \dots, x_n; \hat{\theta}_0)} \tag{114}$$

其中 $\hat{\theta}$ 和 $\hat{\theta}_0$ 分别为参数空间 Θ 和 Θ_0 上的最大似然估计。

定义7.4.2 似然比检验：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

$$H_0 : \theta \in \Theta_0 \quad \text{vs} \quad H_a : \theta \in \Theta - \Theta_0 \tag{115}$$

其似然比统计量

$$\Lambda(x_1, \dots, x_n) = \frac{p(x_1, \dots, x_n; \hat{\theta})}{p(x_1, \dots, x_n; \hat{\theta}_0)} \tag{116}$$

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 λ_0 满足对于任意 $\theta \in \Theta_0$, 成立

$$P_\theta(\Lambda(x_1, \dots, x_n) \geq \lambda_0) \leq \alpha \tag{117}$$

那么称此检验为显著性水平为 α 的似然比检验，简记为LRT。

7.4.2 分布数据的 χ^2 拟合优度检验

定理7.4.1：总体被分为 r 类 A_1, \dots, A_r , 考虑假设检验

$$H_0: A_k \text{ 所占的比率为 } p_k, \quad k = 1, \dots, r \quad (118)$$

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

$$\chi^2 = \sum_{k=1}^r \frac{(n_k - np_k)^2}{np_k} \quad (119)$$

那么当 H_0 成立时, 成立

$$\chi^2 \xrightarrow{L} \chi^2(r-1) \quad (120)$$

因此对于显著性水平 α , 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

$$\chi^2 = \sum_{k=1}^r \frac{(n_k - n\hat{p}_k)^2}{n\hat{p}_k} \xrightarrow{L} \chi^2(r-s-1) \quad (121)$$

7.4.3 分布的 χ^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

$$H_0: F(x) = F_0(x) \quad (122)$$

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 为至多可数个数 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

$$H_0: A_k \text{ 所占的比率为 } p_k, \quad k = 1, \dots, r \quad (123)$$

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

$$p_k = P(X \in A_k) = F_0(a_k) - F_0(a_{k-1}), \quad k = 1, \dots, r \quad (124)$$

于是原假设转化为

$$H_0: A_k \text{ 所占的比率为 } p_k, \quad k = 1, \dots, r \quad (125)$$

7.4.4 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 和 p_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

$$H_0: p_{ij} = p_i p_j, \quad i = 1, \dots, r; j = 1, \dots, s \quad (126)$$

H_0 成立时 p_{ij} 的最大似然估计为

$$\hat{p}_{ij} = \frac{1}{n^2} \sum_{k=1}^r n_{kj} \sum_{l=1}^s n_{il} \quad (127)$$

那么检验统计量为

$$\chi^2 = \sum_{i,j} \frac{(n_{ij} - n\hat{p}_{ij})^2}{n\hat{p}_{ij}} \xrightarrow{L} \chi^2((r-1)(s-1)) \quad (128)$$

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

7.5 正态性检验

7.5.1 正态概率纸

- 对于给定的样本观测值 x_1, \dots, x_n , 做点

$$\left(x_{(k)}, \frac{k - 0.375}{n + 0.25}\right), \quad k = 1, \dots, n \quad (129)$$

- 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。

7.5.2 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

$$W = \frac{\sum_{k=1}^n (w_k - \bar{w})^2 (x_{(k)} - \bar{x})^2}{\sum_{k=1}^n (w_k - \bar{w})^2 \sum_{k=1}^n (x_{(k)} - \bar{x})^2} = \frac{\sum_{k=1}^{\lfloor \frac{n}{2} \rfloor} w_k^2 (x_{(k)} - x_{(n+1-k)})^2}{\sum_{k=1}^n (x_{(k)} - \bar{x})^2} \quad (130)$$

其中

$$\mathbf{e} = \begin{pmatrix} E\left(\frac{x_{(1)} - \mu}{\sigma}\right) \\ \vdots \\ E\left(\frac{x_{(n)} - \mu}{\sigma}\right) \end{pmatrix}, \quad \mathbf{C} = \begin{pmatrix} \text{Cov}\left(\frac{x_{(1)} - \mu}{\sigma}, \frac{x_{(1)} - \mu}{\sigma}\right) & \cdots & \text{Cov}\left(\frac{x_{(1)} - \mu}{\sigma}, \frac{x_{(n)} - \mu}{\sigma}\right) \\ \vdots & \ddots & \vdots \\ \text{Cov}\left(\frac{x_{(n)} - \mu}{\sigma}, \frac{x_{(1)} - \mu}{\sigma}\right) & \cdots & \text{Cov}\left(\frac{x_{(n)} - \mu}{\sigma}, \frac{x_{(n)} - \mu}{\sigma}\right) \end{pmatrix} \quad (131)$$

$$\begin{pmatrix} w_1 \\ \vdots \\ w_n \end{pmatrix} = \frac{\mathbf{C}^{-1} \mathbf{e}}{\sqrt{\mathbf{e}^T (\mathbf{C}^{-1})^2 \mathbf{e}}} \quad (132)$$

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

7.5.3 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

$$T_{\text{EP}} = 1 + \frac{n}{\sqrt{3}} + \frac{2}{n} \sum_{i=2}^n \sum_{j=1}^{i-1} e^{-\frac{(x_j - x_i)^2}{2 \frac{i-1}{n} s^2}} - \sqrt{2} \sum_{k=1}^n e^{-\frac{(x_k - \bar{x})^2}{4 \frac{n-1}{n} s^2}} \quad (133)$$

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

7.6 非参数检验

7.6.1 游程检验

游程: 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

$$y_k = \begin{cases} 1, & x_k \geq m_e \\ 0, & x_k < m_e \end{cases} \quad (134)$$

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

$$H_0: \text{样本序列符合随机抽取的原则} \quad (135)$$

R 的分布如下

$$P(R = 2k) = \frac{2 \binom{n_1-1}{k-1} \binom{n_2-1}{k-1}}{\binom{n_1+n_2}{n_1}}, \quad k = 1, \dots, \left[\frac{n_1+n_2}{2}\right] \quad (136)$$

$$P(R = 2k+1) = \frac{\binom{n_1-1}{k-1} \binom{n_2-1}{k} + \binom{n_1-1}{k} \binom{n_2-1}{k-1}}{\binom{n_1+n_2}{n_1}}, \quad k = 1, \dots, \left[\frac{n_1+n_2-1}{2}\right] \quad (137)$$

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

7.6.2 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

$$S^+ = \sum_{k=1}^n I_{x_k > x_0} \quad (138)$$

7.6.3 秩和检验

定义7.6.1 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7.6.2 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。

定义7.6.3 符号秩和统计量: 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 θ 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

$$W^+ = \sum_{k=1}^n R_k I_{x_k > 0} \sim W^+(n) \quad (139)$$

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

$$R = (Q_1, \dots, Q_m, R_1, \dots, R_n) \quad (140)$$

那么秩和统计量为

$$W = \sum_{k=1}^n R_k \sim W(m, n) \quad (141)$$

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

第八章：方差分析与回归分析

8.1 方差分析

8.1.1 问题的提出

因子：A

水平： A_1, \dots, A_r

结果： y_{ij} , 其中 $i = 1, \dots, r$

8.1.2 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:

- 每一个总体均为正态分布, 记为 $N(\mu_k, \sigma_k^2)$, 其中 $k = 1, \dots, r$ 。
- 各总体的方差相同, 记为 $\sigma_1^2 = \dots = \sigma_r^2 = \sigma^2$ 。
- 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:

$$H_0: \mu_1 = \dots = \mu_r \quad \text{vs} \quad H_a: \mu_1, \dots, \mu_r \text{不全相等} \quad (142)$$

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

$$\varepsilon_{ij} = y_{ij} - \mu_i \quad (143)$$

那么试验结果 y_{ij} 的数据结构式为

$$y_{ij} = \mu_i + \varepsilon_{ij} \quad (144)$$

单因子方差分析的统计模型为

$$\begin{cases} y_{ij} = \mu_i + \varepsilon_{ij} \\ \varepsilon_{ij} \text{相互独立} \\ \varepsilon_{ij} \sim N(0, \sigma^2) \end{cases} \quad (145)$$

总均值

$$\mu = \frac{1}{r} \sum_{k=1}^r \mu_k \quad (146)$$

因子A的第 k 个水平的主效应

$$a_i = \mu_i - \mu \quad (147)$$

容易知道

$$\sum_{k=1}^r a_k = 0 \quad (148)$$

$$\mu_k = \mu + a_k \quad (149)$$

于是统计模型改写为

$$\begin{cases} y_{ij} = \mu + a_i + \varepsilon_{ij} \\ \sum_{i=1}^r a_i = 0 \\ \varepsilon_{ij} \text{相互独立} \\ \varepsilon_{ij} \sim N(0, \sigma^2) \end{cases} \quad (150)$$

统计假设改写为

$$H_0: a_1 = \dots = a_r = 0 \quad \text{vs} \quad H_a: a_1, \dots, a_r \text{不全为0} \quad (151)$$

8.1.3 平方和分解

一、实验数据

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

$$y_{ij} - \bar{y} = (y_{ij} - \bar{y}_i) + (\bar{y}_i - \bar{y}) \quad (152)$$

$$\bar{\varepsilon}_i = \frac{1}{m} \sum_{j=1}^m \varepsilon_{ij} \quad (153)$$

$$\bar{\varepsilon} = \frac{1}{r} \sum_{j=1}^m \bar{\varepsilon}_i = \frac{1}{n} \sum_{i,j} \varepsilon_{ij} \quad (154)$$

组内偏差为

$$y_{ij} - \bar{y}_i = \varepsilon_{ij} - \bar{\varepsilon}_i \quad (155)$$

组间偏差为

$$\bar{y}_i - \bar{y} = a_i + \bar{\varepsilon}_i - \bar{\varepsilon} \quad (156)$$

三、偏差平方和及其自由度

偏差平方和

$$Q = \sum_{k=1}^n (y_k - \bar{y})^2 \quad (157)$$

自由度

$$f_Q = n - 1 \quad (158)$$

四、总平方和分解公式

总偏差平方和

$$S_T = \sum_{i,j} (y_{ij} - \bar{y})^2 = \sum_{i,j} y_{ij}^2 - \frac{T^2}{n}, \quad f_T = n - 1 \quad (159)$$

组内偏差平方和（因子A的偏差平方和）

$$S_A = m \sum_{i=1}^r (\bar{y}_i - \bar{y})^2 = \frac{1}{m} \sum_{i=1}^r T_i^2 - \frac{T^2}{n}, \quad f_A = r - 1 \quad (160)$$

组内偏差平方和（误差偏差平方和）

$$S_e = \sum_{i,j} (y_{ij} - \bar{y}_i)^2 = \sum_{i,j} y_{ij}^2 - \frac{1}{m} \sum_{i=1}^r T_i^2, \quad f_e = n - r \quad (161)$$

定理8.1.1 总平方和分解式：

$$S_T = S_A + S_e \quad (162)$$

8.1.4 检验方法

均方

$$MS = \frac{Q}{f_Q} \tag{163}$$

因子均方和误差均方

$$MS_A = \frac{S_A}{f_A}, \quad MS_e = \frac{S_e}{f_e} \tag{164}$$

定理8.1.2:

•

$$\frac{S_e}{\sigma^2} \sim \chi^2(n-r), \quad E(S_e) = (n-r)\sigma^2 \tag{165}$$

•

$$E(S_A) = (r-1)\sigma^2 + m \sum_{i=1}^r a_i^2 \tag{166}$$

若 H_0 成立, 那么

$$\frac{S_A}{\sigma^2} \sim \chi^2(r-1), \quad E(S_e) = (r-1)\sigma^2 \tag{167}$$

- S_A 与 S_e 相互独立。

检验统计量

$$F = \frac{MS_A}{MS_e} \sim F(r-1, n-r) \tag{168}$$

拒绝域

$$W = \{F \geq F_{1-\alpha}\} \tag{169}$$

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

检验的 p 值为

$$p = P(F \geq F_0) \tag{170}$$

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	$f_T = n - 1$			

8.1.5 参数估计

一、点估计

•

$$y_{ij} \sim N(\mu + a_i, \sigma^2) \tag{171}$$

- μ 的最大似然估计为

$$\hat{\mu} = \bar{y} \tag{172}$$

- a_i 的最大似然估计为

$$\hat{a}_i = \bar{y}_i - \bar{y} \tag{173}$$

- σ^2 的最大似然估计为

$$\hat{\sigma}^2 = MS_e \tag{174}$$

二、置信区间

由于

$$\bar{y}_i \sim N(\mu_i, \frac{\sigma^2}{m}), \quad \frac{S_e}{\sigma^2} \sim \chi^2(n-r) \quad (175)$$

因此

$$\frac{\sqrt{m}(\bar{y}_i - \mu_i)}{\sqrt{\hat{\sigma}^2}} \sim T(f_e) \quad (176)$$

进而 μ_i 的 $1 - \alpha$ 的置信区间为

$$\left[\bar{y}_i - t_{1-\frac{\alpha}{2}} \frac{\hat{\sigma}}{\sqrt{m}}, \quad \bar{y}_i + t_{1-\frac{\alpha}{2}} \frac{\hat{\sigma}}{\sqrt{m}} \right] \quad (177)$$

8.1.6 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

$$\begin{cases} y_{ij} = \mu_i + \varepsilon_{ij} \\ \varepsilon_{ij} \text{相互独立} \\ \varepsilon_{ij} \sim N(0, \sigma^2) \end{cases} \quad (178)$$

二、总均值

加权均值:

$$\mu = \frac{1}{n} \sum_{i=1}^r m_i \mu_i \quad (179)$$

水平效应:

$$a_i = \mu_i - \mu \quad (180)$$

统计模型为

$$\begin{cases} y_{ij} = \mu + a_i + \varepsilon_{ij} \\ \sum_{i=1}^r m_i a_i = 0 \\ \varepsilon_{ij} \text{相互独立} \\ \varepsilon_{ij} \sim N(0, \sigma^2) \end{cases} \quad (181)$$

四、各平方和的计算

$$T_i = \sum_{j=1}^{m_i} y_{ij}, \quad \bar{y}_i = \frac{T_i}{m_i} \quad (182)$$

$$T = \sum_{i,j} y_{ij} = \sum_{i=1}^r T_i, \quad \bar{y} = \frac{T}{n} \quad (183)$$

$$S_T = \sum_{i,j} (y_{ij} - \bar{y})^2 = \sum_{i,j} y_{ij}^2 - \frac{T^2}{n}, \quad f_T = n - 1 \quad (184)$$

$$S_A = \sum_{i=1}^r m_i (\bar{y}_i - \bar{y})^2 = \sum_{i=1}^r \frac{T_i^2}{m_i} - \frac{T^2}{n}, \quad f_A = r - 1 \quad (185)$$

$$S_e = \sum_{i,j} (y_{ij} - \bar{y}_i)^2 = \sum_{i,j} y_{ij}^2 - \sum_{i=1}^r \frac{T_i^2}{m_i}, \quad f_e = n - r \quad (186)$$

8.2 多重比较

8.2.1 水平均值差的置信区间

检验问题

$$H_0: \mu_i - \mu_j = 0 \quad \text{vs} \quad H_a: \mu_i - \mu_j \neq 0 \quad (187)$$

由于

$$\bar{y}_i - \bar{y}_j \sim N\left(\mu_i - \mu_j, \left(\frac{1}{m_i} + \frac{1}{m_j}\right)\sigma^2\right) \quad (188)$$

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

$$\frac{(\bar{y}_i - \bar{y}_j) - (\mu_i - \mu_j)}{\sqrt{\left(\frac{1}{m_i} + \frac{1}{m_j}\right)\hat{\sigma}^2}} \sim T(n-r) \quad (189)$$

那么置信水平为 $1 - \alpha$ 的置信区间为

$$\left[\bar{y}_i - \bar{y}_j - t_{1-\frac{\alpha}{2}} \sqrt{\left(\frac{1}{m_i} + \frac{1}{m_j}\right)\hat{\sigma}^2}, \quad \bar{y}_i - \bar{y}_j + t_{1-\frac{\alpha}{2}} \sqrt{\left(\frac{1}{m_i} + \frac{1}{m_j}\right)\hat{\sigma}^2} \right] \quad (190)$$

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 μ_i 和 μ_j 无显著差异; 反之拒绝原假设, 认为 μ_i 和 μ_j 存在显著差异。

8.2.2 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

$$H_0^{ij} : \mu_i = \mu_j, \quad 1 \leq i < j \leq r \quad (191)$$

拒绝域

$$W = \bigcup_{1 \leq i < j \leq r} \{|\bar{y}_i - \bar{y}_j| \geq c_{ij}\} \quad (192)$$

8.2.3 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

$$\frac{\sqrt{m}(\bar{y}_i - \mu_i)}{\hat{\sigma}} \sim T(n-r) \quad (193)$$

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

$$P(W) = P\left(q(r, n-r) \geq \frac{c\sqrt{m}}{\hat{\sigma}}\right) \quad (194)$$

其中 t 化极差统计量为

$$q(r, n-r) = \max_{1 \leq i \leq r} \frac{\sqrt{m}(\bar{y}_i - \mu_i)}{\hat{\sigma}} - \min_{1 \leq j \leq r} \frac{\sqrt{m}(\bar{y}_j - \mu_j)}{\hat{\sigma}} \quad (195)$$

仅与 n 和 r 有关。由 $P(W) = \alpha$, 可知

$$c = q_{1-\alpha} \frac{\hat{\sigma}}{\sqrt{m}} \quad (196)$$

因此, 如果

$$|\bar{y}_i - \bar{y}_j| \geq q_{1-\alpha} \frac{\hat{\sigma}}{\sqrt{m}} \quad (197)$$

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。

8.2.4 重复数不等场合的S法

由于

$$\frac{(\bar{y}_i - \bar{y}_j) - (\mu_i - \mu_j)}{\sqrt{\left(\frac{1}{m_i} + \frac{1}{m_j}\right)\hat{\sigma}^2}} \sim T(n-r) \quad (198)$$

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

$$\frac{(\bar{y}_i - \bar{y}_j)^2}{\left(\frac{1}{m_i} + \frac{1}{m_j}\right)\hat{\sigma}^2} \sim F(1, n-r) \quad (199)$$

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

$$P(W) = P\left(\max_{1 \leq i < j \leq r} \frac{(\bar{y}_i - \bar{y}_j)^2}{\left(\frac{1}{m_i} + \frac{1}{m_j}\right)\hat{\sigma}^2} \geq \frac{c^2}{\hat{\sigma}^2}\right) \quad (200)$$

其中

$$\frac{\max_{1 \leq i < j \leq r} \frac{(\bar{y}_i - \bar{y}_j)^2}{\left(\frac{1}{m_i} + \frac{1}{m_j}\right)\hat{\sigma}^2}}{r-1} \sim F(r-1, n-r) \quad (201)$$

由 $P(W) = \alpha$, 可知

$$\frac{c^2}{\hat{\sigma}^2} = (r-1)f_{1-\alpha} \quad (202)$$

即

$$c_{ij} = \sqrt{(r-1)f_{1-\alpha}\hat{\sigma}^2\left(\frac{1}{m_i} + \frac{1}{m_j}\right)} \quad (203)$$

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

$$|\bar{y}_i - \bar{y}_j| \geq \sqrt{(r-1)f_{1-\alpha}\hat{\sigma}^2\left(\frac{1}{m_i} + \frac{1}{m_j}\right)} \quad (204)$$

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。

8.3 方差齐性检验

方差齐性检验:

$$H_0: \sigma_1^2 = \cdots = \sigma_r^2 \quad (205)$$

8.3.1 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

$$s_i^2 = \frac{1}{m-1} \sum_{k=1}^m (y_{ik} - \bar{y}_i)^2 \quad (206)$$

检验统计量为

$$H = \frac{\max\{s_i^2\}}{\min\{s_i^2\}} \sim H(r, m-1) \quad (207)$$

拒绝域为

$$W = \{H \geq H_{1-\alpha}\} \quad (208)$$

8.3.2 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

$$s_i^2 = \frac{1}{m_i-1} \sum_{k=1}^{m_i} (y_{ik} - \bar{y}_i)^2 = \frac{Q_i}{f_i} \quad (209)$$

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

$$MS_e = \frac{1}{f_e} \sum_{i=1}^r Q_i = \sum_{i=1}^r \frac{f_i}{f_e} s_i^2 \quad (210)$$

其加权几何平均为

$$GMS_e = \left(\prod_{i=1}^r (s_i^2)^{f_i} \right)^{\frac{1}{f_e}} \quad (211)$$

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

$$\text{MS}_e \geq \text{GMS}_e \quad (212)$$

当且仅当 $s_1^2 = \cdots = s_r^2$ 时等号成立。而

$$B = \frac{f_e}{C} \ln \frac{\text{MS}_e}{\text{GMS}_e} \sim \chi^2(r-1) \quad (213)$$

其中

$$C = 1 + \frac{1}{3(r-1)} \left(\sum_{i=1}^r \frac{1}{f_i} - \frac{1}{f_e} \right) \quad (214)$$

因此拒绝域为

$$W = \{B \geq \chi_{1-\alpha}^2\} \quad (215)$$

8.3.3 修正的Bartlett检验

修正的检验统计量

$$B' = \frac{f_2 BC}{f_1(A - BC)} \sim F(f_1, f_2) \quad (216)$$

其中

$$f_1 = r - 1, \quad f_2 = \frac{r+1}{(C-1)^2}, \quad A = \frac{f_2}{2 - C + \frac{2}{f_2^0}} \quad (217)$$

拒绝域为

$$W = \{B' \geq F_{1-\alpha}\} \quad (218)$$

8.4 一元线性回归

8.4.1 变量间的两类关系

确定性关系

相关关系

8.4.2 一元线性回归模型

第一类回归问题

$$f(x) = E(Y | x) = \int_{-\infty}^{\infty} yp(y | x) dx \quad (219)$$

第二类回归问题

$$y = f(x) + \varepsilon \quad (220)$$

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

一元回归模型：

$$\begin{cases} y_i = \beta_0 + \beta x_i + \varepsilon_i \\ \varepsilon_i \text{相互独立} \\ \varepsilon_i \sim N(0, \sigma^2) \end{cases} \quad (221)$$

由数据 (x_i, y_i) 得到的 β_0 和 β 的估计 $\hat{\beta}_0$ 和 $\hat{\beta}$, 称

$$\hat{y} = \hat{\beta}_0 + \hat{\beta}x \quad (222)$$

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。

8.4.3 回归系数的最小二乘估计

$$l_{xy} = \sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y}) = \sum_{i=1}^n x_i y_i - \frac{1}{n} \sum_{i=1}^n x_i \sum_{i=1}^n y_i = \sum_{i=1}^n x_i y_i - n\bar{x}\bar{y} \quad (223)$$

$$l_{xx} = \sum_{i=1}^n (x_i - \bar{x})^2 = \sum_{i=1}^n x_i^2 - \frac{1}{n} \left(\sum_{i=1}^n x_i \right)^2 = \sum_{i=1}^n x_i^2 - n\bar{x}^2 \quad (224)$$

$$l_{yy} = \sum_{i=1}^n (y_i - \bar{y})^2 = \sum_{i=1}^n x_i y_i - \frac{1}{n} \left(\sum_{i=1}^n y_i \right)^2 = \sum_{i=1}^n y_i^2 - n\bar{y}^2 \quad (225)$$

β_0 和 β 的最小二乘估计 (LSE) $\hat{\beta}_0$ 和 $\hat{\beta}$ 为

$$\hat{\beta} = \frac{l_{xy}}{l_{xx}} \quad (226)$$

$$\hat{\beta}_0 = \bar{y} - \hat{\beta}\bar{x} \quad (227)$$

定理8.4.1: 在如下模型下, 成立

$$\begin{cases} y_i = \beta_0 + \beta x_i + \varepsilon_i \\ \varepsilon_i \text{相互独立} \\ \varepsilon_i \sim N(0, \sigma^2) \end{cases} \quad (228)$$

•
$$\hat{\beta}_0 \sim N\left(\beta_0, \left(\frac{1}{n} + \frac{\bar{x}^2}{l_{xx}}\right)\sigma^2\right), \quad \hat{\beta} \sim N\left(\beta, \frac{\sigma^2}{l_{xx}}\right) \quad (229)$$

•
$$\text{Cov}(\hat{\beta}_0, \hat{\beta}) = -\frac{\bar{x}}{l_{xx}}\sigma^2 \quad (230)$$

• 给定 x_0

$$\hat{y}_0 = \hat{\beta}_0 + \hat{\beta}x_0 \sim N\left(\beta_0 + \beta x_0, \left(\frac{1}{n} + \frac{(\bar{x} - x_0)^2}{l_{xx}}\right)\sigma^2\right) \quad (231)$$

8.4.4 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:

$$H_0: \beta = 0 \quad \text{vs} \quad H_a: \beta \neq 0 \quad (232)$$

方差分析表

来源	平方和	自由度	均方	F比
回归	S_R	$f_R = 1$	$MS_R = \frac{S_R}{f_R}$	$F = \frac{MS_R}{MS_e}$
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:

$$S_T = \sum_{i=1}^n (y_i - \bar{y})^2 = l_{yy} \quad (233)$$

回归平方和:

$$S_R = \sum_{i=1}^n (\hat{y}_i - \bar{y})^2 = \frac{l_{xy}^2}{l_{xx}} \quad (234)$$

残差平方和:

$$S_e = \sum_{i=1}^n (y_i - \hat{y}_i)^2 \quad (235)$$

平方和分解：

$$S_T = S_R + S_e \quad (236)$$

定理8.4.2：在如下模型下，成立

$$\begin{cases} y_i = \beta_0 + \beta x_i + \varepsilon_i \\ \varepsilon_i \text{相互独立} \\ \varepsilon_i \sim N(0, \sigma^2) \end{cases} \quad (237)$$

$$\bullet \quad E(S_R) = \sigma^2 + \hat{\beta}l_{xx}, \quad E(S_e) = (n-2)\sigma^2 \quad (238)$$

$$\bullet \quad \frac{S_e}{\sigma^2} \sim \chi^2(n-2) \quad (239)$$

• 如果 H_0 成立，那么

$$\frac{S_R}{\sigma^2} \sim \chi^2(1) \quad (240)$$

• S_R 与 S_e 、 \bar{y} 独立。

统计检验量：

$$F = \frac{(n-2)S_R}{S_e} \sim F(1, n-2) \quad (241)$$

拒绝域为 $W = \{F \geq F_{1-\alpha}\}$ 。

二、 T 检验

检验统计量：

$$T = \frac{\sqrt{(n-2)l_{xx}}\hat{\beta}}{\sqrt{S_e}} \sim T(n-2) \quad (242)$$

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：

$$H_0 : \rho = 0 \quad \text{vs} \quad H_a : \rho \neq 0 \quad (243)$$

检验统计量：相关系数

$$r = \frac{l_{xy}}{\sqrt{l_{xx}l_{yy}}} = \sqrt{\frac{F}{F + (n-2)}} \sim r(n-2) = \sqrt{\frac{F(1, n-2)}{F(1, n-2) + (n-2)}} \quad (244)$$

- $|r| = 1$: (x_i, y_i) 共线。
- $r > 0$: (x_i, y_i) 正相关。
- $r < 0$: (x_i, y_i) 负相关。
- $r = 0$: (x_i, y_i) 不相关。

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

$$r_{1-\alpha} = \sqrt{\frac{F_{1-\alpha}}{F_{1-\alpha} + (n-2)}} \quad (245)$$

8.4.5 估计与预测

一、 $E(y_0)$ 的置信区间

枢轴量为

$$\frac{\hat{y}_0 - E(y_0)}{\sqrt{\frac{S_e}{n-2}} \sqrt{\frac{1}{n} + \frac{(x_0 - \bar{x})^2}{l_{xx}}}} \sim T(n-2) \quad (246)$$

$1 - \alpha$ 的置信区间为

$$\left[\hat{y}_0 - t_{1-\frac{\alpha}{2}} \sqrt{\frac{S_e}{n-2}} \sqrt{\frac{1}{n} + \frac{(x_0 - \bar{x})^2}{l_{xx}}}, \quad \hat{y}_0 + t_{1-\frac{\alpha}{2}} \sqrt{\frac{S_e}{n-2}} \sqrt{\frac{1}{n} + \frac{(x_0 - \bar{x})^2}{l_{xx}}} \right] \quad (247)$$

二、 y_0 的预测区间

枢轴量为

$$\frac{y_0 - \hat{y}_0}{\sqrt{\frac{S_e}{n-2}} \sqrt{1 + \frac{1}{n} + \frac{(x_0 - \bar{x})^2}{l_{xx}}}} \sim T(n-2) \quad (248)$$

预测区间为

$$\left[\hat{y}_0 - t_{1-\frac{\alpha}{2}} \sqrt{\frac{S_e}{n-2}} \sqrt{1 + \frac{1}{n} + \frac{(x_0 - \bar{x})^2}{l_{xx}}}, \quad \hat{y}_0 + t_{1-\frac{\alpha}{2}} \sqrt{\frac{S_e}{n-2}} \sqrt{1 + \frac{1}{n} + \frac{(x_0 - \bar{x})^2}{l_{xx}}} \right] \quad (249)$$

8.5 一元非线性回归

8.5.1 确定可能的函数形式

8.5.2 参数估计

8.5.3 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。

$$R^2 = 1 - \frac{\sum (y_i - \hat{y}_i)^2}{\sum (y_i - \bar{y})^2} \quad (250)$$

剩余标准差：越小，回归曲线拟合越好。

$$s = \sqrt{\frac{\sum (y_i - \hat{y}_i)^2}{n-2}} \quad (251)$$

附录：概率模型

概率模型	概率分布 $p(x)$	数学期望 $E\zeta$	方差 $D\zeta$	特征函数 $f(t)$
退化分布 $L(x)$	$p(x) = \begin{cases} 1, & x = c \\ 0, & x \neq c \end{cases}$	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 $B(n, p)$	$b(k; n, p) = \binom{n}{k} p^k (1-p)^{n-k}$ $k = 0, \cdots, n; 0 < p < 1$	np	$np(1-p)$	$(pe^{it} + 1 - p)^n$
Poisson分布 $P(\lambda)$	$p(k; \lambda) = \frac{\lambda^k}{k!} e^{-\lambda}$ $k = 0, 1, \cdots; \lambda > 0$	λ	λ	$e^{\lambda(e^{it}-1)}$
几何分布	$g(k; p) = p(1-p)^{k-1}$ $k = 1, 2, \cdots; 0 < p < 1$	$\frac{1}{p}$	$\frac{1}{p^2}$	$\frac{pe^{it}}{1-(1-p)e^{it}}$
超几何分布	$p_k = \frac{\binom{M}{k} \binom{N-M}{n-k}}{\binom{N}{n}}$ $M, n \leq N; M, N, n \in N^+$ $k = 0, \cdots, \min(M, n)$	$\frac{nM}{N}$	$\frac{nM(N-M)(N-n)}{N^2(N-1)}$	$\sum_{k=0}^n \frac{\binom{M}{k} \binom{N-M}{n-k}}{\binom{N}{n}} e^{ikt}$
Pascal分布	$p_k = \binom{k-1}{r-1} p^r (1-p)^{k-r}$ $k = r, r+1, \cdots; 0 < p < 1; r \in N^+$	$\frac{r}{p}$	$\frac{r(1-p)}{p^2}$	$\left(\frac{(1-p)e^{it}}{1-(1-p)e^{it}}\right)^r$
负二项分布	$p_k = \binom{r-1}{k} p^r (p-1)^k$ $k = 0, 1, \cdots; 0 < p < 1; r > 0$	$\frac{r(1-p)}{p}$	$\frac{r(1-p)}{p^2}$	$\left(\frac{1-(1-p)e^{it}}{1-(1-p)e^{it}}\right)^r$
正态分布 $N(\mu, \sigma^2)$	$p(x) = \frac{1}{\sqrt{2\pi\sigma}} e^{-\frac{(x-\mu)^2}{2\sigma^2}}$ $-\infty < x < \infty$	μ	σ^2	$e^{i\mu t - \frac{1}{2}\sigma^2 t^2}$
均匀分布 $U[a, b]$	$p(x) = \begin{cases} \frac{1}{b-a}, & a \leq x \leq b \\ 0, & \text{其他} \end{cases}$ $a < b$	$\frac{a+b}{2}$	$\frac{(b-a)^2}{12}$	$\frac{e^{ibt}-e^{iat}}{ib(a-b)}$
指数分布 $\text{Exp}(\lambda)$	$p(x) = \begin{cases} \lambda e^{-\lambda x}, & x \geq 0 \\ 0, & x < 0 \end{cases}$ $\lambda > 0$	λ^{-1}	λ^{-2}	$(1 - \frac{it}{\lambda})^{-1}$
χ^2 分布	$p(x) = \begin{cases} \frac{1}{2^{n/2}\Gamma(\frac{n}{2})} x^{\frac{n}{2}-1} e^{-\frac{x}{2}}, & x \geq 0 \\ 0, & x < 0 \end{cases}$ $n \in N^+$	n	$2n$	$(1 - 2it)^{-\frac{n}{2}}$
Γ 分布 $\Gamma(r, \lambda)$ ($r \in N^+$ 即为Erlang分布)	$p(x) = \begin{cases} \frac{\lambda^r}{\Gamma(r)} x^{r-1} e^{-\lambda x}, & x \geq 0 \\ 0, & x < 0 \end{cases}$ $r, \lambda > 0$	$\frac{r}{\lambda}$	$\frac{r}{\lambda^2}$	$(1 - \frac{it}{\lambda})^{-r}$
Cauchy分布	$p(x) = \frac{1}{\pi} \frac{\lambda}{\lambda^2 + (x-\mu)^2}$ $-\infty < x, \mu < \infty; \lambda > 0$	不存在	不存在	$e^{i\mu t - \lambda t }$
t 分布	$p(x) = \frac{\Gamma(\frac{n+1}{2})}{\sqrt{n\pi} \Gamma(\frac{n}{2})} (1 + \frac{x^2}{n})^{-\frac{n+1}{2}}$ $-\infty < x < \infty; n \in N^+$	$0(n > 1)$	$-\frac{n}{n-2}(n > 2)$	
Pareto分布	$p(x) = \begin{cases} rA^r x^{-r-1}, & x \geq A \\ 0, & x < A \end{cases}$ $r, A > 0$	($r > 1$ 时存在)	($r > 2$ 时存在)	
F 分布	$p(x) = \begin{cases} \frac{\Gamma(\frac{m+n}{2})}{\Gamma(\frac{m}{2})\Gamma(\frac{n}{2})} m^n n^{\frac{m}{2}} \frac{x^{\frac{m}{2}-1}}{(n+mx)^{\frac{m+n}{2}}}, & x \geq 0 \\ 0, & x < 0 \end{cases}$ $m, n \in N^+$	$\frac{n}{m-2}(n > 2)$	$\frac{2n^2(m+n-2)}{m(m-2)(n-4)}(n > 4)$	
β 分布	$p(x) = \begin{cases} \frac{\Gamma(p+q)}{\Gamma(p)\Gamma(q)} x^{p-1}(1-x)^{q-1}, & 0 < x < 1 \\ 0, & \text{其他} \end{cases}$ $p, q > 0$	$\frac{p}{p+q}$	$\frac{pq}{(p+q)^2(p+q+1)}$	$\frac{\Gamma(p+q)}{\Gamma(p)\Gamma(q)} \sum_{k=0}^{\infty} \frac{\Gamma(p-1)(q)^k}{\Gamma(p+k)\Gamma(q+1-k)}$
对数正态分布	$p(x) = \begin{cases} -\frac{1}{\sqrt{2\pi\sigma^2}} \frac{\ln x + \frac{\sigma^2}{2}}{\sigma^2 x}, & x > 0 \\ 0, & x \leq 0 \end{cases}$ $\sigma, \sigma > 0$	$e^{1+\frac{\sigma^2}{2}}$	$e^{2\sigma^2+1}(e^{\sigma^2}-1)$	
Weibull分布	$p(x) = \begin{cases} \alpha \lambda x^{\alpha-1} e^{-\lambda x^\alpha}, & x > 0 \\ 0, & x \leq 0 \end{cases}$ $\lambda, \alpha > 0$	$\Gamma(\frac{1}{\alpha} + 1) \lambda^{-\frac{1}{\alpha}}$	$\lambda^{-\frac{1}{\alpha}} (\Gamma(\frac{1}{\alpha} + 1) - (\Gamma(\frac{1}{\alpha} + 1))^2)$	
Rayleigh分布	$p(x) = \begin{cases} 2xe^{-x^2}, & x \geq 0 \\ 0, & x < 0 \end{cases}$	$\sqrt{\frac{\pi}{2}}$	$2 - \frac{\pi}{2}$	
Laplace分布	$p(x) = \frac{1}{2\alpha} e^{-\frac{ x }{\alpha}}$ $-\infty < x < \infty; \alpha > 0$	0	$2\alpha^2$	$\frac{1}{1+\alpha^2 t^2}$