

HW4

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:

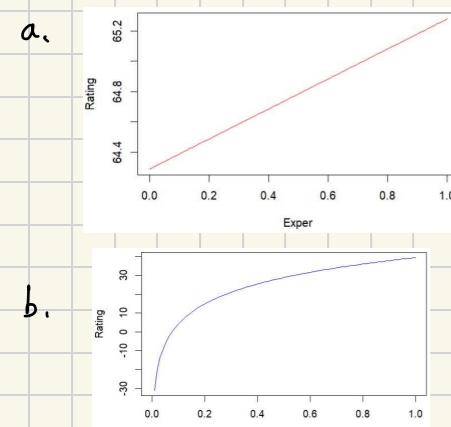
Model 1:

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

Model 2:

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

- 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.
- 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.
- 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.
- e. Which of the two models fits the data better? Estimation of Model 1 using just the technical artists with some experience yields $R^2 = 0.4858$.
- f. Do you find Model 1 or Model 2 more reasonable, or plausible, based on economic reasoning? Explain.



Since the domain of nature logarithm is $(0, \infty)$, the four artists with 0 experience can not be used in the regression model.

C. the margin effect on Rating is

$$\frac{d \text{Rating}}{d \text{Exper}} = 0.99$$

(i) the artists with 10 years Exper: 0.99

(ii) the artists with 20 years Exper: 0.99

D. the margin effect on Rating is

$$\frac{d \text{Rating}}{d \text{Exper}} = \frac{15.312}{\text{Exper}}$$

(i) the artists with 10 years Exper: $\frac{15.312}{10} = 1.5312$

(ii) the artists with 20 years Exper: $\frac{15.312}{20} = 0.7656$

e. since R^2 of mode 1 = 0.3793 < 0.6414 (the R^2 of Model 2)

\Rightarrow Model 2 fits the data better

f. Model 2 is reasonable since it has larger R^2 , and the rating will rapidly increase with less experience.

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 Northamptonshire, consider the following four equations:

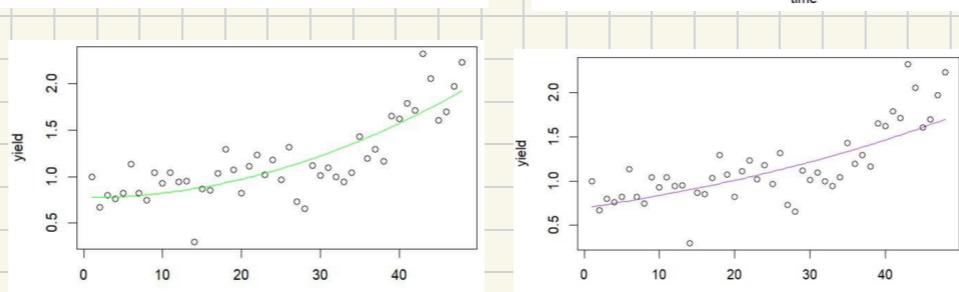
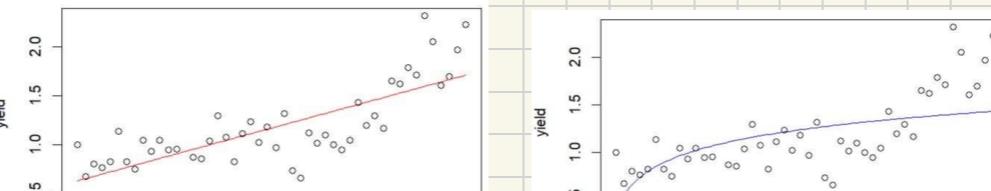
$$\begin{aligned} YIELD_t &= \beta_0 + \beta_1 TIME + e_t \\ YIELD_t &= \alpha_0 + \alpha_1 \ln(TIME) + e_t \\ YIELD_t &= \gamma_0 + \gamma_1 TIME^2 + e_t \\ \ln(YIELD_t) &= \phi_0 + \phi_1 TIME + e_t \end{aligned}$$

- 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.
- b. Interpret the coefficient of the time-related variable in your chosen specification.
- c. Using your chosen specification, identify any unusual observations, based on the studentized residuals, *LEVERAGE*, *DFBETAS*, and *DFITS*.
- 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?

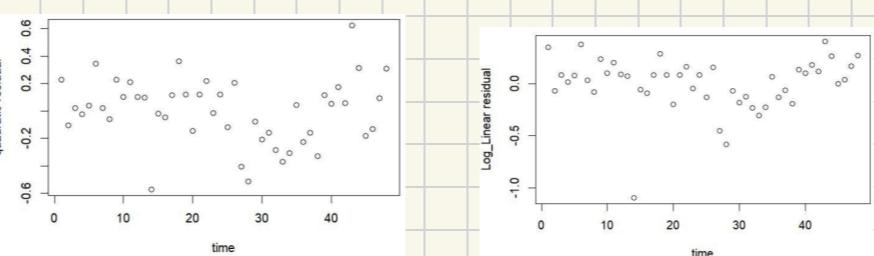
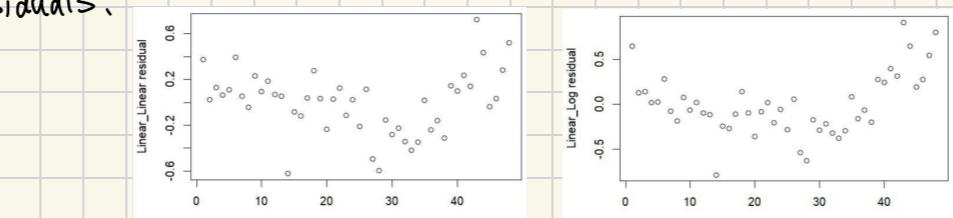
a. the estimation of four equations:

```
> theCoefficientB1
[1] 0.6032451 0.3509620 0.7736655 -0.3639376
> theCoefficientB2
[1] 0.0230779201 0.2790085294 0.0004986181 0.0186323476
```

Fitted equations:



Residuals:



error normality tests:

```
> jbTest1
Jarque Bera Test
data: Bresidual
X-squared = 0.13257, df = 2, p-value = 0.9359
> jbTest2
Jarque Bera Test
data: Aresidual
X-squared = 2.7629, df = 2, p-value = 0.2512
> jbTest3
Jarque Bera Test
data: Cresidual
X-squared = 0.32406, df = 2, p-value = 0.8504
> jbTest4
Jarque Bera Test
data: Dresidual
X-squared = 83.874, df = 2, p-value < 2.2e-16
```

R^2 :

```
> rsq1
[1] 0.5778369
> rsq2
[1] 0.3385733
> rsq3
[1] 0.6890101
> rsq4
[1] 0.5073566
```

According to the above data, I prefer to use model 3. Since it fits the data better, its residual does not have unusual pattern and its R^2 is the nearest to 1. With error normality tests, since the $p\text{-value}=0.8504 > 0.05$, we will not reject that its residual has normal distribution.

b. the estimation:

```
> theCoefficientB1
[1] 0.6032451 0.3509620 0.7736655 -0.3639376
> theCoefficientB2
[1] 0.0230779201 0.2790085294 0.0004986181 0.0186323476
```

In model 3: when time=0, the yield is approximately 0.7736.

Since $C_2 = 0.00049$, when time increase 1 unit, the yield will increase $2 \times C_2 = 0.00098$.

d. the true value: 2.2318

the Interval:

```
> lower
[1] 1.382799
> upper
[1] 2.379424
```

Yes, the true value is contained in the 95% prediction interval.

```

2 yield<-c(wa_wheat$northampton)
3 time<-c(wa_wheat$time)
4 data<-data.frame(time, yield)
5
6 #a
7 modlinear_linear<-lm(yield~time, data=data)
8 modlinear_log<-lm(yield~log(time), data=data)
9 modquadratic<-lm(yield~I(time^2), data=data)
10 modlog_linear<-lm(log(yield)~time, data=data)
11 b1<-coef(modlinear_linear)[[1]]
12 b2<-coef(modlinear_linear)[[2]]
13 Bresidual<-resid(modlinear_linear)
14 a1<-coef(modlinear_log)[[1]]
15 a2<-coef(modlinear_log)[[2]]
16 Aresidual<-resid(modlinear_log)
17 c1<-coef(modquadratic)[[1]]
18 c2<-coef(modquadratic)[[2]]
19 Cresidual<-resid(modquadratic)
20 d1<-coef(modlog_linear)[[1]]
21 d2<-coef(modlog_linear)[[2]]
22 Dresidual<-resid(modlog_linear)
23 theCoefficientB1<-c(b1, a1, c1, d1)
24 theCoefficientB2<-c(b2, a2, c2, d2)
25 theCoefficientB1
26 theCoefficientB2
27 #plot
28 plot(time, yield, xlab="time", ylab="yield", col="black")
29 curve(b1+b2*x, add=TRUE, col="red")
30 curve(a1+a2*log(x), add=TRUE, col="blue")
31 curve(c1+c2*x^2, add=TRUE, col="green")
32 curve(exp(d1+d2*x), add=TRUE, col="purple")
33
34 plot(time, Bresidual, xlab="time", ylab="Linear_Linear residual", col="black")
35 plot(time, Aresidual, xlab="time", ylab="Linear_Log residual", col="black")
36 plot(time, Cresidual, xlab="time", ylab="quadratic residual", col="black")
37 plot(time, Dresidual, xlab="time", ylab="Log_Linear residual", col="black")
38
39 library(tseries)
40 jbTest1<-jarque.bera.test(Bresidual)
41 jbTest2<-jarque.bera.test(Aresidual)
42 jbTest3<-jarque.bera.test(Cresidual)
43 jbTest4<-jarque.bera.test(Dresidual)
44 jbTest1
45 jbTest2
46 jbTest3
47 jbTest4
48
49 rsq1<-summary(modlinear_linear)$r.squared
50 rsq2<-summary(modlinear_log)$r.squared
51 rsq3<-summary(modquadratic)$r.squared
52 rsq4<-summary(modlog_linear)$r.squared
53 rsq1
54 rsq2
55 rsq3
56 rsq4
57

```

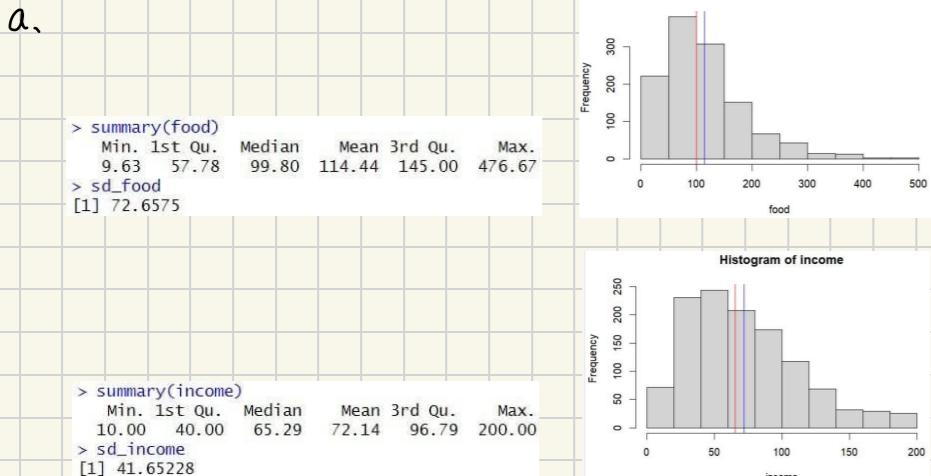
```

67 #d
68 newyield<-c(wa_wheat$northampton[1:47])
69 newtime<-c(wa_wheat$time[1:47])
70 newdata1<-data.frame(newtime, newyield)
71 newmod<-lm(newyield~I(newtime^2), data=newdata1)
72 newb1<-coef(newmod)[[1]]
73 newb2<-coef(newmod)[[2]]
74 newresidual<-resid(newmod)
75 yield_1997<-newb1+(newb2)*(48^2)
76 yield_1997
77 t<-qt(0.975, 45)
78 se<-summary(newmod)$sigma
79 #leverage
80 f<-1+(1/47)+(((48-mean(newtime))^2)/sum((newtime-mean(newtime))^2))
81 lower<-yield_1997-t*se*(sqrt(f))
82 upper<-yield_1997+t*se*(sqrt(f))
83 lower
84 upper

```

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. What do you conclude about the normality of the regression errors in this model?
- 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?
- 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.
- Obtain the least squares residuals from the linear-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?
- Based on this exercise, do you prefer the linear relationship model, or the log-log model or the linear-log model? Explain your reasoning.



Since the two histograms are positive skewness, neither of them are symmetrical nor bell-shape curve.

The sample mean is larger than the median.

```
> jbTest_food
Jarque Bera Test
data: food
X-squared = 648.65, df = 2, p-value < 2.2e-16
```

```
> jbTest_income
Jarque Bera Test
data: income
X-squared = 148.21, df = 2, p-value < 2.2e-16
```

If $p\text{-value} > 0.05$, we will not reject $H_0 = \text{normal distribution}$. Neither of them are normal distribution.

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?

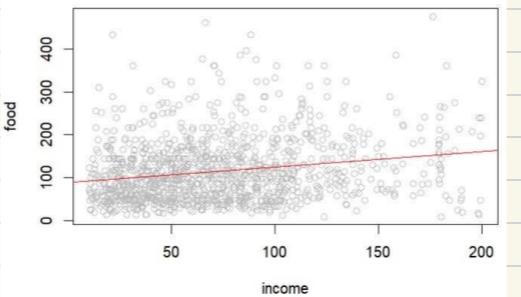
- Construct a point and 95% interval estimate of the elasticity for the log-log model. Is the elasticity of food expenditure from the log-log model similar to that in part (d), or dissimilar? Provide statistical evidence for your claim.
- Obtain the least squares residuals from the log-log model and 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?
- 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?
- 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.
- Obtain the least squares residuals from the linear-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?
- Based on this exercise, do you prefer the linear relationship model, or the log-log model or the linear-log model? Explain your reasoning.

b. Estimation:

```
> b1
[1] 88.5665
> b2
[1] 0.3586867
```

> se_b2
[1] 0.04932104

Plot:

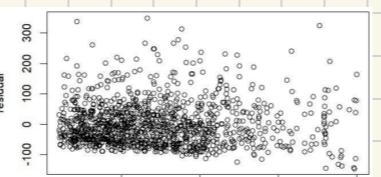


Interval:

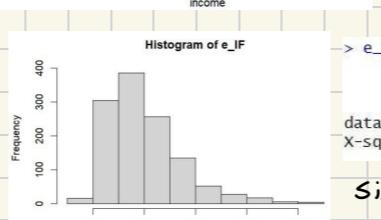
```
> lower
[1] 0.2619215
> upper
[1] 0.455452
```

since the interval is narrow and standard error of β_2 is small \Rightarrow precise.

c.



It does not have any unusual pattern.



```
> e_jbTest
Jarque Bera Test
data: e_IF
X-squared = 624.19, df = 2, p-value < 2.2e-16
```

Since its $p\text{-value} < 0.05$, it is not normal distribution.

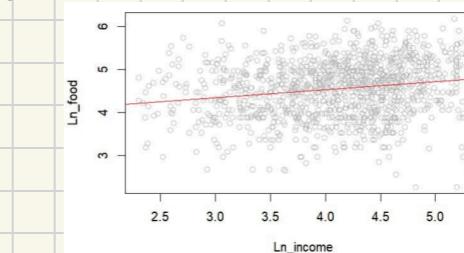
It is more important that random error is normal distribution. If e is normal distribution, the statistic inference (ex: hypothesis test) are valid.

d.

```
> theIncome_p
[1] 19 65 160
> elasticity_p
[1] 0.07145038 0.20838756 0.39319883
> lower_elastic
[1] 0.05217475 0.15216951 0.28712305
> upper_elastic
[1] 0.09072601 0.26460562 0.49927462
```

As the income increase, the portion of income spent on food decrease \rightarrow elasticity decrease (based on economic principle)

e.



```
> G_rsq_linear
[1] 0.04235027
> G_rsq_loglog
[1] 0.03328369
```

since R^2 of linear model is larger, linear model fits better.

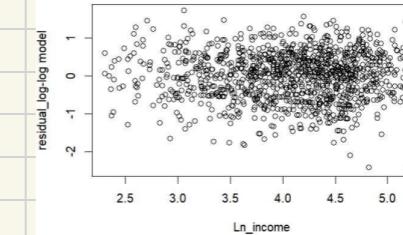
f.

```
> b2_loglog
[1] 0.1863054
> lower_elastic_loglog
[1] 0.1293432
> upper_elastic_loglog
[1] 0.2432675
```

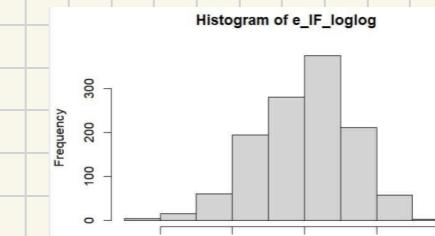
It is dissimilar to part(d).

Since the elasticity of log-log model equal to β_2 , which is a constant.

g.



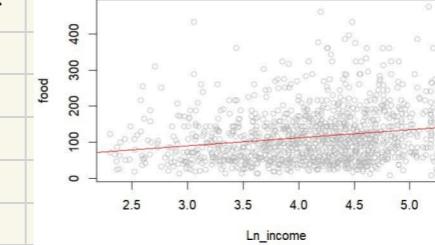
It does not have unusual pattern.



```
> e_jbTest_loglog
Jarque Bera Test
data: residuals(mod_loglog)
X-squared = 25.85, df = 2, p-value = 2.436e-06
```

Since its $p\text{-value} < 0.05$, it is not normal distribution.

h.



```
> rsq_linearlog
[1] 0.03799984
```

Compare with the R^2 in (e), the R^2 of linear model is greatest.

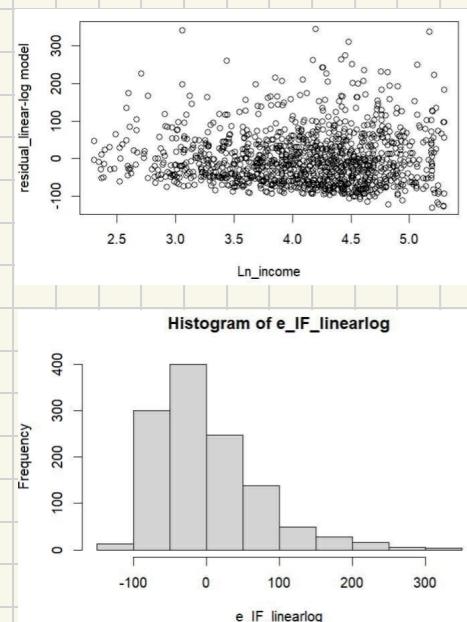
Thus, linear model fits the data better.

The elasticities are dissimilar.
the interval is not overlap.

i.
 19 65 160
 > elasticity_p_linearlog
 [1] 0.2495828 0.1909624 0.1629349
 > lower_elastic_linearlog
 [1] 0.1784009 0.1364992 0.1164652
 > upper_elastic_linearlog
 [1] 0.3207648 0.2454256 0.2094046

As the income increase, the elasticity decrease.

Hence, it is dissimilar to others.



It does not have unusual pattern.

Its p-value < 0.05, so it is not normal distribution.

h. I prefer log log model.

Based on economic principal, as the income increase, the elasticity decrease. However, in Linear model, its elasticity increase as the income increase. Thus, I will not choose

linear model. About the other two model, log log model has

higher R². Therefore, I prefer log log model.

```

2 food<-c(cex5_small$food)
3 income<-c(cex5_small$income)
4 sd_food<-sd(food)
5 sd_income<-sd(income)
6
7 #a
8 summary(food)
9 sd_food
10 summary(income)
11 sd_income
12 #histogram
13 hist(food)
14 abline(v=mean(food), col="blue")
15 abline(v=median(food), col="red")
16 hist(income)
17 abline(v=mean(income), col="blue")
18 abline(v=median(income), col="red")
19 #jbTest
20 jbTest_food<-jarque.bera.test(food)
21 jbTest_income<-jarque.bera.test(income)
22 jbTest_food
23 jbTest_income
24
25 #b
26 data=data.frame(income, food)
27 modIF<-lm(food~income, data=data)
28 b1<-coef(modIF)[[1]]
29 b2<-coef(modIF)[[2]]
30 b1
31 b2
32 #plot
33 plot(income, food, xlab="income", ylab="food", col="grey")
34 abline(b1, b2, col="red")
35 #95% interval
36 summaryModIF<-summary(modIF)
37 se_b2<-summaryModIF$coefficients[2, 2]
38 n<-length(food)
39 t<-qt(0.975, n-2)
40 lower<-b2-t*(se_b2)
41 upper<-b2+t*(se_b2)
42 lower
43 upper
44
45 #c
46 e_IF<-resid(modIF)
47 plot(income, e_IF, xlab="income", ylab="residual")
48 #histogram of residual
49 hist(e_IF)
50 #jbtest
51 e_jbTest<-jarque.bera.test(e_IF)
52 e_jbTest
53
54 #d
55 #point
56 theIncome_p<-c(19, 65, 160)
57 theFood_p<-b1+b2*theIncome_p
58 margin_effect<-b2
59 elasticity_p<-margin_effect*(theIncome_p/theFood_p)
60 elasticity<-margin_effect*(income/food)
61 n_elastic<-length(income)
62 t_elastic<-qt(0.975, n_elastic-2)
63 lower_elastic<-elasticity_p-t_elastic*(theIncome_p/theFood_p)*se_b2
64 upper_elastic<-elasticity_p+t_elastic*(theIncome_p/theFood_p)*se_b2
65 theIncome_p
66 elasticity_p
67 lower_elastic
68 upper_elastic

```

```

70 #e
71 mod_loglog<-lm(log(food)~log(income), data=data)
72 plot(log(income), log(food), xlab="Ln_income", ylab="Ln_food", col="grey")
73 b1_loglog<-coef(mod_loglog)[[1]]
74 b2_loglog<-coef(mod_loglog)[[2]]
75 se_b2_loglog<-summary(mod_loglog)$coefficients[2, 2]
76 abline(b1_loglog, b2_loglog, col="red")
77 #r square
78 rsq_linear<-summary(modIF)$r.squared#linear
79 rsq_loglog<-summary(mod_loglog)$r.squared#log-log
80 N<-nrow(data)
81 p1<-length(coef(modIF))-1
82 p2<-length(coef(mod_loglog))-1
83 G_rsq_linear<-1-exp((log(1-rsq_linear)*N)/(N-p1-1))
84 G_rsq_loglog<-1-exp((log(1-rsq_loglog)*N)/(N-p2-1))
85 G_rsq_linear
86 G_rsq_loglog
87
88 #f
89 #95%interval
90 n_elastic_loglog<-length(elasticity_loglog)
91 t_elastic_loglog<-qt(0.975, n_elastic_loglog-2)
92 lower_elastic_loglog<-b2_loglog-t_elastic_loglog*se_b2_loglog
93 upper_elastic_loglog<-b2_loglog+t_elastic_loglog*se_b2_loglog
94 b2_loglog
95 lower_elastic_loglog
96 upper_elastic_loglog
97
98 #g
99 e_IF_loglog<-resid(mod_loglog)
100 plot(log(income), e_IF_loglog, xlab="Ln_income", ylab="residual_log-log model")
101 hist(e_IF_loglog)
102 e_jbTest_loglog<-jarque.bera.test(residuals(mod_loglog))
103 e_jbTest_loglog

105 #h
106 mod_linearlog<-lm(food~log(income), data=data)
107 plot(log(income), food, xlab="Ln_income", ylab="food", col="grey")
108 b1_linearlog<-coef(mod_linearlog)[[1]]
109 b2_linearlog<-coef(mod_linearlog)[[2]]
110 se_b2_linearlog<-summary(mod_linearlog)$coefficient[2, 2]
111 abline(b1_linearlog, b2_linearlog, col="red")
112 #r square
113 rsq_linearlog<-summary(mod_linearlog)$r.squared
114 rsq_linearlog
115
116 #i
117 #95%interval
118 theIncome_p<-c(19, 65, 160)
119 theFood_p_linearlog<-b1_linearlog+b2_linearlog*log(theIncome_p)
120 margin_effect_linearlog<-b2_linearlog/theIncome_p
121 elasticity_p_linearlog<-margin_effect_linearlog*(theIncome_p/theFood_p_linearlog)
122 n_elastic_linearlog<-length(income)
123 t_elastic_linearlog<-qt(0.975, n_elastic_linearlog-2)
124 lower_elastic_linearlog<-elasticity_p_linearlog-t_elastic_linearlog*(se_b2_linearlog)/theFood_p_linearlog
125 upper_elastic_linearlog<-elasticity_p_linearlog+t_elastic_linearlog*(se_b2_linearlog)/theFood_p_linearlog
126 elasticity_p_linearlog
127 lower_elastic_linearlog
128 upper_elastic_linearlog
129
130 #j
131 e_IF_linearlog<-resid(mod_linearlog)
132 plot(log(income), e_IF_linearlog, xlab="Ln_income", ylab="residual_linear-log model")
133 hist(e_IF_linearlog)
134 e_jbTest_linearlog<-jarque.bera.test(residuals(mod_linearlog))
135 e_jbTest_linearlog

```