

回归分析 (01714601)

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

第七讲 简单线性回归模型（续）

2020.3.11

$$R^2 = \frac{VAR(\hat{y})}{VAR(y)}, \text{ } VAR: \text{方差或样本方差}$$

总体复相关系数平方

$$R^2_{\text{总体}} = \frac{\text{var}\left(\overset{\vee}{y}\right)}{\text{var}(y)} = \frac{\text{var}(E(y|x))}{\text{var}(y)}$$

正交分解: $y = E(y|\mathbf{x}) + y - E(y|\mathbf{x})$

方差分解: $\text{var}(y) = \text{var}(E(y|\mathbf{x})) + E(\text{var}(y|\mathbf{x}))$

\mathbf{x} 解释 y 的比例: $R^2 = \text{var}(E(y|\mathbf{x}))/\text{var}(y)$

简单线性模型/不相关化(总体版本): $E(y|x) = a + bx$

正交分解: $y = a + bx + \varepsilon$, 【 ε 与 x 独立】

方差分解: $\sigma_y^2 = \text{var}(y) = \text{var}(a + bx) + \sigma^2 = b^2\sigma_x^2 + \sigma^2 = \rho^2\sigma_y^2 + \sigma^2$

$$R^2 = \frac{\text{var}(a + bx)}{\text{var}(y)} = \rho_{xy}^2 = \frac{\sigma_{xy}^2}{\sigma_{xx}\sigma_{yy}} = \frac{\sigma_{yx}\sigma_{xx}^{-1}\sigma_{xy}}{\sigma_{yy}} = \frac{\sigma_{xy}\sigma_{yy}^{-1}\sigma_{yx}}{\sigma_{xx}} \text{ (参见第4讲P6)}$$

样本复相关系数平方

简单线性模型平方和分解(样本版本):

正交分解: $y_i = \hat{y}_i + e_i = \hat{a} + \hat{b}x_i + e_i, \quad i = 1, \dots, n \quad \left[\sum \hat{y}_i e_i = 0 \right]$

两边求平方和: $s_{yy} = SS_{\text{总}} = \sum (y_i - \bar{y})^2 = \sum (\hat{y}_i - \bar{y})^2 + \sum e_i^2 = s_{\hat{y}\hat{y}} + s_{ee}$
 $\equiv SS_{\text{回}} + \text{RSS}$

注: 样本方差 $s_y^2 = s_{yy} / (n-1), s_{\hat{y}}^2 = s_{\hat{y}\hat{y}} / (n-1)$

由 $\hat{b} = s_{xy} / s_{xx}, \quad \hat{a} = \bar{y} - \hat{b}\bar{x} \Rightarrow$

$$SS_{\text{回}} = \sum (\hat{y}_i - \bar{y})^2 = \sum (\hat{a} + \hat{b}x_i - \bar{y})^2 = \sum (\bar{y} - \hat{b}\bar{x} + \hat{b}x_i - \bar{y})^2 = \hat{b}^2 s_{xx}$$

所以

$$\begin{aligned}SS_{\text{回}} &= \hat{b}^2 s_{xx} = s_{xy}^2 / s_{xx} = r_{xy}^2 s_{yy} \\RSS &= s_{yy} - s_{xy}^2 / s_{xx} = s_{yy} - s_{yx} s_{xx}^{-1} s_{xy} \hat{=} s_{yy \bullet x} = (1 - r_{xy}^2) s_{yy}\end{aligned}$$

$$R^2 = \frac{s_{\hat{y}}^2}{s_y^2} = \frac{s_{\hat{y}\hat{y}}}{s_{yy}} = \frac{SS_{\text{回}}}{SS_{\text{总}}} = \frac{\sum (\hat{y}_i - \bar{y})^2}{\sum (y_i - \bar{y})^2} = r_{xy}^2$$

$$R^2_{\text{样本}} = \frac{\hat{y} \text{的样本方差}}{y \text{的样本方差}}$$

$$R^2 = r_{xy}^2 = \frac{s_{xy}^2}{s_{xx} s_{yy}} = \frac{s_{yx} s_{xx}^{-1} s_{xy}}{s_{yy}} = \frac{s_{xy} s_{yy}^{-1} s_{yx}}{s_{xx}}$$

误差方差估计

因为 $\sigma^2 = \text{var}(\varepsilon_i) = E(\varepsilon_i^2)$, 其中 $\varepsilon_i = y_i - a - bx_i$

而 $e_i = y_i - \hat{a} - \hat{b}x_i$ 可看作是 ε_i 的预测（不是估计）

σ^2 的 “LS” 估计取为:

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

- 虽然 $\hat{\sigma}^2$ 不是由最小二乘法（极小化误差平方和）直接得到的，但通常也称它为LS估计。

- 为什么RSS除以 $n-2$? 样本量为 n , 估计了两个参数。
这样定义的 $\hat{\sigma}^2$ 是 σ^2 的无偏估计。

- $\hat{\sigma}^2 = \frac{n-1}{n-2} (1-r_{xy}^2) s_y^2$ 【回忆误差方差 $\sigma^2 = (1-\rho_{xy}^2)\sigma_y^2$ 】

$$\begin{aligned}\hat{\sigma}^2 &= \frac{1}{n-2} RSS = \frac{1}{n-2} (1-r_{xy}^2) s_{yy} = \frac{n-1}{n-2} (1-r_{xy}^2) s_y^2 \\ &\approx (1-r_{xy}^2) s_y^2\end{aligned}$$

其中 $s_y^2 = s_{yy} / (n-1)$ 为 y_1, \dots, y_n 的样本方差。

总结：简单回归模型的参数及其LS估计

模型(总体): $y = a + bx + \varepsilon$,
 $\varepsilon \sim (0, \sigma^2), \varepsilon$ 与 x 独立

模型(样本): $y_i = a + bx_i + \varepsilon_i$,
 $\varepsilon_i \sim (0, \sigma^2), \varepsilon_i$ 与 x_i 独立

参数	估计
(1) $b = \text{cov}(x, y) / \text{var}(x) = \rho \sigma_y / \sigma_x$ (2) $a = \mu_y - b \mu_x$, (3) $\sigma^2 = (1 - \rho^2) \sigma_y^2$	(1) $\hat{b} = s_{xy} / s_{xx} = r_{xy} s_y / s_x$ (2) $\hat{a} = \bar{y} - \hat{b} \bar{x}$ (3) $\hat{\sigma}^2 = (1 - r_{xy}^2) s_y^2 \times (n - 1) / (n - 2)$
回归函数: $a + bx$ 误差: $\varepsilon = y - (a + bx)$ ε 与 x 独立	拟合值: $\hat{y}_i = \hat{a} + \hat{b} x_i$ 残差 $e_i = y_i - (\hat{a} + \hat{b} x_i)$ $(e_1, \dots, e_n)^T \perp (x_1, \dots, x_n)^T$
$R^2 = \frac{\text{var}(a + bx)}{\text{var}(y)} = \rho^2$	$R^2 = \frac{'\text{var}'(\hat{y})}{'\text{var}'(y)} = \frac{s_{\hat{y}\hat{y}}}{s_{yy}} = r^2$

最小二乘估计的性质

命题3 (无偏性). $E(\hat{b}) = b, E(\hat{a}) = a$.

证明: $\hat{b} = s_{xy} / s_{xx}$

$$\begin{aligned} E(\hat{b} | \mathbf{x}) &= E\left(\frac{\sum (x_i - \bar{x}) y_i}{\sum (x_i - \bar{x})^2} \mid \mathbf{x}\right) = \frac{\sum (x_i - \bar{x}) E(y_i | x_i)}{\sum (x_i - \bar{x})^2} \\ &= \frac{\sum (x_i - \bar{x})(a + bx_i + E(\varepsilon_i | x_i))}{\sum (x_i - \bar{x})^2} = \frac{\sum (x_i - \bar{x})(bx_i)}{\sum (x_i - \bar{x})^2} = b \Rightarrow E(\hat{b}) = b \end{aligned}$$

由 ε 与 \mathbf{x} 独立, $E(\varepsilon_i | x_i) = E(\varepsilon_i) = 0$

另外, $E(\hat{a} | \mathbf{x}) = E(\bar{y} - \hat{b}\bar{x} | \mathbf{x}) = a + b\bar{x} - E(\hat{b} | \mathbf{x})\bar{x} = a$

注: 无偏性需要条件 $E(\varepsilon) = 0$ 以及 ε 与 \mathbf{x} 的独立性, 但不需要方差齐性条件。

命题4(估计的方差和方差的估计).

(1) 给定所有自变量 \mathbf{x} 条件下, $\text{var}(\hat{\mathbf{b}} | \mathbf{x}) = \frac{\sigma^2}{s_{xx}}$.

(2) $E(\hat{\sigma}^2) = \sigma^2$.

证明: $\hat{b} = s_{xy} / s_{xx}$, 因为 $\text{var}(y_i | x_i) = \sigma^2$,

$$\begin{aligned}\text{var}(\hat{b} | \mathbf{x}) &= \text{var}\left(\frac{\sum (x_i - \bar{x}) y_i}{\sum (x_i - \bar{x})^2} \mid \mathbf{x}\right) \\ &= \frac{\sum (x_i - \bar{x})^2 \text{var}(y_i | x_i)}{\left[\sum (x_i - \bar{x})^2\right]^2} = \frac{\sigma^2}{\sum (x_i - \bar{x})^2} = \frac{\sigma^2}{s_{xx}}\end{aligned}$$

下面的计算都是在给定 x_1, \dots, x_n 的条件下:

(2) 下面证明 $E(RSS) = (n-2)\sigma^2$. 因为 $SS_{\square} = \hat{b}^2 s_{xx}$, $RSS = s_{yy} - \hat{b}^2 s_{xx}$

$$(i) E(\hat{b}^2) = \text{var}(\hat{b}) + (E(\hat{b}))^2 = \sigma^2 / s_{xx} + b^2$$

$$(ii) E(s_{yy}) = E \sum_{i=1}^n (y_i - \bar{y})^2 = \sum_{i=1}^n E y_i^2 - E(n\bar{y})^2$$

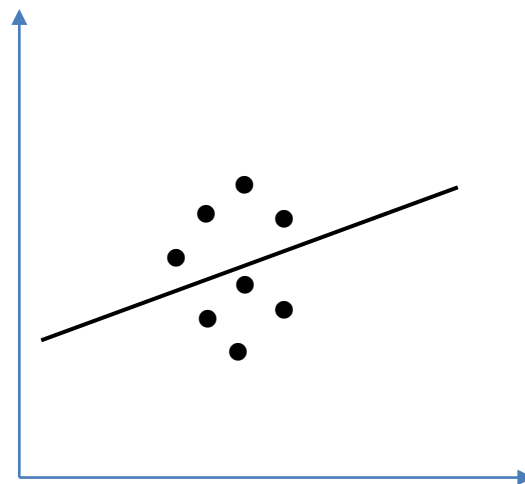
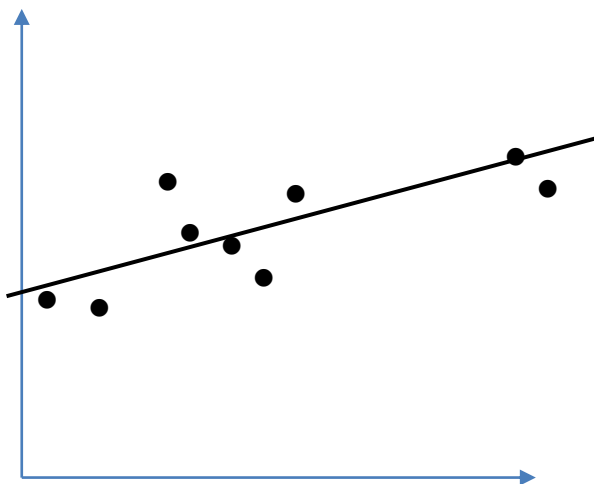
$$= \sum_{i=1}^n (\text{var}(y_i) + (E y_i)^2) - n(\text{var}(\bar{y}) + (E \bar{y})^2)$$

$$= n\sigma^2 + \sum_{i=1}^n (a + b x_i)^2 - n(\sigma^2 / n) - n(a + b \bar{x})^2 = (n-1)\sigma^2 + b^2 s_{xx}$$

$$\Rightarrow E(RSS) = E(s_{yy}) - E(\hat{b}^2 s_{xx}) = (n-1)\sigma^2 + b^2 s_{xx} - s_{xx}(\sigma^2 / s_{xx} + b^2) = (n-2)\sigma^2$$

$$\text{所以 } E(\hat{\sigma}^2 | \mathbf{x}) = \sigma^2, E(\hat{\sigma}^2) = \sigma^2$$

命题4表明, 自变量的样本方差 $s_x^2 = s_{xx} / (n-1)$ 越大, 斜率估计的方差越小。



右图自变量方差比左图小,
回归直线不容易估计准确(方差大)。

LS估计的方差的估计

$$\text{var}(\hat{b} | \mathbf{x}) = \frac{\sigma^2}{\sum (x_i - \bar{x})^2}$$



“plug in” 未知参数 σ^2 的估计

$$\widehat{\text{var}(\hat{b} | \mathbf{x})} = \frac{\hat{\sigma}^2}{\sum (x_i - \bar{x})^2}$$

$$\widehat{se(\hat{b})} = \widehat{sd(\hat{b} | \mathbf{x})} = \frac{\hat{\sigma}}{\sqrt{\sum (x_i - \bar{x})^2}}$$

命题5. (Gauss - Markov定理).

所有 b 的线性无偏估计中，LS估计的方差最小。

证明：设 $\mathbf{u}^\top \mathbf{y}$ 是 b 的任一线性无偏估计，即

$$b = E(\mathbf{u}^\top \mathbf{y} | \mathbf{x}) = \mathbf{u}^\top (\mathbf{1}a + \mathbf{x}b) = a\mathbf{u}^\top \mathbf{1} + b\mathbf{u}^\top \mathbf{x},$$

对任意 a, b 成立。所以 $\mathbf{u}^\top \mathbf{1} = 0$, $\mathbf{u}^\top \mathbf{x} = 1$

$$\begin{aligned}\mathbf{1} &= (1, 1, \dots, 1)^\top, \\ \mathbf{x} &= (x_1, \dots, x_n)^\top, \\ \mathbf{y} &= (y_1, \dots, y_n)^\top \\ \mathbf{y} &= \mathbf{1}a + \mathbf{x}b + \boldsymbol{\varepsilon}\end{aligned}$$

我们要证： $\text{var}(\mathbf{u}^\top \mathbf{y} | \mathbf{x}) = \sigma^2 \mathbf{u}^\top \mathbf{u} \geq \text{var}(\hat{b} | \mathbf{x}) = \sigma^2 / s_{xx}$,

即要证： $\mathbf{u}^\top \mathbf{u} \cdot s_{xx} \geq 1$

由Cauchy不等式：

$$\mathbf{u}^\top \mathbf{u} \cdot s_{xx} = \mathbf{u}^\top \mathbf{u} \cdot (\mathbf{x} - \mathbf{1}\bar{x})^\top (\mathbf{x} - \mathbf{1}\bar{x}) \geq \left(\mathbf{u}^\top (\mathbf{x} - \mathbf{1}\bar{x}) \right)^2 = (\mathbf{u}^\top \mathbf{x})^2 = 1$$

注：注意到自变量的样本方差为 $s_x^2 = s_{xx} / (n-1)$, 对任一(线性)无偏估计 \tilde{b} ,

$$\text{var}(\tilde{b} \mid \mathbf{x}) \geq \sigma^2 / s_{xx},$$

$$\text{var}(\tilde{b} \mid \mathbf{x}) \times s_x^2 \geq \frac{1}{n-1} \sigma^2$$

这说明 b 的估计和 x 的测量不可能同时非常精确
(类似于物理学的测不准原理/不确定性准则)。