

- 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:

$$\widehat{\text{RATING}} = 64.289 + 0.990 \text{EXPER} \quad N = 50 \quad R^2 = 0.3793$$

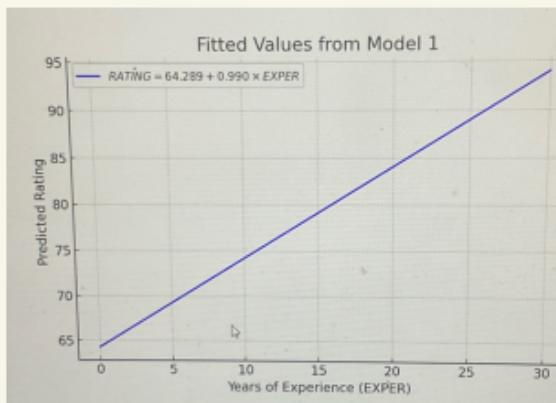
(se) (2.422) (0.183)

Model 2:

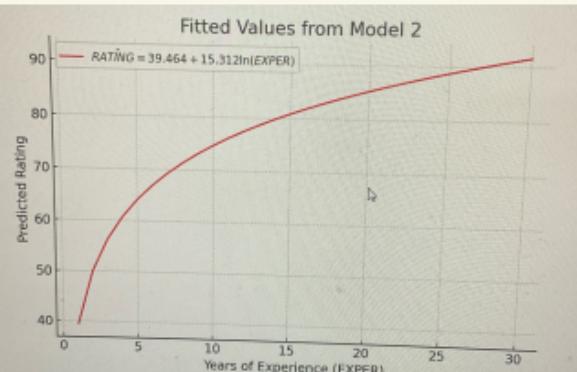
$$\widehat{\text{RATING}} = 39.464 + 15.312 \ln(\text{EXPER}) \quad N = 46 \quad R^2 = 0.6414$$

(se) (4.198) (1.727)

- a. Sketch the fitted values from Model 1 for *EXPER* = 0 to 30 years.



- b. 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.



$\because \ln(0)$ doesn't exist
 \therefore experience = 0 can't be used

- c. 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.

$$\text{marginal effect} = \frac{\partial(\text{RATING})}{\partial(\text{EXPER})} = 0.990$$

!-. No matter what value experience is,
increase one year of EXPER, Rating increase
0.99 point

- d. 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.

$$\text{marginal effect} = \frac{\partial(\text{RATING})}{\partial(\text{EXPER})} = \frac{15.312}{\text{EXPER}}$$

(i) $\text{EXPER} = 10$, marginal effect = $\frac{15.312}{10} = 1.5312$
 (ii) $\text{EXPER} = 20$, marginal effect = $\frac{15.312}{20} = 0.7656$

當藝術家有10年經驗，額外1年經驗提升1.5312分評分
2011-2012 - - - - - 0.7165分評分

- e. Which of the two models fits the data better? Estimation of Model 1 using just the technical artists with some experience yields $R^2 = 0.4858$.

Model 1: $R^2 = 0.3793$

Model 1 (with some experience yield): $R^2 = 0.4858$

Model 2: $R^2 = 0.6414$

\therefore Model 2 has larger R^2 \therefore it can fit data better

- f. Do you find Model 1 or Model 2 more reasonable, or plausible, based on economic reasoning? Explain.

Model 2 更符合經濟學原理，在現實情況中，職場經驗在前幾年帶來更大的回報，但到了一定階段，額外的經驗不會有太大增值，符合邊際報酬遞減現象

- 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:

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

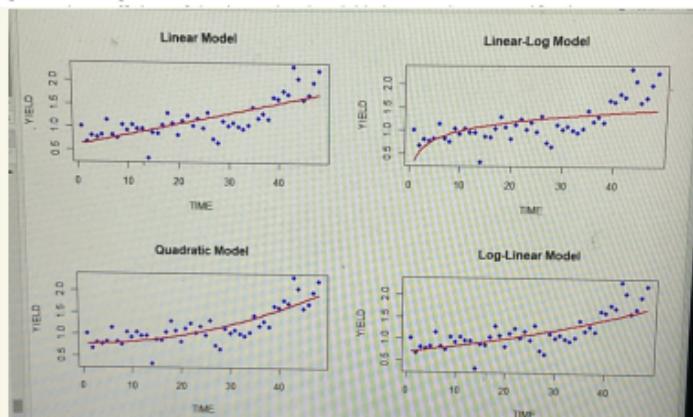
$$YIELD_t = \alpha_0 + \alpha_1 \ln(TIME) + e_t$$

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

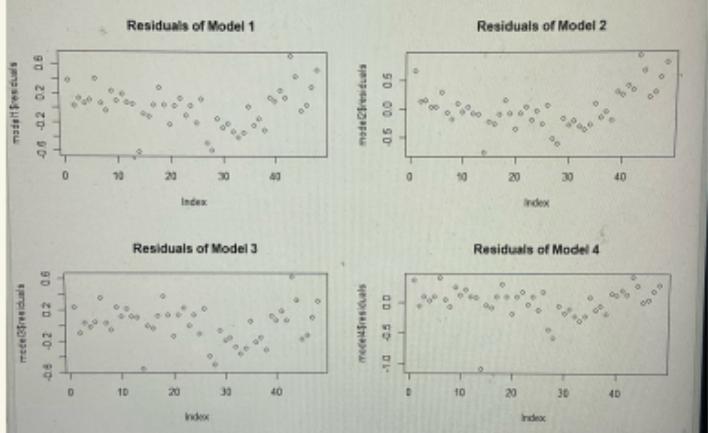
$$\ln(YIELD_t) = \phi_0 + \phi_1 TIME + e_t$$

- a. 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.

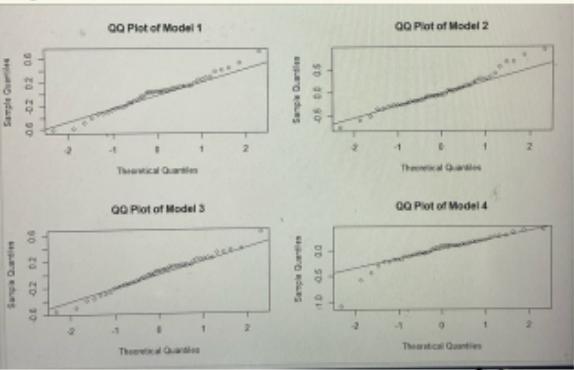
(i)



(ii)



(iii)



```
> print(result_table)
```

Model	Test	V	p_value
1 Model 1	Shapiro-Wilk Normality Test	0.9823582	0.6792
2 Model 2	Shapiro-Wilk Normality Test	0.9686594	0.1856
3 Model 3	Shapiro-Wilk Normality Test	0.9858917	0.8266
4 Model 4	Shapiro-Wilk Normality Test	0.8689369	7.205e-05

根據 Shapiro-Wilk 正態性檢驗，如果 P 值 > 0.05 ，說明沒有顯著偏離正態分布。
∴ model 1, 2, 3 的殘差可以視為正態分佈。

(iv)

Model 3 adj-R² 最大，且殘差為正態，因此是最優模型。

```
> names(adj_r2_values) <- c("Model 1", "Model 2", "Model 3", "Model 4")
> print(adj_r2_values)
Model 1 Model 2 Model 3 Model 4
0.5686594 0.3241945 0.6822494 0.4966469
```

- b. Interpret the coefficient of the time-related variable in your chosen specification.

```
最後模型為: Model 3
> #(b)
> summary(model3)

Call:
lm(formula = YIELD ~ TIME2, data = northampton_data)

Residuals:
    Min      1Q  Median      3Q     Max 
-0.56899 -0.14970  0.03119  0.12176  0.62049 

Coefficients:
            Estimate Std. Error t value Pr(>|t|)    
(Intercept) 7.737e-01 5.223e-02 14.82 < 2e-16 ***
TIME2       4.908e-04 4.939e-05 10.10 3.01e-13 ***  
---
Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

Residual standard error: 0.2396 on 46 degrees of freedom
Multiple R-squared:  0.489, Adjusted R-squared:  0.4822 
F-statistic: 101.9 on 1 and 46 DF, p-value: 3.080e-13
```

$\gamma_1 = 0.000486$, 表示 TIME 與 YIELD 之間是正相關

t value = 10.10, p -value = 3.01×10^{-13} 表示 TIME² 對於 YIELD 有很強的統計顯著性

- c. Using your chosen specification, identify any unusual observations, based on the studentized residuals, LEVERAGE, DFBETAS, and DFFITS.

residuals	leverage	dfbetas.(Intercept)	dfbetas.TIME2	dffits
-2.5606825	0.0393775	-0.48450177	0.320519995	-0.4944002
-2.2468473	0.020832817	-0.219808713	0.003822742	-0.3277591
2.8894474	0.06329921	-0.202525494	0.652179762	0.7823199
1.3786286	0.07843579	-0.116349648	0.338316939	0.3988881
-0.7794852	0.08542511	0.077302871	-0.207150911	-0.2382270
-0.5604893	0.09531255	0.065201030	-0.163400687	-0.1845466
0.4111600	0.10614453	0.053605470	0.127022369	0.1416857
1.3884647	0.11796846	-0.203926321	0.460766675	0.50777802

studentized threshold = ± 2
unusual observation:
14, 28, 43

outlier_residuals	outlier_leverage	outlier_dfbetas.outlier_dfbetas.TIME2
14	TRUE	FALSE
28	TRUE	FALSE
43	TRUE	FALSE
44	FALSE	TRUE
45	FALSE	TRUE
46	FALSE	TRUE
47	FALSE	TRUE
48	FALSE	TRUE

Leverage threshold = $\frac{2P}{n} = 0.0833$
unusual observation:
45, 46, 47, 48

dfbetas threshold = $\pm \frac{2}{\sqrt{n}} = \pm 0.88$; dffits threshold = $\pm \frac{2}{\sqrt{P}} = 0.408$
unusual observation:
14, 43, 44, 48

unusual observation:
14, 43, 48

- d. 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?

95% PI

$$=[1368, 2382]$$

22318有落在這個區間

```
> print(pred_1997)
  fit    lwr     upr
1 1.875113 1.367898 2.382328
> # 1997 年的實際值
> actual_1997 <- wa_wheat$northampton[wa_wheat$time == 48]
> # 比較
> print(actual_1997)
[1] 2.2318
```

4.28R

```
1 data("wa_wheat")
2 northampton_data <- data.frame(
3   YIELD = wa_wheat$northampton,
4   TIME = wa_wheat$time
5 )
6 n(a)
7 model1 <- lm(YIELD ~ TIME, data = northampton_data)
8 res1 <- rstudent(model1)$residuals
9 model1 <- lm(YIELD ~ TIME, data = northampton_data)
10 northampton_data$TIME <- northampton_data$TIME + 1
11 model2 <- lm(YIELD ~ TIME, data = northampton_data)
12 res2 <- rstudent(model2)$residuals
13 model3 <- lm(YIELD ~ TIME, data = northampton_data)
14 res3 <- rstudent(model3)$residuals
15 model4 <- lm(YIELD ~ TIME, data = northampton_data)
16 res4 <- rstudent(model4)$residuals
17 plot(northampton_data$TIME, fitted(model1), col = "#FF0000", pch = 16)
18 lines(northampton_data$TIME, fitted(model1), col = "#FF0000", lwd = 2)
19 lines(northampton_data$TIME, fitted(model2), col = "#0000FF", lwd = 2)
20 lines(northampton_data$TIME, fitted(model3), col = "#008000", lwd = 2)
21 lines(northampton_data$TIME, fitted(model4), col = "#00008B", lwd = 2)
22 plot(northampton_data$TIME, northampton_data$YIELD, main = "Linear Model",
23      xlab = "TIME", ylab = "YIELD", col = "#000000", pch = 16)
24 plot(northampton_data$TIME, northampton_data$YIELD, main = "Quadratic Model",
25      xlab = "TIME", ylab = "YIELD", col = "#000000", pch = 16)
26 plot(northampton_data$TIME, northampton_data$YIELD, main = "Logarithmic Model",
27      xlab = "TIME", ylab = "YIELD", col = "#000000", pch = 16)
28 plot(northampton_data$TIME, exp(fitted(model1)), col = "#FF0000", lwd = 2)
29 lines(northampton_data$TIME, exp(fitted(model2)), col = "#0000FF", lwd = 2)
30 lines(northampton_data$TIME, exp(fitted(model3)), col = "#008000", lwd = 2)
31 lines(northampton_data$TIME, exp(fitted(model4)), col = "#00008B", lwd = 2)
32 n(a)
```

```
66 n(b)
67 summary(model2)
68 n(a)
69 wa_wheat$residuals <- rstudent(model2) # Standardized residuals
70 wa_wheat$leverage <- hatvalues(model2) # Leverage 高於 2p/n
71 wa_wheat$dfBeta <- dfBeta(model2) # DFBETAS
72 wa_wheat$dfFits <- dfFits(model2) # DFFITS
73 n(a)
74 a <- xvar(wa_wheat) # 判斷直線性
75 p <- length(asf(model2)) # 判斷斜率的個數
76 n(a)
77 # 常用的簡單判斷標準
78 threshold_residual <- 2 # Studentized residuals 級別大於 2
79 threshold_leverage <- 2 * p / n # Leverage 高於 2p/n
80 threshold_dfbeta <- 2 * sqrt(p) * DFBETAS # 級別大於 2 * sqrt(p)/n
81 threshold_dffits <- 2 * sqrt(p / n) * DFFITS # 級別大於 2 * sqrt(p/n)
82 n(a)
83 # 檢查異常值
84 wa_wheat$outlier_residuals <- abs(wa_wheat$residuals) > threshold_residual
85 wa_wheat$outlier_leverage <- wa_wheat$leverage > threshold_leverage
86 wa_wheat$outlier_dfbetas <- wa_wheat$dfBeta > threshold_dfbeta
87 wa_wheat$outlier_dffits <- wa_wheat$dfFits > threshold_dffits
88 n(a)
89 # 通過異常值判斷
90 outliers <- wa_wheat$wa_wheat$outlier_residuals |
91   wa_wheat$wa_wheat$leverage |
92   wa_wheat$wa_wheat$dfBeta |
93   wa_wheat$wa_wheat$dfFits |
94   wa_wheat$wa_wheat$dfBeta_time2 .
95 print(outliers)
96 print(model2)
```

```
33 n(c)
34 par(mfrow = c(2,2))
35 openc(model1$residuals, main="QQ Plot of Model 1"), qqline(model1$residuals)
36 qqplot(model2$residuals, main="QQ Plot of Model 2"), qqline(model2$residuals)
37 qqplot(model3$residuals, main="QQ Plot of Model 3"), qqline(model3$residuals)
38 qqplot(model4$residuals, main="QQ Plot of Model 4"), qqline(model4$residuals)
39 n(c)
40 model1$stat <- shapiro.test(model1$residuals)
41 model2$stat <- shapiro.test(model2$residuals)
42 model3$stat <- shapiro.test(model3$residuals)
43 model4$stat <- shapiro.test(model4$residuals)
44 result_table <- data.frame( # 整理一個表格來匯總每個模型的統計 data frame
45   Model = c("Model 1", "Model 2", "Model 3", "Model 4"),
46   Test = c("Shapiro-Wilk Normality Test", "Shapiro-Wilk Normality Test"),
47   T_value = c(model1$stat$p.value, model2$stat$p.value, model3$stat$p.value, model4$stat$p.value),
48   P_value = c(model1$stat$p.value, model2$stat$p.value, model3$stat$p.value, model4$stat$p.value),
49   format(model1$stat$p.value, digits = 4),
50   format(model2$stat$p.value, digits = 4),
51   format(model3$stat$p.value, digits = 4),
52   format(model4$stat$p.value, digits = 4),
53   )
54 n(c)
55 print(result_table)
56 n(a)
57 wa(a)
58 df1$value <- a
59 wa(a)
60 wa(a)
61 wa(a)
62 wa(a)
63 wa(a)
64 wa(a)
65 wa(a)
66 wa(a)
67 wa(a)
68 wa(a)
69 wa(a)
70 wa(a)
71 wa(a)
72 wa(a)
73 wa(a)
74 wa(a)
75 wa(a)
76 wa(a)
77 wa(a)
78 wa(a)
79 wa(a)
80 wa(a)
81 wa(a)
82 wa(a)
83 wa(a)
84 wa(a)
85 wa(a)
86 wa(a)
87 wa(a)
88 wa(a)
89 wa(a)
90 wa(a)
91 wa(a)
92 wa(a)
93 wa(a)
94 wa(a)
95 wa(a)
96 wa(a)
97 wa(a)
98 wa(a)
99 data_1996 <- subset(wa_wheat, time < 47)
100 model3_1996 <- lm(northampton ~ I(time^2), data = wa_wheat)
101 # 挑出 1997 年的 time 值
102 time_1997 <- wa_wheat$time[wa_wheat$time == 47]
103 # 對 1997 年進行預測
104 pred_1997 <- predict(model3_1996,
105   newdata = data.frame(time = time_1997),
106   interval = "prediction", level = 0.95)
107 print(pred_1997)
108 # 1997 年的實際值
109 actual_1997 <- wa_wheat$northampton[wa_wheat$time == 48]
110 print(actual_1997)
111
```

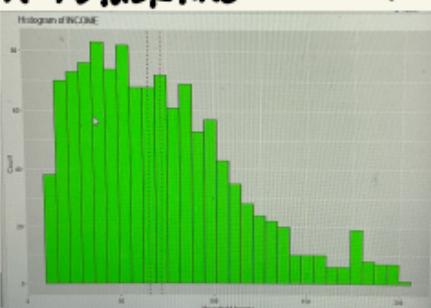
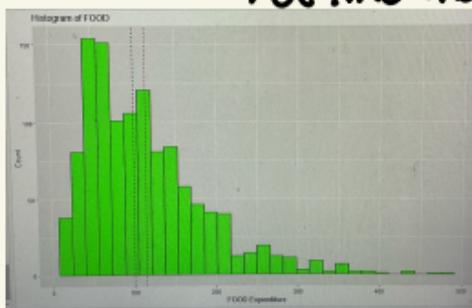
```
98 n(d)
99 data_1996 <- subset(wa_wheat, time < 47)
100 model3_1996 <- lm(northampton ~ I(time^2), data = wa_wheat)
101 # 挑出 1997 年的 time 值
102 time_1997 <- wa_wheat$time[wa_wheat$time == 47]
103 # 對 1997 年進行預測
104 pred_1997 <- predict(model3_1996,
105   newdata = data.frame(time = time_1997),
106   interval = "prediction", level = 0.95)
107 print(pred_1997)
108 # 1997 年的實際值
109 actual_1997 <- wa_wheat$northampton[wa_wheat$time == 48]
110 print(actual_1997)
111
```

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.

- a. 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.

```
> print(summary_stats)
  Variable      Mean Median   Min   Max Std_Dev
1    FOOD 114.44311 99.80 9.63 476.67 72.65750
2  INCOME  72.14264 65.29 10.00 200.00 41.65228
```

red line: mean, blackline: median



兩個分佈 mean > median ; 都是右偏分佈
而不是鍾型或對稱分佈

```
Jarque Bera Test
data: df$food
X-squared = 648.65, df = 2, p-value < 2.2e-16
> print(jb_income)
Jarque Bera Test
data: df$income
X-squared = 148.21, df = 2, p-value < 2.2e-16
```

兩個檢驗之 p-Value
都非常小 ; reject H₀
都不符合常態分佈

- b. 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?

```
Call:
lm(formula = food ~ income, data = df)
```

Residuals:

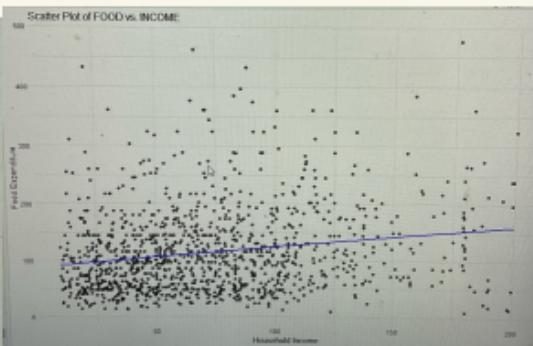
Min	1Q	Median	3Q	Max
-145.37	-51.48	-13.52	35.50	349.81

Coefficients:

	Estimate	Std. Error	t value	Pr(> t)
(Intercept)	88.56650	4.10819	21.559	< 2e-16 ***
income	0.35869	0.04922	7.272	6.36e-13 ***

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

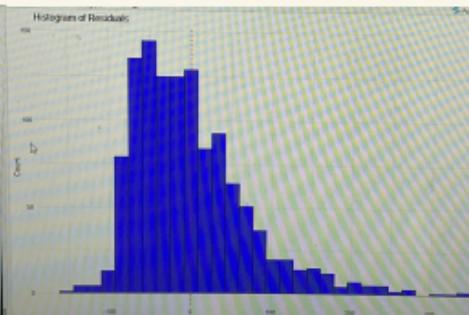
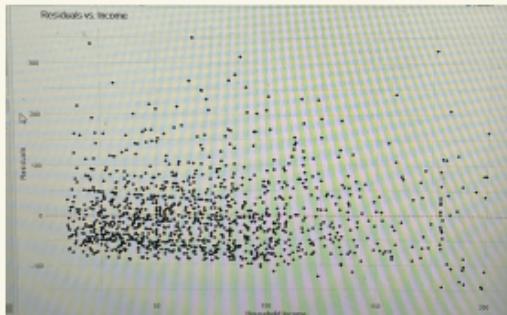
Residual standard error: 71.13 on 1198 degrees of freedom
 Multiple R-squared: 0.04228, Adjusted R-squared: 0.04148
 F -statistic: 52.89 on 1 and 1198 DF, p-value: 6.357e-13



95% interval estimate
 $[0.2619, 0.4555]$

the effect of changing income on food expenditure may not be estimated very precisely, as the low R^2 value indicates that the model explain only a small proportion of the variation in food expenditure.

- c. 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.



Jarque-Bera Test

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

p-value < 0.05 ; reject H₀

不符合常態分布

殘差大致分佈在0附近，沒有嚴重偏移，但在高收入 (Income > 150) 時，殘差的變異數變大，可能有異質變異數問題。

It is more important for the random error to be normally distributed, because the validity of hypothesis tests in regression analysis relies on the normality of the residuals.

- d. 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?

$$E = \beta_2 \times \frac{\text{INCOME}}{\text{Food}} \quad \beta_2 = 0.3587$$

```
> print(elasticity_df)
  Income Predicted_FOOD Elasticity Lower_95CI Upper_95CI
1     19      95.38155 0.07145302 0.05217047 0.09073558
2     65     111.88114 0.20839527 0.15215702 0.26463353
3    160     145.95638 0.39321338 0.28709948 0.49932727
> |
```

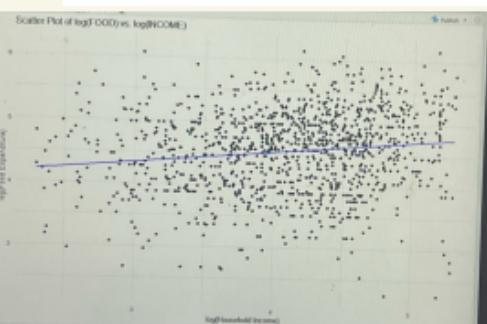
the estimated elasticities are dissimilar.
INCOME increase lead to a higher elasticity.

the interval estimated don't overlap(不重疊).

Based on Economics principles, food is necessity good (必需品), so elasticity must less than 1, and with income increase, the elasticity will decrease.

- e. 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

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?



```
Call:  
lm(formula = log(food) ~ log(income), data = df)  
  
Residuals:  
    Min      1Q  Median      3Q     Max  
-2.48115 -0.45497  0.06151  0.46063  1.72315  
  
Coefficients:  
            Estimate Std. Error t value Pr(>|t|)  
(Intercept) 3.77892  0.12035 31.400  <2e-16 ***  
log(income) 0.18631  0.02903  6.417  2e-10 ***  
  
Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1  
  
Residual standard error: 0.6418 on 1198 degrees of freedom  
Multiple R-squared:  0.03322, Adjusted R-squared:  0.03242  
F-statistic: 41.18 on 1 and 1198 DF, p-value: 1.999e-10
```

```
> print(paste("Generalized R^2 for Log-Log Model: ", generalized_r2_log_log))  
[1] "Generalized R^2 for Log-Log Model: 0.0396516120669934"  
>
```

generalized log-log model $R^2 = 0.0396 < 0.0423 = \text{linear model } R^2$

i. linear model fit the data better

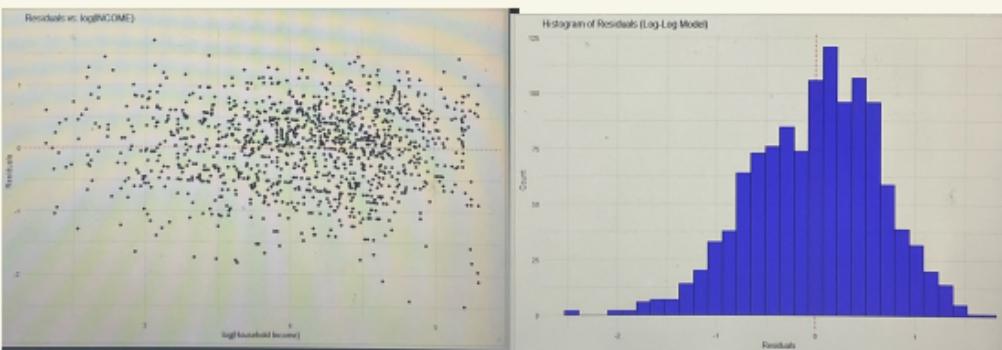
- f. Construct a point and 95% interval estimate of the elasticity for the log-log model. Is the elasticity of food expenditure from the log-log model similar to that in part (d), or dissimilar? Provide statistical evidence for your claim.

```
> confint(model_loglog, level = 0.95)  
    2.5 %   97.5 %  
(Intercept) 3.5428135 4.0150507  
log(income) 0.1293432 0.2432675  
>
```

95% interval estimate
 $= [0.1293, 0.2432]$

$E = r_2 = 0.1863$, it is a constant so dissimilar to the elasticity of linear model.

- g. Obtain the least squares residuals from the log-log model and plot them against $\ln(INCOME)$. 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?

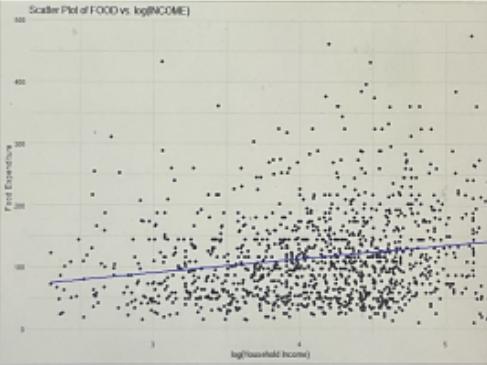


Jarque Bera Test

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

residuals seems no any strong pattern
In Jarque Bera Test, p-value < 0.05.
 \therefore reject H_0 , 不符合常態分布

- h. For expenditures on food, estimate the linear-log relationship $FOOD = \alpha_1 + \alpha_2 \ln(INCOME) + e$. Create a scatter plot for $FOOD$ versus $\ln(INCOME)$ 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?



```

Call:
lm(formula = food ~ log(income), data = df)

Residuals:
    Min      1Q  Median      3Q     Max 
-129.18 -51.47 -13.98  35.05 345.54 

Coefficients:
            Estimate Std. Error t value Pr(>|t|)    
(Intercept) 23.568     13.370   1.763   0.0782 .  
log(income) 22.187      3.225   6.879 9.68e-12 *** 
---
Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

Residual standard error: 71.29 on 1198 degrees of freedom
Multiple R-squared:  0.038,   Adjusted R-squared:  0.0372 
F-statistic: 47.32 on 1 and 1198 DF,  p-value: 9.681e-12

```

$$R^2 \text{ of linear model} = 0.0423$$

$$\text{Generalized } R^2 \text{ of log-log model} = 0.0397$$

$$R^2 \text{ of linear-log model} = 0.038$$

By R^2 , linear model fit the data better

- Construct a point and 95% interval estimate of the elasticity for the linear-log model at $INCOME = 19, 65,$ 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.

```

> confint(model_linlog, level = 0.95)
        2.5 % 97.5 %
(Intercept) -2.661958 49.79892
log(income) 15.859458 28.51531
>

```

```

> print(income_values)
  income predicted_food elasticity elasticity_lower elasticity_upper
1     19     88.89627  0.2495830     0.1784339     0.3206771
2     65    116.18513  0.1909625     0.1365665     0.2453584
3    160    136.17088  0.1629350     0.1165227     0.2093473
> #(1)

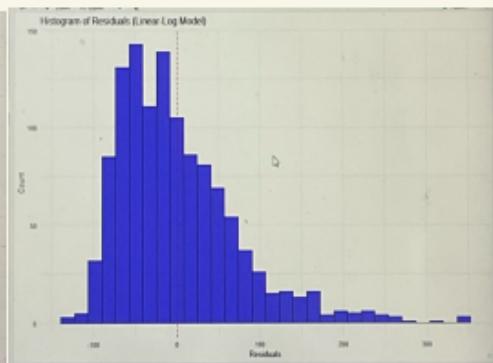
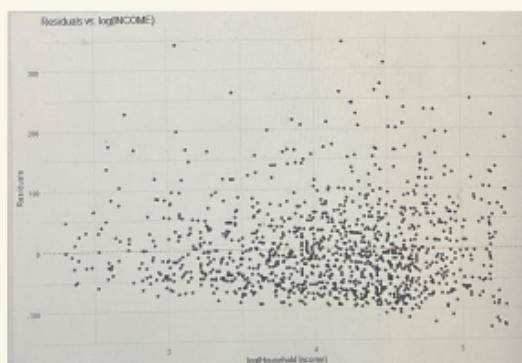
```

95% interval estimate
 $= [15.859, 28.515]$

$$E = \frac{\alpha_2}{INCOME}$$

Income increase lead to a decrease on elasticity, so it is dissimilar to other models

- 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?



Jarque-Bera Test

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

residuals seems no any strong pattern
In Jarque-Bera Test, $p\text{-value} < 0.05$.
 $\therefore \text{reject } H_0$, 不符合常態分布

- 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.

Linear Model 的所得彈性隨收入增加而上升，不符合經濟學原理首先排除，log-log model 相比之下 R² 更高且殘差更接近隨機分佈，因此選 log-log model

4.29 R

```
1  data(cx5_small)
2  df <- cx5_small
3  #(a)
4  summary_stats <- data.frame(
5    Variable = c("FOOD", "INCOME"),
6    Mean = c(mean(df$food), mean(df$income)),
7    Median = c(median(df$food), median(df$income)),
8    Min = c(min(df$food), min(df$income)),
9    Max = c(max(df$food), max(df$income)),
10   Std_Dev = c(sd(df$food), sd(df$income)))
11 )
12 print(summary_stats)
13 ggplot(df, aes(x = food)) +
14   geom_histogram(color = "#00AEEF", fill = "#4CAF50", bins = 30, alpha = 0.7) +
15   geom_vline(aes(xintercept = mean(food)), color = "#E91E63", linetype = "dashed") +
16   geom_vline(aes(xintercept = median(food)), color = "#00AEEF", linetype = "dashed") +
17   labs(title = "Histogram of FOOD", x = "Food Expenditure", y = "Count")
18 ggplot(df, aes(x = income)) +
19   geom_histogram(color = "#00AEEF", fill = "#4CAF50", bins = 30, alpha = 0.7) +
20   geom_vline(aes(xintercept = mean(income)), color = "#E91E63", linetype = "dashed") +
21   geom_vline(aes(xintercept = median(income)), color = "#00AEEF", linetype = "dashed") +
22   labs(title = "Histogram of INCOME", x = "Household Income", y = "Count")
23 jb_food <- jarque.bera.test(df$food)
24 jb_income <- jarque.bera.test(df$income)
25 print(jb_food)
26 print(jb_income)
27
28 u(h)
29 model_linear <- lm(food ~ income, data = df)
30 summary(model_linear)
31 ggplot(df, aes(x = income, y = food)) +
32   geom_point(alpha = 0.6) + # Scatter points
33   geom_smooth(method = "lm", color = "#00AEEF", se = FALSE) + # Regression line
34   labs(title = "Scatter Plot of FOOD vs. INCOME",
35        x = "Household Income",
36        y = "Food Expenditure") +
37   theme_minimal()
38 confint(model_linear, level = 0.95)
39
40 #(c)
41 df_residuals <- residuals(model_linear)
42 ggplot(df, aes(x = income, y = residuals)) +
43   geom_point(alpha = 2) +
44   geom_hline(yintercept = 0, color = "#E91E63", linetype = "dashed") +
45   labs(title = "Residuals vs. Income",
46        x = "Household Income",
47        y = "Residuals") +
48   theme_minimal()
49 ggplot(df, aes(x = residuals)) +
50   geom_histogram(color = "#00AEEF", fill = "#4CAF50", bins = 30, alpha = 0.7) +
51   geom_vline(aes(xintercept = mean(residuals)), color = "#E91E63", linetype = "dashed") +
52   labs(title = "Histogram of Residuals", x = "Residuals", y = "Count") +
53   theme_minimal()
54 jb_test_residuals <- jarque.bera.test(df$residuals)
55 print(jb_test_residuals)
56
```

```

57 # (d)
58 income_values <- data.frame(income = c(19, 65, 160))
59 predicted_food <- predict(model_linear, newdata = income_values)
60 elasticity <- 0.3587 * (income_values$income / predicted_food)
61 elasticity_df <- data.frame(income_values$income, predicted_food, elasticity)
62 colnames(elasticity_df) <- c("Income", "Predicted FOOD", "Elasticity")
63 elasticity_lower <- 0.2619 * (income_values$income / predicted_food)
64 elasticity_upper <- 0.4555 * (income_values$income / predicted_food)
65 elasticity_df$Lower_95CI <- elasticity_lower
66 elasticity_df$Upper_95CI <- elasticity_upper
67 print(elasticity_df)
68
69 # (e)
70 model_loglog <- lm(log(food) ~ log(income), data = df)
71 ggplot(df, aes(x = log(income), y = log(food))) +
72   geom_point(alpha = 2) + # Scatter points
73   geom_smooth(method = "lo", color = "#FF0000", se = FALSE) + # Regression line
74   labs(title = "Scatter Plot of log(FOOD) vs. log(INCOME)", 
75        x = "log(Household Income)", 
76        y = "log(Food Expenditure)") +
77   theme_minimal()
78 summary(model_loglog)
79 # 由圖可知Food NO income 都大於 0
80 data_log <- subset(df, food > 0 & income > 0)
81 # 計算 log 值數
82 data_log$ln_food <- log(data_log$food)
83 data_log$ln_income <- log(data_log$income)
84 model_loglog <- lm(ln_food ~ ln_income, data = data_log)
85
86 # 採用 Log-Log 模型的準則值
87 data_log$food_pred <- exp(predict(model_loglog)) # 需要將預測值反log，回到原始單位
88
89 # 1) Generalized R2 的計算
90 generalized_r2_crry <- function(actual, predicted) {
91   r_xr <- cor(actual, predicted)^2 # 2) 計算總體平方
92   return(r_xr)
93 }
94 # 2) Log-Log Model 的 Generalized R2
95 generalized_r2_log_log <- generalized_r2_crry(data_log$food, data_log$food_pred)
96 print(paste("Generalized R2 for Log-Log Model: ", generalized_r2_log_log))
97
98 # (f)
99 elasticity_loglog <- coef(model_loglog)[["log(income)"]]
100 print(elasticity_loglog)
101 confint(model_loglog, level = 0.95)
102
103 # (g)
104 df$residuals_loglog <- residuals(model_loglog)
105 ggplot(df, aes(x = log(income), y = residuals_loglog)) +
106   geom_point(alpha = 0.5) + # Scatter points
107   geom_smooth(method = "lo", color = "#FF0000", linetype = "dashed") +
108   labs(title = "Residuals vs. log(INCOME)", 
109        x = "log(Household Income)", 
110        y = "Residuals") +
111   theme_minimal()
112 ggplot(df, aes(x = log(income), y = residuals_loglog)) +
113   geom_histogram(color = "#FF0000", fill = "#FF0000", bins=30, alpha=0.3) +
114   geom_vline(aes(xintercept = mean(residuals_loglog)), color="#FF0000", linetype="dashed") +
115   labs(title="Histogram of Residuals (Log-Log Model)", x="Residuals", y="Count") +
116   theme_minimal()
117 jb_test_residuals_loglog <- jarque.bera.test(df$residuals_loglog)
118 print(jb_test_residuals_loglog)
119
120 # (h)
121 model_llinlog <- lm(food ~ log(income), data = df)
122 income_values <- data.frame(income = c(19, 65, 160))
123 ggplot(df, aes(x = log(income), y = food)) +
124   geom_point(alpha = 2) + # Scatter points
125   geom_smooth(method = "lo", color = "#FF0000", se = FALSE) + # Regression line
126   labs(title = "Scatter Plot of FOOD vs. log(INCOME)", 
127        x = "log(Household Income)", 
128        y = "Food Expenditure") +
129   theme_minimal()
130
131 # (i)
132 confint(model_llinlog, level = 0.95)
133 income_values <- data.frame(income = c(19, 65, 160))
134 income_values$predicted_food <- 23.568 + 22.187 * log(income_values$income)
135 income_values$elasticity <- 22.187 / income_values$predicted_food
136 income_values$elasticity_lower <- 15.867 / income_values$predicted_food
137 income_values$elasticity_upper <- 28.507 / income_values$predicted_food
138 print(income_values)
139
140 # (j)
141 df$residuals_llinlog <- residuals(model_llinlog)
142 ggplot(df, aes(x = log(income), y = residuals_llinlog)) +
143   geom_point(alpha = 0.5) +
144   geom_hline(intercept = 0, color = "#FF0000", linetype = "dashed") +
145   labs(title = "Residuals vs. log(INCOME)", 
146        x = "log(Household Income)", 
147        y = "Residuals") +
148   theme_minimal()
149 ggplot(df, aes(x = residuals_llinlog)) +
150   geom_histogram(color="#FF0000", fill="#FF0000", bins=30, alpha=0.3) +
151   geom_vline(aes(xintercept = mean(residuals_llinlog)), color="#FF0000", linetype="dashed") +
152   labs(title="Histogram of Residuals (Linear-Log Model)", x="Residuals", y="Count") +
153   theme_minimal()
154 jb_test_residuals_llinlog <- jarque.bera.test(df$residuals_llinlog)
155 print(jb_test_residuals_llinlog)
156

```