

4.4 The general manager of a large engineering firm wants to know whether the experience of technical artists influences their work quality. A random sample of 50 artists is selected. Using years of work experience (*EXPER*) and a performance rating (*RATING*, on a 100-point scale), two models are estimated by least squares. The estimates and standard errors are as follows:

Model 1:

$$\widehat{RATING} = 64.289 + 0.990EXPER \quad N = 50 \quad R^2 = 0.3793 \\ (\text{se}) \quad (2.422) \quad (0.183)$$

Model 2:

$$\widehat{RATING} = 39.464 + 15.312 \ln(EXPER) \quad N = 46 \quad R^2 = 0.6414 \\ (\text{se}) \quad (4.198) \quad (1.727)$$

CHAPTER 4 Prediction, Goodness-of-Fit, and Modeling Issues

- Sketch the fitted values from Model 1 for *EXPER* = 0 to 30 years.
- Sketch the fitted values from Model 2 against *EXPER* = 1 to 30 years. Explain why the four artists with no experience are not used in the estimation of Model 2.
- Using Model 1, compute the marginal effect on *RATING* of another year of experience for (i) an artist with 10 years of experience and (ii) an artist with 20 years of experience.
- Using Model 2, compute the marginal effect on *RATING* of another year of experience for (i) an artist with 10 years of experience and (ii) an artist with 20 years of experience.
- Which of the two models fits the data better? Estimation of Model 1 using just the technical artists with some experience yields $R^2 = 0.4858$.
- Do you find Model 1 or Model 2 more reasonable, or plausible, based on economic reasoning? Explain.

d.

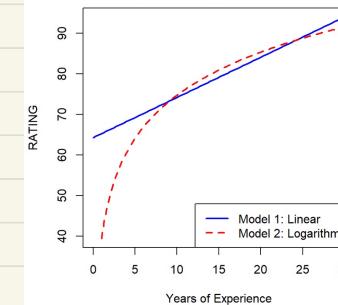
$$\frac{\partial RATING}{\partial EXPER} = 15.312 \times \frac{1}{EXPER}$$

$$EXPER = 10 \quad 15.312 / 10 = 1.5312$$

$$EXPER = 20 \quad 15.312 / 20 = 0.7656$$

a. b.

Comparison of Fitted Models



$\ln(EXPER)$ 要求 $EXPER > 0$

$\ln(0)$ 沒有定義

⇒ 有 4 名不包含在內

c.

$$EXPER = 10 \quad \text{marginal effect} = \frac{\partial RATING}{\partial EXPER} = 0.99$$

$$EXPER = 20 \quad \text{marginal effect} = 0.99$$

e. model 2 $R^2 >$ model 1 R^2

model 2 better

f. model 2 較合理，考慮邊際報酬遞減現象

- 4.28 The file *wa-wheat.dat* contains observations on wheat yield in Western Australian shires. There are 48 annual observations for the years 1950–1997. For the Northampton shire, consider the following four equations:

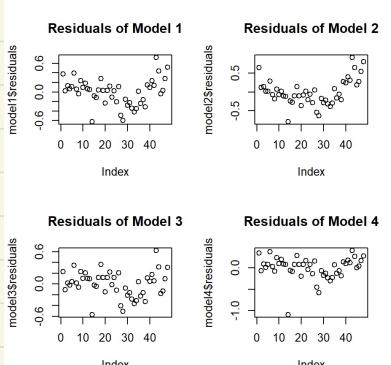
$$YIELD_t = \beta_0 + \beta_1 TIME + e_t$$

$$YIELD_t = \alpha_0 + \alpha_1 \ln(TIME) + e_t$$

$$YIELD_t = \gamma_0 + \gamma_1 TIME^2 + e_t$$

$$\ln(YIELD_t) = \phi_0 + \phi_1 TIME + e_t$$

- Estimate each of the four equations. Taking into consideration (i) plots of the fitted equations, (ii) plots of the residuals, (iii) error normality tests, and (iv) values for R^2 , which equation do you think is preferable? Explain.
- Interpret the coefficient of the time-related variable in your chosen specification.
- Using your chosen specification, identify any unusual observations, based on the studentized residuals, *LEVERAGE*, *DFBETAS*, and *DFFITS*.
- Using your chosen specification, use the observations up to 1996 to estimate the model. Construct a 95% prediction interval for *YIELD* in 1997. Does your interval contain the true value?

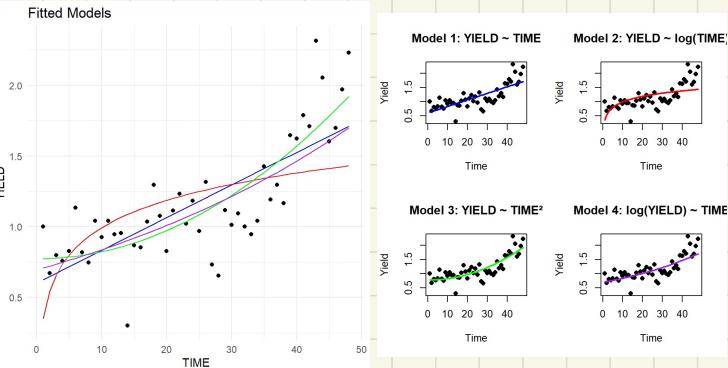


```
> shapiro.test(model1$residuals)
Shapiro-Wilk normality test
data: model1$residuals
W = 0.98236, p-value = 0.6792

> shapiro.test(model2$residuals)
Shapiro-Wilk normality test
data: model2$residuals
W = 0.96657, p-value = 0.1856

> shapiro.test(model3$residuals)
Shapiro-Wilk normality test
data: model3$residuals
W = 0.98589, p-value = 0.8266

> shapiro.test(model4$residuals)
Shapiro-Wilk normality test
data: model4$residuals
W = 0.86894, p-value = 7.205e-05
```



Model	R^2
Linear	0.5778369
Linear-Log	0.3385733
Quadratic	0.6890101
Log-Linear	0.5073566

model 3, R^2 最大且不拒绝常態假設

b. $time^2$ 增加 1 單位, yield 增加 4.986×10^{-4}

C. Studentized Residuals (学生化殘差)

Unusual observations: $| \text{std-res} | > 2$

$\Rightarrow 14, 28, 43$

k : 變數數量 (包含截距)

n : 樣本數

Leverage (樁桿值, 影響回歸的程度)

? $\text{lev} > 2 \times k/n$

$\Rightarrow 45, 41, 47, 49$

DFBETAS (單個變數影響回歸系數的變化程度)

$| \text{dfb} | > 2 / \sqrt{n}$

$\Rightarrow 14, 43, 44, 48$ (對 the 变数) 6.14 (對截距)

DFFITS (刪除該觀測值後, 模型擬合值的變化程度)

$| \text{dff} | > 2 \times \sqrt{k/n}$

d. $[1.372403, 2.389819]$, Yes, the value = 2.2319

$\Rightarrow 14, 43, 48$

4.29 Consider a model for household expenditure as a function of household income using the 2013 data from the Consumer Expenditure Survey, *cex5_small*. The data file *cex5* contains more observations. Our attention is restricted to three-person households, consisting of a husband, a wife, plus one other. In this exercise, we examine expenditures on a staple item, food. In this extended example, you are asked to compare the linear, log-log, and linear-log specifications.

- Calculate summary statistics for the variables: *FOOD* and *INCOME*. Report for each the sample mean, median, minimum, maximum, and standard deviation. Construct histograms for both variables. Locate the variable mean and median on each histogram. Are the histograms symmetrical and “bell-shaped” curves? Is the sample mean larger than the median, or vice versa? Carry out the Jarque–Bera test for the normality of each variable.
- Estimate the linear relationship $FOOD = \beta_1 + \beta_2 INCOME + e$. Create a scatter plot *FOOD* versus *INCOME* and include the fitted least squares line. Construct a 95% interval estimate for β_2 . Have we estimated the effect of changing income on average *FOOD* relatively precisely, or not?
- Obtain the least squares residuals from the regression in (b) and plot them against *INCOME*. Do you observe any patterns? Construct a residual histogram and carry out the Jarque–Bera test for normality. Is it more important for the variables *FOOD* and *INCOME* to be normally distributed, or that the random error e be normally distributed? Explain your reasoning.
- Calculate both a point estimate and a 95% interval estimate of the elasticity of food expenditure with respect to income at *INCOME* = 19, 65, and 160, and the corresponding points on the fitted line, which you may treat as not random. Are the estimated elasticities similar or dissimilar? Do the interval estimates overlap or not? As *INCOME* increases should the income elasticity for food increase or decrease, based on Economics principles?
- For expenditures on food, estimate the log-log relationship $\ln(FOOD) = \gamma_1 + \gamma_2 \ln(INCOME) + e$. Create a scatter plot for $\ln(FOOD)$ versus $\ln(INCOME)$ and include the fitted least squares line. Compare this to the plot in (b). Is the relationship more or less well-defined for the log-log model

```
> jarque.bera.test(cex5_small$food)
    Jarque Bera Test
data: cex5_small$food
X-squared = 648.65, df = 2, p-value < 2.2e-16

> jarque.bera.test(cex5_small$income)
    Jarque Bera Test
data: cex5_small$income
X-squared = 148.21, df = 2, p-value < 2.2e-16
```

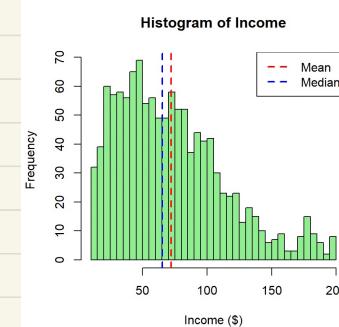
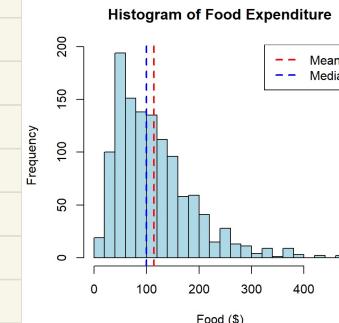
$p\text{-value} < 0.05$

→ 不符合常態分布

a.

```
> print(stats)
```

Variable	Mean	Median	Min	Max	SD
Food	114.44311	99.80	9.63	476.67	72.65750
Income	72.14264	65.29	10.00	200.00	41.65228

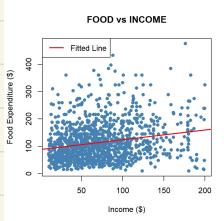


Food & Income
 $\rightarrow \text{mean} > \text{median}$
 right-skewed
 they're not symmetrical
 and "bell-shaped"

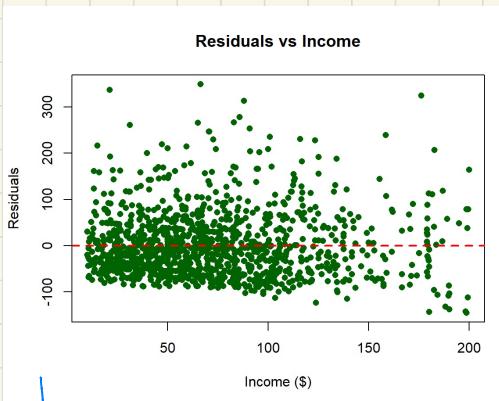
b.

$$CI: [0.2519215, 0.455452]$$

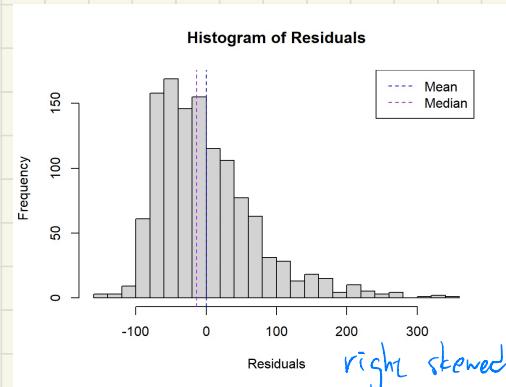
N, R^2 only 0.04228



c



未呈隨機分布



right skewed

Jarque Bera Test

```
data: residuals
X-squared = 624.19, df = 2, p-value < 2.2e-16
```

p-value < 0.05

→ 不符合常態分布

E的常態性較重要 → $E(e)=0$ 、若無、同負
才符合假定要求

d.

	Income	Fitted_Food	Elasticity	Lower_CI	Upper_CI
(Intercept)	19	95.38155	0.07145038	0.05217475	0.09072601
(Intercept)1	65	111.88114	0.20838756	0.15216951	0.26460562
(Intercept)2	160	145.95638	0.39319883	0.28712305	0.49927462

dissimilar

not overlap

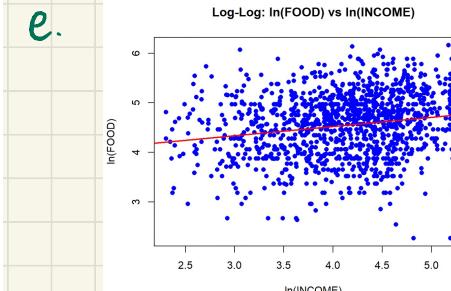
根據經濟原理，food應隨著 income 上升
花費比例減少 (E 下降)

relative to the linear specification? Calculate the generalized R^2 for the log-log model and compare it to the R^2 from the linear model. Which of the models seems to fit the data better?

- f. Construct a point and 95% interval estimate of the elasticity for the log-log model. Is the elasticity of food expenditure from the log-log model similar to that in part (d), or dissimilar? Provide statistical evidence for your claim.
- g. Obtain the least squares residuals from the log-log model and plot them against $\ln(\text{INCOME})$. Do you observe any patterns? Construct a residual histogram and carry out the Jarque–Bera test for normality. What do you conclude about the normality of the regression errors in this model?
- h. For expenditures on food, estimate the linear-log relationship $\text{FOOD} = \alpha_1 + \alpha_2 \ln(\text{INCOME}) + e$. Create a scatter plot for FOOD versus $\ln(\text{INCOME})$ and include the fitted least squares line. Compare this to the plots in (b) and (e). Is this relationship more well-defined compared to the others? Compare the R^2 values. Which of the models seems to fit the data better?
- i. Construct a point and 95% interval estimate of the elasticity for the linear-log model at $\text{INCOME} = 19, 65$, and 160 , and the corresponding points on the fitted line, which you may treat as not random. Is the elasticity of food expenditure similar to those from the other models, or dissimilar? Provide statistical evidence for your claim.
- j. Obtain the least squares residuals from the linear-log model and plot them against $\ln(\text{INCOME})$. Do you observe any patterns? Construct a residual histogram and carry out the Jarque–Bera test for normality. What do you conclude about the normality of the regression errors in this model?
- k. Based on this exercise, do you prefer the linear relationship model, or the log-log model or the linear-log model? Explain your reasoning.

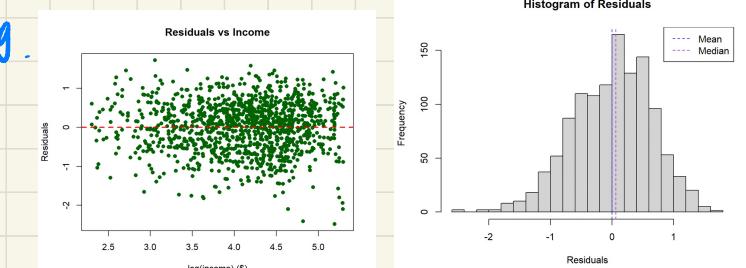
	Income	Fitted_Food	Elasticity	Lower_CI	Upper_CI
(Intercept)	19	75.75443	0.1863054	0.1293432	0.2432675
(Intercept)1	65	95.26313	0.1863054	0.1293432	0.2432675
(Intercept)2	160	112.67012	0.1863054	0.1293432	0.2432675

log-log model 的 ϵ 是固定的



Model	R_squared
1 Linear	0.04228120
2 Log-Log	0.03322915

R^2 都很低
都沒有很好 fit the data



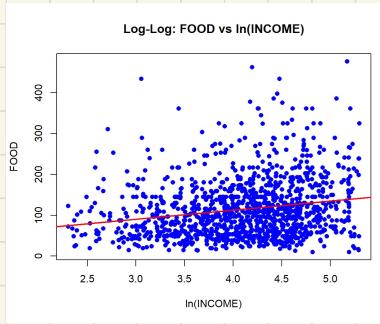
Jarque Bera Test

```
data: residuals
X-squared = 25.85, df = 2, p-value = 2.436e-06
```

p-value < 0.05

→ 不符合常態分布

h.



Model	R_squared
1 Linear	0.04228120
2 Log-Log	0.03322915
3 Linear-Log	0.03799984

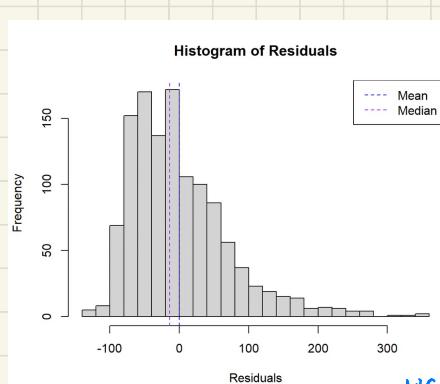
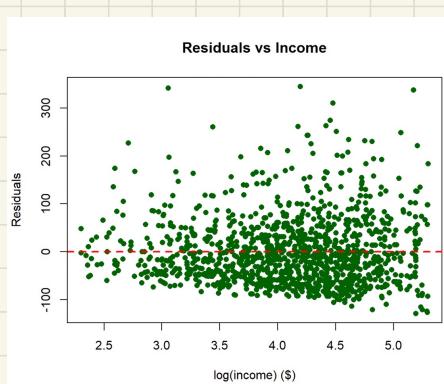
← 在兩個中間 (都不對)
Linear R² 最高

i.

	Income	Fitted_Food	Elasticity	Lower_CI	Upper_CI
(Intercept)	19	88.89788	0.2495828	0.1784009	0.3207648
(Intercept)1	65	116.18722	0.1909624	0.1364992	0.2454256
(Intercept)2	160	136.17332	0.1629349	0.1164652	0.2094046

Income↑, E↓
⇒ 貧富經濟兩型

j.



Jarque Bera Test

```
data: residuals
X-squared = 628.07, df = 2, p-value < 2.2e-16
```

不具分常態

right-skewed

k. log-log model, residual 較 random,