

STAT 151-A HW 1

Cao jilin

$$1. a) \hat{\beta}_1 = \frac{\sum_{i=1}^n (X_i - \bar{X}) y_i}{\sum_{i=1}^n (X_i - \bar{X})^2}$$

$$\hat{\beta}_0 = \bar{y} - \hat{\beta}_1 \bar{x}$$

$$b) \hat{\alpha}_1 = \frac{\sum_{i=1}^n (y_i - \bar{y}) x_i}{\sum_{i=1}^n (y_i - \bar{y})^2}$$

$$\hat{\alpha}_0 = \bar{x} - \hat{\alpha}_1 \bar{y}$$

$$c) \text{ No } \frac{1}{\hat{\beta}_1} = \frac{\sum_{i=1}^n (X_i - \bar{X})^2}{\sum_{i=1}^n (X_i - \bar{X}) y_i} \neq \hat{\alpha}_1$$

d) stored as a picture on the next page

$$2. (a) \hat{y} = -0.31886 \hat{\beta}_1 x + 32.14271 \hat{\beta}_0$$

$$SD(\hat{\beta}_1) = 0.03484$$

$$SD(\hat{\beta}_0) = 0.99758$$

(b) our model shows a negative relationship, so the assumption was wrong. It turns out the lunch program has a negative effect on student performance.

3. a) Yes. More money spent means that a student can get more resources, thus increasing grade. In R we get $\beta_0 = 1.34 \times 10^1$ $\beta_1 = 2.456 \times 10^{-3}$ which also suggests a positive relationship

b) percentage change in math10 $\Delta \text{math10} = \beta_1 \Delta \log(\text{expend})$ We observe that change in $\log(\text{expend})$
 model: $\text{math10} = \beta_0 + \beta_1 \log(\text{expend}) + e$ $\Delta \log(\text{expend}) \approx \frac{1}{100} (\% \Delta \text{expend})$, so if $\% \Delta \text{expend} = 10$
 $\Delta \text{math10} = \frac{1}{10} \beta_1$

c) $\widehat{\text{math10}} = -69.341 + 11.164 \log(\text{expend})$
 $\beta_0 \quad \beta_1$
 $SD(\beta_0) = 26.53$
 $SD(\beta_1) = 3.169$

d) $\Delta \text{expend} = 10\%$ $\Delta \text{math10} = 11.164 \times 10\% = 1.1164$, so $\widehat{\text{math10}}$ increases by 1.1164

e) inspect the data set, we found $\max(\text{math10}) = 66.7$, $\max(\text{expend}) = 7419$
 $66.7 < 100$ and if we plug 7419 into $\widehat{\text{math10}} = -69.34 + 11.164 \times \log(7419) = 30.15 < 100$
 the largest fitted value

so we don't need to worry that.

4. a) (3)

$$\begin{aligned} \hat{\beta}_0 + \hat{\beta}_1 \alpha &= \left(\frac{1}{n} \sum_{i=1}^n y_i - \frac{\sum_{i=1}^n (X_i - \bar{X}) y_i}{\sum_{i=1}^n (X_i - \bar{X})^2} \cdot \frac{1}{n} \sum_{i=1}^n X_i \right) + \frac{\sum_{i=1}^n (X_i - \bar{X}) y_i}{\sum_{i=1}^n (X_i - \bar{X})^2} \alpha \\ &= \frac{\sum_{i=1}^n (X_i - \bar{X}) y_i}{\sum_{i=1}^n (X_i - \bar{X})^2} (\alpha - \bar{X}) + \frac{1}{n} \sum_{i=1}^n y_i \\ &= \sum_{i=1}^n y_i \left\{ \frac{1}{n} + \frac{(\alpha - \bar{X})(X_i - \bar{X})}{\sum_{i=1}^n (X_i - \bar{X})^2} \right\} \end{aligned}$$

b) we know $\text{Var}_X(\hat{\beta}_1) = \frac{\sigma^2}{\sum_{i=1}^n (X_i - \bar{X})^2}$ and $\bar{y} = \hat{\beta}_0 + \hat{\beta}_1 \bar{X} \Rightarrow \hat{\beta}_0 = \bar{y} - \hat{\beta}_1 \bar{X}$

$$\text{Var}_X(\bar{y} - \hat{\beta}_1 \bar{X} + \hat{\beta}_1 \alpha) = \text{Var}_X(\bar{y} + \hat{\beta}_1 (\alpha - \bar{X})) = \text{Var}_X(\bar{y}) + (\alpha - \bar{X})^2 \text{Var}_X(\hat{\beta}_1) + 2(\alpha - \bar{X}) \text{Cov}(\bar{y}, \hat{\beta}_1)$$

we have $\text{Cov}(\hat{\beta}_0, \hat{\beta}_1) = \frac{-\sigma^2 n \bar{X}}{n^2 \text{Var}(X)} = \text{Cov}(\bar{y} - \hat{\beta}_1 \bar{X}, \hat{\beta}_1) = \text{Cov}(\bar{y}, \hat{\beta}_1) - \bar{X} \text{Cov}(\hat{\beta}_1, \hat{\beta}_1) \Rightarrow$

$$\text{Cov}(\bar{y}, \hat{\beta}_1) = 0, \quad \text{So } \text{Var}_X(\hat{\beta}_0 + \hat{\beta}_1 \alpha) = \frac{\sigma^2}{n} + \frac{\sigma^2 (\alpha - \bar{X})^2}{\sum_{i=1}^n (X_i - \bar{X})^2}$$

c) given the formula in part cb), we observe that $\frac{\sigma^2 (\alpha - \bar{X})^2}{\sum_{i=1}^n (X_i - \bar{X})^2} \geq 0$ and when $\alpha = \bar{X}$

it equals to 0, which minimizes the $\text{Var}(\hat{\beta}_0 + \hat{\beta}_1 \alpha | X_1, \dots, X_n)$ to $\frac{\sigma^2}{n}$

6. a) given $E(y|x) = \beta_1 x$ and assume $E(e_i|x_i) = 0$

We can write $y_i = \beta_1 x_i + e_i$

LSE of β_1 is value of b_1 that minimize $Q(b_1) = \sum_{i=1}^n (y_i - b_1 x_i)^2$

$$\frac{\partial Q}{\partial b_1} = -2 \sum_{i=1}^n x_i (y_i - b_1 x_i) = 0 \Rightarrow b_1 = \hat{\beta}_1 = \frac{\sum_{i=1}^n x_i y_i}{\sum_{i=1}^n x_i^2}$$

$$b) E(\hat{\beta}_1 | X) = E\left(\frac{\sum_{i=1}^n x_i y_i}{\sum_{i=1}^n x_i^2} \mid X\right) = \frac{\sum_{i=1}^n x_i E(y_i | X)}{\sum_{i=1}^n x_i^2} = \frac{\sum_{i=1}^n x_i (\beta_1 x_i)}{\sum_{i=1}^n x_i^2} = \beta_1$$

and by law of iterated expectations $E(\hat{\beta}_1) = E(E(\hat{\beta}_1 | X)) = E(\beta_1) = \beta_1$

Therefore, $\hat{\beta}_1$ is unbiased

$$c) \text{Var}_X(\hat{\beta}_1) = \text{Var}_X\left(\frac{\sum_{i=1}^n x_i y_i}{\sum_{i=1}^n x_i^2}\right) = \frac{\sum_{i=1}^n x_i^2 \text{Var}_X(y_i)}{\sum_{i=1}^n x_i^2} = \frac{\sigma^2 \sum_{i=1}^n x_i^2}{\sum_{i=1}^n x_i^2}$$

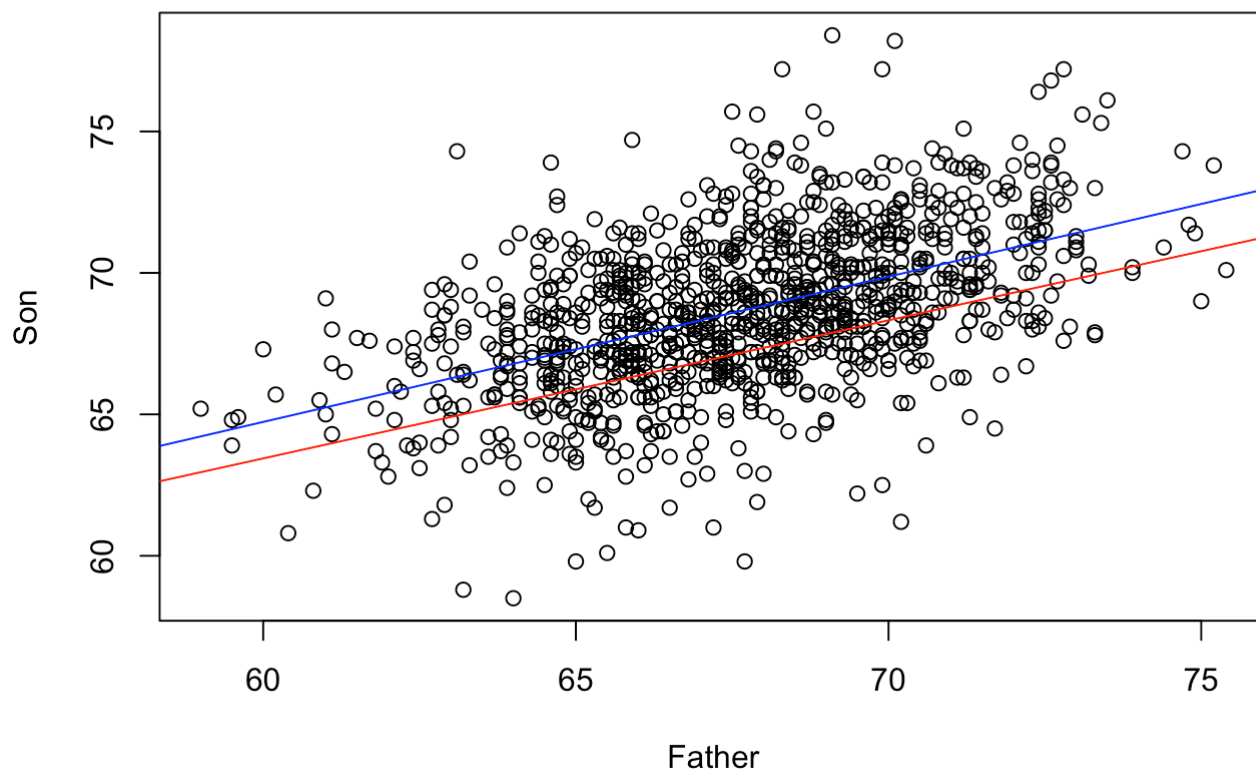
HW01

caojilin

9/5/2018

Problem 1d)

```
dat = read.table("PearsonHeightData.txt", header = T)
lmod1 = lm(Son ~ Father, data = dat)
lmod2 = lm(Father ~ Son, data = dat)
plot(Son ~ Father, data = dat)
abline(lmod1, col = "blue")
abline(lmod2, col = "red")
```



Problem 5

```

library(datasets)
a <- anscombe
par(mfrow=c(2,2))

lm1 = lm(a$y1 ~ a$x1)

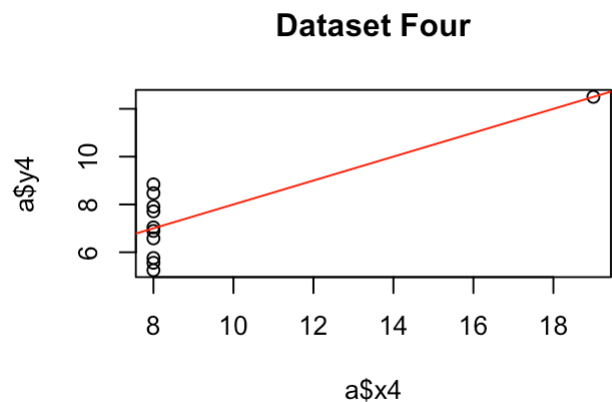
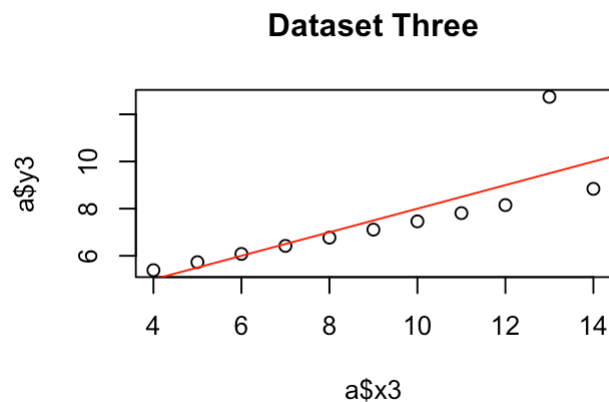
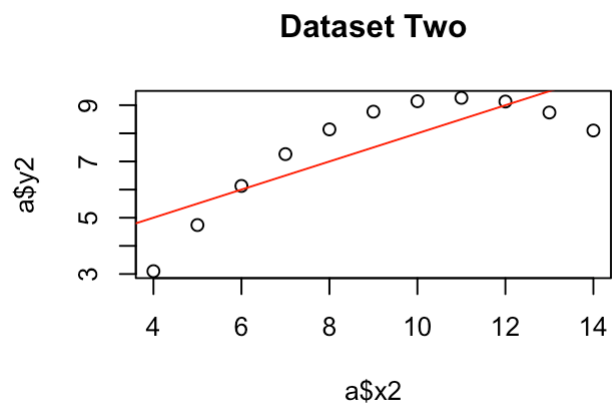
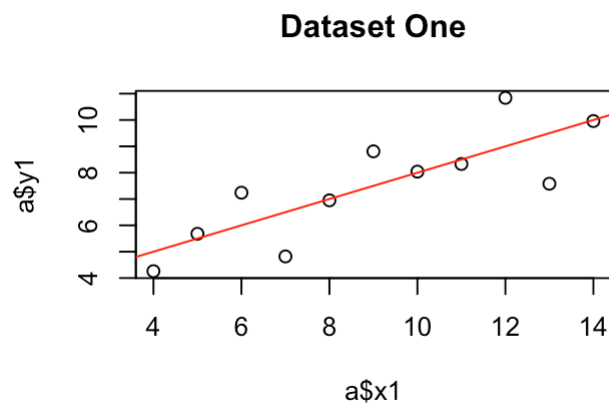
lm2 = lm(a$y2 ~ a$x2)

lm3 = lm(a$y3 ~ a$x3)

lm4 = lm(a$y4 ~ a$x4)

plot(a$x1,a$y1, main=paste("Dataset One"))
abline(lm1,col="red")
plot(a$x2,a$y2, main=paste("Dataset Two"))
abline(lm2,col="red")
plot(a$x3,a$y3, main=paste("Dataset Three"))
abline(lm3,col="red")
plot(a$x4,a$y4, main=paste("Dataset Four"))
abline(lm4,col="red")

```



dataset 1, and 3 make sense. Although dataset 3 has a outlier, it doesn't matter.

dataset 2 looks like a parabola, linear model is not accurate as dataset gets larger, but it's ok for small dataset like this.

dataset 4 doesn't make sense because most of the data are clustered on the line $x = 8$

predictions are:

```
lm1$coefficients[1] + 10 * lm1$coefficients[2]
```

```
## (Intercept)
##          8.001
```

```
lm2$coefficients[1] + 10 * lm2$coefficients[2]
```

```
## (Intercept)
##          8.000909
```

```
lm3$coefficients[1] + 10 * lm3$coefficients[2]
```

```
## (Intercept)
##          7.999727
```

```
lm4$coefficients[1] + 10 * lm4$coefficients[2]
```

```
## (Intercept)
##          8.000818
```

again, predictions for dataset 1,2,3 are close to real data values. So they make sense. but dataset 4 doesn't make sense.

Problem 7

a)

we assume simple linear regression model for

$$(x_1, y_1), (x_2, y_2) \dots (x_n, y_n)$$

to be

$$y_i = \beta_0 + \beta_1 x_i + e_i$$

where

$$\mathbb{E}(e_i | X_i) = 0$$

That means we have to make sure

$$\mathbb{E}(e_i | X_i) = 0$$

in order to assume a linear model.

To verify that: for

$$x_i \leq 65$$

$$y_i$$

can be written as

$$y_i = N(\beta_0 + \beta_1 x_i, 25) + N(0, 25)$$

where

$$e_i = N(0, 25)$$

this satisfies that

$$\mathbb{E}(e_i | X_i) = 0$$

similarly, for

$$65 < x_i \leq 70$$

,

$$y_i = \beta_0 + \beta_1 x_i + e_i$$

where

$$e_i = 10T_i$$

whose mean is 0

for

$$x_i > 70$$

,

$$y_i = \beta_0 + \beta_1 x_i + e_i$$

where

$$e_i = U_i$$

, and

$$\mathbb{E}(U_i) = 0$$

Therefore, the condition

$$\mathbb{E}(e_i | X_i) = 0$$

for simple linear regression model is satisfied. Since we have proved in the class and in the homework that least square estimators are unbiased. we can conclude that

$$\hat{\beta}_0$$

and

$$\hat{\beta}_1$$

are unbiased.

b)

```
M = 1:10000

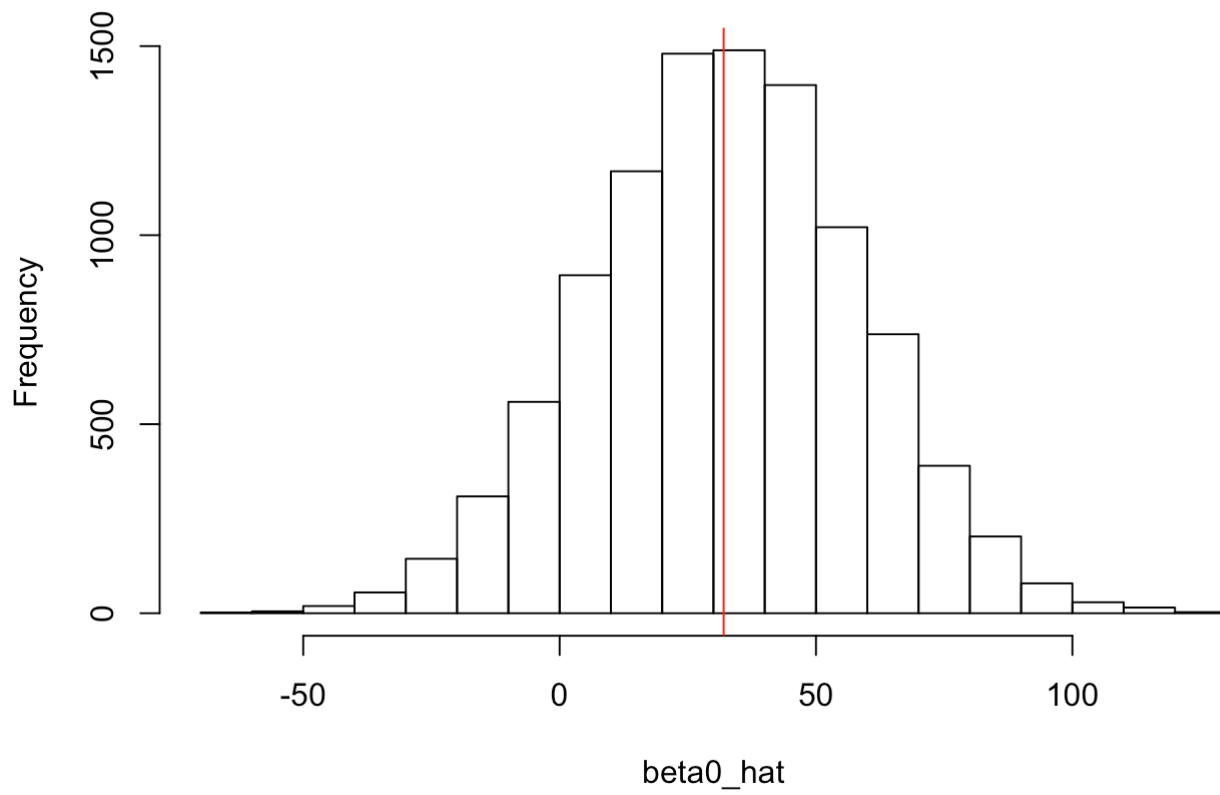
samp = function(){
  x = seq(59, 76, length.out = 100)
  x1=x[x<=65]
  y1 = rnorm(length(x1),mean = 32+0.5*x1,sd = 25)
  x2=x[x>65 & x<=70]
  y2 = 32+0.5*x2+10*rt(n = length(x2),df = 3)
  x3=x[x>70]
  y3=32+0.5*x3+runif(length(x3),min = -8,max = 8)
  y = c(y1,y2,y3)
  return(list(x,y))
}

beta0_hat = c()
beta1_hat = c()
sd_beta0 = c()
sd_beta1 = c()

for (i in M) {
  sample = samp()
  x = sample[[1]]
  y = sample[[2]]
  lmod = lm(y ~ x)
  beta0_hat = c(beta0_hat, lmod$coefficients[1])
  beta1_hat = c(beta1_hat, lmod$coefficients[2])
  sd_beta0 = c(sd_beta0,summary(lmod)$coefficients[3])
  sd_beta1 = c(sd_beta1,summary(lmod)$coefficients[4])
}

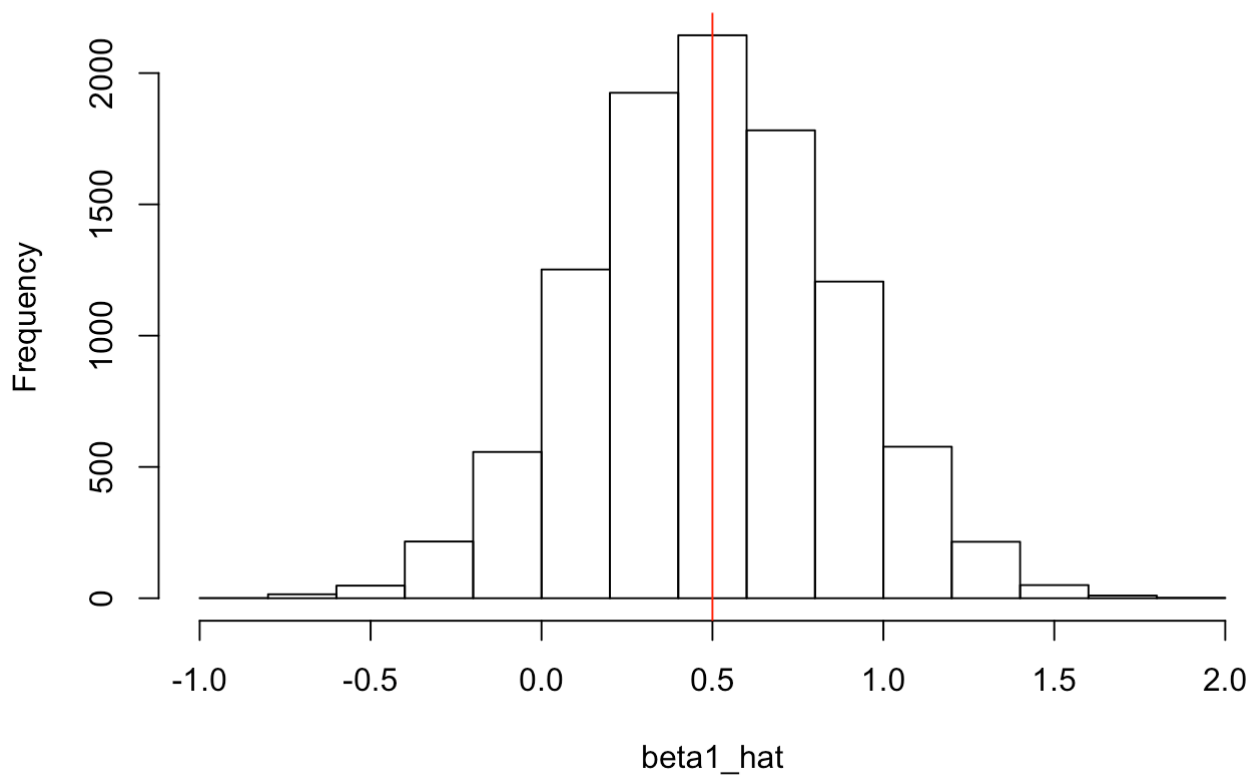
hist(beta0_hat)
abline(v=32,col="red")
```

Histogram of beta0_hat



```
hist(beta1_hat)
abline(v=0.5,col="red")
```

Histogram of beta1_hat



```
#bias of beta 0
sum(beta0_hat -32)/10000
```

```
## [1] 0.2965467
```

```
#bias of beta 1
sum(beta1_hat -0.5)/10000
```

```
## [1] -0.004216398
```

From histogram, the estimates of the bias are close to 0 and we found that the center for both graph are close to real

$$\beta_0 = 32$$

and

$$\beta_1 = 0.5$$

This verifies unbiasedness

c)

homoskedasticity means that same variance of the errors,

$$\text{Var}(e_i | X) = \sigma^2$$

for each i . However, in this problem, we see that

$$\text{Var}(e_i | X) = \sigma^2$$

are not same for each i . Thus, homoskedasticity is not valid here.

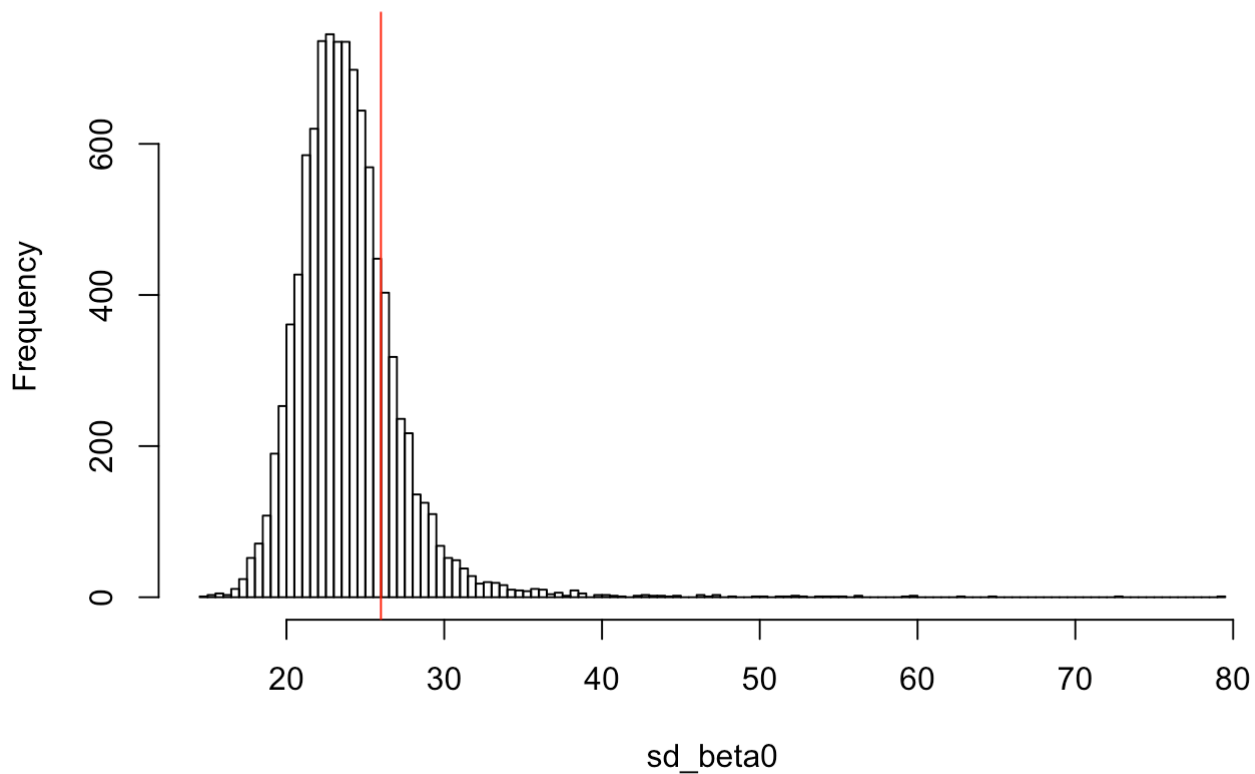
d)

we would like to check unbiasedness of

$$\mathbb{E}(\text{Var}(\hat{\beta})) = \text{Var}(\hat{\beta})$$

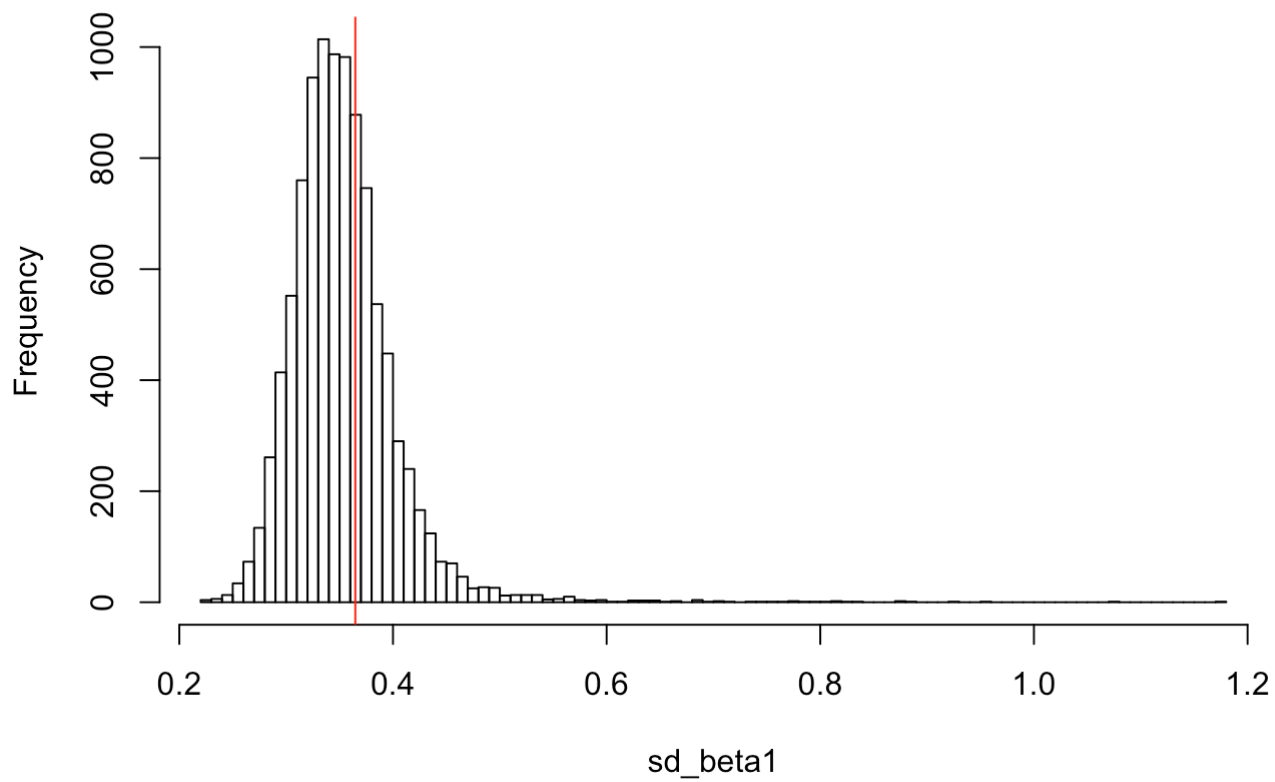
```
sd0 = sd(beta0_hat)
sd1 = sd(beta1_hat)
hist(sd_beta0, breaks = 100)
abline(v=sd0, col="red")
```

Histogram of sd_beta0



```
hist(sd_beta1, breaks=100)
abline(v=sd1, col="red")
```

Histogram of sd_beta1



As we see from histograms, the centers are not close to sd0 and sd1, meaning that

$$\mathbb{E}(\hat{Var}(\hat{\beta})) \neq Var(\hat{\beta})$$

this verifies that homoskedasticity is not valid