Let
$$k=2$$
, then $y_i = \beta_1 + \beta_2 \chi_i + e_i$

$$\Rightarrow \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{bmatrix} = \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_n \end{bmatrix} \begin{bmatrix} \beta_1 \\ \beta_2 \end{bmatrix} + \begin{bmatrix} e_1 \\ e_2 \\ \vdots \\ e_n \end{bmatrix}$$

$$\Rightarrow \chi'\chi = \begin{bmatrix} 1 & 1 & \cdots & 1 \\ \chi_1 & \chi_2 & \cdots & \chi_n \end{bmatrix} \begin{bmatrix} 1 & \chi_1 \\ \chi_2 & \cdots & \chi_n \end{bmatrix} = \begin{bmatrix} n & \sum \chi_{\bar{i}} \\ \sum \chi_{\bar{i}} & \sum \chi_{\bar{i}} \end{bmatrix}_{2,\chi_2}$$

$$X' Y = \begin{bmatrix} 1 & \cdots & 1 \\ \chi_1 \chi_2 & \cdots & \chi_n \end{bmatrix} \begin{bmatrix} \gamma_1 \\ \gamma_2 \\ \vdots \\ \gamma_n \end{bmatrix} = \begin{bmatrix} \Sigma \gamma_2 \\ \Sigma \chi_1 \gamma_2 \end{bmatrix}_{2 \times 1}$$

$$\exists \left(\chi'\chi\right)^{-1} = \frac{1}{n \sum \chi_{i}^{2} - \left(\sum \chi_{i}\right)^{2}} \left[\sum \chi_{i}^{2} - \sum \chi_{i}^{2}\right]$$

$$\Rightarrow b = (X'X)^{-1}(X'Y) = \frac{1}{n \sum \chi_{i}^{2} - (\sum \chi_{i})^{2}} \left(\sum \chi_{i}^{2} - \sum \chi_{i}\right) \left(\sum \chi_{i} \chi_{i}^{2}\right)$$

$$=\frac{1}{h \sum \chi_i^2 - (\sum \chi_i)^2} \cdot \left[\sum \chi_i^2 \sum \chi_i \sum \chi_i \sum \chi_i \chi_i - \sum \chi_i \sum \chi_i \chi_i \right]$$

Set
$$\overline{X} = \frac{1}{h} \geq \overline{\chi}_i$$
, $\overline{Y} = \frac{1}{h} \sum \overline{\chi}_i$

then
$$b_{2} = \frac{-\sum \chi_{i} \sum \gamma_{i} + n \sum \chi_{i} \gamma_{i}}{n \sum \chi_{i}^{2} - (\sum \chi_{i})^{2}} = \frac{\sum \chi_{i} \gamma_{i} - n (\frac{1}{n} \sum \chi_{i}) \cdot (\frac{1}{n} \sum \chi_{i})}{\sum \chi_{i}^{2} - n \cdot (\frac{1}{n} \sum \chi_{i})^{2}}$$

$$= \frac{\sum \chi_{i} \gamma_{i} - n \overline{\chi}}{\sum \chi_{i}^{2} - n \overline{\chi}^{2}} = \frac{\sum (\chi_{i} - \overline{\chi})(\gamma_{i} - \overline{\gamma})}{\sum (\chi_{i} - \overline{\chi})^{2}}$$

$$= \frac{\sum \chi_{i} \gamma_{i} - n \overline{\chi}^{2}}{\sum \chi_{i}^{2} - n \overline{\chi}^{2}} = \frac{\sum (\chi_{i} - \overline{\chi})(\gamma_{i} - \overline{\gamma})}{\sum (\chi_{i} - \overline{\chi})^{2}}$$

and
$$b_1 = \overline{y} - b_2 \overline{x}$$

$$V_{nr}(b) = \sigma^{2} \left(\chi' \chi \right)^{-1} = \frac{\sigma^{2}}{n \sum \pi_{i}^{2} - \left(\sum \pi_{i} \right)^{2}} \left(-\sum \pi_{i} \right) = \frac{\left(v_{n}(b_{1},b_{2}|\pi) \right) \left(v_{n}(b_{2}|\pi) \right)}{\left(c_{n}(b_{1},b_{2}|\pi) \right) \left(c_{n}(b_{1},b_{2}|\pi) \right)}$$

$$Var(b_{1}|X) = G^{2} \left[\sum_{i} \chi_{i}^{2} \right]$$

$$= \int Var(b_{2}|X) = G^{2} \left[\sum_{i} (\chi_{i} - \overline{\chi})^{2} \right]$$

$$= \int Var(b_{2}|X) = G^{2} \left[\sum_{i} (\chi_{i} - \overline{\chi})^{2} \right]$$

$$= \int Var(b_{1}|X) = G^{2} \left[\sum_{i} (\chi_{i} - \overline{\chi})^{2} \right]$$

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

Dependent Variable: WALC

Output for Exercise 5.3

Included observations: 1200				
Variable	Coefficient	Std. Error	t-Statistic	Prob.
C	1.4515	2.2019		0.5099
ln(TOTEXP)	2.7648		5.7103	0.0000
NK		0.3695	-3.9376	0.0001
AGE	-0.1503	0.0235	-6.4019	0.0000
R-squared	Mean dependent var			6.19434
S.E. of regression	S.D. dependent var			6.39547

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

46221.62

- i. The *t*-statistic for b_1 .
- ii. The standard error for b_2 .
- iii. The estimate b_3 .
- iv. R^2 .

Sum squared resid

v.
$$\hat{\sigma}$$
.

i.

$$t = \frac{h_1}{se(b_1)} = \frac{1.4515}{2.2019} \approx 0.6592$$

ii.

$$SE = \frac{b_2}{t} = \frac{2.7148}{5.7103} \approx 0.4842$$

iii.

$$b_3 = t \times JE = -3.4376 \times 0.3695 \approx -1.4549$$

iv.

$$R^2 = 1 - \frac{55E}{517} = 1 - \frac{46521.12}{(639547)^2 \times (1041)} \approx 0.0515$$

$$V_1 = \frac{1}{104} = \frac{1.4515}{104} \approx 0.2116$$

$$V_2 = \frac{1}{104} = \frac{1.4515}{104} \approx 0.2116$$

b. Interpret each of the estimates b_2 , b_3 , and b_4 .

b2=2.7648, if the household's total spending goes up by 1%, the alcohol budget share goes up by about 2.76 percentage points.

b3=0.3695, Each extra child in the family increases the alcohol budget

share by about 0.37 percentage points.

b4=-0.1503, if the head of the household gets one year older, the alcohol budget share goes down by about 0.15 percentage points.

c. Compute a 95% interval estimate for β_4 . What does this interval tell you?

b4 ± 1.96 × 5 E = -0.1503 ± 1.96 × 0.0235 = -0.1503 ± 0.0461

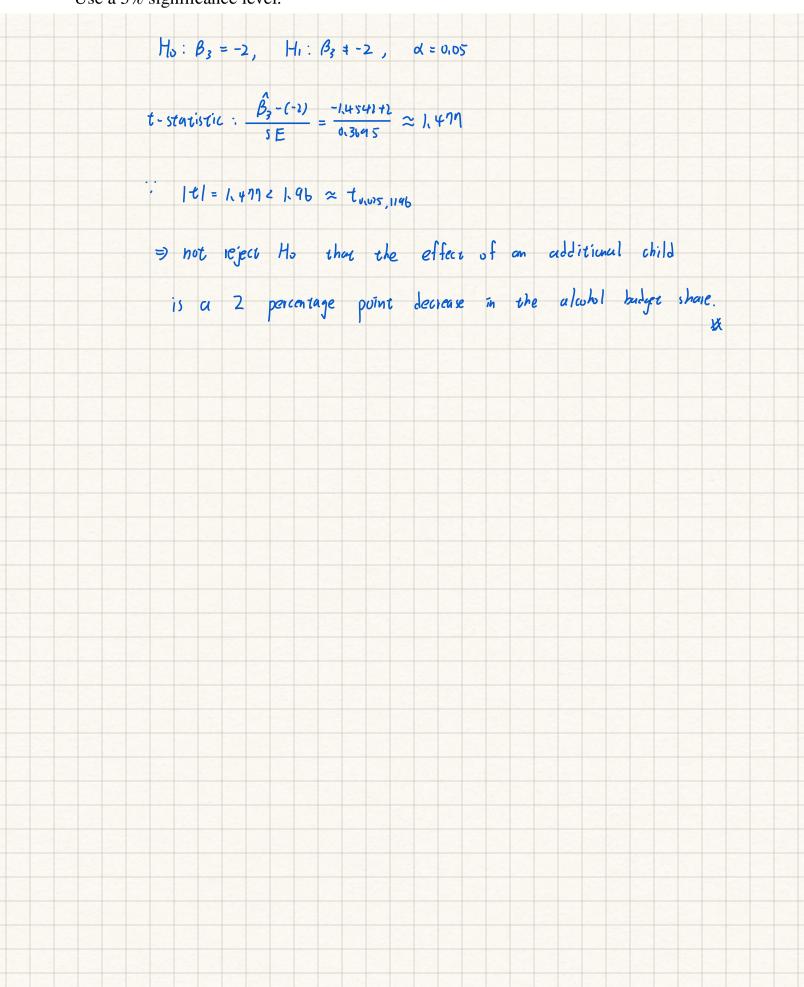
⇒ [-0.1964, -0.1042]

★

d. Are each of the coefficient estimates significant at a 5% level? Why?

At the 5% significance level, all coefficient estimates except for the intercept are startiscally significant. This is because their products are all less than aus, in dicorting strong evidence against the null hyphothesis than the coefficients are zero. The intercept, with a produce of association, is not startiscally significant at the 5% level.

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.



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 Premia 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

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

a. What signs would you expect on the coefficients β_2 , β_3 , and β_4 ?

```
> model <- lm(price ~ quant + qual + trend, data = cocaine)</p>
                                                                                B, =-0,0599720
> # 查看結果摘要
> summary(model)
lm(formula = price ~ quant + qual + trend, data = cocaine)
                                                                      7
Residuals:
   Min
             1Q Median
-43.479 -12.014 -3.743 13.969 43.753
Coefficients:
            Estimate Std. Error t value Pr(>|t|)
                         8.58025 10.588 1.39e-14 ***
(Intercept) 90.84669
            -0.05997
                         0.01018 -5.892 2.85e-07 ***
guant
                                  0.572
             0.11621
                                            0.5700
                         0.20326
gual
                         1.38612 -1.699
                                            0.0954
            -2.35458
trend
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
F-statistic: 18.02 on 3 and 52 DF, p-value: 3.806e-08
```

b. Use your computer software to estimate the equation. Report the results and interpret the coefficient estimates. Have the signs turned out as you expected?

The signs of the coefficients are consistent with expectations, though

not all are statistically significant.

c. What proportion of variation in cocaine price is explained jointly by variation in quantity, quality, and time?

d. 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

e. 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.

f. What is the average annual change in the cocaine price? Can you suggest why price might be changing in this direction?

```
From the regression output at (a) => the coefficient on the trend variable is -2.355 gs.

There is evidence that the price is declining over time, possibly due to increasing supply, improved distribution efficiency, or shifting market dynamics
```