

<u>Course</u> > <u>Section 2: Linear M...</u> > <u>2.2: Least Squares ...</u> > Assessment: Least ...

# Assessment: Least Squares Estimates, part 1

# Question 1

1/1 point (graded)

The following code was used in the video to plot RSS with  $\beta_0 = 25$ .

In a model for sons' heights vs fathers' heights, what is the least squares estimate (LSE) for  $\beta_1$  if we assume  $\hat{\beta}_0$  is 36?

Hint: modify the code above to do your analysis.

○ 0.65○ 0.5

O.2

O 12



Correct: Correct. You can tell from a plot of RSS vs  $eta_1$  that the minimum estimate is 0.5

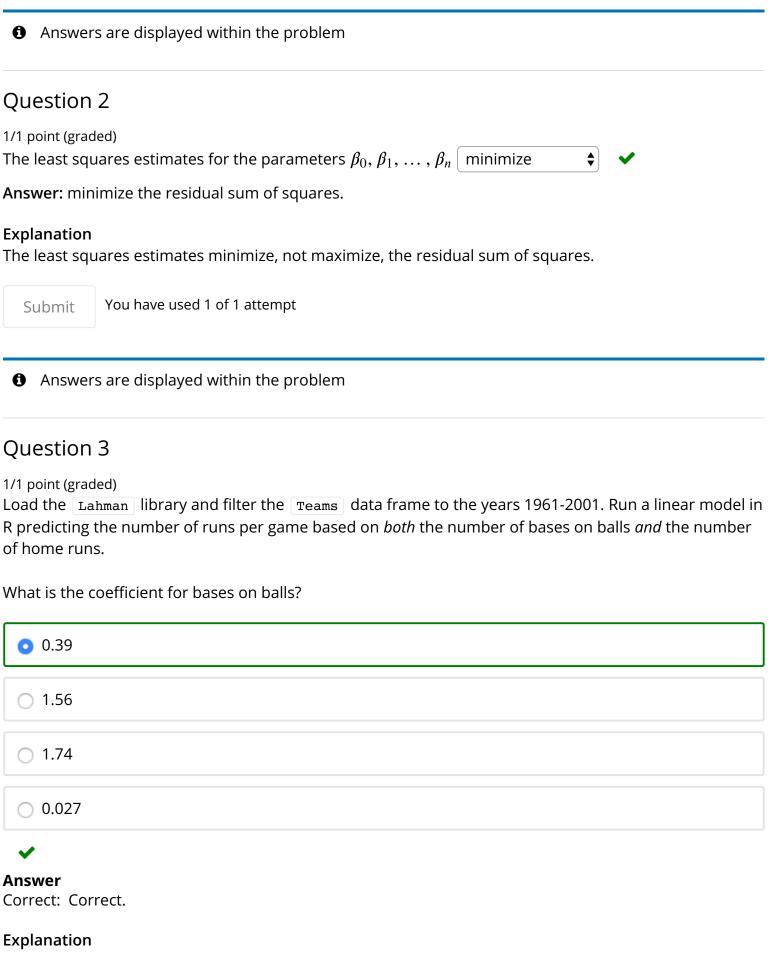
### **Explanation**

Using the code from the video, you can plot RSS vs  $\beta_1$  to find the value for  $\beta_1$  that minimizes the RSS. In this case, that value is 0.5 when we assume that  $\hat{\beta}_0$  is 36.

When we assumed that  $\hat{\beta}_0$  was 25, as in the sample code, the LSE for  $\beta_1$  was 0.65.

Submit

You have used 1 of 2 attempts



The coefficient for bases on balls is 0.39; the coefficient for home runs is 1.56; the intercept is 1.74; the standard error for the BB coefficient is 0.027.

Submit

You have used 1 of 2 attempts

**1** Answers are displayed within the problem

## Question 4

1/1 point (graded)

We run a Monte Carlo simulation where we repeatedly take samples of N = 100 from the Galton heights data and compute the regression slope coefficients for each sample:

```
B <- 1000
N <- 100
lse <- replicate(B, {
    sample_n(galton_heights, N, replace = TRUE) %>%
    lm(son ~ father, data = .) %>% .$coef
})
lse <- data.frame(beta_0 = lse[1,], beta_1 = lse[2,])</pre>
```

What does the central limit theorem tell us about the variables beta\_0 and beta\_1? Select ALL that apply.

- They are approximately normally distributed.
- ightharpoonup The expected value of each is the true value of  $eta_0$  and  $eta_1$  (assuming the Galton heights data is a complete population).
- The central limit theorem does not apply in this situation.



#### **Answer**

Correct:

Correct. With a large enough N, the distributions of both beta\_0 and beta\_1 are approximately normal.

## **Explanation**

With a large enough N, the central limit theorem applies and tells us that the distributions of both beta\_0 and beta\_1 are approximately normal. The expected values of beta\_0 and beta\_1 are the true values of  $\beta_0$  and  $\beta_1$ , assuming that the Galton heights data are a complete population.

For hypothesis testing, we assume that the errors in the model are normally distributed.

Submit

You have used 1 of 2 attempts

• Answers are displayed within the problem

## Question 5

1/1 point (graded)

In an earlier video, we ran the following linear model and looked at a summary of the results.

```
> mod <- lm(son ~ father, data = galton_heights)</pre>
> summary(mod)
Call:
lm(formula = son ~ father, data = galton heights)
Residuals:
  Min 1Q Median 3Q
                            Max
-5.902 -1.405 0.092 1.342 8.092
Coefficients:
              Estimate Std. Error t value
                                            Pr(>|t|)
(Intercept) 35.7125 4.5174 7.91 2.8e-13 ***
                        0.0653
father
              0.5028
                                   7.70 9.5e-13 ***
Signif. codes: 0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
```

What null hypothesis is the second p-value (the one in the father row) testing?

- $\bigcirc$   $eta_1=1$ , where  $eta_1$  is the coefficient for the variable "father."
- $\bigcirc$   $eta_1=0.503$ , where  $eta_1$  is the coefficient for the variable "father."
- $oldsymbol{\circ}$   $eta_1=0$ , where  $eta_1$  is the coefficient for the variable "father."



#### **Explanation**

The p-value for "father" tests the null hypothesis that  $\beta_1=0$ , i.e., the fathers' heights are not associated with the sons' heights, where  $\beta_1$  is the coefficient for the variable father.

Submit

You have used 1 of 1 attempt

**1** Answers are displayed within the problem

1/1 point (graded)

Which R code(s) below would properly plot the predictions and confidence intervals for our linear model of sons' heights?

Select ALL that apply.

```
galton_heights %>% ggplot(aes(father, son)) +
    geom_point() +
    geom_smooth()
```

```
galton_heights %>% ggplot(aes(father, son)) +
    geom_point() +
    geom_smooth(method = "lm")
```

```
model <- lm(son ~ father, data = galton_heights)
predictions <- predict(model, interval = c("confidence"), level = 0.95)
data <- as.tibble(predictions) %>% bind_cols(father = galton_heights$father)

ggplot(data, aes(x = father, y = fit)) +
    geom_line(color = "blue", size = 1) +
    geom_ribbon(aes(ymin=lwr, ymax=upr), alpha=0.2) +
    geom_point(data = galton_heights, aes(x = father, y = son))
```

```
model <- lm(son ~ father, data = galton_heights)
predictions <- predict(model)
data <- as.tibble(predictions) %>% bind_cols(father = galton_heights$father)

ggplot(data, aes(x = father, y = fit)) +
    geom_line(color = "blue", size = 1) +
    geom_point(data = galton_heights, aes(x = father, y = son))
```



#### **Answer**

#### Correct:

Correct. This is one way to plot predictions and confidence intervals for a linear model of sons' heights vs. fathers' heights. This is one of two correct answers.

Correct. This code uses the <code>predict</code> command to generate predictions and 95% confidence intervals for the linear model of sons' heights vs. fathers' heights. This is one of two correct answers.

#### **Explanation**

If using the <code>geom\_smooth</code> command, you need to specify that <code>method = "lm"</code> in your <code>geom\_smooth</code> command, otherwise the smooth line is a loess smooth and not a linear model.

If using the <code>predict</code> command, you need to include the confidence intervals on your figure by first specifying that you want confidence intervals in the <code>predict</code> command, and then adding them to your figure as a <code>geom\_ribbon</code>.

Submit

You have used 2 of 2 attempts

**1** Answers are displayed within the problem

© All Rights Reserved