

回归分析 (01714601)

课程主页: <http://staff.ustc.edu.cn/~ynyang/lm2020>

第八讲 简单线性模型的统计推断

2020.3.13

$$\frac{(n-2)\hat{\sigma}^2}{\sigma^2} \sim \chi_{n-2}^2, \text{ 与 } (\hat{a}, \hat{b}) \text{ 独立}$$

预备：正态模型下的统计推断

例1. 假设 $y_1, \dots, y_n \text{ iid} \sim (a, \sigma^2)$, 可表示为线性模型的形式

$$y_i = a + \varepsilon_i, \quad \varepsilon_i, \quad i = 1, \dots, n \text{ iid} \sim (0, \sigma^2)$$

是最简单的线性模型.

$$\text{最小二乘: } \min \sum (y_i - a)^2,$$

$$\text{对 } a \text{ 求导} \Rightarrow \sum (y_i - a) = 0 \Rightarrow LS \text{ 估计 } \hat{a} = \bar{y} \text{ (样本均值)}$$

$$RSS = \sum (y_i - \hat{a})^2 = \sum (y_i - \bar{y})^2$$

$$\hat{\sigma}^2 = RSS / (n-1) = s^2 \text{ (样本方差)}$$

正态假设下, 假设 $y_1, \dots, y_n \text{ iid} \sim (a, \sigma^2)$, 则

(1) $\bar{y} \sim N(a, \sigma^2 / n)$, (2) $(n-1)s^2 / \sigma^2 \sim \chi_{n-1}^2$, (3) \bar{y} 与 s^2 独立.

(4) $\sqrt{n}(\bar{y} - a) / s \sim t_{n-1}$

证明(2)–(3):

$$(n-1)s^2 = \sum (y_i - \bar{y})^2 = \sum (y_i - a - (\bar{y} - a))^2 = \sum \varepsilon_i^2 - n\bar{\varepsilon}^2$$

$$\text{记 } \boldsymbol{\varepsilon} = (\varepsilon_1, \dots, \varepsilon_n)^\top \sim N(0, \sigma^2 I_n)$$

令 A 为一个 n 阶正交矩阵, $A^\top A = AA^\top = I_n$, 其第一行为

$$(1/\sqrt{n}, \dots, 1/\sqrt{n})$$

令 $\mathbf{z} = A\boldsymbol{\varepsilon}$, 则 $\mathbf{z} = A\boldsymbol{\varepsilon} \sim N(0, \sigma^2 I_n)$, 且 $\|\mathbf{z}\|^2 = \|\boldsymbol{\varepsilon}\|^2$,

$$z_1 = (\varepsilon_1 + \dots + \varepsilon_n) / \sqrt{n} = \sqrt{n}\bar{\varepsilon}.$$

$$(n-1)s^2 = \sum_{i=1}^n \varepsilon_i^2 - n\bar{\varepsilon}^2 = \sum_{i=1}^n z_i^2 - z_1^2 = \sum_{i=2}^n z_i^2 \sim \sigma^2 \chi_{n-1}^2,$$

s^2 只与 z_2, \dots, z_n 有关, 而 $\bar{y} = a + \bar{\varepsilon}$ 只与 z_1 有关, 故 \bar{y} 与 s^2 独立.

简单正态线性模型的统计推断

大多数软件假设模型中误差项服从正态分布，以后称之为正态线性模型.

我们仅考虑 b 的推断

命题6. 假设模型 $y_i = a + bx_i + \varepsilon_i, \varepsilon_1, \dots, \varepsilon_n \text{ iid} \sim N(0, \sigma^2)$, 则
给定 $\mathbf{x} = (x_1, \dots, x_n)^\top$ 的条件下

$$(1) \quad \hat{b} \sim N(b, \sigma^2 / s_{xx})$$

$$(2) \quad \frac{(n-2)\hat{\sigma}^2}{\sigma^2} \sim \chi_{n-2}^2, \text{ 且 } \hat{\sigma}^2 \text{ 与 } (\hat{a}, \hat{b}) \text{ 独立}$$

$$(3) \quad \frac{\hat{b} - b}{\sqrt{\widehat{\text{var}}(\hat{b})}} = \frac{\hat{b} - b}{\sqrt{\hat{\sigma}^2 / s_{xx}}} \sim t_{n-2}$$

证明: (1) $y_i | x_i \sim N(a + bx_i, \sigma^2) \Rightarrow$ 给定 x_1, \dots, x_n 时,

$$\hat{b} = \frac{\sum (x_i - \bar{x}) y_i}{\sum (x_i - \bar{x})^2} \sim N(b, \sigma^2 / s_{xx}).$$

(2) $RSS = s_{yy} - s_{xy}^2 / s_{xx} = s_{\varepsilon\varepsilon} - s_{x\varepsilon}^2 / s_{xx}$ (后页验证)

$$\begin{aligned} &= \sum_{i=1}^n (\varepsilon_i - \bar{\varepsilon})^2 - \left(\sum_{i=1}^n (x_i - \bar{x})(\varepsilon_i - \bar{\varepsilon}) \right)^2 / s_{xx} \\ &= \sum_{i=1}^n \varepsilon_i^2 - n\bar{\varepsilon}^2 - \left(\sum_{i=1}^n \left(\frac{x_i - \bar{x}}{\sqrt{s_{xx}}} \right) \varepsilon_i \right)^2 \\ &= \|\boldsymbol{\varepsilon}\|^2 - (\mathbf{u}^\top \boldsymbol{\varepsilon})^2 - (\mathbf{v}^\top \boldsymbol{\varepsilon})^2, \end{aligned}$$

其中 $\mathbf{u}^\top = (1/\sqrt{n}, \dots, 1/\sqrt{n})$,
 $\mathbf{v}^\top = ((x_1 - \bar{x})/\sqrt{s_{xx}}, \dots, (x_n - \bar{x})/\sqrt{s_{xx}})$

$$RSS = \|\boldsymbol{\varepsilon}\|^2 - (\mathbf{u}^\top \boldsymbol{\varepsilon})^2 - (\mathbf{v}^\top \boldsymbol{\varepsilon})^2$$

令 $\boldsymbol{\varepsilon} = (\varepsilon_1, \dots, \varepsilon_n)^\top \sim N(0, \sigma^2 I_n)$, 令 $\mathbf{z} = A\boldsymbol{\varepsilon} \sim N(0, \sigma^2 I_n)$,

其中 A 是一个正交矩阵, 其第一二行分别为 $\mathbf{u}^\top, \mathbf{v}^\top$:

$$A = \begin{pmatrix} \mathbf{u}^\top \\ \mathbf{v}^\top \\ * \end{pmatrix} = \begin{pmatrix} \frac{1}{\sqrt{n}} & \frac{1}{\sqrt{n}} & \cdots & \frac{1}{\sqrt{n}} \\ \frac{x_1 - \bar{x}}{\sqrt{s_{xx}}} & \frac{x_2 - \bar{x}}{\sqrt{s_{xx}}} & \cdots & \frac{x_n - \bar{x}}{\sqrt{s_{xx}}} \\ * & * & \cdots & * \end{pmatrix}$$

$$\text{则 } z_1 = \mathbf{u}^\top \boldsymbol{\varepsilon} = \sqrt{n}\bar{\varepsilon}, \quad z_2 = \mathbf{v}^\top \boldsymbol{\varepsilon} = \sum_{i=1}^n (x_i - \bar{x})\varepsilon_i / \sqrt{s_{xx}},$$

因为 A 是正交矩阵, $\|\mathbf{z}\|^2 = \|\boldsymbol{\varepsilon}\|^2$,

$$RSS = \|\boldsymbol{\varepsilon}\|^2 - (\mathbf{u}^\top \boldsymbol{\varepsilon})^2 - (\mathbf{v}^\top \boldsymbol{\varepsilon})^2 = \|\mathbf{z}\|^2 - z_1^2 - z_2^2 = \sum_{i=3}^n z_i^2$$

$$(n-2)\hat{\sigma}^2 / \sigma^2 = RSS / \sigma^2 = \sum_{i=3}^n z_i^2 / \sigma^2 \sim \chi_{n-2}^2, \text{ 且与 } z_1, z_2 \text{ 独立.}$$

而 \hat{a}, \hat{b} 仅与 z_1, z_2 有关:

$$\hat{b} = b + \frac{\sum (x_i - \bar{x})\varepsilon_i}{\sum (x_i - \bar{x})^2} = b + z_2 / \sqrt{s_{xx}},$$

$$\hat{a} = \bar{y} - \hat{b}\bar{x} = a + b\bar{x} + \bar{\varepsilon} - \hat{b}\bar{x} = a + z_1 / \sqrt{n} - z_2\bar{x} / \sqrt{s_{xx}},$$

所以 $\hat{\sigma}^2$ 与 (\hat{a}, \hat{b}) 独立。

验证: $RSS = s_{\varepsilon\varepsilon} - s_{x\varepsilon}^2 / s_{xx}$

$$(a) \quad y_i - \bar{y} = a + bx_i + \varepsilon_i - (a + b\bar{x} + \bar{\varepsilon}) = b(x_i - \bar{x}) + (\varepsilon_i - \bar{\varepsilon}),$$

$$(b) \quad \text{由 } \hat{a} = \bar{y} - \hat{b}\bar{x} \Rightarrow \hat{y}_i - \bar{y} = \hat{a} + \hat{b}x_i - \bar{y} = \hat{b}(x_i - \bar{x})$$

$$\text{而 } \hat{b} = \frac{\sum (x_i - \bar{x})y_i}{\sum (x_i - \bar{x})^2} = \frac{\sum (x_i - \bar{x})(a + bx_i + \varepsilon_i)}{\sum (x_i - \bar{x})^2} = b + s_{x\varepsilon} / s_{xx},$$

$$\Rightarrow \hat{y}_i - \bar{y} = \hat{b}(x_i - \bar{x}) = b(x_i - \bar{x}) + (x_i - \bar{x})s_{x\varepsilon} / s_{xx}$$

$$\text{所以 } RSS = \sum_{i=1}^n (y_i - \hat{y}_i)^2 = \sum_{i=1}^n (y_i - \bar{y} - (\hat{y}_i - \bar{y}))^2$$

$$= \sum_{i=1}^n ((\varepsilon_i - \bar{\varepsilon}) - (x_i - \bar{x})s_{x\varepsilon} / s_{xx})^2 = \sum_{i=1}^n (\varepsilon_i - \bar{\varepsilon})^2 - s_{x\varepsilon}^2 / s_{xx}$$

(3). 由 (1),(2)知:

$$\frac{\sqrt{s_{xx}}(\hat{b}-b)}{\sigma} \sim N(0,1), \quad \frac{(n-2)\hat{\sigma}^2}{\sigma^2} \sim \chi_{n-2}^2,$$

且两者独立, 所以

$$\frac{\frac{\sqrt{s_{xx}}(\hat{b}-b)}{\sigma}}{\sqrt{\frac{(n-2)\hat{\sigma}^2}{\sigma^2} / (n-2)}} = \frac{\sqrt{s_{xx}}(\hat{b}-b)}{\hat{\sigma}} \sim t_{n-2}$$

plug-in(代入)方差估计

$$\frac{\hat{b} - b}{\sqrt{\text{var}(\hat{b})}} = \frac{\hat{b} - b}{\sigma / \sqrt{s_{xx}}} \sim N(0,1)$$

代入 σ^2 的估计

Wald test

$$\frac{\hat{b} - b}{\sqrt{\widehat{\text{var}}(\hat{b})}} = \frac{\hat{b} - b}{\hat{\sigma} / \sqrt{s_{xx}}} \sim t_{n-2}$$

置信区间

基于事实: $\frac{\sqrt{s_{xx}}(\hat{b} - b)}{\hat{\sigma}} \sim t_{n-2}$

$$\begin{aligned} 1 - \alpha &= P\left(\left|\frac{\sqrt{s_{xx}}(\hat{b} - b)}{\hat{\sigma}}\right| \leq t_{n-2}(\alpha/2)\right) \\ &= P\left(\hat{b} - \frac{\hat{\sigma}}{\sqrt{s_{xx}}}t_{n-2}(\alpha/2) \leq b \leq \hat{b} + \frac{\hat{\sigma}}{\sqrt{s_{xx}}}t_{n-2}(\alpha/2)\right) \end{aligned}$$

(1) b 的 $(1 - \alpha)100\%$ 置信区间: $\left[\hat{b} \mp \frac{\hat{\sigma}}{\sqrt{s_{xx}}}t_{n-2}(\alpha/2) \right]$

(2) $H_0: b = b_0$ (b_0 为已知常数, 通常 $b_0 = 0$)

t -检验统计量: $t = \frac{\sqrt{s_{xx}}(\hat{b} - b_0)}{\hat{\sigma}} \stackrel{H_0}{\sim} t_{n-2}$

回归系数的显著性 t -检验:

$$H_0 : b = 0, \quad t = \frac{\hat{b}}{\sqrt{\hat{\sigma}^2 / s_{xx}}} = \sqrt{n-2} \frac{r}{\sqrt{1-r^2}} \stackrel{H_0}{\sim} t_{n-2}$$

$|t| \geq t_{n-2}(\alpha/2)$ 时拒绝原假设。

验证: 由 $\hat{b} = s_{xy} / s_{xx}$, $\hat{\sigma}^2 = \frac{1}{n-2} RSS = \frac{1}{n-2} (s_{yy} - s_{xy}^2 / s_{xx})$,

$$t = \frac{\hat{b}}{\sqrt{\hat{\sigma}^2 / s_{xx}}} = \frac{s_{xy} / s_{xx}}{\sqrt{\frac{1}{n-2} \frac{s_{yy} - s_{xy}^2 / s_{xx}}{s_{xx}}}} = \sqrt{n-2} \frac{r}{\sqrt{1-r^2}}$$

两样本 t -检验是回归系数显著性检验的特殊情形

$$\begin{aligned} y_1, \dots, y_{n_1} & \text{ iid } \sim N(\mu_1, \sigma^2) & \leftarrow x_1, \dots, x_{n_1} = 1 \\ y_{n_1+1}, \dots, y_{n_1+n_2} & \text{ iid } \sim N(\mu_2, \sigma^2) & \leftarrow x_{n_1+1}, \dots, x_{n_1+n_2} = 0 \end{aligned}$$

应用线性模型: $y_i = a + bx_i + \varepsilon_i$ ($a = \mu_2, b = \mu_1 - \mu_2$)

容易验证:
$$t = \frac{\hat{b}}{\sqrt{\hat{\sigma}^2 / s_{xx}}} = \frac{\bar{y}_1 - \bar{y}_2}{\sqrt{(n_1^{-1} + n_2^{-1}) s^2}}$$

非正态假设下的显著性检验

z -检验:

不假设误差服从正态分布的情形下, 原假设 $H_0: b = 0$ 成立时,

$$t = \frac{\hat{b}}{\sqrt{\hat{\sigma}^2/s_{xx}}} = \sqrt{n-2} \frac{r}{\sqrt{1-r^2}} \approx z = \sqrt{nr} \xrightarrow{d} N(0,1), n \rightarrow \infty$$

$|t| \geq z_{\alpha/2}$, 或 $|z| \geq z_{\alpha/2}$ 时否定原假设。

卡方检验:

卡方检验统计量 $z^2 = nr^2$ 在原假设 $H_0: b = 0$ 下近似服从 χ_1^2 , $n \rightarrow \infty$.

当 $z^2 \geq \chi_1^2(\alpha)$ 时, 在 α 水平下拒绝原假设.

结果的解释：因果还是关联？

$$y = a + bx + \varepsilon, \quad x \text{ 与 } \varepsilon \text{ 是否独立?}$$

- 随机化控制试验或天然试验, x 与 ε 独立

结论：因果, x 每增加一个单位, y 的期望增加 b 个单位。 **Key: 同一个对象。**

- 观察研究, 自变量与误差几乎不可能独立,

结论：关联, 如果一个研究对象的 x 比另外一个大 1 个单位, 则其 y 平均大 b 个单位。 **Key: 不同的对象。**

例如，分析**2001**年人口抽样调查数据，得到妻子教育水平（上学的年数）与丈夫教育水平的回归方程如下：

$$\text{WifeEdLevel} = 5.60 + 0.57 \times \text{HusbandEdLevel} + \text{residual}$$

如果公司送王先生到大学在职培养一年，你是否预期王太太的教育水平会上升**0.57**年？若不是，**0.57**的含义是什么？

这是观察研究而非试验，**结果是关联而不是因果**：

b=0.57的含义是：如果该研究中某人比另外一个人多上一年学，那么这个人的妻子的比另外一人的妻子多上（平均意义上）**0.57**年学。