

$$Y = X\beta + e$$

$$\begin{bmatrix} Y_1 \\ Y_2 \\ \vdots \\ Y_n \end{bmatrix}_{n \times 1} = \begin{bmatrix} 1 & X_{12} & \dots & X_{1n} \\ 1 & X_{22} & \vdots & \vdots \\ \vdots & \vdots & \ddots & \vdots \\ 1 & X_{n2} & \dots & X_{nn} \end{bmatrix}_{n \times n} \begin{bmatrix} \beta_1 \\ \beta_2 \\ \vdots \\ \beta_k \end{bmatrix}_{n \times 1} + \begin{bmatrix} e_1 \\ \vdots \\ e_n \end{bmatrix}_{n \times 1}$$

$$SSE(\beta) = (Y - X\beta)'(Y - X\beta)$$

$$\frac{\partial SSE(\beta)}{\partial \beta} = -2X'(Y - X\beta) = 0$$

$$X'Y - X'X\beta = 0$$

$$X'Y = X'X\beta, \beta = (X'X)^{-1}X'Y$$

$$X'X = \begin{bmatrix} 1 & \dots & 1 \\ X_{12} & \dots & X_{1n} \\ \vdots & \ddots & \vdots \\ X_{n2} & \dots & X_{nn} \end{bmatrix} \begin{bmatrix} 1 & X_{12} & \dots & X_{1n} \\ \vdots & \vdots & \ddots & \vdots \\ 1 & X_{n2} & \dots & X_{nn} \end{bmatrix} = \begin{bmatrix} n & \sum_{i=1}^n X_{i2} & \dots & \sum_{i=1}^n X_{in} \\ \sum_{i=1}^n X_{i2} & \sum_{i=1}^n X_{i2}^2 & \dots & \sum_{i=1}^n X_{i2}X_{in} \\ \vdots & \vdots & \ddots & \vdots \\ \sum_{i=1}^n X_{in} & \sum_{i=1}^n X_{in}X_{i2} & \dots & \sum_{i=1}^n X_{in}^2 \end{bmatrix}$$

$k=2$

$$X'X = \begin{bmatrix} 1 & \dots & 1 \\ X_1 & \dots & X_n \end{bmatrix} \begin{bmatrix} 1 \\ \vdots \\ X_n \end{bmatrix} = \begin{bmatrix} n & \sum_{i=1}^n X_i \\ \sum_{i=1}^n X_i & \sum_{i=1}^n X_i^2 \end{bmatrix}$$

$$(X'X)^{-1} = \frac{1}{n \sum_{i=1}^n X_i^2 - (\sum_{i=1}^n X_i)^2} \begin{bmatrix} \sum_{i=1}^n X_i^2 & -\sum_{i=1}^n X_i \\ -\sum_{i=1}^n X_i & n \end{bmatrix}$$

$$X'Y = \begin{bmatrix} 1 & \dots & 1 \\ X_1 & \dots & X_n \end{bmatrix} \begin{bmatrix} Y_1 \\ Y_2 \\ \vdots \\ Y_n \end{bmatrix} = \begin{bmatrix} \sum_{i=1}^n Y_i \\ \sum_{i=1}^n X_i Y_i \end{bmatrix}$$

$$\begin{bmatrix} b_1 \\ b_2 \end{bmatrix} = \frac{1}{\sum_{i=1}^n X_i^2 - (\sum_{i=1}^n X_i)^2} \begin{bmatrix} \sum_{i=1}^n X_i^2 & -\sum_{i=1}^n X_i \\ -\sum_{i=1}^n X_i & n \end{bmatrix} \begin{bmatrix} \sum_{i=1}^n Y_i \\ \sum_{i=1}^n X_i Y_i \end{bmatrix}$$

$$= \frac{1}{n \sum_{i=1}^n X_i^2 - (\sum_{i=1}^n X_i)^2} \begin{bmatrix} \sum_{i=1}^n X_i^2 \sum_{i=1}^n Y_i - \sum_{i=1}^n X_i \sum_{i=1}^n X_i Y_i \\ -\sum_{i=1}^n X_i \sum_{i=1}^n Y_i + n \sum_{i=1}^n X_i Y_i \end{bmatrix}$$

$$= \frac{1}{\frac{1}{n} [\sum_{i=1}^n (X_i - \bar{X})^2]} \begin{bmatrix} \frac{1}{n} [\sum_{i=1}^n (Y_i - \bar{Y}) \sum_{i=1}^n (X_i - \bar{X})] \\ \frac{1}{n} [\sum_{i=1}^n (X_i - \bar{X})(Y_i - \bar{Y})] \end{bmatrix}$$

$$= \begin{bmatrix} \bar{Y} - \frac{\sum_{i=1}^n (X_i - \bar{X})(Y_i - \bar{Y})}{\sum_{i=1}^n (X_i - \bar{X})^2} \bar{X} \\ \frac{\sum_{i=1}^n (X_i - \bar{X})(Y_i - \bar{Y})}{\sum_{i=1}^n (X_i - \bar{X})^2} \end{bmatrix} \Rightarrow \begin{cases} b_1 = \bar{Y} - b_2 \bar{X} \\ b_2 = \frac{\sum_{i=1}^n (X_i - \bar{X})(Y_i - \bar{Y})}{\sum_{i=1}^n (X_i - \bar{X})^2} \end{cases}$$

The Ordinary Least Squares (OLS) Estimators

$$b_2 = \frac{\sum (x_i - \bar{x})(y_i - \bar{y})}{\sum (x_i - \bar{x})^2} \quad (2.7)$$

$$b_1 = \bar{y} - b_2 \bar{x} \quad (2.8)$$

where $\bar{y} = \sum y_i / N$ and $\bar{x} = \sum x_i / N$ are the sample means of the observations on y and x .

$$X'Y = \begin{bmatrix} 1 & \dots & 1 \\ X_{12} & \dots & X_{1n} \\ \vdots & \ddots & \vdots \\ X_{n2} & \dots & X_{nn} \end{bmatrix} \begin{bmatrix} Y_1 \\ Y_2 \\ \vdots \\ Y_n \end{bmatrix}$$

(note)

$$\begin{aligned} ① \sum (X - \bar{X})^2 &= \sum (X^2 - 2X\bar{X} + \bar{X}^2) \\ &= \sum X^2 - 2\bar{X} \sum X + n\bar{X}^2 \\ &= \sum X^2 - n\bar{X}^2 = \frac{1}{n} [n \sum X^2 - (\sum X)^2] \end{aligned}$$

$$\begin{aligned} ② \sum (X - \bar{X})(Y - \bar{Y}) &= \sum (XY - \bar{X}Y - X\bar{Y} + \bar{X}\bar{Y}) \\ &= \sum XY - \bar{X} \sum Y - \bar{Y} \sum X + n\bar{X}\bar{Y} \\ &= \sum XY - n\bar{X}\bar{Y} - n\bar{X}\bar{Y} + n\bar{X}\bar{Y} \\ &= \sum XY - n\bar{X}\bar{Y} \\ &= \sum XY - \frac{1}{n} \sum X \sum Y = \frac{1}{n} [n \sum XY - \sum X \sum Y] \end{aligned}$$

$$\begin{aligned} ③ \bar{Y} - b_2 \bar{X} &= \bar{Y} - \frac{\sum XY - n\bar{X}\bar{Y}}{\sum X^2 - n\bar{X}^2} \bar{X} \\ &= \bar{Y} - \frac{\bar{X} \sum XY - n\bar{X}^2 \bar{Y}}{\sum X^2 - n\bar{X}^2} \\ &= \frac{\bar{Y} \sum X^2 - n\bar{X}^2 \bar{Y} - \bar{X} \sum XY + n\bar{X}^2 \bar{Y}}{\sum X^2 - n\bar{X}^2} \\ &= \frac{\bar{Y} \sum X^2 - \bar{X} \sum XY}{\sum X^2 - n\bar{X}^2} \end{aligned}$$

$$y \sim N(x\beta, \sigma^2 I)$$

$$\text{var}(b) = \text{var}[(X'X)^{-1}X'y]$$

$$= (X'X)^{-1}X' \text{var}(y) [(X'X)^{-1}X']'$$

$$= (X'X)^{-1}X' \cdot \sigma^2 I \cdot X(X'X)^{-1}$$

$$= \sigma^2 (X'X)^{-1}$$

$k=2$

$$\text{由上題可得 } (X'X)^{-1} = (X'X)^{-1} = \frac{1}{n \sum_{i=1}^n x_i^2 - (\sum_{i=1}^n x_i)^2} \begin{bmatrix} \sum_{i=1}^n x_i^2 & -\sum_{i=1}^n x_i \\ -\sum_{i=1}^n x_i & n \end{bmatrix}$$

$$= \begin{bmatrix} \frac{\sum_{i=1}^n x_i^2}{n \sum_{i=1}^n (x_i - \bar{x})^2} & \frac{-\sum_{i=1}^n x_i}{n \sum_{i=1}^n (x_i - \bar{x})^2} \\ \frac{-\sum_{i=1}^n x_i}{n \sum_{i=1}^n (x_i - \bar{x})^2} & \frac{n}{n \sum_{i=1}^n (x_i - \bar{x})^2} \end{bmatrix}$$

$$\text{var}(b) = \begin{bmatrix} \frac{\sigma^2 \sum x_i^2}{n \sum (x_i - \bar{x})^2} & \frac{-\bar{x} \sigma^2}{\sum (x_i - \bar{x})^2} \\ \frac{-\bar{x} \sigma^2}{\sum (x_i - \bar{x})^2} & \frac{\sigma^2}{\sum (x_i - \bar{x})^2} \end{bmatrix}$$

$\text{var}(b_1)$ $\text{cov}(b_1, b_2)$ $\text{var}(b_2)$

$$b_1 | x \sim N\left(\beta_1, \frac{\sigma^2 \sum x_i^2}{N \sum (x_i - \bar{x})^2}\right) \quad (2.17)$$

$$b_2 | x \sim N\left(\beta_2, \frac{\sigma^2}{\sum (x_i - \bar{x})^2}\right) \quad (2.18)$$

5.3. Consider the following model that relates the percentage of a household's budget spent on alcohol $WALC$ to total expenditure $TOTEXP$, age of the household head AGE , and the number of children in the household NK .

$$WALC = \beta_1 + \beta_2 \ln(TOTEXP) + \beta_3 NK + \beta_4 AGE + e$$

This model was estimated using 1200 observations from London. An incomplete version of this output is provided in Table 5.6.

TABLE 5.6 Output for Exercise 5.3

Dependent Variable: WALC				
Included observations: 1200				
Variable	Coefficient	Std. Error	t-Statistic	Prob.
C	1.4515	2.2019	0.6592	0.5099
$\ln(TOTEXP)$	2.7648	0.9892	5.7103	0.0000
NK	-1.4549	0.3695	-3.9376	0.0001
AGE	-0.1503	0.0235	-6.4019	0.0000
R-squared	0.0575	Mean dependent var		6.19434
S.E. of regression	6.2167	S.D. dependent var		6.39547
Sum squared resid	46221.62			

a. Fill in the following blank spaces that appear in this table.

- The t-statistic for b_1 .
 - The standard error for b_2 .
 - The estimate b_3 .
 - R^2 .
 - $\hat{\sigma}$.
- b. Interpret each of the estimates b_2 , b_3 , and b_4 .
- c. Compute a 95% interval estimate for β_4 . What does this interval tell you?
- d. Are each of the coefficient estimates significant at a 5% level? Why?
- e. Test the hypothesis that the addition of an extra child decreases the mean budget share of alcohol by 2 percentage points against the alternative that the decrease is not equal to 2 percentage points. Use a 5% significance level.

(a.)

$$\langle i \rangle t = \frac{b_1}{SE(b_1)} = \frac{1.4515}{2.2019} = 0.6592$$

$$\langle ii \rangle SE(b_2) = \frac{2.7648}{5.7103} = 0.4842$$

$$\langle iii \rangle b_3 = 0.3695 \times (-3.9376) = -1.4549$$

$$\langle v \rangle SSE = 46221.62$$

$$MSE = \frac{SSE}{n-k} = \frac{46221.62}{1200-4} = 38.6468$$

$$\hat{\sigma} = \sqrt{MSE} = \sqrt{38.6468} = 6.2167$$

$$\langle iv \rangle SSTO = S_y^2 \times (n-1) = 6.39547^2 \times 1199$$

$$R^2 = 1 - \frac{SSE}{SSTO} = 1 - \frac{46221.62}{6.39547^2 \times 1199} = 0.0575$$

(b) b_2 : 家庭 total expenditure 每增加 1%, $WALC$ 增加 2.7648 百分比 (控制其他条件)
 b_4 : 家庭 head 每多 1 岁, $WALC$ 下降 0.1503 百分比
 b_3 : 家庭 children 每多一位, $WALC$ 下降 1.4549 百分比
 不变下.

(c) c.i. of $b_4 \Rightarrow b_4 \pm Z_{0.025} \cdot SE(b_4) = [-0.1503 - 1.96 \cdot 0.0235, -0.1503 + 1.96 \cdot 0.0235]$
 $= [-0.1964, -0.1042]$

0 在 $[-0.1964, -0.1042]$, 代表年龄对于 $WALC$, 在 95% 信心水准下, 有显著负向影响。

(d) b_1 p-value: $0.5099 > 0.05 \rightarrow b_1$ is insignificant at a 5% level.
 b_2 p-value: $0.0000 < 0.05$
 b_3 p-value: $0.0001 < 0.05$
 b_4 p-value: $0.0000 < 0.05$ } b_2, b_3, b_4 are significant at a 5% level.

(e) $\begin{cases} H_0: \beta_3 = -2 \\ H_1: \beta_3 \neq -2 \end{cases}$
 $\alpha = 0.05$

$$\text{test statistic: } t = \frac{b_3 - (-2)}{SE(b_3)} \sim t_{(1196)} \xrightarrow{CLT} Z$$

$$RR = \{|t| \geq Z_{0.025} = 1.96\}$$

$t = \frac{-1.4549 + 2}{0.3695} = 1.4752$ & $RR \Rightarrow$ do not reject H_0 , the addition of an extra child decrease the mean budget share of alcohol is not significantly differ from 2 percentage.

5.23 The file *cocaine* contains 56 observations on variables related to sales of cocaine powder in northeastern California over the period 1984–1991. The data are a subset of those used in the study Caulkins, J. P. and R. Padman (1993), "Quantity Discounts and Quality Premium for Illicit Drugs," *Journal of the American Statistical Association*, 88, 748–757. The variables are

PRICE = price per gram in dollars for a cocaine sale
 QUANT = number of grams of cocaine in a given sale
 QUAL = quality of the cocaine expressed as percentage purity
 TREND = a time variable with 1984 = 1 up to 1991 = 8
 Consider the regression model

$$PRICE = \beta_1 + \beta_2 QUANT + \beta_3 QUAL + \beta_4 TREND + e$$

- What signs would you expect on the coefficients β_2 , β_3 , and β_4 ?
- Use your computer software to estimate the equation. Report the results and interpret the coefficient estimates. Have the signs turned out as you expected?
- What proportion of variation in cocaine price is explained jointly by variation in quantity, quality, and time?
- It is claimed that the greater the number of sales, the higher the risk of getting caught. Thus, sellers are willing to accept a lower price if they can make sales in larger quantities. Set up H_0 and H_1 that would be appropriate to test this hypothesis. Carry out the hypothesis test.
- Test the hypothesis that the quality of cocaine has no influence on expected price against the alternative that a premium is paid for better-quality cocaine.
- What is the average annual change in the cocaine price? Can you suggest why price might be changing in this direction?

(a) I expect the coefficient β_2 will be negative, the coefficient β_3 will be positive, and the coefficient β_4 will be uncertain.

(b) $PRICE = 90.84669 - 0.05997 QUANT + 0.11621 QUAL - 2.35458 TREND$
 Call:
 lm(formula = PRICE ~ QUANT + QUAL + TREND, data = DATA)

Residuals:
 Min 1Q Median 3Q Max
 -43.479 -12.014 -3.743 13.969 43.753

Coefficients:
 Estimate Std. Error t value Pr(>|t|)
 (Intercept) 90.84669 8.58025 10.588 1.39e-14 ***
 QUANT -0.05997 0.01018 -5.892 2.85e-07 ***
 QUAL 0.11621 0.20326 0.572 0.5700
 TREND -2.35458 1.38612 -1.699 0.0954 .

 Signif. codes: 0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

Residual standard error: 20.06 on 52 degrees of freedom
 Multiple R-squared: 0.5097, Adjusted R-squared: 0.4814
 F-statistic: 18.02 on 3 and 52 DF, p-value: 3.806e-08

Holding other variables constant, for every increase in QUANT, the PRICE drops by 0.05997, as expected.
 For every increase in QUAL, the PRICE increases by 0.11621, as expected.
 For every increase in TREND, the PRICE drops by -2.35458.

(c) $R^2 = \frac{55R}{55T0} = 50.97\%$

(d) $\begin{cases} H_0: \beta_2 \geq 0 \\ H_1: \beta_2 < 0 \text{ (left tail)} \end{cases}$

$\alpha = 0.05$

test statistic: $t = \frac{b_2}{SE(b_2)} \sim t(56-4)$

$RR = \{t \leq -t_{0.05}(52) = -1.675\}$

$t = \frac{-0.05997}{0.01018} = -5.891 \in RR$, reject H_0 , accept H_1 , there is a significant evidence to support the claim that the greater number of sales, the lower price they made.

(e) $\begin{cases} H_0: \beta_3 = 0 \\ H_1: \beta_3 > 0 \text{ (right tail)} \end{cases}$

$\alpha = 0.05$

test statistic: $t = \frac{b_3}{SE(b_3)} \sim t(52)$

p-value = 0.57 > 0.05, we do not reject H_0 , there isn't a significant evidence to support the claim that a premium is paid for better-quality cocaine.

(f) Holding other variables constant, every increase in the trend results in a drop of 2.3548, which I think it may be due to the increase in supply, changes in demand, and changes in policy or law enforcement.