

1 公式

$$\text{Var}(X) = \mathbb{E}(X^2) - (\mathbb{E}(X))^2$$

$$\Sigma_x = \mathbb{E}[(x - \mathbb{E}[x])(x - \mathbb{E}[x])^\top]$$

$$y = Ax \Rightarrow \Sigma_y = A\Sigma_x A^\top$$

$$A \perp B \Rightarrow \sigma_{A+B}^2 = \sigma_A^2 + \sigma_B^2$$

$$X = \begin{bmatrix} 1 & 1 & \cdots & 1 \\ x_1 & x_2 & \cdots & x_n \end{bmatrix}^\top, X^\top X = \begin{bmatrix} 1 & \sum x_i \\ \sum x_i & \sum x_i^2 \end{bmatrix}$$

2 Simple Linear Regression

2.1 Model Fitting

Residual: $e_i = y_i - \hat{y}_i$

$$y_i = \beta_0 + \beta_1 x_i + \epsilon_i, \mathbb{E}[\epsilon_i] = 0, \text{Var}(\epsilon_i) = \sigma^2$$

$$\begin{cases} S_{xx} = \sum_{i=1}^n (x_i - \bar{x})^2 \\ S_{xy} = \sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y}) \\ S_{yy} = \sum_{i=1}^n (y_i - \bar{y})^2 \end{cases}$$

$$\hat{\beta}_0 = \bar{y} - \hat{\beta}_1 \bar{x}, \hat{\beta}_1 = \frac{S_{xy}}{S_{xx}}$$

Gauss-Markov Theorem (quality of LSE): Among all estimates that are linear combination of y_1, \dots, y_n and unbiased, the LSE has the smallest variance.

$$\sum_{i=1}^n e_i = 0, \sum_{i=1}^n x_i e_i = 0, \sum_{i=1}^n y_i = \sum_{i=1}^n \hat{y}_i$$

2.2 Statistic Inference & Model Test

Error sum of squares: $\text{SSE} = \sum_{i=1}^n (y_i - \hat{y}_i)^2$

$$s^2 = \text{MSE} = \frac{\text{SSE}}{n-2}$$

Regression sum of squares: $\text{SSR} = \sum_{i=1}^n (\hat{y}_i - \bar{y})^2$

$$S_{yy} = \text{SSR} + \text{SSE}$$

$$\hat{\beta}_0 = \mathbf{l}^\top \mathbf{y}, \hat{\beta}_1 = \mathbf{k}^\top \mathbf{y}, k_i = \frac{x_i - \bar{x}}{S_{xx}}, l_i = \frac{1}{n} - k_i \bar{x}$$

$$\sum_{i=1}^n l_i = 0, \sum_{i=1}^n l_i x_i = 0$$

$$\frac{\hat{\beta}_1 - \beta_1}{\sqrt{\sigma^2/S_{xx}}} \sim \mathcal{N}(0, 1), \frac{\hat{\beta}_0 - \beta_0}{\sqrt{(\sigma^2 \sum_{i=1}^n x_i^2)/(nS_{xx})}} \sim \mathcal{N}(0, 1)$$

$$\frac{\hat{\beta}_1 - \beta_1}{\sqrt{s^2/S_{xx}}} \sim t_{n-2}, \frac{\hat{\beta}_0 - \beta_0}{\sqrt{(s^2 \sum_{i=1}^n x_i^2)/(nS_{xx})}} \sim t_{n-2}$$

$$\frac{\hat{y}_0 - \mathbb{E}[y_0]}{s\sqrt{\frac{1}{n} + \frac{(x_0 - \bar{x})^2}{S_{xx}}}} \sim t_{n-2}$$

$$\frac{y_{\text{new}} - \hat{y}_{\text{new}}}{s\sqrt{1 + \frac{1}{n} + \frac{(x_0 - \bar{x})^2}{S_{xx}}}} \sim t_{n-2}$$

$$\mathcal{R}^2 = \frac{\text{SSR}}{S_{yy}} = \frac{S_{xy}^2}{S_{xx}S_{yy}} = r_{xy}^2$$

ANOVA for SLR:

Source of variation	df	Sum of squares (SS)	MS	F	p-value
Regression	1	$\text{SSR} = \sum (\hat{y}_i - \bar{y})^2$	$\text{MS}_{\text{Reg}} = \frac{\text{SSR}}{1}$	$\frac{\text{MS}_{\text{Reg}}}{s^2}$	
Residual	$n-2$	$\text{SSE} = \sum (y_i - \hat{y}_i)^2$	s^2		
Total	$n-1$	$S_{yy} = \sum (y_i - \bar{y})^2$			

$$\text{MS} = \frac{\text{SS}}{\text{df}}$$

$$\text{SSR} \sim \sigma^2 \chi_1^2, \text{SSE} \sim \sigma^2 \chi_{n-2}^2$$

$$F = \frac{\text{MS}_{\text{reg}}}{s^2} \sim F(1, n-2)$$

$$\mathbb{E}[\text{MS}_{\text{reg}}] = \sigma^2 + \beta_1^2 S_{xx}$$

F 越大说明 MS_{reg} 越大, 也就是 SSR 的贡献更大, 更能说明 $\beta_1 \neq 0$, X effective in explaining variation in Y

2.3 Model Diagnostics

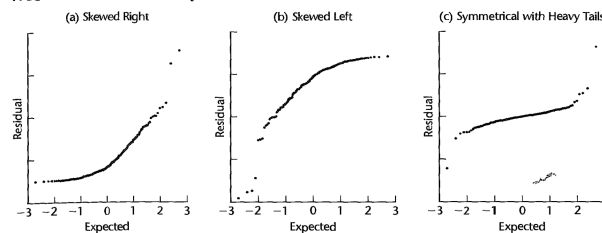
Standardized (semistudentized) residual: $e_i^* = \frac{e_i}{\sqrt{\text{MSE}}}$

Rule of Thumb: $|e_i^*| > 3 \Leftrightarrow \text{outliers}$

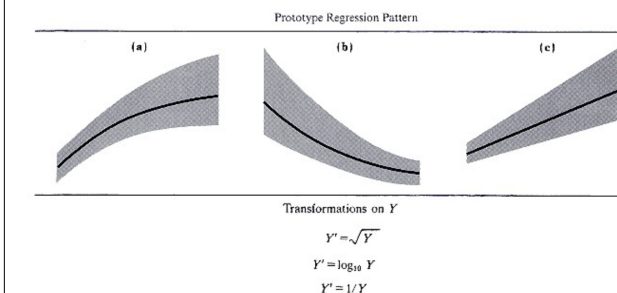
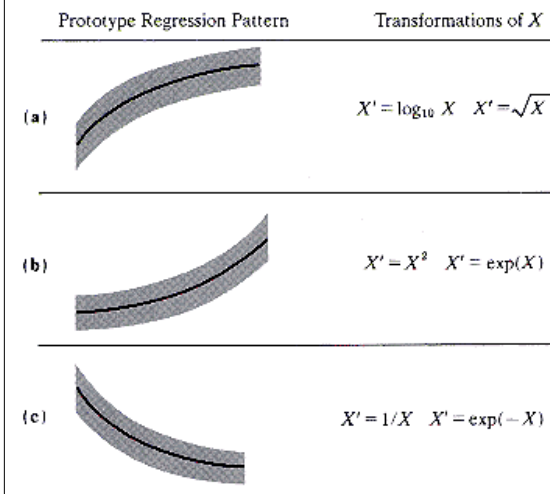
QQ-Plot: k -th smallest, $\sqrt{\text{MSE}} [z(\frac{k-0.375}{n+0.25})]$

residual	expected residual
e_{\min}	$\sqrt{\text{MSE}} \left[z\left(\frac{1-0.375}{n+0.25}\right) \right]$
$e_{2\text{nd smallest}}$	$\sqrt{\text{MSE}} \left[z\left(\frac{2-0.375}{n+0.25}\right) \right]$
\vdots	\vdots
e_{\max}	$\sqrt{\text{MSE}} \left[z\left(\frac{n-0.375}{n+0.25}\right) \right]$

FIGURE 3.9 Normal Probability Plots when Error Term Distribution Is Not Normal.



左边: 往下尾巴大; 右边: 往上尾巴大



3 Calculator

$$\sigma^2_{\mathbf{x}} = \frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})^2$$

$$s^2_{\mathbf{x}} = \sqrt{\frac{1}{n-1} \sum_{i=1}^n (x_i - \bar{x})^2}$$

$$S_{xx} = n \cdot \sigma^2_{\mathbf{x}}, S_{yy} = n \cdot \sigma^2_{\mathbf{y}}$$

$$S_{xy} = \Sigma xy - n \cdot \bar{x} \cdot \bar{y}$$