

AI 中的数学

第二十四讲

方聪，概率统计部分参考章复熹和张原老师课件

2024 年秋季

① 概率论分册

② 参数估计

③ 假设检验

④ 回归分析

2. 设 X_1, X_2, \dots 是随机变量序列, 且 X_n 取值 0 或 n^2 ,
 $P(X_n = n^2) = \frac{1}{n^2} = 1 - P(X_n = 0)$ ($n = 1, 2, \dots$), 易知
 $E(X_n) = 1$ (一切 n)。试证: 随机变量序列 X_1, X_2, \dots 不服从大数律。

2. 设 X_1, X_2, \dots 是随机变量序列, 且 X_n 取值 0 或 n^2 , $P(X_n = n^2) = \frac{1}{n^2} = 1 - P(X_n = 0)$ ($n = 1, 2, \dots$), 易知 $E(X_n) = 1$ (一切 n)。试证: 随机变量序列 X_1, X_2, \dots 不服从大数律。

解: 事件 $D_n: \forall 2 \leq i \leq n, X_i = 0$, 其概率

$$P(D_n) = \prod_{i=2}^n \left(1 - \frac{1}{i^2}\right) = \prod_{i=2}^n \frac{(i-1)(i+1)}{i^2} = \frac{1}{2} \left(1 + \frac{1}{n}\right)$$

这意味着均值在 0 附近的概率接近 0.5, 故它不依概率收敛于期望。若 $n \rightarrow \infty$,

$$\lim_{n \rightarrow \infty} P\left(\frac{S_n - E(S_n)}{n} = -1\right) \geq \lim_{n \rightarrow \infty} P(D_n) = \frac{1}{2}$$

可见 $\frac{S_n - E(S_n)}{n}$ 不依概率收敛于 0, 不服从大数定律。

3. 设随机变量序列 ξ_1, ξ_2, \dots 满足

$$\xi_n \xrightarrow{w} 0 \quad (n \rightarrow \infty),$$

试证:

$$\xi_n \xrightarrow{P} 0 \quad (n \rightarrow \infty).$$

3. 设随机变量序列 ξ_1, ξ_2, \dots 满足

$$\xi_n \xrightarrow{w} 0 \quad (n \rightarrow \infty),$$

试证：

$$\xi_n \xrightarrow{P} 0 \quad (n \rightarrow \infty).$$

证明：记 $F_n(x) = P(\xi_n < x)$ ，由于 $\xi_n \xrightarrow{w} 0 \quad (n \rightarrow \infty)$ ，令 $\varepsilon > 0$ ，则 $\lim_{n \rightarrow \infty} F_n(\varepsilon) = 0$ ， $\lim_{n \rightarrow \infty} F_n(-\varepsilon) = 1$ 。

$$\begin{aligned} \lim_{n \rightarrow \infty} P(|\xi_n| > \varepsilon) &= \lim_{n \rightarrow \infty} P(\xi_n > \varepsilon) + \lim_{n \rightarrow \infty} P(\xi_n < -\varepsilon) \\ &= \lim_{n \rightarrow \infty} F_n(\varepsilon) + 1 - \lim_{n \rightarrow \infty} F_n(-\varepsilon) = 0 \end{aligned}$$

因此 $\xi_n \xrightarrow{P} 0$ 。

6. 设 X_1, X_2, \dots 是独立同分布的随机变量序列, 共同分布是区间 $[0, a]$ 上的均匀分布 ($a > 0$), $\xi_n = \max\{X_1, \dots, X_n\}$ ($n = 1, 2, \dots$), 试证:

$$\xi_n \xrightarrow{P} a \quad (n \rightarrow \infty).$$

证明: $\xi_n = \max\{X_1, X_2, \dots, X_n\}$ 的分布函数为:

$$F_{\xi_n}(x) = P(\max\{X_1, X_2, \dots, X_n\} \leq x) = P(X_1 \leq x) \cdots P(X_n \leq x) = \left(\frac{x}{a}\right)^n$$

对于任意 $\epsilon > 0$, 当 $n \rightarrow \infty$ 时, 我们有:

$$P(\xi_n > a - \epsilon) = 1 - \left(\frac{a - \epsilon}{a}\right)^n \rightarrow 1, \quad P(\xi_n < a + \epsilon) = 1$$

$$P(a - \epsilon < \xi_n < a + \epsilon) \rightarrow 1$$

因此

$$\xi_n \xrightarrow{P} a \quad (n \rightarrow \infty).$$

9. 试证下列条件对应的各个相互独立的随机变量序列服从大数律:

$$(1) P(X_k = \sqrt{\ln k}) = P(X_k = -\sqrt{\ln k}) = \frac{1}{2} \quad (k = 2, 3, \cdots);$$

$$(2) P(X_k = \frac{2^n}{n^2}) = \frac{1}{2^n} \quad (k = 1, 2, \cdots; n = 1, 2, \cdots);$$

$$(3) P(X_k = n) = \frac{c}{n^2 \ln^2 n} \quad (k = 1, 2, \cdots; n = 2, 3, \cdots), \text{ 其中}$$

$$c = \left(\sum_{n=2}^{\infty} \frac{1}{n^2 \ln^2 n} \right)^{-1}.$$

证明: (1) 随机变量的期望和方差为 $E(X_k) = 0, \text{var}(X_k) = \ln k$, 序列和的方差为 $\text{var}(S_n) = \sum_{k=2}^n \ln k$, 由切比雪夫不等式,

$$P\left(\left|\frac{S_n}{n}\right| \geq \varepsilon\right) \leq \frac{1}{\varepsilon^2} \frac{\sum_{k=2}^n \ln k}{n^2} < \frac{\ln n}{\varepsilon^2 n}$$

对于任意 $\varepsilon > 0$, 有 $\lim_{n \rightarrow \infty} \frac{\ln n}{\varepsilon^2 n} = 0$, 因此 $\frac{S_n}{n}$ 依概率收敛于 0, 服从大数定律。

(2) 独立同分布随机变量的期望为 $E(X_k) = \sum_{n=1}^{\infty} \frac{1}{2^n} \frac{2^n}{n^2} = \frac{\pi^2}{6}$ 存在, 由 Kolmogorov's SLLN 知 $\frac{1}{n} S_n \xrightarrow{\text{a.s.}} \frac{\pi^2}{6}$, 服从大数定律。

(3) 独立同分布随机变量的期望为 $E(X_k) = \sum_{n=1}^{\infty} \frac{c}{n \ln^2 n}$, 该序列收敛, 因此期望 $E(X_k), k = 1, 2, \dots$ 存在, 由 Kolmogorov' s SLLN 知 $\frac{1}{n} S_n \xrightarrow{\text{a.s.}} E(X_1)$, 服从大数定律。

12. 计算机在进行加法运算时，对每个加数取整，设所有的取整误差相互独立且都服从 $[-0.5, 0.5]$ 上的均匀分布。

(1) 若将 1500 个数相加，问：误差总和的绝对值超过 15 的概率是多少？

(2) 多少个数相加在一起可使得误差总和的绝对值小于 10 的概率为 0.90？

12. 计算机在进行加法运算时，对每个加数取整，设所有的取整误差相互独立且都服从 $[-0.5, 0.5]$ 上的均匀分布。

(1) 若将 1500 个数相加，问：误差总和的绝对值超过 15 的概率是多少？

(2) 多少个数相加在一起可使得误差总和的绝对值小于 10 的概率为 0.90？

解：(1) 设 $\xi_n \sim U[-0.5, 0.5]$, $S_n = \sum_{i=0}^n \xi_i$ ，显然，

$$E(\xi_n) = 0, \quad \text{var}(\xi_n) = \frac{1}{12}, \quad E(S_n) = 0, \quad \text{var}(S_n) = \frac{n}{12}$$

因此

$$S_n^* = \frac{S_n}{\sqrt{\frac{n}{12}}} = \frac{S_n}{\sqrt{125}} \xrightarrow{\omega} N(0, 1)$$

$$P(|S_n| > 15) = P\left(|S_n^*| > \frac{15}{\sqrt{125}}\right) = 2\left(1 - \Phi\left(\frac{3}{\sqrt{5}}\right)\right) \approx 0.18$$

(2)

$$P(|S_n| < 10) = P\left(|S_n^*| < 10\sqrt{\frac{12}{n}}\right) = 2\left(1 - \Phi\left(10\sqrt{\frac{12}{n}}\right)\right) = 0.9$$

解得 $n \approx 441$ 。

15. 对足够多的选民进行民意调查，以确定赞成某一候选人的百分比。假设选民中有未知的百分比 p 的人赞成该候选人，并且选民彼此是独立行动的，问：为了有 95% 的把握预测 p 的值在 0.045 的误差幅度内，应该调查多少人？

解：设 $X_i = \begin{cases} 1, & \text{选民 } i \text{ 赞成,} \\ 0, & \text{选民 } i \text{ 不赞成,} \end{cases}$, $S_n = \sum_{i=1}^n X_i$ ，显然

$$E(S_n) = np, \quad \text{var}(S_n) = n(p - p^2)$$

当 n 足够大时：

$$S_n^* = \frac{S_n - np}{\sqrt{n(p - p^2)}} \sim N(0, 1)$$

$$P\left(\left|\frac{\sum_{i=1}^n X_i - np}{n}\right| < 0.045\right) = P\left(\left|\frac{\sqrt{n(p - p^2)} S_n^*}{n}\right| < 0.045\right) \geq 0.95$$

在 $p = 0.5$ 时为最坏情况，此时有 $2\Phi(0.09\sqrt{n}) - 1 \geq 0.95$ ，解得 n 最小为 475.

① 概率论分册

② 参数估计

③ 假设检验

④ 回归分析

1. 设 X 的分布为几何分布：

$$P(X = k) = p(1 - p)^{k-1} \quad (k = 1, 2, \dots),$$

其中 $p \in (0, 1)$. 这个分布的实际背景为独立同分布试验序列，其中 p 为一次试验成功的概率， X 为试验序列中取得第一次成功所需的试验次数。设 X_1, \dots, X_n 为来自总体 X 的样本。

- (1) 写出这个模型的似然函数；
- (2) 求出参数 p 的 ML 估计；
- (3) 求出参数 p 的矩估计。

1. 设 X 的分布为几何分布：

$$P(X = k) = p(1 - p)^{k-1} \quad (k = 1, 2, \cdots),$$

其中 $p \in (0, 1)$. 这个分布的实际背景为独立同分布试验序列，其中 p 为一次试验成功的概率， X 为试验序列中取得第一次成功所需的试验次数。设 X_1, \cdots, X_n 为来自总体 X 的样本。

- (1) 写出这个模型的似然函数；
- (2) 求出参数 p 的 ML 估计；
- (3) 求出参数 p 的矩估计。

解：(1) 似然函数为

$$\begin{aligned} L(p; x_1, x_2, \cdots, x_n) &= \prod_{i=1}^n P(X_i = x_i) \\ &= \prod_{i=1}^n p(1 - p)^{x_i-1} = p^n (1 - p)^{\sum_{i=1}^n x_i - n} \end{aligned}$$

(2) 对数似然函数:

$$\begin{aligned}\ell(p; x_1, x_2, \dots, x_n) &= \log L(p; x_1, x_2, \dots, x_n) \\ &= n \log p + \left(\sum_{i=1}^n x_i - n \right) \log(1 - p)\end{aligned}$$

求导并令导数为零:

$$\frac{d\ell}{dp} = \frac{n}{p} - \frac{\sum_{i=1}^n x_i - n}{1 - p} = 0$$

解出 p : $p = \frac{n}{\sum_{i=1}^n x_i}$ 。令样本均值 $\bar{x} = \frac{\sum_{i=1}^n x_i}{n}$, 因此, 最大似然估计为:

$$\hat{p}_{\text{MLE}} = \frac{n}{\sum_{i=1}^n x_i} = \frac{1}{\bar{x}}$$

(3) 几何分布期望为 $E(X) = \frac{1}{p}$, 由矩估计, $\hat{p} = \frac{1}{\bar{x}}$ 。

3. 设 X 具有分布密度

$$p(x, \theta) = \begin{cases} \exp\{-(x - \theta)\}, & x \geq \theta, \\ 0, & x < \theta \end{cases} \quad (\theta \in (-\infty, +\infty)).$$

X_1, \dots, X_n 为来自总体 X 的样本, 求 θ 的 ML 估计 T_1 .

3. 设 X 具有分布密度

$$p(x, \theta) = \begin{cases} \exp\{-(x - \theta)\}, & x \geq \theta, \\ 0, & x < \theta \end{cases} \quad (\theta \in (-\infty, +\infty)).$$

X_1, \dots, X_n 为来自总体 X 的样本, 求 θ 的 ML 估计 T_1 .

解: 似然函数为:

$$L(\theta; x_1, x_2, \dots, x_n) = \prod_{i=1}^n p(x_i, \theta) = \prod_{i=1}^n \exp\{-(x_i - \theta)\} \cdot I(x_i \geq \theta)$$

似然函数在 $x_i < \theta$ 时为 0。当 $x_i \geq \theta$ 时, 似然函数随着 θ 增大而减小。然而, θ 必须满足 $x_i \geq \theta$ 对所有 i 成立, 因此, θ 的最大似然估计为样本中的最小值:

$$\hat{\theta}_{\text{MLE}} = \min_{1 \leq i \leq n} x_i$$

4. 设 $X_1, \dots, X_n \sim \text{iid} B(1, p)$, 即
 $P(X_i = 1) = p, P(X_i = 0) = 1 - p \ (i = 1, \dots, n)$.

4. 设 $X_1, \dots, X_n \sim \text{iid} B(1, p)$, 即
 $P(X_i = 1) = p, P(X_i = 0) = 1 - p \ (i = 1, \dots, n).$

- (1) 计算 $\text{var}(X_1)$;
- (2) 求 $\text{var}(X_1)$ 的 ML 估计 $T(X_1, \dots, X_n)$;
- (3) 求 $E(T(X_1, \dots, X_n))$.

解: (1) 随机变量 X_i 的期望和方差分别为:

$$E(X_i) = p, \quad \text{var}(X_i) = p(1 - p)$$

(2) 似然函数:

$$L(p; x_1, x_2, \dots, x_n) = \prod_{i=1}^n p^{x_i} (1 - p)^{1-x_i}$$

对数似然函数:

$$\ell(p; x_1, x_2, \dots, x_n) = \sum_{i=1}^n (x_i \log p + (1 - x_i) \log(1 - p))$$

求导并令导数为零：

$$\frac{d\ell}{dp} = \sum_{i=1}^n \left(\frac{x_i}{p} - \frac{1-x_i}{1-p} \right) = 0$$

解得 p 的最大似然估计为：

$$\hat{p}_{MLE} = \bar{x}$$

$\text{var}(X_1)$ 的 MLE 为：

$$T(X_1, \dots, X_n) = \hat{p}_{MLE}(1 - \hat{p}_{MLE}) = \bar{x}(1 - \bar{x})$$

(3) \bar{x} 的期望和方差分别为：

$$E(\bar{x}) = p, \quad \text{var}(\bar{x}) = \frac{p(1-p)}{n}$$

$$\begin{aligned} E(T(X_1, \dots, X_n)) &= E(\bar{x}(1 - \bar{x})) = E(\bar{x}) - \text{var}(\bar{x}) - [E(\bar{x})]^2 \\ &= p - \left(\frac{p(1-p)}{n} + p^2 \right) = p(1-p) \frac{n-1}{n} \end{aligned}$$

5. 在例 1.2 中 (二项分布), 求:

- (1) $\text{var}(\hat{p})$, 其中 \hat{p} 为参数 p 的 ML 估计;
- (2) $\text{var}(\hat{p})$ 的 ML 估计.

5. 在例 1.2 中 (二项分布), 求:

- (1) $\text{var}(\hat{p})$, 其中 \hat{p} 为参数 p 的 ML 估计;
- (2) $\text{var}(\hat{p})$ 的 ML 估计.

解: (1) 总体 $X \sim B(1, p), p \in [0, 1]$. 参数 p 的最大似然估计:

$$\hat{p} = \bar{x} = \frac{1}{n} \sum_{i=1}^n X_i$$

期望和方差分别为:

$$E(\hat{p}) = p, \quad \text{var}(\hat{p}) = \frac{p(1-p)}{n}$$

(2) 将 p 的最大似然估计 \hat{p} 代入得 $\text{var}(\hat{p})$ 的 ML 估计:

$$\text{var}(\hat{p})_{ML} = \frac{\hat{p}(1-\hat{p})}{n} = \frac{\bar{x}(1-\bar{x})}{n}$$

6. 在第 3 题中, 求:

(1) θ 的矩估计 $T_2(X_1, \dots, X_n)$;

(2) θ 的 ML 估计 T_1 和矩估计 T_2 的均方误差 $E_\theta[(T_1 - \theta)^2]$ 和 $E_\theta[(T_2 - \theta)^2]$.

解: (1) 注意到随机变量 $Z = X - \theta$ 服从参数为 1 的指数分布, 因此 X 的期望为:

$$E(X) = E(Z) + \theta = \theta + 1$$

因此, θ 的矩估计为:

$$T_2(X_1, \dots, X_n) = \bar{X} - 1$$

(2) 已经知道 θ 的最大似然估计 (MLE) 为:

$$T_1(X_1, \dots, X_n) = \min_{1 \leq i \leq n} X_i$$

其分布函数为：

$$F_{T_1}(x) = 1 - P(X_1 > x, X_2 > x, \dots, X_n > x) = 1 - [\exp\{-(x - \theta)\}]^n$$

密度函数为：

$$f_{T_1}(x) = \frac{d}{dx} F_{T_1}(x) = n \exp\{-n(x - \theta)\} \quad \text{for } x \geq \theta$$

令 $Y = T_1 - \theta$, Y 服从参数为 n 的指数分布, T_1 的期望和方差：

$$E(T_1) = E(Y) + \theta = \frac{1}{n} + \theta, \quad \text{var}(T_1) = \text{var}(Y) = \frac{1}{n^2}$$

因此, T_1 的均方误差为：

$$E_{\theta}[(T_1 - \theta)^2] = \text{var}(T_1) + [E(T_1 - \theta)]^2 = \frac{2}{n^2}$$

对于矩估计,

$$E(T_2) = E(\bar{X} - 1) = E(\bar{X}) - 1 = \theta + 1 - 1 = \theta$$

由于 $Z = X - \theta$ 服从指数分布,

$$\text{var}(T_2) = \text{var}(\bar{X} - 1) = \text{var}(\bar{X}) = \frac{\text{var}(X)}{n} = \frac{\text{var}(Z)}{n} = \frac{1}{n}$$

因此, T_2 的均方误差为:

$$E_{\theta}[(T_2 - \theta)^2] = \text{var}(T_2) + [E(T_2 - \theta)]^2 = \frac{1}{n}$$

7. 在例 1.4 (均匀分布参数估计) 中, 求出 $\hat{\theta} = \max_{1 \leq i \leq n} \{X_i\}$ 的
(1) $E_{\theta}(\hat{\theta})$; (2) 分布.

7. 在例 1.4 (均匀分布参数估计) 中, 求出 $\hat{\theta} = \max_{1 \leq i \leq n} \{X_i\}$ 的
(1) $E_{\theta}(\hat{\theta})$; (2) 分布.

解: 由独立同分布, $P(\hat{\theta} \leq t) = \left(\frac{t}{\theta}\right)^n, (0 \leq t \leq \theta)$, 因此分布为

$$p_{\hat{\theta}}(t) = \frac{dF_{\hat{\theta}}(t)}{dt} = \begin{cases} \frac{n}{\theta} \left(\frac{t}{\theta}\right)^{n-1}, & (0 \leq t \leq \theta), \\ 0, & \text{其他,} \end{cases}$$

期望为

$$E_{\theta}(\hat{\theta}) = \int_{-\infty}^{\infty} p_{\hat{\theta}}(t) t dt = \int_0^{\theta} \frac{n}{\theta} \left(\frac{t}{\theta}\right)^{n-1} t dt = \frac{n\theta}{n+1}$$

8. 设 X_1, \dots, X_n 为来自总体 X 的一个样本, 又总体 X 的分布密度为

$$p(x, \theta) = \begin{cases} \theta x^{\theta-1}, & 0 \leq x \leq 1, \\ 0, & \text{其他} \end{cases} \quad (\theta > 0).$$

(1) 求 θ 的矩估计; (2) 求 θ 的 ML 估计.

8. 设 X_1, \dots, X_n 为来自总体 X 的一个样本, 又总体 X 的分布密度为

$$p(x, \theta) = \begin{cases} \theta x^{\theta-1}, & 0 \leq x \leq 1, \\ 0, & \text{其他} \end{cases} (\theta > 0).$$

(1) 求 θ 的矩估计; (2) 求 θ 的 ML 估计.

解: (1) X 的期望为:

$$E[X] = \int_0^1 x \cdot \theta x^{\theta-1} dx = \theta \int_0^1 x^{\theta} dx = \frac{\theta}{\theta + 1}.$$

因此 θ 的矩估计量 $\hat{\theta}$ 为

$$\hat{\theta} = \frac{\bar{x}}{1 - \bar{x}}.$$

(2) 对数似然函数:

$$l(\theta) = \ln(L(\theta)) = \ln \prod_{i=1}^n \theta X_i^{\theta-1} = n \ln(\theta) + (\theta - 1) \sum_{i=1}^n \ln(X_i).$$

求导并令导数等于零:

$$\frac{d}{d\theta} l(\theta) = \frac{n}{\theta} + \sum_{i=1}^n \ln(X_i) = 0.$$

解这个方程可得 ML 估计为

$$\hat{\theta}_{\text{ML}} = -\frac{n}{\sum_{i=1}^n \ln(X_i)}.$$

10. 设 $\xi \sim N(0, 1)$, $\eta \sim \chi^2(n)$, 且 ξ 和 η 相互独立, 证明:
 $\xi/\sqrt{\eta/n}$ 服从 $t(n)$ 分布.

10. 设 $\xi \sim N(0, 1)$, $\eta \sim \chi^2(n)$, 且 ξ 和 η 相互独立, 证明:
 $\xi/\sqrt{\eta/n}$ 服从 $t(n)$ 分布.

不用证明, 当作定义。

13. 在 (4.6) 式中将 e 指数上的表达式写成

$$\left(\bar{x}, \bar{y}, \sum_{i=1}^n (x_i - \bar{x})^2, \sum_{i=1}^n (y_i - \bar{y})^2, \sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y}) \right)$$

的函数.

解: 原始的指数部分为:

$$-\frac{1}{2(1-\rho^2)} \sum_{i=1}^n \left[\left(\frac{x_i - \mu_1}{\sigma_1} \right)^2 - \frac{2\rho(x_i - \mu_1)(y_i - \mu_2)}{\sigma_1\sigma_2} + \left(\frac{y_i - \mu_2}{\sigma_2} \right)^2 \right].$$

其中

$$\sum_{i=1}^n \left(\frac{x_i - \mu_1}{\sigma_1} \right)^2 = \frac{1}{\sigma_1^2} \sum_{i=1}^n (x_i - \mu_1)^2 = \frac{1}{\sigma_1^2} \left[\sum_{i=1}^n (x_i - \bar{x})^2 + n(\bar{x} - \mu_1)^2 \right].$$

$$\begin{aligned} & \sum_{i=1}^n (x_i - \mu_1)(y_i - \mu_2) \\ &= \sum_{i=1}^n ((x_i - \bar{x}) + (\bar{x} - \mu_1)) ((y_i - \bar{y}) + (\bar{y} - \mu_2)) \\ &= S_{xy} + n(\bar{x} - \mu_1)(\bar{y} - \mu_2). \end{aligned}$$

$$\sum_{i=1}^n \left(\frac{y_i - \mu_2}{\sigma_2} \right)^2 = \frac{1}{\sigma_2^2} \left[\sum_{i=1}^n (y_i - \bar{y})^2 + n(\bar{y} - \mu_2)^2 \right].$$

合并所有项，e 指数上的表达式可以写成

$$\begin{aligned} & - \frac{1}{2(1 - \rho^2)} \left(\frac{S_x^2}{\sigma_1^2} + \frac{S_y^2}{\sigma_2^2} - \frac{2\rho S_{xy}}{\sigma_1 \sigma_2} + n \left[\frac{(\bar{x} - \mu_1)^2}{\sigma_1^2} - \frac{2\rho(\bar{x} - \mu_1)(\bar{y} - \mu_2)}{\sigma_1 \sigma_2} \right. \right. \\ & \left. \left. + \frac{(\bar{y} - \mu_2)^2}{\sigma_2^2} \right] \right). \end{aligned}$$

14. 设 $X \sim B(n, \theta)$, 即 X 的分布由下式给出:

$$P_{\theta}(X = k) = C_n^k \theta^k (1 - \theta)^{n-k} \quad (k = 0, 1, \dots, n).$$

(1) 求 θ 和 $\theta(1 - \theta)$ 的无偏估计; (2) 求 θ 的无偏估计 $\hat{\theta}$ 的均方误差.

解: (1) 注意到 $n\theta$ 和 $n\theta(1 - \theta)$ 分别为分布的期望和方差, 因此只需考虑对均值和方差的无偏估计, 对于均值, X 的期望为:

$$E[X] = n\theta.$$

这表明 $\frac{X}{n}$ 是 θ 的无偏估计。

代入 $\theta(1 - \theta)$, 得到 $\hat{\theta}(1 - \hat{\theta}) = \frac{X}{n} (1 - \frac{X}{n})$, 其期望:

$$\begin{aligned} E \left[\frac{X}{n} \left(1 - \frac{X}{n} \right) \right] &= \frac{1}{n} E[X] - \frac{1}{n^2} E[X^2] = \frac{1}{n} (n\theta) - \frac{1}{n^2} (n\theta(1 - \theta) + n^2\theta^2) \\ &= \frac{\theta(1 - \theta)(n - 1)}{n}. \end{aligned}$$

据此修正这个估计: $\frac{X}{n} \left(1 - \frac{X}{n}\right) \cdot \frac{n}{n-1}$. 这样, 得到无偏估计

$$\frac{x}{n} \left(1 - \frac{x}{n}\right) \cdot \frac{n}{n-1} = \frac{x(n-x)}{n(n-1)}.$$

(2) 由于 \bar{X} 是无偏估计, 因此均方误差等于方差, 样本均值 \bar{X} 的方差为:

$$\text{Var}(\bar{X}) = \text{Var}\left(\frac{1}{n} \sum_{i=1}^n X_i\right) = \frac{1}{n^2} \sum_{i=1}^n \text{Var}(X_i) = \frac{1}{n^2} \cdot n\theta(1-\theta) = \frac{\theta(1-\theta)}{n}.$$

因此, \bar{X} 的均方误差为 $\frac{\theta(1-\theta)}{n}$.

16. 设 $X_1, \dots, X_n \sim \text{iid} N(\mu, 1)$, $\mu \in (-\infty, +\infty)$, 求 μ, μ^2, μ^3 的 UMVU 估计.

解: 根据指数分布性质与 \bar{X} 的无偏性, \bar{X} 是 μ 的 UMVU 估计量。

由于 \bar{X} 是 μ 的完全充分统计量, 可以构造一个依赖于 \bar{X} 的函数作为 μ^2 的估计。

$$E[\bar{X}^2] = \text{Var}(\bar{X}) + [E(\bar{X})]^2 = \frac{1}{n} + \mu^2.$$

因此, μ^2 的 UMVU 估计为:

$$\hat{\mu}^2 = \bar{X}^2 - \frac{1}{n}.$$

由于:

$$E[\bar{X}^3] = \mu^3 + 3\mu \cdot \frac{1}{n} + 3\mu \cdot \frac{1}{n} = \mu^3 + \frac{3\mu}{n}.$$

因此, μ^3 的 UMVU 估计为: $\hat{\mu}^3 = \bar{X}^3 - \frac{3\bar{X}}{n}.$

17. 设 $X_1, \dots, X_n \sim \text{iid} N(\mu, \sigma^2)$, $\mu \in (-\infty, +\infty)$, $\sigma^2 > 0$, 证明:

$$T(X_1, \dots, X_n) = \bar{X}^2 - \frac{\sum_{i=1}^n (X_i - \bar{X})^2}{n(n-1)}$$

是 μ^2 的 UMVU 估计.

17. 设 $X_1, \dots, X_n \sim \text{iid} N(\mu, \sigma^2)$, $\mu \in (-\infty, +\infty)$, $\sigma^2 > 0$, 证明:

$$T(X_1, \dots, X_n) = \bar{X}^2 - \frac{\sum_{i=1}^n (X_i - \bar{X})^2}{n(n-1)}$$

是 μ^2 的 UMVU 估计.

解: 样本均值 \bar{X} 也是正态分布的, 且其均值为 μ , 均值的方差为 $\frac{\sigma^2}{n}$. 因此,

$$E[\bar{X}^2] = \text{Var}(\bar{X}) + [E(\bar{X})]^2 = \frac{\sigma^2}{n} + \mu^2.$$

$$E\left[\frac{\sum_{i=1}^n (X_i - \bar{X})^2}{n(n-1)}\right] = \frac{(n-1)\sigma^2}{n(n-1)} = \frac{\sigma^2}{n}.$$

估计是无偏的:

$$E[T] = E\left[\bar{X}^2 - \frac{\sum_{i=1}^n (X_i - \bar{X})^2}{n(n-1)}\right] = \left(\frac{\sigma^2}{n} + \mu^2\right) - \frac{\sigma^2}{n} = \mu^2.$$

对于正态分布 $N(\mu, \sigma^2)$, (\bar{X}, S^2) 是 (μ, σ^2) 的充分统计量。

20. 设 $X_1, \dots, X_n \sim \text{iid} B(1, p), 0 < p < 1$, 证明 (X_1, \dots, X_n) 的分布为指数族分布, 并找出其完全充分统计量.

20. 设 $X_1, \dots, X_n \sim \text{iid} B(1, p)$, $0 < p < 1$, 证明 (X_1, \dots, X_n) 的分布为指数族分布, 并找出其完全充分统计量.

解: 考虑 X_1, \dots, X_n 的联合概率分布:

$$\begin{aligned} P(X_1 = x_1, \dots, X_n = x_n) &= \prod_{i=1}^n P(X_i = x_i) = \prod_{i=1}^n p^{x_i} (1-p)^{1-x_i} = p^{\sum_{i=1}^n x_i} (1-p)^{n - \sum_{i=1}^n x_i} \\ &= \exp \left(\sum_{i=1}^n x_i \log(p) + (n - \sum_{i=1}^n x_i) \log(1-p) \right) \\ &= \exp \left(\sum_{i=1}^n x_i (\log(p) - \log(1-p)) + n \log(1-p) \right) \end{aligned}$$

(X_1, \dots, X_n) 的联合分布具备指数族分布的形式。

由指数族分布的形式知, $T(x) = \sum_{i=1}^n x_i$ 是完全充分统计量。

21. 设 $X_1, \dots, X_n \sim \text{iid} p(x, \theta), 0 < \theta < 1$, 其中 $p(x, \theta)$ 为离散型随机变量的分布列:

$$p(x, \theta) = \theta(1 - \theta)^{x-1} \quad (x = 1, 2, \dots),$$

证明 (X_1, \dots, X_n) 的分布为指数族分布, 并找出完全充分统计量.

21. 设 $X_1, \dots, X_n \sim \text{iid} p(x, \theta), 0 < \theta < 1$, 其中 $p(x, \theta)$ 为离散型随机变量的分布列:

$$p(x, \theta) = \theta(1 - \theta)^{x-1} \quad (x = 1, 2, \dots),$$

证明 (X_1, \dots, X_n) 的分布为指数族分布, 并找出完全充分统计量.

$$\begin{aligned} P(X_1 = x_1, \dots, X_n = x_n) &= \prod_{i=1}^n p(x_i, \theta) = \prod_{i=1}^n \theta(1 - \theta)^{x_i-1} = \theta^n (1 - \theta)^{\sum_{i=1}^n (x_i-1)} \\ &= \exp \left(n \log(\theta) + \sum_{i=1}^n (x_i - 1) \log(1 - \theta) \right) \\ &= \exp \left(\log(1 - \theta) \sum_{i=1}^n x_i + n(\log(\theta) - \log(1 - \theta)) \right) \end{aligned}$$

$T(x) = \sum_{i=1}^n x_i$ 是完全充分统计量

26. 证明：例 6.1 中 σ^2 的 ML 估计 $\hat{\sigma}^2$ 与 σ^2 的 UMVU 估计 S_n 具有相同的渐近分布。

26. 证明：例 6.1 中 σ^2 的 ML 估计 $\hat{\sigma}^2$ 与 σ^2 的 UMVU 估计 S_n 具有相同的渐近分布。

证明：已经知道， σ^2 的 MLE 是：

$$\hat{\sigma}^2 = \frac{1}{n} \sum_{i=1}^n (X_i - \bar{X})^2$$

UMVUE S_n 是样本方差：

$$S_n^2 = \frac{1}{n-1} \sum_{i=1}^n (X_i - \bar{X})^2$$

$(n-1)S_n^2/\sigma^2$ 服从自由度为 $n-1$ 的卡方分布。设 $\xi_1, \dots, \xi_{n-1} \sim \text{iid} N(0, 1)$ ，则 $\sum_{i=1}^{n-1} \xi_i^2$ 的分布是自由度为 $n-1$ 的 χ^2 分布。这样 $(n-1)S_n^2 - (n-1)\sigma^2$ 与 $\sum_{i=1}^{n-1} (\xi_i^2 - 1)\sigma^2$ 具有相同的分布，或 $\sqrt{n-1}S_n^2 - \sqrt{(n-1)}\sigma^2$ 与 $\sqrt{(n-1)} \left[\sum_{i=1}^{n-1} (\xi_i^2 - 1)\sigma^2 \right] / (n-1)$ 具有相同的分布。

由于 $\text{var}(\xi_i^2) = 2$ ，利用中心极限定理可知，当 $n \rightarrow \infty$ 时，

$$\sigma^2 \sqrt{n-1} \left[\sum_{i=1}^{n-1} (\xi_i^2 - 1) \right] / (n-1) \xrightarrow{w} N(0, 2\sigma^4).$$

因此， S_n 的渐近分布可以写成：

$$\sqrt{n}(S_n - \sigma^2) \xrightarrow{d} N(0, 2\sigma^4)$$

同样的，对于 MLE $\hat{\sigma}^2 = \frac{n-1}{n} S_n$ ，有

$$\sqrt{n} \left(\frac{n}{n-1} \hat{\sigma}^2 - \sigma^2 \right) \xrightarrow{d} N(0, 2\sigma^4)$$

因此，当 $n \rightarrow \infty$ 时， $\hat{\sigma}^2$ 的渐近分布可以写成：

$$\sqrt{n}(\hat{\sigma}^2 - \sigma^2) \xrightarrow{d} N(0, 2\sigma^4)$$

即 ML 估计 $\hat{\sigma}^2$ 与 σ^2 的 UMVU 估计 S_n 具有相同的渐近分布。

33. 已知某统计工作者对某种面值纸币的长度（单位：mm）进行测量，得数据：

156.2, 155.3, 155.5, 155.1, 155.3, 154.5, 154.9, 155.1, 154.7, 154.7.

(1) 求出该种纸币长度均值的置信度为 0.95 的置信区间；(2) 求出该种纸币长度标准差的置信度为 0.95 的置信上限。

33. 已知某统计工作者对某种面值纸币的长度 (单位: mm) 进行测量, 得数据:

156.2, 155.3, 155.5, 155.1, 155.3, 154.5, 154.9, 155.1, 154.7, 154.7.

(1) 求出该种纸币长度均值的置信度为 0.95 的置信区间; (2) 求出该种纸币长度标准差的置信度为 0.95 的置信上限。

解: (1) 样本均值和标准差 \bar{x} 和 s 为:

$$\bar{x} = 155.04$$

$$s = \sqrt{\frac{1}{n-1} \sum_{i=1}^n (x_i - \bar{x})^2} = \sqrt{\frac{2.242}{9}} \approx 0.4991$$

t 分布的临界值 $t_{\alpha/2, df}$ 可以从 t 表中查找, 其中 $df = n - 1 = 9$, 对于 0.95 的置信水平, $\alpha = 0.05$, 查表得到 $t_{0.025, 9} \approx 2.262$ 。

$$CI = \left(\bar{x} - t_{0.025, 9} \cdot \frac{s}{\sqrt{n}}, \bar{x} + t_{0.025, 9} \cdot \frac{s}{\sqrt{n}} \right) \approx (154.78, 155.48)$$

(2) 选取 $\frac{1}{\sigma^2}(n-1)s^2$ 为枢轴量, 对于 0.95 置信水平, 自由度 $df = n - 1 = 9$, 查表得到 $\chi_{0.05,9}^2 \approx 3.325$ 。置信上限为:

$$\sqrt{\frac{(n-1)s^2}{\chi_{0.05,9}^2}} \approx 0.8062$$

34. 设 $\mathbf{X} = (X_1, X_2)^\top$, $X_1, X_2 \sim \text{iid}N(\mu, \sigma^2)$ 。(1) 求 $X_1 + X_2$ 与 $X_1 - X_2$ 的联合分布密度; (2) 证明: $X_1 + X_2$ 与 $X_1 - X_2$ 相互独立。

34. 设 $\mathbf{X} = (X_1, X_2)^\top$, $X_1, X_2 \sim \text{iid} N(\mu, \sigma^2)$ 。(1) 求 $X_1 + X_2$ 与 $X_1 - X_2$ 的联合分布密度; (2) 证明: $X_1 + X_2$ 与 $X_1 - X_2$ 相互独立。

解: 设: $\mathbf{Y} = \mathbf{A}\mathbf{X}$, 其中 \mathbf{A} 是一个转换矩阵:

$$\mathbf{A} = \begin{pmatrix} 1 & 1 \\ 1 & -1 \end{pmatrix}$$

随机变量 \mathbf{Y} 的联合分布也是一个二维正态分布, 其均值向量和协方差矩阵:

$$\mu_{\mathbf{Y}} = \mathbf{A}\mu_{\mathbf{X}} = \begin{pmatrix} 1 & 1 \\ 1 & -1 \end{pmatrix} \begin{pmatrix} \mu \\ \mu \end{pmatrix} = \begin{pmatrix} 2\mu \\ 0 \end{pmatrix} \quad \Sigma_{\mathbf{Y}} = \mathbf{A}\Sigma_{\mathbf{X}}\mathbf{A}^\top = \sigma^2 \begin{pmatrix} 1 & 1 \\ 1 & -1 \end{pmatrix}$$

即 $X_1 + X_2$ 与 $X_1 - X_2$ 服从二维正态分布 $N\left(\begin{pmatrix} 2\mu \\ 0 \end{pmatrix}, 2\sigma^2\mathbf{I}\right)$, 独立

36. 设 $\mathbf{X} = (X_1, X_2)^\top \sim N(\mu, \mathbf{M})$, 其中 μ 为 2 维向量, \mathbf{M} 为 2×2 正定矩阵, 求系数 b , 使得 X_1 与 $X_2 - bX_1$ 相互独立。

36. 设 $\mathbf{X} = (X_1, X_2)^\top \sim N(\mu, \mathbf{M})$, 其中 μ 为 2 维向量, \mathbf{M} 为 2×2 正定矩阵, 求系数 b , 使得 X_1 与 $X_2 - bX_1$ 相互独立。

解: X_1 与 $X_2 - bX_1$ 的协方差:

$$\text{Cov}(X_1, X_2 - bX_1) = \text{Cov}(X_1, X_2) - b \cdot \text{Cov}(X_1, X_1) = m_{12} - bm_{11}$$

令协方差为 0, 解得 b :

$$b = \frac{m_{12}}{m_{11}}$$

37. 设 $\mathbf{X} = (X_1, X_2, X_3)^\top \sim N(\mu, \mathbf{I}_3)$, 其中 $\mu^\top = (1, 0, 3)$, 又设

$$\mathbf{d} = \begin{bmatrix} 2 \\ 1 \end{bmatrix}, \quad \mathbf{A} = \begin{bmatrix} 2 & 1 & 0 \\ 0 & 1 & 1 \end{bmatrix},$$

求 $\mathbf{Y} = \mathbf{A}\mathbf{X} + \mathbf{d}$ 的分布.

解: \mathbf{Y} 的均值向量 μ_Y 和协方差矩阵 Σ_Y :

$$\mu_Y = \mathbf{A}\mu + \mathbf{d} = \begin{bmatrix} 2 & 1 & 0 \\ 0 & 1 & 1 \end{bmatrix} \begin{bmatrix} 1 \\ 0 \\ 3 \end{bmatrix} + \begin{bmatrix} 2 \\ 1 \end{bmatrix} = \begin{bmatrix} 4 \\ 4 \end{bmatrix}$$

$$\Sigma_Y = \mathbf{A}\mathbf{I}_3\mathbf{A}^\top = \begin{bmatrix} 2 & 1 & 0 \\ 0 & 1 & 1 \end{bmatrix} \begin{bmatrix} 2 & 0 \\ 1 & 1 \\ 0 & 1 \end{bmatrix} = \begin{bmatrix} 5 & 1 \\ 1 & 2 \end{bmatrix}$$

\mathbf{Y} 服从参数为 μ_Y 和 Σ_Y 的正态分布。

① 概率论分册

② 参数估计

③ 假设检验

④ 回归分析

1. 设 $X_1, \dots, X_n \sim \text{iid} N(\mu, \sigma_0^2)$, σ_0^2 为已知, 假设检验问题为

$$H_0 : \mu \geq \mu_0 \leftrightarrow H_1 : \mu < \mu_0,$$

求出它的水平为 α 的 UMP 否定域。

1. 设 $X_1, \dots, X_n \sim \text{iid} N(\mu, \sigma_0^2)$, σ_0^2 为已知, 假设检验问题为

$$H_0: \mu \geq \mu_0 \leftrightarrow H_1: \mu < \mu_0,$$

求出它的水平为 α 的 UMP 否定域。

解: 跟单参指数族性质, 使用样本均值 \bar{X} 作为检验统计量。样本均值 \bar{X} 也服从正态分布:

$$\bar{X} \sim N\left(\mu, \frac{\sigma_0^2}{n}\right)$$

标准化后的检验统计量 Z 为:

$$Z = \frac{\bar{X} - \mu_0}{\sigma_0 / \sqrt{n}}$$

在原假设 $H_0: \mu \geq \mu_0$ 下, Z 服从标准正态分布 $N(0, 1)$ 。设 z_α 是标准正态分布的 α 分位数, 即 $P(Z \leq z_\alpha) = \alpha$ 。因此, 水平为 α 的 UMP 否定域可以表示为: $\{z: z \leq z_\alpha\}$ 。

将 Z 表达式代入上述不等式中，得到水平为 α 的 UMP 否定域：

$$\{\mathbf{x} : \frac{\bar{\mathbf{x}} - \mu_0}{\sigma_0/\sqrt{n}} \leq z_\alpha\}$$

2. 设某接收站收到的信号为 X ，当对方发信号时， X 的分布为 $U(0, 2)$ ；当对方不发信号时， X 的分布为 $U(-1, 1)$ 。考虑如下的假设检验问题：

$$H_0 : X \sim U(-1, 1) \leftrightarrow H_1 : X \sim U(0, 2).$$

求出此假设检验问题依赖于观察值 $X = x$ 的水平为 α 的 UMP 否定域。

解：设 $f_0(x)$ 和 $f_1(x)$ 分别是原假设 H_0 和备择假设 H_1 下的密度函数，似然比 $\Lambda(x)$ 定义为：

$$\Lambda(x) = \frac{f_1(x)}{f_0(x)} = \begin{cases} 0, & -1 \leq x \leq 0, \\ 1, & 0 \leq x \leq 1, \\ \infty, & 1 \leq x \leq 2, \end{cases}$$

因此，否定域应该是 $\{x : x \geq c\}, (0 < c < 1)$ 。

需要调整否定域以确保其概率等于 α ，即 $P(X > c|H_0) = \alpha$ 。

$$P(X > c|H_0) = \frac{1 - c}{2} = \alpha$$

因此，水平为 α 的 UMP 否定域是：

$$\{x : x > 1 - 2\alpha\}$$

3. 设 X 可能来自两个不同的总体，它们的分布密度分别为 $f_0(x)$ 和 $f_1(x)$ ，其中 $f_0(x)$ 为区间 $(0, 1)$ 上均匀分布 $U(0, 1)$ 的分布密度， $f_1(x) = 3x^2, x \in (0, 1)$ 。相应的假设检验问题为

$$H_0 : f = f_0(x) \leftrightarrow H_1 : f = f_1(x).$$

求出相应的依赖于观察值 $X = x$ 的水平为 α 的 UMP 否定域。

3. 设 X 可能来自两个不同的总体，它们的分布密度分别为 $f_0(x)$ 和 $f_1(x)$ ，其中 $f_0(x)$ 为区间 $(0, 1)$ 上均匀分布 $U(0, 1)$ 的分布密度， $f_1(x) = 3x^2, x \in (0, 1)$ 。相应的假设检验问题为

$$H_0 : f = f_0(x) \leftrightarrow H_1 : f = f_1(x).$$

求出相应的依赖于观察值 $X = x$ 的水平为 α 的 UMP 否定域。

解：似然比 $\Lambda(x)$ 定义为：

$$\Lambda(x) = \frac{f_1(x)}{f_0(x)} = \frac{3x^2}{1} = 3x^2$$

设似然比的临界值 c ，使得当 $3x^2 > c$ 时，拒绝 H_0 。由于 $x \in (0, 1)$ ，所以否定域 $\{x : x > \sqrt{\frac{c}{3}}\}$ 。

令否定域的大小等于 α ，即

$$P(X > \sqrt{\frac{c}{3}} | H_0) = 1 - \sqrt{\frac{c}{3}} = \alpha$$

解得：

$$c = 3(1 - \alpha)^2$$

因此，水平为 α 的 UMP 否定域是：

$$\{x : x > 1 - \alpha\}$$

6. 已知矿井中瓦斯的含量（浓度）为随机变量，其分布为 $N(\mu, \sigma^2)$, $\sigma^2 > 0$ 。按规定， $\mu \geq \mu_0$ 为危险浓度。为了保证安全，矿里决定设立 10 个监测点。为了通过监测值监测矿上的安全状况，采用假设检验的方法。假设检验问题有两种提法：(1) $H_0: \mu \geq \mu_0 \leftrightarrow H_1: \mu < \mu_0$; (2) $H_0: \mu \leq \mu_0 \leftrightarrow H_1: \mu > \mu_0$. 你认为应采用哪一种提法？并说明理由。

6. 已知矿井中瓦斯的含量（浓度）为随机变量，其分布为 $N(\mu, \sigma^2)$, $\sigma^2 > 0$ 。按规定， $\mu \geq \mu_0$ 为危险浓度。为了保证安全，矿里决定设立 10 个监测点。为了通过监测值监测矿上的安全状况，采用假设检验的方法。假设检验问题有两种提法：(1) $H_0: \mu \geq \mu_0 \leftrightarrow H_1: \mu < \mu_0$; (2) $H_0: \mu \leq \mu_0 \leftrightarrow H_1: \mu > \mu_0$. 你认为应采用哪一种提法？并说明理由。

解：应采用提法 (1)。在提法 (1) 中，第一类错误（错误地拒绝 H_0 ）意味着我们在实际上瓦斯浓度不安全的情况下错误地认为它是安全的，这将导致严重的安全风险。假设检验问题提法 (1) 能够尽可能避免这种错误。

10. 设总体 $X \sim N(\mu_0, \sigma^2)$, μ_0 为已知, $\mathbf{X} = (X_1, \dots, X_n)$ 为来自 X 的一个样本, 假设检验问题为

$$H_0 : \sigma^2 \leq \sigma_0^2 \leftrightarrow H_1 : \sigma^2 > \sigma_0^2.$$

(1) 利用单参数指数族中的方法求出该假设检验问题的水平为 α 的否定域; (2) 利用广义似然比方法求出该假设检验问题的水平为 α 的否定域。

10. 设总体 $X \sim N(\mu_0, \sigma^2)$, μ_0 为已知, $\mathbf{X} = (X_1, \dots, X_n)$ 为来自 X 的一个样本, 假设检验问题为

$$H_0 : \sigma^2 \leq \sigma_0^2 \leftrightarrow H_1 : \sigma^2 > \sigma_0^2.$$

(1) 利用单参数指数族中的方法求出该假设检验问题的水平为 α 的否定域; (2) 利用广义似然比方法求出该假设检验问题的水平为 α 的否定域。

解: (1) 正态分布总体为单参数指数族, $T(\mathbf{X}) = (x - \mu_0)^2$, UMP 否定域形如

$$\mathcal{W} = \{\mathbf{x} : \sum_{i=1}^n (x_i - \mu_0)^2 > c\}$$

在方差 σ_0^2 条件下 $\frac{(x_i - \mu_0)}{\sigma_0} \sim N(0, 1)$, 故

$$\frac{1}{\sigma_0^2} \sum_{i=1}^n (x_i - \mu_0)^2 \sim \chi^2(n).$$

$$P_{\sigma_0^2}(\mathbf{X} \in \mathcal{W}) = P\left(\frac{1}{\sigma_0^2} \sum_{i=1}^n (x_i - \mu_0)^2 > \frac{c}{\sigma_0^2}\right) = \alpha.$$

卡方分布的 $1 - \alpha$ 分位数记为 $\chi_{1-\alpha}^2(n)$ ，因此 $c = \sigma_0^2 \chi_{1-\alpha}^2(n)$ 。
否定域为

$$\left\{ \mathbf{x} : \sum_{i=1}^n (x_i - \mu_0)^2 > \sigma_0^2 \chi_{1-\alpha}^2(n) \right\}$$

卡方分布的 $1 - \alpha$ 分位数记为 $\chi_{1-\alpha}^2(n)$ ，因此 $c = \sigma_0^2 \chi_{1-\alpha}^2(n)$ 。
否定域为

$$\left\{ \mathbf{x} : \sum_{i=1}^n (x_i - \mu_0)^2 > \sigma_0^2 \chi_{1-\alpha}^2(n) \right\}$$

(2) 似然函数为：

$$L(\sigma^2) = \prod_{i=1}^n \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{(X_i - \mu_0)^2}{2\sigma^2}\right)$$

设 $U = \frac{1}{\sigma_0^2} \sum_{i=1}^n (X_i - \mu_0)^2$ 。对似然函数求导得，方差的最大似然估计为

$$\hat{\sigma}^2 = \frac{U\sigma_0^2}{n}$$

在原假设 $H_0: \sigma^2 \leq \sigma_0^2$ 下, 最大似然估计为:

$$\hat{\sigma}_0^2 = \min \left(\frac{U\sigma_0^2}{n}, \sigma_0^2 \right)$$

因此, 广义似然比为:

$$\Lambda = \frac{L(\hat{\sigma}^2)}{L(\hat{\sigma}_0^2)} = \frac{\left(\frac{1}{\sqrt{2\pi\hat{\sigma}^2}} \right)^n \exp \left(-\frac{U\sigma_0^2}{2\hat{\sigma}^2} \right)}{\left(\frac{1}{\sqrt{2\pi\hat{\sigma}_0^2}} \right)^n \exp \left(-\frac{U\sigma_0^2}{2\hat{\sigma}_0^2} \right)} = \left(\frac{\hat{\sigma}^2}{\hat{\sigma}_0^2} \right)^{-\frac{n}{2}} \exp \left(\frac{U\sigma_0^2}{2} \left(\frac{1}{\hat{\sigma}_0^2} - \frac{1}{\hat{\sigma}^2} \right) \right)$$

当 $\frac{U}{n} \leq 1$ 时, $\Lambda = 1$; 当 $\frac{U}{n} > 1$ 时, $\Lambda = \left(\frac{U}{n} \right)^{-\frac{n}{2}} \exp \left(\frac{n}{2} \left(\frac{U}{n} - 1 \right) \right)$ 。

否定域形如

$$\left\{ \mathbf{x} : \frac{U}{n} > 1, \quad \left(\frac{U}{n} \right)^{-\frac{n}{2}} \exp \left(\frac{n}{2} \left(\frac{U}{n} - 1 \right) \right) > c \right\}$$

$$= \left\{ \mathbf{x} : \frac{U}{n} > 1, \quad U^{-\frac{n}{2}} \exp \left(\frac{U}{2} \right) > \tilde{c} \right\}$$

当 $\frac{U}{n} > 1$ 时, $U^{-\frac{n}{2}} \exp \left(\frac{U}{2} \right)$ 关于 U 单调减, 当 $\frac{U}{n} > 1$ 时,
 $U^{-\frac{n}{2}} \exp \left(\frac{U}{2} \right)$ 关于 U 单调减, 因此否定域可以写为

$$\{\mathbf{x} : U < c'\}$$

$\forall \sigma > \sigma_0^2, U \geq \frac{\sum_{i=1}^n (X_i - \mu_0)^2}{\sigma^2} \sim \chi^2(n)$. 在 $\sigma^2 = \sigma_0^2$ 时等号成立, 因此取 $c' = \chi_{\alpha}^2(n)$ 即可满足

$$\max_{\sigma^2 \geq \sigma_0^2} P_{\sigma^2}(U < c) = \alpha$$

所求否定域为

$$\left\{ \mathbf{x} : \frac{1}{\sigma_0^2} \sum_{i=1}^n (x_i - \mu_0)^2 > \chi_{1-\alpha}^2(n) \right\}$$

11. 设总体 $X \sim N(\mu, \sigma_0^2)$, σ_0^2 为已知, $\mathbf{X} = (X_1, \dots, X_n)$ 为来自 X 的一个样本, 假设检验问题为

$$H_0 : \mu \leq \mu_0 \leftrightarrow H_1 : \mu > \mu_0,$$

利用广义似然比方法求出该假设检验问题的水平为 α 的否定域。

解：在 $H_0: \mu \leq \mu_0$ 下，MLE 为： $\hat{\mu}_0 = \min(\bar{X}, \mu_0)$ ，在无约束下，MLE 为： $\hat{\mu} = \bar{X}$ 。
因此，似然比为：

$$\begin{aligned}\Lambda &= \frac{L(\hat{\mu})}{L(\hat{\mu}_0)} = \frac{\exp\left(-\frac{1}{2\sigma_0^2} \sum_{i=1}^n (X_i - \hat{\mu})^2\right)}{\exp\left(-\frac{1}{2\sigma_0^2} \sum_{i=1}^n (X_i - \hat{\mu}_0)^2\right)} \\ &= \exp\left(\frac{1}{2\sigma_0^2} \left(\sum_{i=1}^n (X_i - \hat{\mu}_0)^2 - \sum_{i=1}^n (X_i - \hat{\mu})^2\right)\right) \\ &= \exp\left(\frac{n}{2\sigma_0^2} (\hat{\mu}_0 - \hat{\mu})^2\right) = \exp\left(\frac{n}{2\sigma_0^2} (\min(\bar{X}, \mu_0) - \bar{X})^2\right) \\ &= \begin{cases} 1, & \bar{X} \leq \mu_0 \\ \exp\left(\frac{n}{2\sigma_0^2} (\mu_0 - \bar{X})^2\right), & \bar{X} > \mu_0, \end{cases}\end{aligned}$$

由正态分布性质, $\frac{\bar{X}-\mu_0}{\sigma_0/\sqrt{n}} \sim N(0,1)$, 因此取标准正态分布的 $1-\alpha$ 分位数 $z_{1-\alpha}$ 。我们拒绝 H_0 当且仅当:

$$\frac{\bar{X} - \mu_0}{\sigma_0/\sqrt{n}} > z_{1-\alpha}$$

即:

$$\bar{X} > \mu_0 + \frac{\sigma_0}{\sqrt{n}} z_{1-\alpha}$$

因此, 水平为 α 的广义似然比检验的否定域是:

$$\left\{ \mathbf{x} : \bar{x} > \mu_0 + \frac{\sigma_0}{\sqrt{n}} z_{1-\alpha} \right\}$$

① 概率论分册

② 参数估计

③ 假设检验

④ 回归分析

1. 设 b_0 和 b 是一元线性回归模型 (1.9) 中的截距和回归系数, 而 \hat{b}_0 和 \hat{b} 是相应的最小二乘估计, 记 $\hat{y}_i = \hat{b}_0 + \hat{b}x_i$, 证明:

$$\sum_{i=1}^n (y_i - \hat{y}_i) = 0, \quad \sum_{i=1}^n x_i (y_i - \hat{y}_i) = 0.$$

1. 设 b_0 和 b 是一元线性回归模型 (1.9) 中的截距和回归系数, 而 \hat{b}_0 和 \hat{b} 是相应的最小二乘估计, 记 $\hat{y}_i = \hat{b}_0 + \hat{b}x_i$, 证明:

$$\sum_{i=1}^n (y_i - \hat{y}_i) = 0, \quad \sum_{i=1}^n x_i (y_i - \hat{y}_i) = 0.$$

解: 最小二乘估计需要使得残差平方和 (RSS) 最小化:

$$RSS = \sum_{i=1}^n (y_i - \hat{y}_i)^2 = \sum_{i=1}^n (y_i - (\hat{b}_0 + \hat{b}x_i))^2$$

对 \hat{b}_0 和 \hat{b} 求偏导数, 并令其等于零。

$$\frac{\partial RSS}{\partial \hat{b}_0} = -2 \sum_{i=1}^n (y_i - \hat{b}_0 - \hat{b}x_i) = 0$$

$$\frac{\partial RSS}{\partial \hat{b}} = -2 \sum_{i=1}^n x_i (y_i - \hat{b}_0 - \hat{b}x_i) = 0$$

简化得到：

$$\sum_{i=1}^n (y_i - \hat{b}_0 - \hat{b}x_i) = \sum_{i=1}^n (y_i - \hat{y}_i) = 0$$

$$\sum_{i=1}^n x_i (y_i - \hat{b}_0 - \hat{b}x_i) = \sum_{i=1}^n x_i (y_i - \hat{y}_i) = 0$$