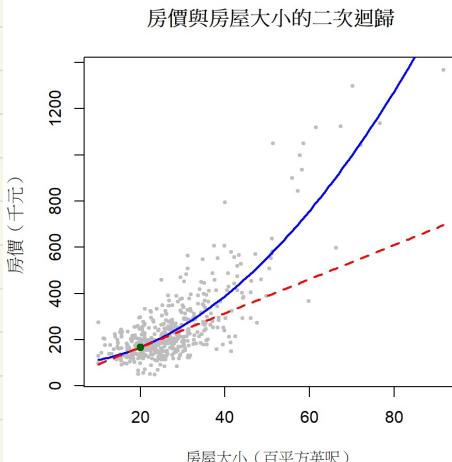


2.17 The data file *collegetown* contains observations on 500 single-family houses sold in Baton Rouge, Louisiana, during 2009–2013. The data include sale price (in thousands of dollars), $PRICE$, and total interior area of the house in hundreds of square feet, $SQFT$.

- Plot house price against house size in a scatter diagram.
- Estimate the linear regression model $PRICE = \beta_1 + \beta_2 SQFT + e$. Interpret the estimates. Draw a sketch of the fitted line.
- Estimate the quadratic regression model $PRICE = \alpha_1 + \alpha_2 SQFT^2 + e$. Compute the marginal effect of an additional 100 square feet of living area in a home with 2000 square feet of living space.
- Graph the fitted curve for the model in part (c). On the graph, sketch the line that is tangent to the curve for a 2000-square-foot house.
- For the model in part (c), compute the elasticity of $PRICE$ with respect to $SQFT$ for a home with 2000 square feet of living space.
- For the regressions in (b) and (c), compute the least squares residuals and plot them against $SQFT$. Do any of our assumptions appear violated?
- One basis for choosing between these two specifications is how well the data are fit by the model. Compare the sum of squared residuals (SSE) from the models in (b) and (c). Which model has a lower SSE ? How does having a lower SSE indicate a “better-fitting” model?

c. d.



$$price = 93.57 + 0.185 SQFT^2$$

marginal effect \Rightarrow

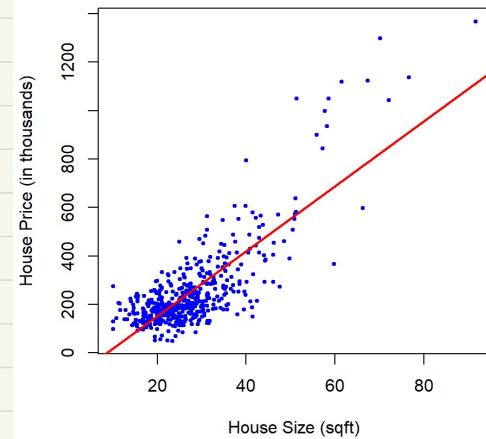
$$\frac{d PRICE}{d SQFT} = 2 \alpha_2 SQFT$$

$$= 2 \times 0.185 \times 20 = 7.38$$

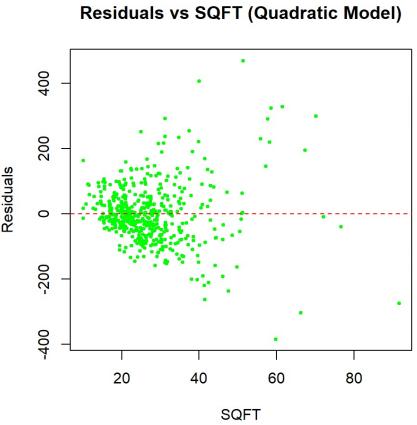
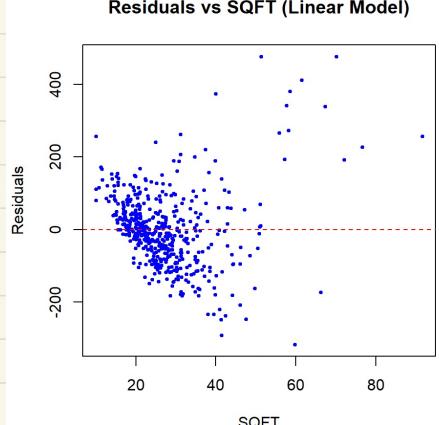
$$e. E = \frac{d PRICE}{d SQFT} \times \frac{SQFT}{PRICE} = 7.38 \times \frac{20}{167.17} = 0.882$$

a. b. $price = -155.42 + 13.4 SQFT$

Scatter Plot of House Price vs. House Size



f.



```
> SSE_linear  
[1] 5262847  
> SSE_quad  
[1] 4222356
```

In the linear model, the residuals exhibit heteroskedasticity, with the variance of residuals increasing as SQFT increases. There also appears to be a systematic pattern in the residuals, suggesting that the linear specification may be inappropriate.

Even in the quadratic model, the residuals appear to fan out as SQFT increases, particularly beyond 40 hundred square feet (4000 sqft). This suggests that the issue of heteroskedasticity is not fully resolved, although it is less severe compared to the linear model.

g

In comparing the two models, we find that the quadratic model in (c) has a lower sum of squared residuals (SSE) than the linear model in (b). A lower SSE indicates that the predicted values from the model are, on average, closer to the observed data points. Therefore, the quadratic model provides a better in-sample fit to the data. However, we should also be cautious of overfitting, as adding complexity to a model will often reduce SSE at the cost of generalizability.

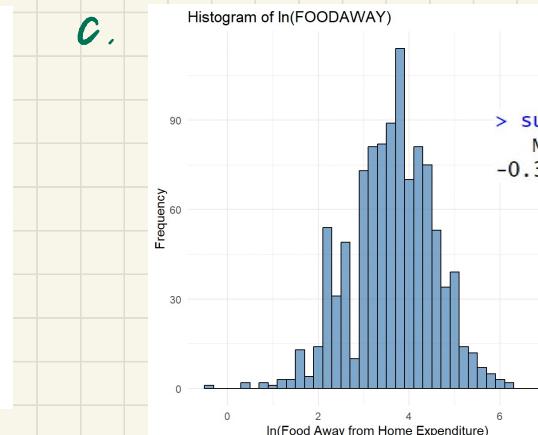
2.25 Consumer expenditure data from 2013 are contained in the file *cex5_small*. [Note: *cex5* is a larger version with more observations and variables.] Data are on three-person households consisting of a husband and wife, plus one other member, with incomes between \$1000 per month to \$20,000 per month. *FOODAWAY* is past quarter's food away from home expenditure per month per person, in dollars, and *INCOME* is household monthly income during past year, in \$100 units.

- Construct a histogram of *FOODAWAY* and its summary statistics. What are the mean and median values? What are the 25th and 75th percentiles?
- What are the mean and median values of *FOODAWAY* for households including a member with an advanced degree? With a college degree member? With no advanced or college degree member?
- Construct a histogram of $\ln(\text{FOODAWAY})$ and its summary statistics. Explain why *FOODAWAY* and $\ln(\text{FOODAWAY})$ have different numbers of observations.
- Estimate the linear regression $\ln(\text{FOODAWAY}) = \beta_1 + \beta_2 \text{INCOME} + e$. Interpret the estimated slope.
- Plot $\ln(\text{FOODAWAY})$ against *INCOME*, and include the fitted line from part (d).
- Calculate the least squares residuals from the estimation in part (d). Plot them vs. *INCOME*. Do you find any unusual patterns, or do they seem completely random?

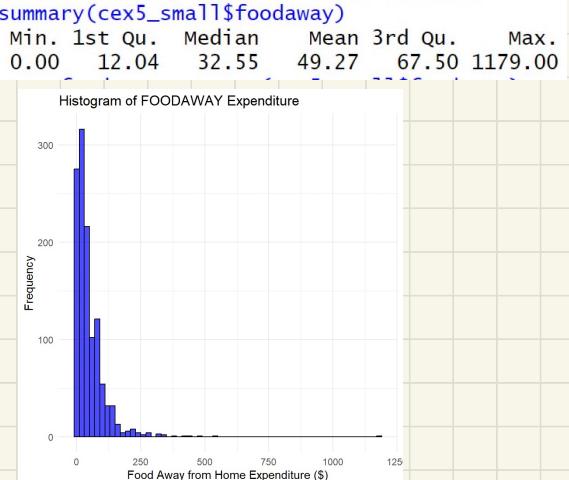
b.

```
> mean(adv_degree$foodaway)
[1] 73.15494
> median(adv_degree$foodaway)
[1] 48.15
>
> mean(college_degree$foodaway)
[1] 48.59718
> median(college_degree$foodaway)
[1] 36.11
>
> mean(no_degree$foodaway)
[1] 39.01017
> median(no_degree$foodaway)
[1] 26.02
```

c.



a.



> summary(cex5_small\$ln_foodaway)

	Min.	1st Qu.	Median	Mean	3rd Qu.	Max.
ln_foodaway	-0.3011	3.0759	3.6865	3.6508	4.2797	7.0724

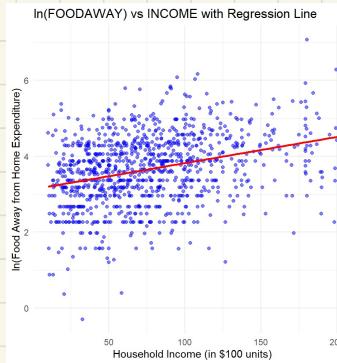
$\ln(\text{FOODAWAY})$ 值减少 178 个。
 ∵ 有 178 个家庭在 foodaway = 0
 $\ln(0)$ 未定义 → 缺失值

d.

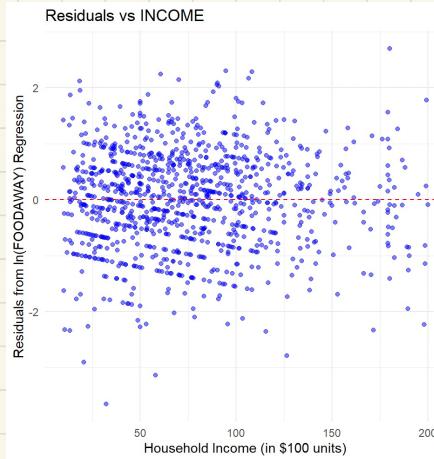
$$\ln(\text{FOODAWAY}) = 3.1293 + 0.0069 \text{ INCOME}$$

At \$100 income, $\log \text{ of food away} = 3.1293 + 0.0069 \times 100 = 3.7979$
 or 79.79%

e.



f.



They seem completely random

2.28 How much does education affect wage rates? The data file *cps5_small* contains 1200 observations on hourly wage rates, education, and other variables from the 2013 Current Population Survey (CPS). [Note: *cps5* is a larger version.]

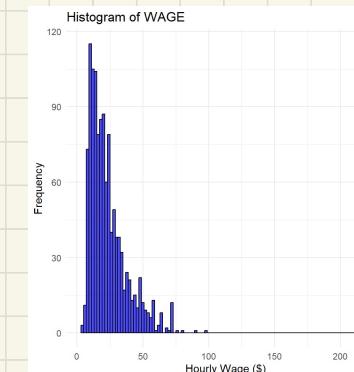
- Obtain the summary statistics and histograms for the variables *WAGE* and *EDUC*. Discuss the data characteristics.
- Estimate the linear regression $WAGE = \beta_1 + \beta_2 EDUC + e$ and discuss the results.
- Calculate the least squares residuals and plot them against *EDUC*. Are any patterns evident? If assumptions SR1–SR5 hold, should any patterns be evident in the least squares residuals?
- Estimate separate regressions for males, females, blacks, and whites. Compare the results.
- Estimate the quadratic regression $WAGE = \alpha_1 + \alpha_2 EDUC^2 + e$ and discuss the results. Estimate the marginal effect of another year of education on wage for a person with 12 years of education and for a person with 16 years of education. Compare these values to the estimated marginal effect of education from the linear regression in part (b).
- Plot the fitted linear model from part (b) and the fitted values from the quadratic model from part (e) in the same graph with the data on *WAGE* and *EDUC*. Which model appears to fit the data better?

a. > summary(cps5_small\$wage)

	Min.	1st Qu.	Median	Mean	3rd Qu.	Max.
	3.94	13.00	19.30	23.64	29.80	221.10

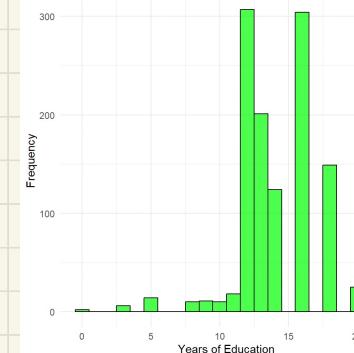
> summary(cps5_small\$educ)

	Min.	1st Qu.	Median	Mean	3rd Qu.	Max.
	0.0	12.0	14.0	14.2	16.0	21.0



WAGE 右偏分布

Histogram of EDUC

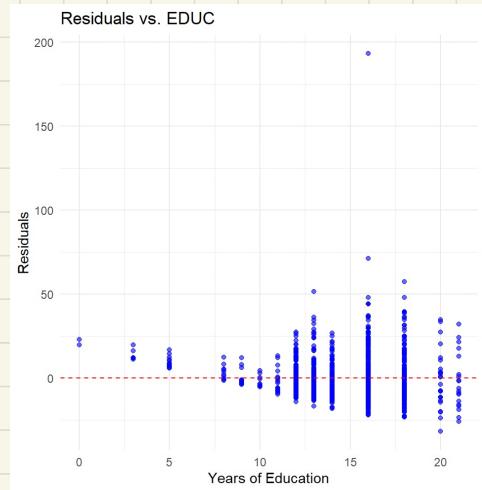


EDUC 大部分为 12~16 年

b. $WAGE = -0.4 + 2.3968 \cdot EDUC$

EDUC 增加 1 年 \rightarrow WAGE 增加 2.3968 \$

C.



$EDUC \uparrow$ residual \uparrow

\Rightarrow 不滿足 Homoscedasticity (SR5)

若 SR1-SR5 成立，殘差應在各個 EDUC 下均勻分佈

d. 男: $wage = -8.2849 + 2.3785 \text{ educ}$

女: $wage = -16.6028 + 2.6595 \text{ educ}$

黑: $wage = -6.2541 + 1.9233 \text{ educ}$

白: $wage = -10.4747 + 2.4178 \text{ educ}$

] 女性 教育对薪资影響較大，但都距低

] 黑人 教育对薪资影響較大，但都距低

e. $\text{WAGE} = 4.9165 + 0.08913 \text{ EDUC}^2$

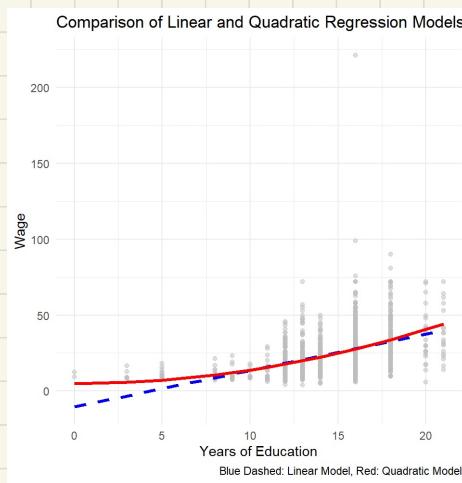
marginal effect = $2 \times 0.08913 \times \text{EDUC}$

$$\text{EDUC} = 12 \Rightarrow 2 \times 0.08913 \times 12 = 2.139$$

$$\text{EDUC} = 16 \Rightarrow 2 \times 0.08913 \times 16 = 2.852$$

二次回歸的 marginal effect 會隨 EDUC 增大而變大

f.



quadratic regression fits the data better