

TA Session 2

Chi-Yuan Fang

March 9, 2021

Contents

1	Introduction	1
1.1	TA Information	1
1.2	TA Sessions Schedule	1
1.3	Reference	2
2	Empirical Exercise 4.1	2
3	Empirical Exercise 4.2	6

1 Introduction

1.1 TA Information

TA: Chi-Yuan Fang

TA sessions: Tuesday 1:20 – 3:10 PM (SS 501)

Email: r09323017@ntu.edu.tw

Office hours: Friday 2:00 – 3:30 PM or by appointments (SS 643)

Class group on Facebook: Statistics (Fall 2020) and Econometrics (Spring 2021)

<https://www.facebook.com/groups/452292659024369/>

Because screens are not clear in SS 501, I will provide the link of live streaming in the group.

1.2 TA Sessions Schedule

Week	TA Sessions	Quiz	Content	Remind
1	02/23: No class			
2	03/02: Class 1		Function, Confidence Interval, T test	03/10 Turn in HW1
3	03/09: Class 2		Loops, Linear Model	03/10 Turn in HW1, 03/16 Quiz 1
4	03/16: Class 3	Quiz 1		03/24 Turn in HW2
5	03/23: Class 4			03/24 Turn in HW2, 03/30 Quiz 2
6	03/30: Class 5	Quiz 2		04/14 Turn in HW3
7	04/06: No class			04/14 Turn in HW3
8	04/13: Class 6			04/14 Turn in HW3, 04/20 Quiz 3
9	04/20: Class 7	Quiz 3		04/28 Midterm
10	04/27: Class 8		Review and Q&A	04/28 Midterm , 05/05 Turn in HW4
11	05/04: Class 9			05/05 Turn in HW4, 05/11 Quiz 4
12	05/11: Class 10	Quiz 4		05/19 Turn in HW5

Week	TA Sessions	Quiz	Content	Remind
13	05/18: Class 11			05/19 Turn in HW5, 05/25 Quiz 5
14	05/25: Class 12	Quiz 5		06/02 Turn in HW6
15	06/01: Class 13			06/02 Turn in HW6, 06/08 Quiz 6
16	06/08: Class 14	Quiz 6	Review and Q&A	06/16 Final Exam
17	06/15: No class			06/16 Final Exam
18	06/22: No class			

1.3 Reference

Introduction to Econometrics with R

<https://www.econometrics-with-r.org>

R for Data Science

<https://r4ds.had.co.nz>

R Markdown

<https://rmarkdown.rstudio.com>

Introduction to R Markdown

<https://rpubs.com/brandonkopp/RMarkdown>

What is a good book on learning R with examples?

<https://www.quora.com/What-is-a-good-book-on-learning-R-with-examples>

2 Empirical Exercise 4.1

On the text website, <http://www.pearsonglobaleditions.com>, you will find the data file **Growth**, which contains data on average growth rates from 1960 through 1995 for 65 countries, along with variables that are potentially related to growth. A detailed description is given in **Growth_Description**, also available on the website. In this exercise, you will investigate the relationship between growth and trade.

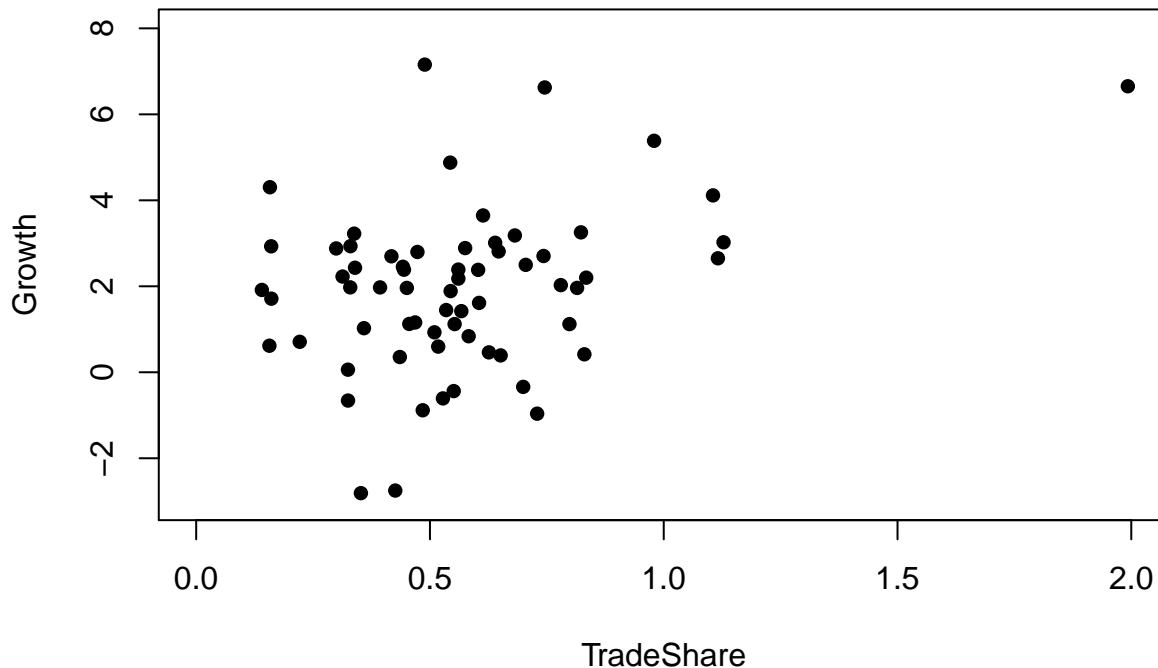
- Construct a scatterplot of average annual growth rate (*Growth*) on the average trade share (*TradeShare*). Does there appear to be a relationship between the variables?

Solution

```
# import data
#install.packages("readxl")
library(readxl)
Growth <- read_xlsx("Growth/Growth.xlsx")

plot(x = Growth$tradeshare,
     y = Growth$growth,
     pch = 16, # filled circle
     col = "black",
     xlim = c(0, 2),
     ylim = c(-3, 8),
     xlab = "TradeShare",
     ylab = "Growth",
     main = "E4.1 (a)")
```

E4.1 (a)



- b. One country, Malta, has a trade share much larger than the other countries. Find Malta on the scatterplot. Does Malta look like an outlier?

Solution

```
Growth[Growth$country_name == "Malta",]
```

```
## # A tibble: 1 x 8
##   country_name growth    oil rgdp60 tradeshare yearsschool rev_coups
##   <chr>         <dbl> <dbl> <dbl>      <dbl>      <dbl>      <dbl>
## 1 Malta          6.65     0  1374      1.99        5.64         0
## # ... with 1 more variable: assassinations <dbl>
```

Malta is the “outlying” observation with a trade share of 1.99.

- c. Using all observations, run a regression of *Growth* on *TradeShare*. What is the estimated slope? What is the estimated intercept? Use the regression to predict the growth rate for a country with a trade share of 0.5 and for another with a trade share equal to 1.0.

Solution

```
# regression
E41c <- lm(formula = growth ~ tradeshare, data = Growth)
```

```
# estimated intercept, estimated slope
summary(E41c)
```

```
##
## Call:
## lm(formula = growth ~ tradeshare, data = Growth)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
```

```
## -4.3739 -0.8864 0.2329 0.9248 5.3889
##
## Coefficients:
##             Estimate Std. Error t value Pr(>|t|)
## (Intercept)  0.6403     0.4900   1.307  0.19606
## tradeshare   2.3064     0.7735   2.982  0.00407 **
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 1.79 on 63 degrees of freedom
## Multiple R-squared:  0.1237, Adjusted R-squared:  0.1098
## F-statistic: 8.892 on 1 and 63 DF,  p-value: 0.00407
```

```
# predict value
E41c_predict <- function(x){
  E41c$coefficients %*% matrix(c(1, x), ncol = 1)
}
```

```
# predict value: tradeshare = 0.5
E41c_predict(0.5)
```

```
##           [,1]
## [1,] 1.793482
```

```
# predict value: tradeshare = 1.0
E41c_predict(1)
```

```
##           [,1]
## [1,] 2.946699
```

- d. Estimate the same regression, excluding the data from Malta. Answer the same questions in (c).

Solution

```
# excluding the data from Malta
Growth_n <- Growth[Growth$country_name != "Malta",]

# regression
E41d <- lm(formula = growth ~ tradeshare, data = Growth_n)

# estimated intercept, estimated slope
summary(E41d)
```

```
##
## Call:
## lm(formula = growth ~ tradeshare, data = Growth_n)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -4.4247 -0.9383  0.2091  0.9265  5.3776
##
## Coefficients:
##             Estimate Std. Error t value Pr(>|t|)
## (Intercept)  0.9574     0.5804   1.650  0.1041
## tradeshare   1.6809     0.9874   1.702  0.0937 .
```

```
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 1.789 on 62 degrees of freedom
## Multiple R-squared:  0.04466,    Adjusted R-squared:  0.02925
## F-statistic: 2.898 on 1 and 62 DF,  p-value: 0.09369

# predict value
E41d_predict <- function(x){
  E41d$coefficients %*% matrix(c(1, x), ncol = 1)
}

# predict value: tradeshare = 0.5
E41d_predict(0.5)

##           [,1]
## [1,] 1.797863

# predict value: tradeshare = 1.0
E41d_predict(1)

##           [,1]
## [1,] 2.638315
```

- e. Plot the estimated regression functions from (c) and (d). Using the scatterplot in (a), explain why the regression function that includes Malta is steeper than the regression function that excludes Malta.

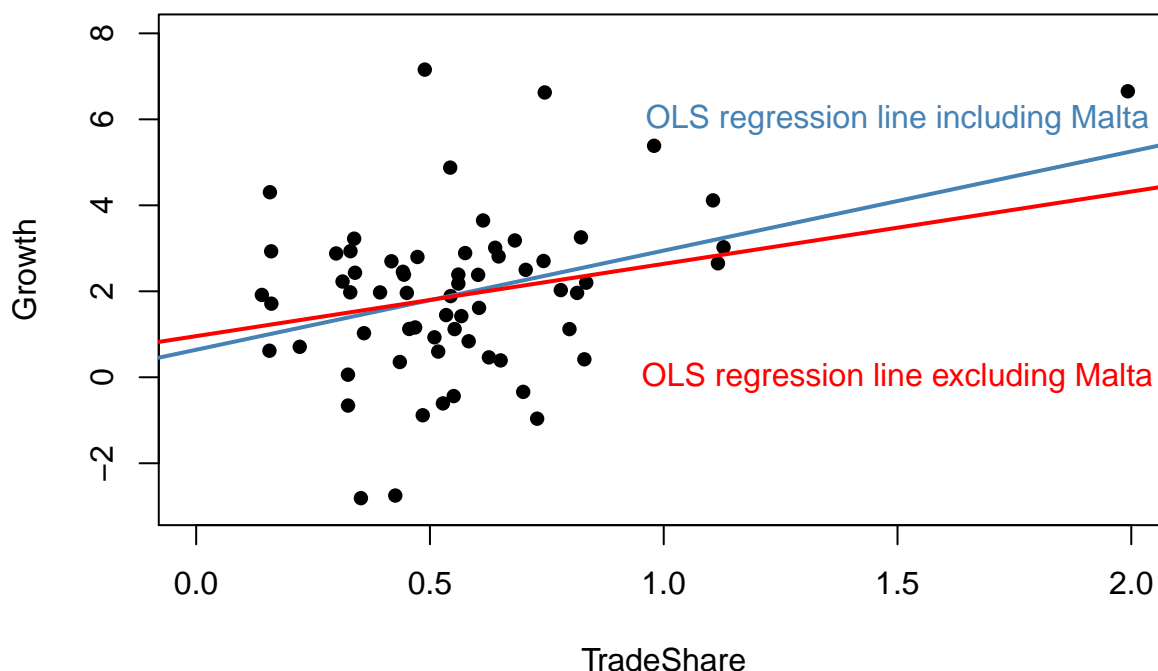
Solution

```
plot(x = Growth$tradeshare,
     y = Growth$growth,
     pch = 16, # filled circle
     col = "black",
     xlim = c(0, 2),
     ylim = c(-3, 8),
     xlab = "TradeShare",
     ylab = "Growth",
     main = "E4.1 (e)")

# with Malta
abline(E41c, lwd = 2, col = "steelblue")
text(1.5, 6, "OLS regression line including Malta", col = "steelblue")

# without Malta
abline(E41d, lwd = 2, col = "red")
text(1.5, 0, "OLS regression line excluding Malta", col = "red")
```

E4.1 (e)



- f. Where is Malta? Why is the Malta trade share so large? Should Malta be included or excluded from the analysis?

Solution

Malta is an island nation in the Mediterranean Sea, south of Sicily.

Malta is a freight transport site, which explains its large “trade share.” Many goods coming into Malta (imports into Malta) and are immediately transported to other countries (as exports from Malta).

Thus, Malta’s imports and exports are unlike the imports and exports of most other countries. Malta should not be included in the analysis.

3 Empirical Exercise 4.2

On the text website, <http://www.pearsonglobaleditions.com>, you will find the data file **Earnings_and_Height**, which contains data on earnings, height, and other characteristics of a random sample of U.S. workers. A detailed description is given in **Earnings_and_Height_Description**, also available on the website. In this exercise, you will investigate the relationship between earnings and height.

- a. What is the median value of height in the sample?

Solution

```
# import data
#install.packages("readxl")
library(readxl)
Earnings_and_Height <- read_xlsx("Earnings_and_Height/Earnings_and_Height.xlsx")

median(Earnings_and_Height$height)
```

```
## [1] 67
```

- b.
 - i. Estimate average earnings for workers whose height is at most 67 inches.
 - ii. Estimate average earnings for workers whose height is greater than 67 inches.
 - iii. On average, do taller workers earn more than shorter workers? How much more? What is a 95% confidence interval for the difference in average earnings?

Solution

```
# create "group" variable
group <- c()

for (i in 1:length(Earnings_and_Height$height)){
  if (Earnings_and_Height$height[i] <= 67){
    # group 0: height <= 67
    group[i] <- c(0)
  } else {
    # group 1: height > 67
    group[i] <- c(1)
  }
}

Earnings_and_Height <- cbind(Earnings_and_Height, group)

E42b <- function(x){
  # sample mean
  mu <- mean(x)

  # sample standard deviation (standard error)
  se <- sd(x)/sqrt(length(x))

  # test
  test <- t.test(x,
                 alternative = c("two.sided"),
                 mu = 0, # H0
                 conf.level = 0.95) # alpha = 0.05

  # 95% confidence interval
  lower <- round(test$conf.int[1], digit = 4)
  upper <- round(test$conf.int[2], digit = 4)
  CI <- paste(lower, "-", upper)

  Table <- data.frame(mu, se, CI)
  colnames(Table) <- c("Mean", "Standard Error", "95% Confidence Interval")

  Table
}

# i. # ii.
tapply(Earnings_and_Height$earnings, Earnings_and_Height$group, E42b)

## $`0`
##      Mean Standard Error 95% Confidence Interval
## 1 44488.44      265.4948 43968.0133 - 45008.8585
##
## $`1`
```

```
##           Mean Standard Error 95% Confidence Interval
## 1 49987.88           305.4062 49389.1973 - 50586.5544

# height <= 67
Earnings_and_Height_i <- Earnings_and_Height[Earnings_and_Height$height <= 67, ]

# height > 67
Earnings_and_Height_ii <- Earnings_and_Height[Earnings_and_Height$height > 67, ]

# iii. 95% CI for difference
t.test(Earnings_and_Height_ii$earnings, Earnings_and_Height_i$earnings,
       alternative = c("two.sided"),
       mu = 0, # H0
       var.equal = FALSE,
       conf.level = 0.95) # alpha = 0.05

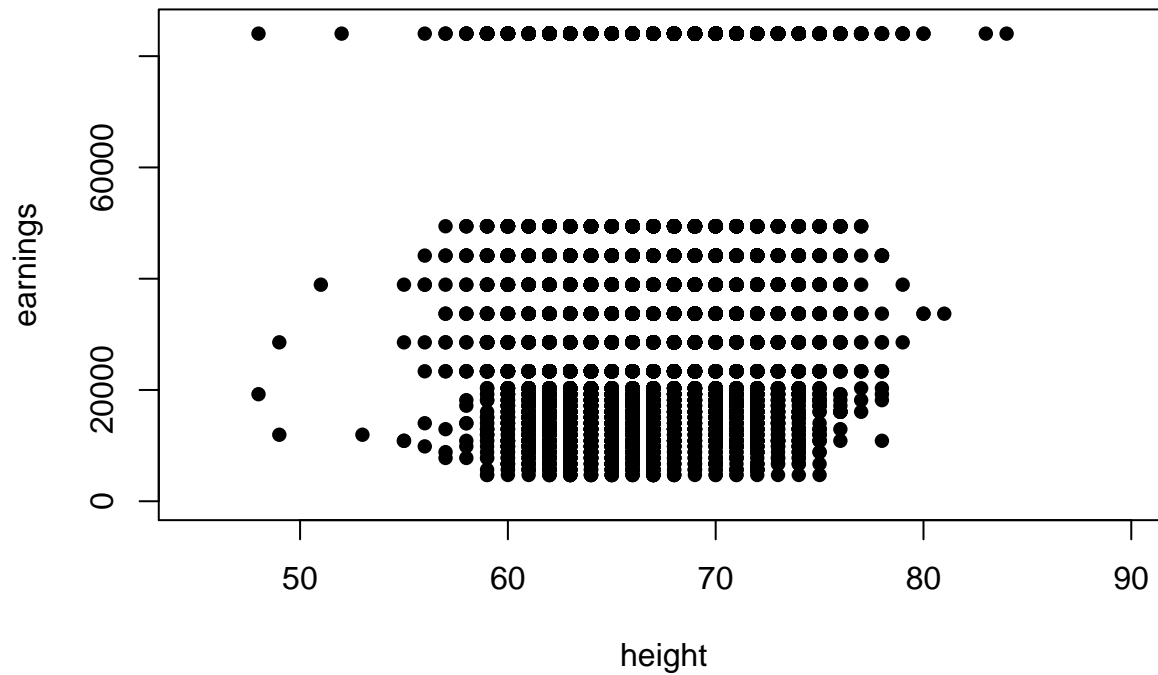
##
## Welch Two Sample t-test
##
## data: Earnings_and_Height_ii$earnings and Earnings_and_Height_i$earnings
## t = 13.59, df = 16624, p-value < 2.2e-16
## alternative hypothesis: true difference in means is not equal to 0
## 95 percent confidence interval:
##  4706.237 6292.643
## sample estimates:
## mean of x mean of y
## 49987.88 44488.44
```

- c. Construct a scatterplot of annual earnings (*Earnings*) on height (*Height*). Notice that the points on the plot fall along horizontal lines. (There are only 23 distinct values of *Earnings*). Why? (Hint: Carefully read the detailed data description.)

Solution

```
plot(x = Earnings_and_Height$height,
     y = Earnings_and_Height$earnings,
     pch = 16, # filled circle
     col = "black",
     xlim = c(45, 90),
     ylim = c(0, 85000),
     xlab = "height",
     ylab = "earnings",
     main = "E4.2 (c)")
```


E4.2 (c)



The data documentation reports that individual earnings were reported in 23 brackets, and a single average value is reported for earnings in the same bracket. Thus, the dataset contains 23 distinct values of earnings.

- d. Run a regression of *Earnings* on *Height*.
 - i. What is the estimated slope?
 - ii. Use the estimated regression to predict earnings for a worker who is 67 inches tall, for a worker who is 70 inches tall, and for a worker who is 65 inches tall.

Solution

```
# regression
E42d <- lm(formula = earnings ~ height, data = Earnings_and_Height)

# i. estimated intercept, estimated slope
summary(E42d)
```

```
##
## Call:
## lm(formula = earnings ~ height, data = Earnings_and_Height)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -47836 -21879  -7976   34323  50599
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept)  -512.73    3386.86  -0.151    0.88
## height         707.67     50.49   14.016 <2e-16 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
```

```
## Residual standard error: 26780 on 17868 degrees of freedom
## Multiple R-squared:  0.01088,    Adjusted R-squared:  0.01082
## F-statistic: 196.5 on 1 and 17868 DF,  p-value: < 2.2e-16
```

```
# predict value
E42d_predict <- function(x){
  E42d$coefficients %*% matrix(c(1, x), ncol = 1)
}
```

```
# ii. predict value: height = 67
E42d_predict(67)
```

```
##           [,1]
## [1,] 46901.26
```

```
# ii. predict value: height = 70
E42d_predict(70)
```

```
##           [,1]
## [1,] 49024.28
```

```
# ii. predict value: height = 65
E42d_predict(65)
```

```
##           [,1]
## [1,] 45485.92
```

- e. Suppose height were measured in centimeters instead of inches. Answer the following questions about the *Earnings* on *Height* (in cm) regression.
- What is the estimated slope of the regression?
 - What is the estimated intercept?
 - What is the R^2 ?
 - What is the standard error of the regression?

Solution

```
# translates from inches to cm
height_cm <- cm(Earnings_and_Height$height)

Earnings_and_Height <- cbind(Earnings_and_Height, height_cm)

# regression
E42e <- lm(formula = earnings ~ height_cm, data = Earnings_and_Height)

# i. estimated slope # ii. estimated intercept
# iii. R^2 # iv. SE
summary(E42e)
```

```
##
## Call:
## lm(formula = earnings ~ height_cm, data = Earnings_and_Height)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -47836 -21879  -7976   34323  50599
##
## Coefficients:
```

```
##           Estimate Std. Error t value Pr(>|t|)
## (Intercept)  -512.73    3386.86  -0.151    0.88
## height_cm     278.61     19.88  14.016 <2e-16 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 26780 on 17868 degrees of freedom
## Multiple R-squared:  0.01088,    Adjusted R-squared:  0.01082
## F-statistic: 196.5 on 1 and 17868 DF,  p-value: < 2.2e-16
```

- f. Run a regression of *Earnings* on *Height*, using data for female workers only.
- What is the estimated slope?
 - A randomly selected woman is 1 inch taller than the average woman in the sample. Would you predict her earnings to be higher or lower than the average earnings for women in the sample? By how much?

Solution

```
# female
Earnings_and_Height_f <- Earnings_and_Height[Earnings_and_Height$sex == 0, ]

# regression
E42f <- lm(formula = earnings ~ height, data = Earnings_and_Height_f)

# i. estimated slope # ii.
summary(E42f)
```

```
##
## Call:
## lm(formula = earnings ~ height, data = Earnings_and_Height_f)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -42748 -22006  -7466   36641  46865
##
## Coefficients:
##           Estimate Std. Error t value Pr(>|t|)
## (Intercept)  12650.9    6383.7   1.982  0.0475 *
## height        511.2      98.9   5.169  2.4e-07 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 26800 on 9972 degrees of freedom
## Multiple R-squared:  0.002672,    Adjusted R-squared:  0.002572
## F-statistic: 26.72 on 1 and 9972 DF,  p-value: 2.396e-07
```

A women who is one inch taller than average is predicted to have earnings that are \$511.2 per year higher than average.

- g. Repeat (f) for male workers.

Solution

```
# male
Earnings_and_Height_g <- Earnings_and_Height[Earnings_and_Height$sex == 1, ]

# regression
E42g <- lm(formula = earnings ~ height, data = Earnings_and_Height_g)
```

```
# i. estimated slope # ii.
summary(E42g)
```

```
##
## Call:
## lm(formula = earnings ~ height, data = Earnings_and_Height_g)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -50158 -22373  -8118   33091   59228
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept) -43130.3      7068.5  -6.102  1.1e-09 ***
## height       1306.9        100.8  12.969 < 2e-16 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 26670 on 7894 degrees of freedom
## Multiple R-squared:  0.02086,    Adjusted R-squared:  0.02074
## F-statistic: 168.2 on 1 and 7894 DF,  p-value: < 2.2e-16
```

A man who is one inch taller than average is predicted to have earnings that are \$1306.9 per year higher than average.

- h. Do you think that height is uncorrelated with other factors that cause earning? That is, do you think that the regression error term, u_i has a conditional mean of 0 given *Height* (X_i)? (You will investigate this more in the Earnings and Height exercises in later chapters.)

Solution

Height may be correlated with other factors that cause earnings. For example, height may be correlated with “strength,” and in some occupations, stronger workers may be more productive. There are many other potential factors that may be correlated with height and cause earnings and we will investigate of these in future exercises.