

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.



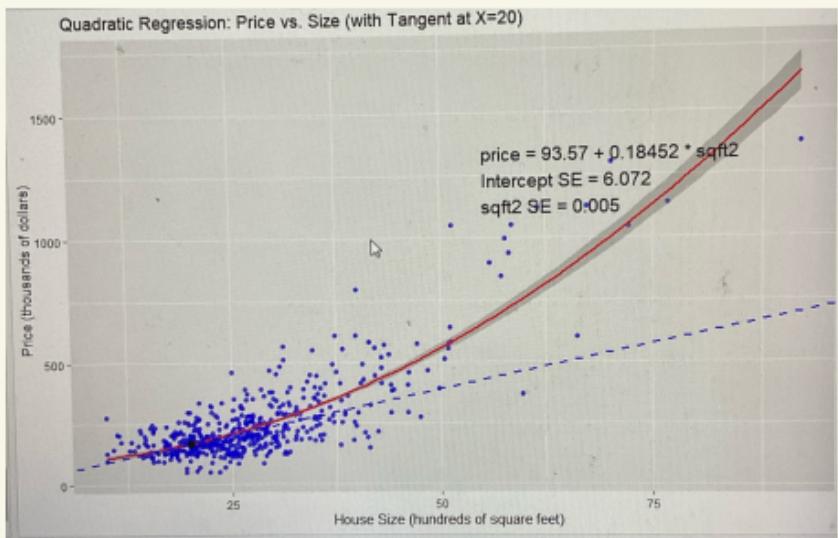
$$\widehat{\text{price}} = -115.42 + 13.4 \text{ sqft}$$

其他條件相同每增加 100 單位的 sqft 會增加  
13.4 thousands 的 home price

$\beta_1 = -115.42$  代表在 0 sqft 時，  
我們預估的 home price 為 -115.42 thousands

- 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.

## Quadratic Regression: Price vs. Size (with Tangent at X=20)



$$\widehat{\text{price}} = 93.57 + 0.18452 * \text{sqft}^2$$

增加 2000 sqft, price 增加  
~~1000~~

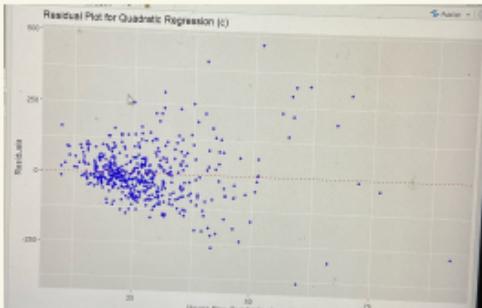
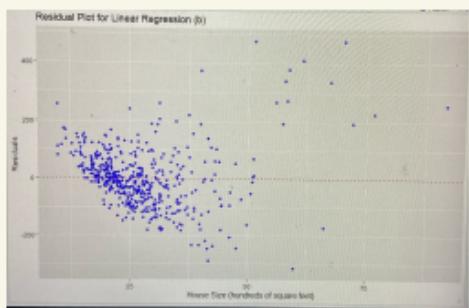
$$0.18452 \times 20^2 = 7.38 \text{ thousands}$$

- e. 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.

$$\begin{aligned} \epsilon &= \frac{\partial \text{price}}{\partial \text{sqft}} \times \frac{\text{sqft}}{\text{price}} = 2B_2 \cdot 20 \cdot \frac{20}{167.37} \\ &= 0.8819 \end{aligned}$$

$$\text{sqft} = 20 \Rightarrow \widehat{\text{price}} = 93.57 + 0.18452 \times 20^2 = 167.37$$

- f. For the regressions in (b) and (c), compute the least squares residuals and plot them against  $SQFT$ . Do any of our assumptions appear violated?



殘差值隨  $SQFT$  而變大，違反  
同質變異數假設

- g. 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?

$$SSE \text{ model (b)} = 5262847$$

$$SSE \text{ model (c)} = 4222356$$

∴ quadratic model has lower SSE

$SSE = \sum (y_i - \hat{y}_i)^2$ , 較小的 SSE 代表  
資料更貼近預估模型

# 2.17 (a)(b) R

```
① Untitled1 × ② Untitled2 × collegetown × model × ④ 2.17ab.R × History ×
Source on Save
1 model <- lm(price ~ sqft, data = collegetown)
2 coefficients <- summary(model)$coefficients
3 coef_intercept <- coefficients[1, 1] # 截距
4 coef_slope <- coefficients[2, 1] # 斜率
5 intercept_se <- coefficients[1, 2] # 截距標準誤
6 slope_se <- coefficients[2, 2] # 斜率標準誤
7 regression_eq <- paste0("price = ", round(coef_intercept, 2), " + ", round(coef_slope, 2), " * sqft")
8 regression_se <- paste0("SE(Intercept) = ", round(intercept_se, 2), ", SE(Slope) = ", round(slope_se, 2))
9 ggplot(collegetown, aes(x = sqft, y = price)) +
10 geom_point(color = "#00AEEF", alpha = 0.6) + # 數點圖
11 geom_smooth(method = "lm", se = FALSE, color = "#00AEEF", size = 1) + # 回歸線
12 labs(
13 title = "Scatter Plot of Price and Size with Regression Line",
14 x = "House Size (Hundreds of square feet)",
15 y = "Price (thousands of dollars)"
16 ) +
17 annotate("text", x = max(collegetown$sqft) * 0.5, y = max(collegetown$price) * 0.8,
18 label = paste(regression_eq, "\n", regression_se), size = 5, color = "#00AEEF", hjust = 0) # 調整大小
19
20
```

# 2.17 (c)(d) R

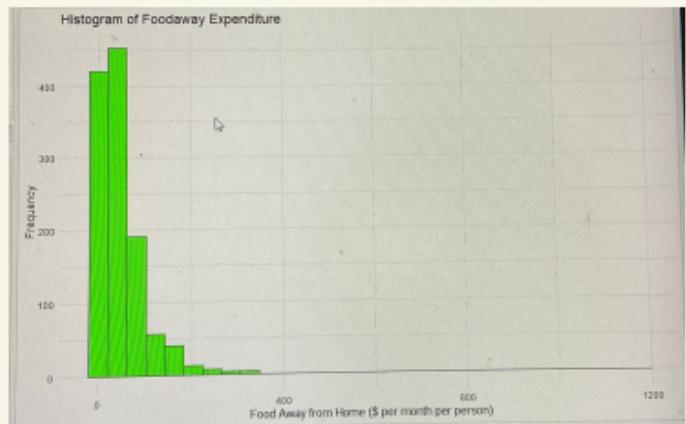
```
① # 創建平方項
2 collegetown$sqft2 <- collegetown$sqft^2
3
4 # 估計線性回歸模型
5 quad_model <- lm(price ~ sqft2, data = collegetown)
6
7
8 coef_intercept <- coef(quad_model)[1]
9 coef_sqft2 <- coef(quad_model)[2]
10
11 # 計算殘差
12 se_intercept <- summary(quad_model)$coefficients[1, 2]
13 se_sqft2 <- summary(quad_model)$coefficients[2, 2]
14
15
16 x0 <- 20
17 y0 <- coef_intercept + coef_sqft2 * x0^2
18 slope_tangent <- 2 * coef_sqft2 * x0 # x=20切線斜率
19
20 # 回歸公式包含標準差
21 regression_eq <- paste0("price = ",
22 round(coef_intercept, 2), " + ",
23 round(coef_sqft2, 5), " * sqft^2\n",
24 "Intercept SE = ", round(se_intercept, 3), "\n",
25 "sqft? SE = ", round(se_sqft2, 3))
26
27
28 ggplot(collegetown, aes(x = sqft, y = price)) +
29 geom_point(color = "#00AEEF", alpha = 0.6) + # 數點圖
30 geom_smooth(method = "lm", formula = y ~ I(x^2), se = TRUE, color = "#00AEEF", size = 1) + # 二次回歸線
31 geom_abline(intercept = y0 - slope_tangent * x0, slope = slope_tangent, color = "#00AEEF", linetype = "dashed", size = 1) + # 切線
32 geom_point(aes(x = x0, y = y0), color = "black", size = 3) +
33 labs(
34 title = "Quadratic Regression: Price vs. Size (with Tangent at X=20)",
35 x = "House Size (Hundreds of square feet)",
36 y = "Price (thousands of dollars)"
37 ) +
38 annotate("text", x = max(collegetown$sqft) * 0.6, y = max(collegetown$price) * 0.9,
39 label = regression_eq, size = 5, color = "#00AEEF", hjust = 0)
40
41
```

# 2.17 (f)(g), R

```
1 # 創建平方項
2 collegetown$sqft2 <- collegetown$sqft^2
3
4
5 linear_model <- lm(price ~ sqft, data = collegetown)
6
7
8 quad_model <- lm(price ~ sqft2, data = collegetown)
9
10 # 計算殘差並存入資料框
11 collegetown$linear_resid <- collegetown$price - predict(linear_model)
12 collegetown$quad_resid <- collegetown$price - predict(quad_model)
13
14
15 p1 <- ggplot(collegetown, aes(x = sqft, y = linear_resid)) +
16   geom_point(color = "blue", alpha = 0.6) +
17   geom_hline(yintercept = 0, color = "red", linetype = "dashed") +
18   labs(
19     title = "Residual Plot for Linear Regression (b)",
20     x = "House Size (hundreds of square feet)",
21     y = "Residuals"
22   )
23
24
25 p2 <- ggplot(collegetown, aes(x = sqft, y = quad_resid)) +
26   geom_point(color = "blue", alpha = 0.6) +
27   geom_hline(yintercept = 0, color = "red", linetype = "dashed") +
28   labs(
29     title = "Residual Plot for Quadratic Regression (c)",
30     x = "House Size (hundreds of square feet)",
31     y = "Residuals"
32   )
33
34
35 print(p1)
36 print(p2)
37 linear_resid <- collegetown$price - predict(linear_model)
38 quad_resid <- collegetown$price - predict(quad_model)
39 SSE_linear <- sum(linear_resid^2)
40 SSE_quad <- sum(quad_resid^2)
41 cat("Linear SSE:", SSE_linear, "\n") #顯示結果
42 cat("Quadratic SSE:", SSE_quad, "\n")
43
```

**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.

- a. Construct a histogram of *FOODAWAY* and its summary statistics. What are the mean and median values? What are the 25th and 75th percentiles?



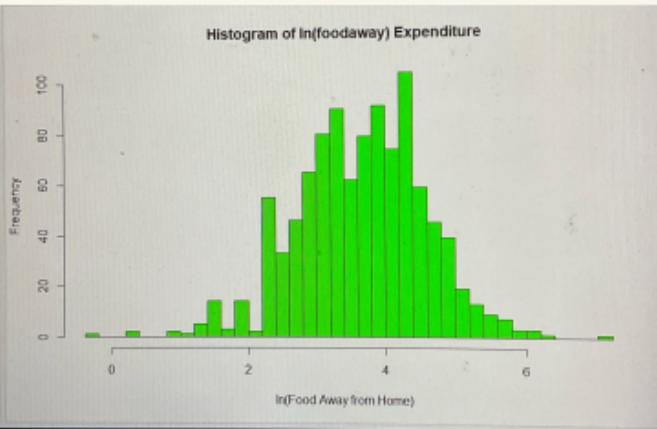
Mean      Median      25th      75th  
49.21085    32.555    12.04    67.5025

- b. 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?

	category	foodaway.Mean	foodaway.Median	N
1	Advanced Degree	73.15494	48.15000	257
2	College Degree	48.59718	36.11000	369
3	No College/Advanced Degree	39.01017	26.02000	574

- c. Construct a histogram of *ln(FOODAWAY)* and its summary statistics. Explain why *FOODAWAY* and *ln(FOODAWAY)* have different numbers of observations.

```
> print(summary_log)
      Mean    Median     Q25     Q75 Count_Log Count_Original
25% 3.650804 3.686499 3.075929 4.279717      1022          1200
      50% 4.000000 4.035696 4.000000 4.035696      1022          1200
      75% 4.350804 4.386499 4.300000 4.386499      1022          1200
      90% 4.650804 4.686499 4.500000 4.686499      1022          1200
      95% 4.850804 4.886499 4.700000 4.886499      1022          1200
      Max 5.000000 5.000000 5.000000 5.000000      1022          1200
```

Histogram of  $\ln(\text{foodaway})$  Expenditure

$$\text{Count\_Original} = 1200$$

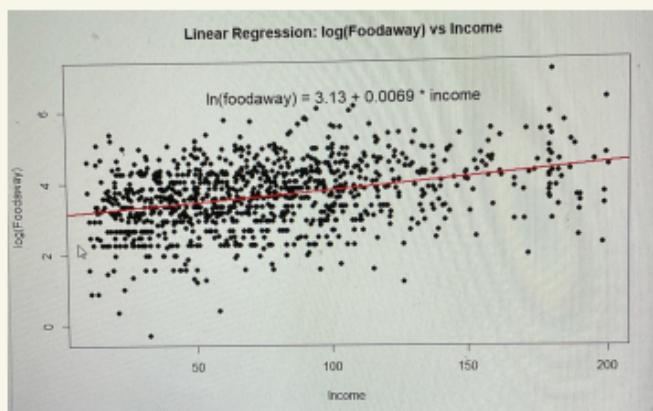
$$\text{Count\_Log} = 1022 \quad 1200 - 1022 = 178$$

有178個家庭花費0元在Foodaway上

$\because \ln(0)$  is not defined

∴ 這178個值無法用在迴歸上

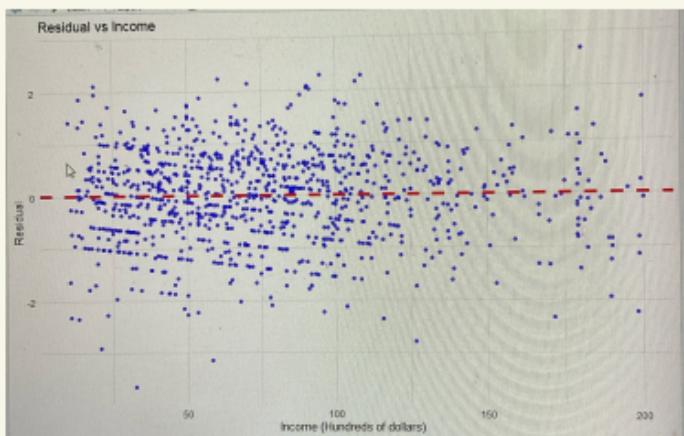
- d. Estimate the linear regression  $\ln(\text{FOODAWAY}) = \beta_1 + \beta_2 \text{INCOME} + e$ . Interpret the estimated slope.
- e. Plot  $\ln(\text{FOODAWAY})$  against  $\text{INCOME}$ , and include the fitted line from part (d).



$$\text{model: } \ln(\text{foodaway}) = 3.13 + 0.0069 \text{ income}$$

其他條件相同時，每增加 100 income，  
會增加 0.69% 的  $\ln(\text{foodaway})$  的花  
費

- f. 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?



殘差圖看起來是隨機沒有特定模  
式的

# 2.25 R

```
1 summary_stats <- data.frame(
2   Mean = mean(cex5_small$foodaway, na.rm = TRUE),
3   Median = median(cex5_small$foodaway, na.rm = TRUE),
4   Q25 = quantile(cex5_small$foodaway, 0.25, na.rm = TRUE),
5   Q75 = quantile(cex5_small$foodaway, 0.75, na.rm = TRUE)
6 )
7 print(summary_stats)
8 ggplot(cex5_small, aes(x = foodaway)) +
9   geom_histogram(color = "black", fill = "#FFFFE0", bins = 30, alpha = 0.7) +
10  labs(
11    title = "Histogram of Foodaway Expenditure",
12    x = "Food Away From Home ($ per month per person)",
13    y = "Frequency"
14  ) +
15  theme_minimal()
16
17
18 # b. 創建category
19 cex5_small$category <- ifelse(cex5_small$advanced == 1, "Advanced Degree",
20                               ifelse(cex5_small$college == 1 & cex5_small$advanced == 0, "College Degree",
21                                     "No College/Advanced Degree"))
22
23 summary_stats <- aggregate(foodaway ~ category, data = cex5_small,
24                             FUN = function(x) c(Mean = mean(x, na.rm = TRUE), Median = median(x, na.rm = TRUE)))
25
26 # 計算每個category的數量N
27 category_counts <- as.data.frame(table(cex5_small$category))
28 colnames(category_counts) <- c("category", "N")
29
30 # 合併數量N到summary_stats
31 summary_stats <- merge(summary_stats, category_counts, by = "category")
32 print(summary_stats)
33
34 #c. #M*log(foodaway) 對於foodaway>0的值進行轉換；其他為NA
35 cex5_small$log_foodaway <- ifelse(cex5_small$log_foodaway > 0, log(cex5_small$foodaway), NA)
36 Mean_log <- mean(cex5_small$log_foodaway, na.rm = TRUE)
37 Median_log <- median(cex5_small$log_foodaway, na.rm = TRUE)
38 Q25_log <- quantile(cex5_small$log_foodaway, 0.25, na.rm = TRUE)
39 Q75_log <- quantile(cex5_small$log_foodaway, 0.75, na.rm = TRUE)
40 Count_Log <- sum(!is.na(cex5_small$log_foodaway))#計算有效值的數量
41 Count_Original <- sum(!is.na(cex5_small$foodaway))
42 summary_log <- data.frame(
43   Mean = Mean_log,
44   Median = Median_log,
45   Q25 = Q25_log,
46   Q75 = Q75_log,
47   Count_Log = Count_Log,
48   Count_Original = Count_Original
49 )
50 print(summary_log)
51 hist(cex5_small$log_foodaway, breaks = 30, col = "GREEN", border = "black",
52       xlab = "ln(Food Away From Home)", ylab = "Frequency",
53       main = "Histogram of ln(foodaway) Expenditure")
54
55 #d. e.
56 cex5_clean <- cex5_small[!is.na(cex5_small$log_foodaway) & !is.na(cex5_small$income), ]#清理數據 去除NA
57 model <- lm(log_foodaway ~ income, data = cex5_clean)
58 summary(model)
59 intercept <- coef(model)[1]
```

```

60 slope <- coef(model)[2]
61 regression_eq <- paste0("ln(foodaway) = ", round(intercept, 2),
62                         " + ", round(slope, 4), " * income")
63 x_pos <- max(cx5_clean$income) - 0.5 * (max(cx5_clean$income) - min(cx5_clean$income)) #調整位置
64 y_pos <- max(cx5_clean$log_foodaway) * 0.9
65 print(regression_eq)
66 plot(cx5_clean$income, cx5_clean$log_foodaway,
67      main = "Linear Regression: log(Foodaway) vs Income",
68      xlab = "Income", ylab = "log(Foodaway)",
69      pch = 19, col = "black")
70 abline(model, col = "red", lwd = 2)
71 text(x_pos, y_pos, labels = regression_eq, col = "black", cex = 1.3)
72
73 #f,
74 cx5_clean$residual <- resid(model)
75 ggplot(cx5_clean, aes(x = income, y = residual)) +
76   geom_point(colour = "blue", alpha = 0.6) +
77   geom_hline(yintercept = 0, color = "red", linetype = "dashed", size = 1.5) + #添加紅色虛線表示0殘差的基準線
78   labs(
79     title = "Residual vs Income",
80     x = "Income (Hundreds of dollars)", y = "Residual",
81   ) +
82   theme_minimal()
83

```

**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.]

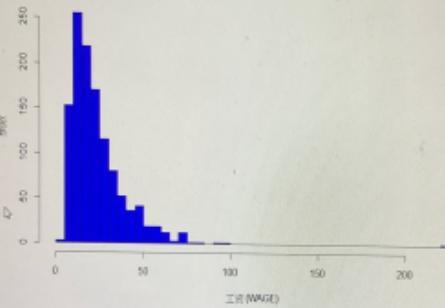
- a. Obtain the summary statistics and histograms for the variables *WAGE* and *EDUC*. Discuss the data characteristics.

```

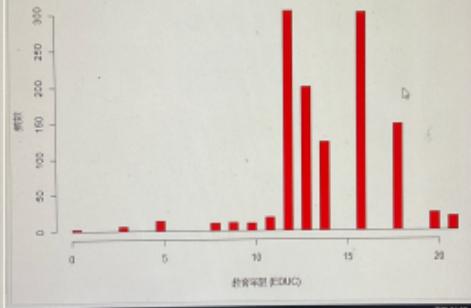
print(summary_stats)
Mean_wage Median_wage SD_wage Min_wage Max_wage Mean_educ Median_educ SD_educ Min_educ Max_educ
23.64004    19.3 15.21655     3.94    221.1   14.2025    14.2.890811     0     21

```

工資分布



教育年限分布



平均工資為 23.6 但中位數僅為 19.3  
工資分佈呈現右偏，並且標準差較大

直方圖顯示大多數人工資較低，少部分工資很高，有工資不均衡的情況

平均教育年限為 14 中位數也是 14，分佈較為對稱，直方圖顯示大部分人受教育程度集中在 12~16 年之間

- b. Estimate the linear regression  $WAGE = \beta_1 + \beta_2 EDUC + e$  and discuss the results.

```
> print(regression_eq) # 输出回归方程  
[1] "Wage = -10.40 + 2.40 * Educ"  
> |
```

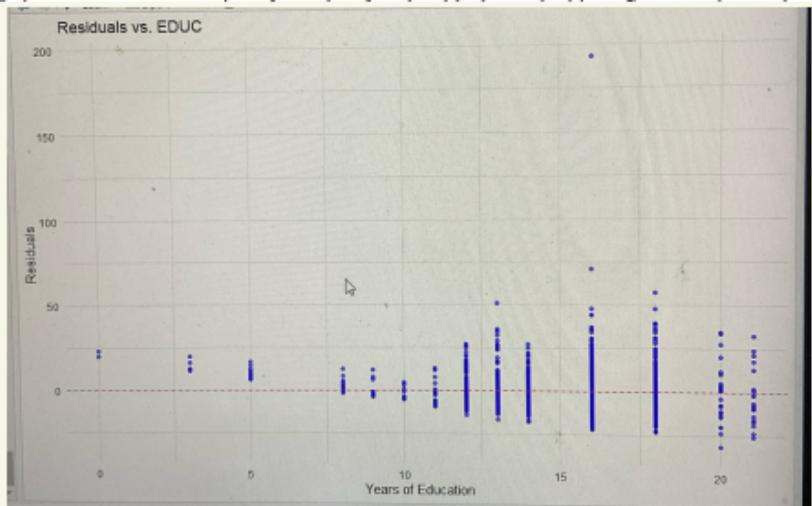
Coefficients:					
	Estimate	Std. Error	t value	Pr(> t )	
(Intercept)	-10.4000	1.9624	-5.3	1.38e-07	***
educ	2.3968	0.1354	17.7	< 2e-16	***
---					

model:  $\widehat{Wage} = -10.4 + 2.4 \text{ Educ}$

$\beta_1 = 2.4$ ：每多一年教育，工資平均增加 2.4

P 值 ( $< 2e-16$ )：educ 對 wage 具有顯著統計關係

- c. 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?



$EDUC \uparrow$ , 殘差也變大, 不滿足同質  
變異數(SR5),  
如果滿足SR1-SR5, 殘差應該在  
各個  $EDUC$  水平下均勻分佈

- d. Estimate separate regressions for males, females, blacks, and whites. Compare the results.

males:  $\widehat{WAGE} = -8.2849 + 2.3785 EDUC$   
 females:  $\widehat{WAGE} = -16.6028 + 2.6595 EDUC$

$$\text{whites: } \widehat{\text{WAGE}} = -10.475 + 2.418 \text{ EDUC}$$

$$\text{blacks: } \widehat{\text{WAGE}} = -6.2541 + 1.9233 \text{ EDUC}$$

女性的教育回報(2.6595) > 男性的教育回報(2.3785). 說明每多一年教育, 女性能獲得更多的薪資增幅, 女性截距低於男性, 說明在低教育水準下, 女性薪資較男性低

白人的教育回報(2.418) > 黑人的教育回報(1.9233). 說明每多一年教育, 白人能獲得更多的薪資增幅, 白人截距低於黑人, 說明在低教育水準下, 黑人薪資較白人高

- e. Estimate the quadratic regression  $\text{WAGE} = \alpha_1 + \alpha_2 \text{EDUC}^2 + e$  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).

## Coefficients:

	Estimate	Std. Error	t value	Pr(> t )
(Intercept)	4.916477	1.091864	4.503	7.36e-06 ***
I(educ^2)	0.089134	0.004858	18.347	< 2e-16 ***

$$\widehat{WAGE} = 4.916477 + 0.089134 \times EDUC^2$$

$$\text{marginal effect} = 2 \times 0.089134 \times EDUC$$

$$EDUC = 12, \text{ margin} = 2 \times 0.089134 \times 12 = 2.1392$$

$$EDUC = 16, \text{ margin} = 2 \times 0.089134 \times 16 = 2.8523$$

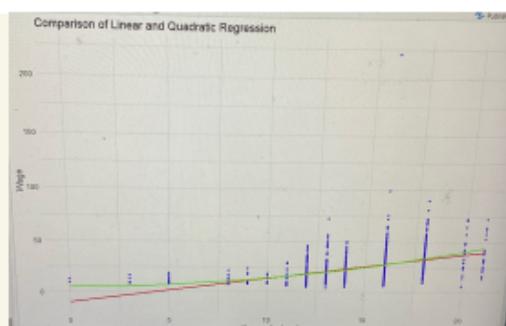
$$\text{from (b)} \quad \widehat{WAGE} = -10.4 + 2.4 \times EDUC$$

$$\text{marginal effect} = 2.4$$

在  $EDUC = 12$  時，線性迴歸模型對工資影響更大

在  $EDUC = 16$  時，二次迴歸模型對工資影響更大

- f. Plot the fitted linear model from part (b) and the fitted values from part (e) in the same graph with the data on WAGE and EDUC. Which model appears to fit the data better?



二次迴歸模型更貼近散點分佈，更合適

# 2.28 R

```
1 #a.
2 head(cps5_small)
3 summary_stats <- data.frame(
4   Mean_wage = mean(cps5_small$wage, na.rm = TRUE),
5   Median_wage = median(cps5_small$wage, na.rm = TRUE),
6   SD_wage = sd(cps5_small$wage, na.rm = TRUE),
7   Min_wage = min(cps5_small$wage, na.rm = TRUE),
8   Max_wage = max(cps5_small$wage, na.rm = TRUE),
9   Mean_educ = mean(cps5_small$educ, na.rm = TRUE),
10  Median_educ = median(cps5_small$educ, na.rm = TRUE),
11  SD_educ = sd(cps5_small$educ, na.rm = TRUE),
12  Min_educ = min(cps5_small$educ, na.rm = TRUE),
13  Max_educ = max(cps5_small$educ, na.rm = TRUE)
14 )
15
16 print(summary_stats)
17 # 工资 (WAGE) 直方图
18 hist(cps5_small$wage, breaks = 50, col = "blue", main = "工资分布", xlab = "工资 (WAGE)", ylab = "频数")
19
20 # 教育年限 (EDUC) 直方图
21 hist(cps5_small$educ, breaks = 50, col = "red", main = "教育年限分布", xlab = "教育年限 (EDUC)", ylab = "频数")
22
23 #d.
24 model <- lm(wage ~ educ, data = cps5_small)
25 summary(model)
26 intercept <- coef(model)[1]
27 slope <- coef(model)[2]
28 regression_eq <- sprintf("Wage = %.2f + %.2f * Educ", intercept, slope)
29 print(regression_eq) # 输出回归方程
30
31 #e.
32 cps5_small$residuals <- residuals(model)
33 ggplot(cps5_small, aes(x = educ, y = residuals)) +
34   geom_point(color = "blue", alpha = 0.5) +
35   geom_hline(yintercept = 0, linetype = "dashed", color = "red") +
36   labs(title = "Residuals vs. EDUC",
37         x = "Years of Education",
38         y = "Residuals") +
39   theme_minimal()
40
41 male_model <- lm(wage ~ educ, data = cps5_small[cps5_small$female == 0, ])
42 female_model <- lm(wage ~ educ, data = cps5_small[cps5_small$female == 1, ])
43 white_model <- lm(wage ~ educ, data = cps5_small[cps5_small$black == 0, ])
44 black_model <- lm(wage ~ educ, data = cps5_small[cps5_small$black == 1, ])
45 summary(male_model)
46 summary(female_model)
47 summary(white_model)
48 summary(black_model)
49
50 #f.
51 quad_model <- lm(wage ~ I(educ^2), data = cps5_small)
52 cps5_small$cps5_small_prdct <- predict(quad_model)
53 summary(quad_model)
54
55 #f.
56 linear_model = lm(wage ~ educ, data = cps5_small)
57 quadratic_model = lm(wage ~ I(educ^2), data = cps5_small)
58 cps5_small$cps5_linear <- predict(linear_model)
59 cps5_small$cps5_quadratic <- predict(quadratic_model)
60
61 ggplot(cps5_small, aes(x = educ, y = wage)) +
62   geom_point(color = "blue", size = 1, alpha = 0.6) +
63   geom_line(aes(y = cps5_linear), color = "red", size = 1) +
64   geom_line(aes(y = cps5_quadratic), color = "green", size = 1) +
65   labs(title = "Comparison of Linear and Quadratic Regression",
66         x = "Years of education",
67         y = "Wage") +
68   theme_minimal()
69
70
71
```