

Homework 8

姓名: 方嘉聪 学号: 2200017849

Problem 1. 给定 $(x_1, y_1), (x_2, y_2), \dots, (x_n, y_n)$, 其中 $y_i = \alpha + \beta x_i + \varepsilon_i$, ε_i 相互独立, 且服从 Laplace 分布, 其概率密度函数 (参考作业三第五题) 满足, 对于任意实数 $x \in \mathbb{R}$,

$$f(x) = \frac{1}{2b} e^{-|x|/b}.$$

这里 α, β 和 $b > 0$ 是未知参数. 证明 α, β 的最大似然估计量为

$$\operatorname{argmin}_{\alpha, \beta} \sum_{i=1}^n |y_i - (\alpha + \beta x_i)|.$$

Solution. 似然函数为

$$L(\alpha, \beta, b) = \prod_{i=1}^n f(\varepsilon_i) = \frac{1}{(2b)^n} \exp \left(-\frac{1}{b} \sum_{i=1}^n |y_i - (\alpha x_i + \beta)| \right)$$

由于 $b > 0$, 那么 $L(\alpha, \beta)$ 关于 $\sum_{i=1}^n |y_i - (\alpha x_i + \beta)|$ 单调递减, 故 α, β 的最大似然估计量为

$$\operatorname{argmin}_{\alpha, \beta} \sum_{i=1}^n |y_i - (\alpha + \beta x_i)|.$$

证毕.

Problem 2. 给定 $(x_1, y_1), (x_2, y_2), \dots, (x_n, y_n)$, 令 $\hat{\alpha}, \hat{\beta}$ 为最小二乘估计量, $\hat{y}_i = \hat{\alpha} + \hat{\beta} x_i$ 为 y_i 的预测值. 令 $\bar{x} = \frac{1}{n} \sum_{i=1}^n x_i$, $\bar{y} = \frac{1}{n} \sum_{i=1}^n y_i$, 证明

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

提示: 利用正规方程, 并证明 $\hat{y}_i = \bar{y} + \hat{\beta}(x_i - \bar{x})$.

Solution. 由正规方程, 我们有

$$\hat{\alpha} = \bar{y} - \hat{\beta} \bar{x}, \quad \hat{\beta} = \frac{\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})}{\sum_{i=1}^n (x_i - \bar{x})^2}.$$

那么

$$\hat{y}_i = \hat{\alpha} + \hat{\beta} x_i = \bar{y} - \hat{\beta} \bar{x} + \hat{\beta} x_i = \bar{y} + \hat{\beta}(x_i - \bar{x}).$$

又有

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

而

$$\begin{aligned}\sum_{i=1}^n (y_i - \hat{y}_i)(\hat{y}_i - \bar{y}) &= \sum_{i=1}^n \hat{\beta}(x_i - \bar{x})(y_i - \bar{y} - \hat{\beta}(x_i - \bar{x})) \\ &= \hat{\beta} \left[\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y}) - \hat{\beta} \sum_{i=1}^n (x_i - \bar{x})^2 \right] = 0.\end{aligned}$$

故有 $\sum_{i=1}^n (y_i - \bar{y})^2 = \sum_{i=1}^n (y_i - \hat{y}_i)^2 + \sum_{i=1}^n (\hat{y}_i - \bar{y})^2$. 证毕. \triangleleft

Problem 3. 给定 $(x_1, y_1), (x_2, y_2), \dots, (x_n, y_n)$, 其中 $y_i = \alpha + \beta x_i + \varepsilon_i$, $\varepsilon_i \stackrel{\text{i.i.d.}}{\sim} \mathcal{N}(0, \sigma^2)$. 沿用第二题中的记号, 并令 $s^2 = \frac{1}{n-2} \sum_{i=1}^n (y_i - \hat{y}_i)^2$, $s_{xx} = \sum_{i=1}^n (x_i - \bar{x})^2$.

(1) 令

$$\begin{aligned}q_1 &= [1/\sqrt{n}, 1/\sqrt{n}, \dots, 1/\sqrt{n}]^\top \in \mathbb{R}^n, \\ q_2 &= \left[\frac{x_1 - \bar{x}}{\sqrt{s_{xx}}}, \frac{x_2 - \bar{x}}{\sqrt{s_{xx}}}, \dots, \frac{x_n - \bar{x}}{\sqrt{s_{xx}}} \right]^\top \in \mathbb{R}^n.\end{aligned}$$

证明: 存在 $q_3, q_4, \dots, q_n \in \mathbb{R}^n$, 使得 $q_1, q_2, q_3, \dots, q_n$ 构成 \mathbb{R}^n 的一组标准正交基.

(2) 将 y 视作 \mathbb{R}^n 中的一个向量. 对于 $1 \leq i \leq n$, 令 $z_i = q_i^\top y$, 也即 $z = Qy \in \mathbb{R}^n$, 其中 $Q \in \mathbb{R}^{n \times n}$ 的第 i 行为 $q_i \in \mathbb{R}^n$. 给出 n 维随机变量 z 服从的分布. 提示: 计算随机向量 y 的数学期望, 并验证其与 q_3, q_4, \dots, q_n 的正交性.

(3) 证明: $z_1 = \sqrt{n}\bar{y}$, $z_2 = \sqrt{s_{xx}}\hat{\beta}$.

(4) 利用第二题中提示和结论, 证明 $\sum_{i=1}^n (\hat{y}_i - \bar{y})^2 = z_2^2$ 及 $(n-2)s^2 = \sum (y_i - \hat{y}_i)^2 = \sum_{i=3}^n z_i^2$.

(5) 给出 $(n-2)s^2/\sigma^2$ 服从的分布, 并证明 s^2 与 $\hat{\alpha}, \hat{\beta}$ 均相互独立.

(6) 当 $\beta = 0$, 给出统计量 $t = \frac{\hat{\beta}}{\sqrt{s^2}/\sqrt{s_{xx}}}$ 服从的分布.

(7) 若 σ^2 未知, 考虑假设检验问题, 原假设 $H_0: \beta = 0$, 备择假设 $H_1: \beta \neq 0$. 拒绝域为

$$W = \{((x_1, y_1), (x_2, y_2), \dots, (x_n, y_n)) \mid |t| \geq c\},$$

其中 c 是待定常数. 若显著性水平为 α , 给出 c 的取值. \blacktriangleleft

Solution. (1) 显然 q_1, q_2 的模长均为 1, 且有

$$q_1^\top q_2 = \frac{1}{\sqrt{n}\sqrt{s_{xx}}} \sum_{i=1}^n (x_i - \bar{x}) = 0.$$

故 q_1, q_2 正交. 由于 q_1, q_2 线性无关, 故存在非零向量 q'_3, q'_4, \dots, q'_n 使得 $q_1, q_2, q'_3, \dots, q'_n$ 构成 \mathbb{R}^n 的一组基. 对其进行 Schmidt 正交化 (并归一化), 即可得到一组标准正交基, 注意到在正交化过程中, q_1, q_2 保持不变. 故存在 q_3, q_4, \dots, q_n 使得 $q_1, q_2, q_3, \dots, q_n$ 构成 \mathbb{R}^n 的一组标准正交基.

(2) 记 $\varepsilon = (\varepsilon_1, \varepsilon_2, \dots, \varepsilon_n) \in \mathbb{R}^n$. 待定系数, 设

$$y = (\alpha + \beta x_1 + \varepsilon_1, \dots, \alpha + \beta x_n + \varepsilon_n) = k_1 q_1 + k_2 q_2 + \varepsilon.$$

解得:

$$y = \sqrt{n}(\alpha + \beta\bar{x}) \cdot q_1 + \beta\sqrt{s_{xx}} \cdot q_2 + \varepsilon.$$

那么有

$$\begin{aligned} z_1 &= q_1^\top y = \sqrt{n}(\alpha + \beta\bar{x}) \cdot \|q_1\|^2 + q_1^\top \varepsilon \sim \mathcal{N}(\sqrt{n}\alpha + \sqrt{n}\beta\bar{x}, \sigma^2), \\ z_2 &= q_2^\top y = \beta\sqrt{s_{xx}} \cdot \|q_2\|^2 + q_2^\top \varepsilon \sim \mathcal{N}(\beta\sqrt{s_{xx}}, \sigma^2). \end{aligned}$$

对于 $i \in \{3, 4, \dots, n\}$, 由正交性有 $q_i^\top q_1 = q_i^\top q_2 = 0$, 故 (注意到 q_i 的模长为 1)

$$z_i = q_i^\top y = q_i^\top \varepsilon \sim \mathcal{N}(0, \sigma^2).$$

记 $z'_1 = z_1 - \sqrt{n}\alpha - \sqrt{n}\beta\bar{x} \sim \mathcal{N}(0, \sigma^2)$, $z'_2 = z_2 - \beta\sqrt{s_{xx}} \sim \mathcal{N}(0, \sigma^2)$. 记

$$z' = (z'_1, z'_2, z_3, \dots, z_n) = (q_1^\top \varepsilon, q_2^\top \varepsilon, q_3^\top \varepsilon, \dots, q_n^\top \varepsilon) = Q\varepsilon.$$

我们来证明随机变量 $z'_1, z'_2, z_3, \dots, z_n$ 是相互独立的. 设 $t = (t_1, t_2, \dots, t_n) \in \mathbb{R}^n$, 有

$$\begin{aligned} f_{z'}(z' = t) &= f_{z'}(Q\varepsilon = t) = f_\varepsilon(\varepsilon = Q^\top t) \\ &= \prod_{i=1}^n f_{\varepsilon_i}(\varepsilon_i = q_i^\top t) = \frac{1}{(2\pi)^{n/2}\sigma^n} \exp\left(-\frac{1}{2\sigma^2}\|Q^\top t\|^2\right). \end{aligned}$$

而

$$\prod_{i=1}^2 f_{z'_i}(z'_i = t_i) \cdot \prod_{i=3}^n f_{z_i}(z_i = t_i) = \frac{1}{(2\pi)^{n/2}\sigma^n} \exp\left(-\frac{1}{2\sigma^2}\|t\|^2\right).$$

由于 Q 是正交矩阵, 故 $\|Q^\top t\|^2 = \|t\|^2$, 故

$$f_{z'}(z' = t) = \prod_{i=1}^2 f_{z'_i}(z'_i = t_i) \cdot \prod_{i=3}^n f_{z_i}(z_i = t_i).$$

故 $z'_1, z'_2, z_3, \dots, z_n$ 是相互独立的. 又我们有

$$z = (z'_1, z'_2, z_3, \dots, z_n)^\top + (\sqrt{n}\alpha + \sqrt{n}\beta\bar{x}, \beta\sqrt{s_{xx}}, 0, \dots, 0)^\top.$$

故 $z \sim \mathcal{N}(\mu, \Sigma)$, 其中 $\mu = (\sqrt{n}\alpha + \sqrt{n}\beta\bar{x}, \beta\sqrt{s_{xx}}, 0, \dots, 0)^\top$, $\Sigma = \sigma^2 I_n$.

(3) 注意到 $\hat{\beta} = s_{xy}/s_{xx}$, 故

$$\begin{aligned} z_1 &= q_1^\top y = \frac{1}{\sqrt{n}} \sum_{i=1}^n y_i = \sqrt{n}\bar{y}, \\ z_2 &= q_2^\top y = \frac{1}{\sqrt{s_{xx}}} \sum_{i=1}^n (x_i - \bar{x})y_i = \frac{1}{\sqrt{s_{xx}}} \sum_{i=1}^n (x_i - \bar{x})(y_i - \hat{y}_i) = \sqrt{s_{xx}}\hat{\beta}. \end{aligned}$$

(4) 第二题中有 $\hat{y}_i = \bar{y} + \hat{\beta}(x_i - \bar{x})$, 那么有

$$\sum_{i=1}^n (\hat{y}_i - \bar{y})^2 = \sum_{i=1}^n \hat{\beta}^2 (x_i - \bar{x})^2 = \hat{\beta}^2 \cdot s_{xx} = z_2^2.$$

由于 Q 为正交矩阵, 故有 $\|z\|^2 = \|Qy\|^2 = \|y\|^2$, 即

$$\sum_{i=1}^n y_i^2 = \sum_{i=1}^n z_i^2 \implies \sum_{i=1}^n (y_i - \hat{y}_i)^2 = \sum_{i=1}^n y_i^2 - n\bar{y}^2 = \sum_{i=1}^n z_i^2 - z_1^2 = \sum_{i=2}^n z_i^2.$$

故我们有

$$(n-2)s^2 = \sum_{i=1}^n (y_i - \hat{y}_i)^2 = \sum_{i=1}^n (y_i - \bar{y})^2 - \sum_{i=1}^n (\hat{y}_i - \bar{y})^2 = \sum_{i=2}^n z_i^2 - z_2^2 = \sum_{i=3}^n z_i^2.$$

证毕.

(5) 在 (2) 中我们证明了对于 $i \in \{3, 4, \dots, n\}$, $z_i \stackrel{\text{i.i.d.}}{\sim} \mathcal{N}(0, \sigma^2) \implies z_i/\sigma \stackrel{\text{i.i.d.}}{\sim} \mathcal{N}(0, 1)$. 故

$$\frac{(n-2)s^2}{\sigma^2} = \sum_{i=3}^n \left(\frac{z_i}{\sigma}\right)^2 \sim \chi^2(n-2).$$

由于 z_1, z_2 与 z_3, z_4, \dots, z_n 相互独立, 而

$$s^2 = \frac{1}{n-2} \sum_{i=3}^n z_i^2, \quad \hat{\beta} = \frac{z_2}{\sqrt{s_{xx}}}, \quad \hat{\alpha} = \bar{y} - \hat{\beta}\bar{x} = \frac{z_1}{\sqrt{n}} - \frac{z_2}{\sqrt{s_{xx}}}\bar{x}.$$

故 s^2 与 $\hat{\alpha}, \hat{\beta}$ 均相互独立.

(6) 当 $\beta = 0$ 时, $z_2 \sim \mathcal{N}(0, \sigma^2) \implies z_2/\sigma \sim \mathcal{N}(0, 1)$, 由 (5) 知 $(n-2)s^2/\sigma^2 \sim \chi^2(n-2)$. 且 z_2/σ 与 $(n-2)s^2/\sigma^2$ 相互独立, 故

$$t = \frac{\hat{\beta}}{\sqrt{s^2}/\sqrt{s_{xx}}} = \frac{z_2/\sqrt{s_{xx}}}{\sqrt{s^2}/\sqrt{s_{xx}}} = \frac{z_2/\sigma}{\sqrt{\frac{(n-2)s^2}{\sigma^2}/(n-2)}} \sim t(n-2).$$

(7) 当 σ^2 未知时, 给定显著性水平 α , 原假设成立时, $t \sim t(n-2)$, 故由 Neyman-Pearson 原则有

$$\mathbb{P}(|t| \geq c) = \alpha \implies c = t_{\alpha/2}(n-2).$$

其中 $t_{\alpha/2}(n-2)$ 表示自由度为 $n-2$ 的 t 分布上侧 $\alpha/2$ 分位点, 即 $\mathbb{P}(t \geq t_{\alpha/2}(n-2)) = \alpha/2$.

◁