

Tutorial Notes 3 of MATH3424

1 Summary of course material

1.1 Covariance, Correlation Coefficient

- Covariance of X and Y is defined by

$$\text{Cov}(X, Y) = \frac{\sum_{i=1}^n (y_i - \bar{y})(x_i - \bar{x})}{n - 1} \quad (1)$$

- Correlation and coefficient between X and Y is given by

$$\text{Cor}(X, Y) = \frac{\text{Cov}(Y, X)}{s_y s_x} \quad (2)$$

1.2 Simple Linear Regression Model

- A simple linear model:

$$Y_i = \beta_0 + \beta_1 X_i + \epsilon_i \quad i=1, \dots, n$$

Core assumption:

$$\epsilon_1, \dots, \epsilon_n \stackrel{i.i.d.}{\sim} N(0, \sigma^2)$$

- Parameter estimation (least square estimates/unbiased/standard error/distribution):

$$\hat{\beta}_1 = \frac{\sum_{i=1}^n (y_i - \bar{y})(x_i - \bar{x})}{\sum_{i=1}^n (x_i - \bar{x})^2} \quad \hat{\beta}_0 = \bar{y} - \hat{\beta}_1 \bar{x}$$

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

- Fitted values and residuals

$$\hat{y}_i = \hat{\beta}_0 + \hat{\beta}_1 x_i, \quad e_i = y_i - \hat{y}_i$$

- Measuring the quality of fit:

$$\text{SST} = \sum (y_i - \bar{y})^2 \quad \text{SSR} = \sum (y_i - \bar{y})^2 \quad \text{SSE} = \sum (y_i - \bar{y})^2$$

$$R^2 = \frac{\text{SSR}}{\text{SST}} = 1 - \frac{\text{SSE}}{\text{SST}}$$

- Hypothesis test and confidence intervals:

$$t = \frac{\hat{\beta}_1 - \beta_1^0}{s.e.(\hat{\beta}_1)}$$

$$\hat{\beta}_1 \pm t_{(n-2, \alpha/2)} \times s.e.(\hat{\beta}_1)$$

$$H_0: \beta_1 = \beta_1^0 \quad \text{v.s.} \quad H_1: \beta_1 \neq \beta_1^0$$

if $|t| > t_{(n-2, \alpha/2)}$, reject H_0 .

- No intercept model: difference

2 Questions

1. Consider a simple linear regression model: $y_i = \beta_0 + \beta_1 x_i + \epsilon_i$ for $i = 1, \dots, n$.

- You are given 5 pairs of (x_i, y_i) where y_4 is missing and the fitted line passes through the point $(3, 1.65)$. Find c and then determine $\sum_{i=1}^5 (y_i - \bar{y})^2$.

(\bar{x}, \bar{y})

x_i	1	2	3	4	5
y_i	0.25	1.75	1.5	c	2.5

- Given the following statistics from $n=25$ pairs of (x_i, y_i) :

$$\bar{x} = 0, \quad \hat{\sigma}^2 = 100, \quad \hat{\beta}_0 = 3$$

determine the length of a 98% confidence interval for β_0 .

$$\alpha = 0.02$$

- Given the following statistics from 10 pairs of (x_i, y_i) :

$$\sum_{i=1}^{10} (x_i - \bar{x})^2 = 400, \quad \sum_{i=1}^{10} (y_i - \bar{y})^2 = 425, \quad \sum_{i=1}^{10} (\hat{y}_i - \bar{y})^2 = 225$$

S_{xx} SST SSR

Calculate the test statistic for testing the hypothesis $H_0: \beta_1 = 0$ versus $H_1: \beta_1 \neq 0$ by t test. Write down your conclusion clearly. Set the significance level at $\alpha = 0.05$

Sol: 1. (a)

$$\bar{x} = \frac{1}{5} \sum_{i=1}^5 x_i = 3$$

$$\bar{y} = 1.65 = \frac{0.25 + 1.75 + 1.5 + 2.5 + c}{5} \Rightarrow c = 1.65 \times 5 - 6 = 2.25$$

$$(b) \quad \hat{\beta}_0 \pm t_{(23, 0.01)} \times s.e.(\hat{\beta}_0)$$

≈ 2.5

$$s.e.(\hat{\beta}_0) = \hat{\sigma} \sqrt{\frac{1}{n} + \frac{\bar{x}^2}{\sum_{i=1}^n (x_i - \bar{x})^2}} = \frac{\hat{\sigma}}{\sqrt{n}} = \frac{10}{5} = 2$$

$$3 \pm 2.5 \times 2$$

$$\text{length of } 98\% \text{ CI} \approx 2 \times 2.5 \times 2 \approx 10$$

$$(c) \quad t = \frac{\hat{\beta}_1}{s.e.(\hat{\beta}_1)}, \quad SSE = SST - SSR = 200$$

$$\hat{\sigma}^2 = \frac{SSE}{n-2} = \frac{200}{8} = 25$$

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

$$s.e.(\hat{\beta}_1) = \frac{\hat{\sigma}}{\sqrt{\sum_{i=1}^n (x_i - \bar{x})^2}} = \frac{\hat{\sigma}}{\sqrt{S_{xx}}} = 2 \frac{5}{20} = \frac{1}{4}$$

$$\begin{aligned}\underline{SSR} &= \sum_{i=1}^n (\hat{y}_i - \bar{y})^2 = \sum_{i=1}^n (\hat{\beta}_0 + \hat{\beta}_1 x_i - \hat{\beta}_0 - \hat{\beta}_1 \bar{x})^2 \\ &= \hat{\beta}_1^2 \sum_{i=1}^n (x_i - \bar{x})^2 = \hat{\beta}_1^2 \underline{S_{xx}}\end{aligned}$$

$$\hat{\beta}_1^2 = \frac{SSR}{S_{xx}} = \frac{225}{400}$$

$$|\hat{\beta}| = \sqrt{\hat{\beta}_1^2} = \frac{15}{20} = \frac{3}{4}$$

$$|t| = \frac{\frac{3}{4}}{\frac{1}{4}} = 3 > t_{(8, 0.025)} = 2.306$$

So we reject H_0 .

2. Consider a linear model, for $i = 1, \dots, 3$ $S(\beta_0, \beta_1) = \sum_{i=1}^n (y_i - \beta_0 - \beta_1 x_i)^2$

$$y_i = \beta_0 + i\beta_1 + \epsilon_i \quad \bar{x} = 2 \quad \begin{matrix} i & 1 & 2 & 3 \\ x_i & 1 & 2 & 3 \end{matrix}$$

where ϵ_i follows independent normal distribution with mean 0 and variance σ^2 .

- Find the least squares estimates of β_0 and β_1 in terms of y_i
- Find the $Var(\hat{\beta}_0)$ and $Var(\hat{\beta}_1)$.

Sol: (a) $\hat{\beta}_1 = \frac{S_{xy}}{S_{xx}} = \frac{\sum_{i=1}^3 (i-2) y_i}{\sum_{i=1}^3 (i-2)^2} = \frac{1}{2} (y_3 - y_1)$

$$\hat{\beta}_0 = \bar{y} - \hat{\beta}_1 \bar{x} = \frac{1}{3} (y_1 + y_2 + y_3) - (y_3 - y_1) = \frac{4}{3} y_1 + \frac{1}{3} y_2 - \frac{2}{3} y_3$$

(b) $Var(\hat{\beta}_0) = Var(\frac{4}{3} y_1 + \frac{1}{3} y_2 - \frac{2}{3} y_3) = \frac{16}{9} Var(y_1) + \frac{1}{9} Var(y_2) + \frac{4}{9} Var(y_3)$

$$= \frac{16}{9} \sigma^2 + \frac{1}{9} \sigma^2 + \frac{4}{9} \sigma^2 = \frac{21}{9} \sigma^2 = \frac{7}{3} \sigma^2$$

$$Var(\hat{\beta}_1) = \frac{1}{4} (3\sigma^2 + \sigma^2) = \sigma^2$$

$$\begin{aligned} \text{Cov}(aX + bY, cZ) &= a \text{Cov}(X, cZ) + b \text{Cov}(Y, cZ) \\ &= ac \text{Cov}(X, Z) + bc \text{Cov}(Y, Z) \end{aligned}$$

3. Consider the simple linear regression model: $y_i = \beta_0 + \beta_1 x_i + \epsilon_i$ for $i = 1, \dots, n$. Show that $\text{Cov}(\bar{y}, \hat{\beta}_1) = 0$ and $\text{Cov}(e_i, \hat{\beta}_1) = 0$.

$$\epsilon_i \stackrel{iid}{\sim} N(0, \sigma^2)$$

population covariance

Proof: $\underline{d_i} = \frac{(x_i - \bar{x})}{\sum_{i=1}^n (x_i - \bar{x})^2}$ $\sum x_i = n\bar{x}$

$$\begin{aligned} \text{Cov}(\bar{y}, \hat{\beta}_1) &= \text{Cov}\left(\frac{1}{n} \sum_{i=1}^n y_i, \frac{\sum_{i=1}^n (x_i - \bar{x}) y_i}{\sum_{i=1}^n (x_i - \bar{x})^2}\right) = \text{Cov}\left(\frac{1}{n} \sum_{i=1}^n y_i, \sum_{i=1}^n d_i y_i\right) \\ &= \frac{1}{n} \sum_{i=1}^n \text{Cov}(y_i, d_i y_i) = \frac{1}{n} \sum_{i=1}^n d_i \underbrace{\text{Cov}(y_i, y_i)}_{\text{Var}(y_i) = \sigma^2} = \frac{\sigma^2}{n} \sum_{i=1}^n d_i = 0 \end{aligned}$$

$$\begin{aligned} \text{Cov}(e_i, \hat{\beta}_1) &= \text{Cov}(y_i - \hat{y}_i, \hat{\beta}_1) = \text{Cov}(y_i - (\hat{\beta}_0 + \hat{\beta}_1 x_i), \hat{\beta}_1) \\ &= \text{Cov}(y_i - (\bar{y} - \hat{\beta}_1 \bar{x}) - \hat{\beta}_1 x_i, \hat{\beta}_1) = \text{Cov}(y_i - (\bar{y} + \hat{\beta}_1 (x_i - \bar{x})), \hat{\beta}_1) \\ &= \text{Cov}(y_i - (\bar{y} + d_i \sum_{j=1}^n (x_j - \bar{x}) y_j), \sum_{j=1}^n d_j y_j) \\ &= \underbrace{\text{Cov}(y_i, \sum_{j=1}^n d_j y_j)}_{d_i \text{Cov}(y_i, y_i)} - \underbrace{\text{Cov}(\bar{y}, \sum_{j=1}^n d_j y_j)}_{0} - \underbrace{\text{Cov}(d_i \sum_{j=1}^n (x_j - \bar{x}) y_j, \sum_{j=1}^n d_j y_j)}_{d_i \sum_{j=1}^n (x_j - \bar{x}) d_j \text{Cov}(y_i, y_j)} \\ &= d_i \sigma^2 \end{aligned}$$

$$d_i \sum_{j=1}^n (x_j - \bar{x}) d_j \text{Cov}(y_i, y_j)$$

$$\underbrace{(d_i \sigma^2)}_{\text{Cov}(y_i, y_i)} \sum_{j=1}^n \frac{(x_j - \bar{x})(x_j - \bar{x})}{\sum_{j=1}^n (x_j - \bar{x})^2}$$

$$= d_i \sigma^2 - d_i \sigma^2 = 0$$

$$= d_i \sigma^2$$