

# Problem Set 2

QTM 200: Applied Regression Analysis

Due: February 10, 2020

## Instructions

- Please show your work! You may lose points by simply writing in the answer. If the problem requires you to execute commands in `R`, please include the code you used to get your answers. Please also include the `.R` file that contains your code. If you are not sure if work needs to be shown for a particular problem, please ask.
- Your homework should be submitted electronically on the course GitHub page in `.pdf` form.
- This problem set is due at the beginning of class on Monday, February 10, 2020. No late assignments will be accepted.
- Total available points for this homework is 100.

## Question 1 (40 points): Political Science

The following table was created using the data from a study run in a major Latin American city.<sup>1</sup> As part of the experimental treatment in the study, one employee of the research team was chosen to make illegal left turns across traffic to draw the attention of the police officers on shift. Two employee drivers were upper class, two were lower class drivers, and the identity of the driver was randomly assigned per encounter. The researchers were interested in whether officers were more or less likely to solicit a bribe from drivers depending on their class (officers use phrases like, “We can solve this the easy way” to draw a bribe). The table below shows the resulting data.

---

<sup>1</sup>Fried, Lagunes, and Venkataramani (2010). “Corruption and Inequality at the Crossroad: A Multi-method Study of Bribery and Discrimination in Latin America. *Latin American Research Review*. 45 (1): 76-97.

	Not Stopped	Bribe requested	Stopped/given warning
Upper class	14	6	7
Lower class	7	7	1

(a) Calculate the  $\chi^2$  test statistic by hand (even better if you can do "by hand" in R).

```

1 X = matrix(c(14, 7, 6, 7, 7, 1), nrow = 2, ncol = 3) #Creating matrix
   from problem
2
3 colnames(X) <- c("Not Stopped", "Bribe requested", "Stopped/given warning")
4 rownames(X) <- c("Upper Class", "Lower Class")
5 X #Creation of Matrix
6
7 no_stop_sum <- sum(X[,1])
8 no_stop_sum #21
9
10 bribe_sum <- sum(X[,2])
11 bribe_sum #13
12
13 stopped_sum <- sum(X[,3])
14 stopped_sum #8
15
16 #getting total for each column
17
18 upper_sum <- sum(X[1,])
19 upper_sum #27
20
21 lower_sum <- sum(X[2,])
22 lower_sum #15
23
24 #Getting total for each row

```

```

25 sum(upper_sum, lower_sum) #Matrix total is 42
26
27 #now remaking matrix, to have column and row totals for observed values
28 new_x = matrix(c(14, 7, 21, 6, 7, 13, 7, 1, 8, 27, 15, 42), nrow = 3,
29               ncol = 4)
30 colnames(new_x) <- c("Not Stopped", "Bribe requested", "Stopped/given
31                   warning", "Column Total")
32 rownames(new_x) <- c("Upper Class", "Lower Class", "Row Total")
33 new_x #Creation of Matrix
34
35 #to calculate expected values for chi squared test = (row total/grand sum
36   ) * column total
37 (matrix(new_x[3,1]) / matrix(new_x[3,4])) * matrix(new_x[1,4]) #13.5
38 (matrix(new_x[3,1]) / matrix(new_x[3,4])) * matrix(new_x[2,4]) #7.5
39 (matrix(new_x[3,2]) / matrix(new_x[3,4])) * matrix(new_x[1,4]) #8.36
40 (matrix(new_x[3,2]) / matrix(new_x[3,4])) * matrix(new_x[2,4]) #4.64
41 (matrix(new_x[3,3]) / matrix(new_x[3,4])) * matrix(new_x[1,4]) #5.14
42 (matrix(new_x[3,3]) / matrix(new_x[3,4])) * matrix(new_x[2,4]) #2.86
43
44 #creating matrix of expected values
45 expected_values = matrix(c(13.5, 7.5, 8.36, 4.64, 5.14, 2.86), nrow = 2,
46                          ncol = 3)
47 expected_values
48
49 #calculating chi squared directly
50 sum(((X - expected_values)^2) / expected_values) #Chi squared values is
51 3.801141
52
53 #Using chi squared function in R to test if values are close, which they
54   are with the function outputting 3.7912
55 chisq.test(X)

```

(b) Now calculate the p-value (in R).<sup>2</sup> What do you conclude if  $\alpha = .1$ ?

```

1 #calculating p-value with chi squared value = 3.801141, df= (2-1)(3-1),
2   and lower tail = false, as we only want the upper end due to
3   distribution type
4 pchisq(3.801141, 2, lower.tail = F)
5 #p-value = 0.1495 meaning we fail to reject the null, that officers would
6   not solicit a bribe depending on the class

```

---

<sup>2</sup>Remember frequency should be  $> 5$  for all cells, but let's calculate the p-value here anyway.

(c) Calculate the standardized residuals for each cell and put them in the table below.

```

1 #calculating standardized residuals and want to see the matrixes again
2 new_x
3 expected_values
4
5 A_one = ((X[1,1] - expected_values[1,1]) / (sqrt(expected_values[1,1]) *
        (1-(27/42)) * (1-(21/42))))
6 A_one
7 A_two = ((X[2,1] - expected_values[2,1]) / (sqrt(expected_values[2,1]) *
        (1-(15/42)) * (1-(21/42))))
8 A_two
9
10 B_one = ((X[1,2] - expected_values[1,2]) / (sqrt(expected_values[1,2]) *
        (1-(27/42)) * (1-(13/42))))
11 B_one
12 B_two = ((X[2,2] - expected_values[2,2]) / (sqrt(expected_values[2,2]) *
        (1-(15/42)) * (1-(13/42))))
13 B_two
14
15 C_one = ((X[1,3] - expected_values[1,3]) / (sqrt(expected_values[1,3]) *
        (1-(27/42)) * (1-(8/42))))
16 C_one
17 C_two = ((X[2,3] - expected_values[2,3]) / (sqrt(expected_values[2,3]) *
        (1-(15/42)) * (1-(8/42))))
18 C_two
19
20 #creating matrix for standardized residuals
21 standardized_matrix = matrix(c(0.762, -0.568, -3.310, 2.468, 2.838,
        -2.113), nrow=2, ncol=3)
22
23 colnames(standardized_matrix) <- c("Not Stopped", "Bribe requested", "
        Stopped/given warning")
24 rownames(standardized_matrix) <- c("Upper Class", "Lower Class")
25 standardized_matrix #Creation of Matrix

```

	Not Stopped	Bribe requested	Stopped/given warning
Upper class			
Lower class			

- (d) How might the standardized residuals help you interpret the results? The standard residual data may help with the interpretation of the results because they show how far off each data point is from a normal projected value. Simply, a the standardized residuals illustrates the difference between the observed and expected results, and show which values are contributing most in either more/less bias or no bias ways.

## Question 2 (20 points): Economics

Chattopadhyay and Duflo were interested in whether women promote different policies than men.<sup>3</sup> Answering this question with observational data is pretty difficult due to potential confounding problems (e.g. the districts that choose female politicians are likely to systematically differ in other aspects too). Hence, they exploit a randomized policy experiment in India, where since the mid-1990s,  $\frac{1}{3}$  of village council heads have been randomly reserved for women. A subset of the data from West Bengal can be found at the following link: <https://raw.githubusercontent.com/kosukeimai/qss/master/PREDICTION/women.csv>

Each observation in the data set represents a village and there are two villages associated with one GP (i.e. a level of government is called "GP"). Figure 1 below shows the names and descriptions of the variables in the dataset. The authors hypothesize that female politicians are more likely to support policies female voters want. Researchers found that more women complain about the quality of drinking water than men. You need to estimate the effect of the reservation policy on the number of new or repaired drinking water facilities in the villages.

Figure 1: Names and description of variables from Chattopadhyay and Duflo (2004).

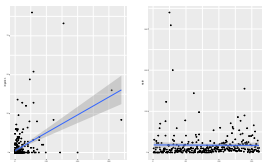
Name	Description
<b>GP</b>	An identifier for the Gram Panchayat (GP)
<b>village</b>	identifier for each village
<b>reserved</b>	binary variable indicating whether the GP was reserved for women leaders or not
<b>female</b>	binary variable indicating whether the GP had a female leader or not
<b>irrigation</b>	variable measuring the number of new or repaired irrigation facilities in the village since the reserve policy started
<b>water</b>	variable measuring the number of new or repaired drinking-water facilities in the village since the reserve policy started

<sup>3</sup>Chattopadhyay and Duflo. (2004). "Women as Policy Makers: Evidence from a Randomized Policy Experiment in India. *Econometrica*. 72 (5), 1409-1443.

- (a) State a null and alternative (two-tailed) hypothesis. Null Hypothesis: There is no effect from the reservation policy on the number of new or repaired drinking water facilities in the villages Alternative Hypothesis: There is an effect from the reservation policy on the number of new or repaired drinking water facilities in the villages

- (b) Run a bivariate regression to test this hypothesis in R (include your code!).

```
1 ggplot(aes(water, irrigation), data = economics_ps) +  
2   geom_point() +  
3   geom_smooth(method = 'lm')  
4  
5 #can also plot against the identifier BP to see in each location how many  
6   water sites there are  
6 ggplot(aes(GP, water), data = economics_ps) +  
7   geom_point() +  
8   geom_smooth(method = 'lm')
```



There seems to be very little change in the water situation as a whole in all GP locations with a near-zero slope, indicating a lack of change per GP. With irrigation however, as water increased, irrigation seemed to also increase. The difficult aspect of this interpretation however, is the lack of data points at the larger end of the scope.

(c) Interpret the coefficient estimate for reservation policy.

```
1 c(round(mean(economics_ps$water), 2), round(sd(economics_ps$water),2)) #x
   bar which is the calculated mean is 17.84, while the standard
   deviation is 33.68.
2
3 standardized_water <- (((economics_ps$water - (mean(economics_ps$water)))
   /sd(economics_ps$water)))
4 round(standardized_water, 2) #calculated standardized difference from the
   observed valued of lifespans
5
6 (1/(322-1)) * (sum(round(standardized_water, 2)))
```

The correlation coefficient is 0.00081. Meaning, the relationship from water as a whole in this data set, only has a relationship that is positive and extremely weak.



### Question 3 (40 points): Biology

There is a physiological cost of reproduction for fruit flies, such that it reduces the lifespan of female fruit flies. Is there a similar cost to male fruit flies? This dataset contains observations from five groups of 25 male fruit flies. The experiment tests if increased reproduction reduces longevity for male fruit flies. The five groups are: males forced to live alone, males assigned to live with one or eight newly pregnant females (non-receptive females), and males assigned to live with one or eight virgin females (interested females). The name of the data set is `fruitfly.csv`.<sup>4</sup>

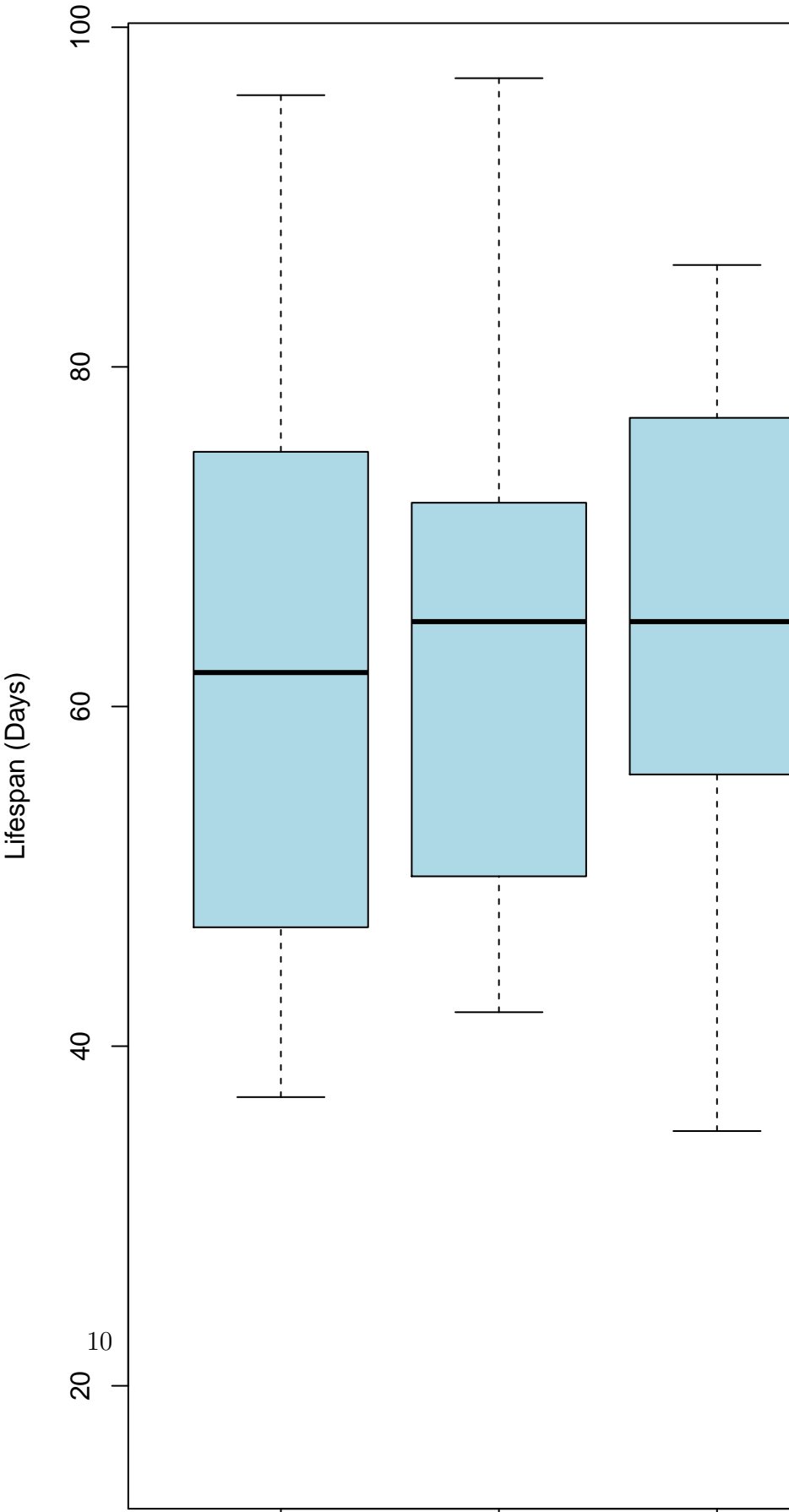
<b>No</b>	serial number (1-25) within each group of 25
<b>type</b>	Type of experimental assignment
	1 = no females
	2 = 1 newly pregnant female
	3 = 8 newly pregnant females
	4 = 1 virgin female
	5 = 8 virgin females
<b>lifespan</b>	lifespan (days)
<b>thorax</b>	length of thorax (mm)
<b>sleep</b>	percentage of each day spent sleeping

1. Import the data set and obtain summary statistics and examine the distribution of the

---

<sup>4</sup>Partridge and Farquhar (1981). "Sexual Activity and the Lifespan of Male Fruitflies". *Nature*. 294, 580-581.

Fruitfly Lifespan

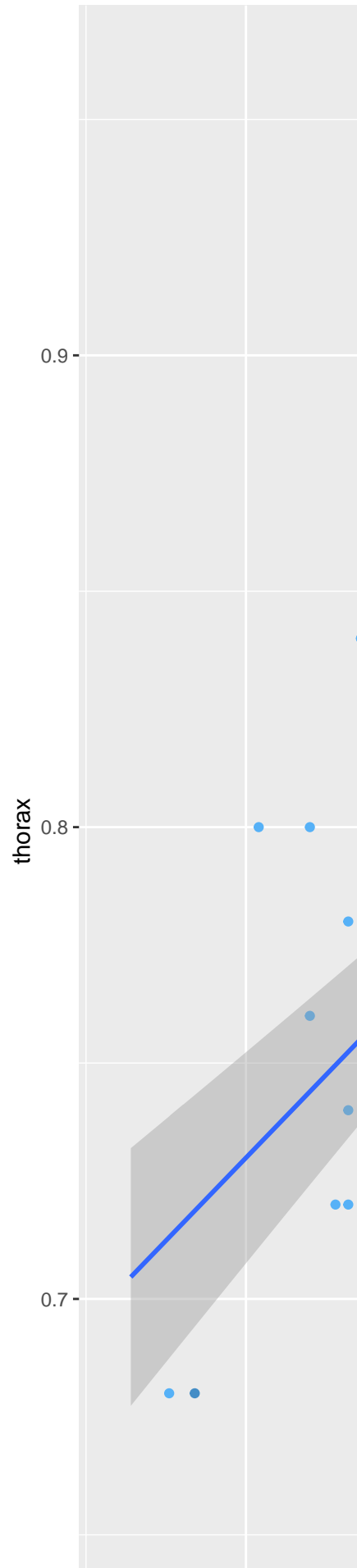


```

1 library(tidyverse)
2 library(dplyr) #inputing library for filter function
3 fruitfly_type1 <- filter(fruitfly , type == 1)
4 fruitfly_type1 #isolating within the fruitfly dataset , each type
5
6 fruitfly_type2 <- filter(fruitfly , type == 2)
7 fruitfly_type2
8
9 fruitfly_type3 <- filter(fruitfly , type == 3)
10 fruitfly_type3
11
12 fruitfly_type4 <- filter(fruitfly , type == 4)
13 fruitfly_type4
14
15 fruitfly_type5 <- filter(fruitfly , type == 5)
16 fruitfly_type5

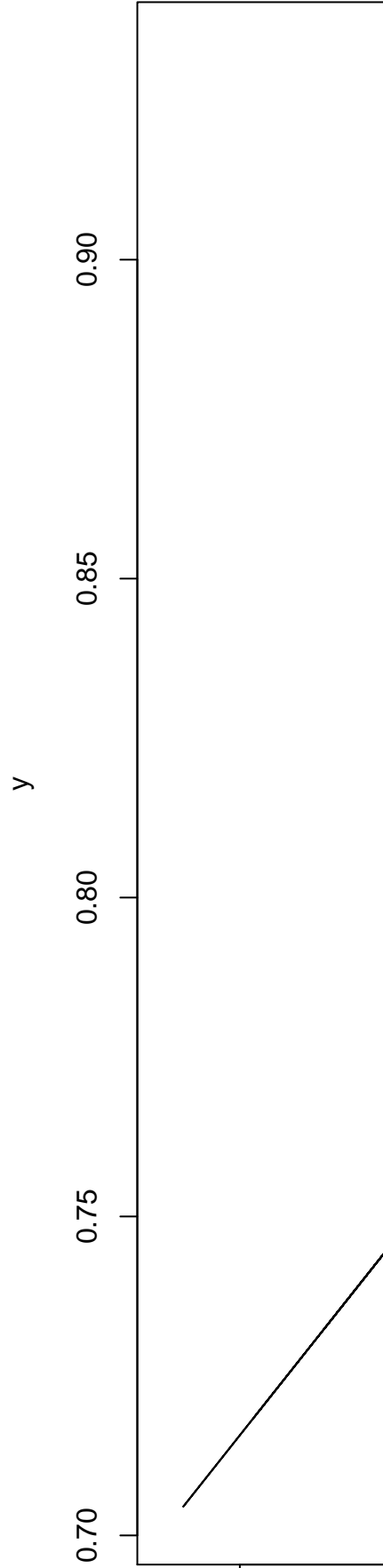
```

2. Plot `lifespan` vs `thorax`. Does it look like there is a linear relationship? Provide the



There does seem to be a positive, linear, relationship between lifespan and thorax length.

```
1 fruitfly
2 c(round(mean(fruitfly$lifespan), 2), round(sd(fruitfly$lifespan), 2)) #x
   bar which is the calculated mean is 57.44, while the standard
   deviation is 17.56.
3
4 standardized_fly <- (((fruitfly$lifespan - (mean(fruitfly$lifespan)))/sd(
   fruitfly$lifespan)))
5 round(standardized_fly, 2) #calculated standardized difference from the
   observed valued of lifespans
6
7 standardized_thorax <- (((fruitfly$thorax - (mean(fruitfly$thorax)))/sd(
   fruitfly$thorax)))
8 round(standardized_thorax, 2) #calculated standardized difference from
   the observed valued of thorax length
9
10 standardized_fly_thorax <- standardized_fly * standardized_thorax #
   multiplying each standardarized values to then calculate correlation
   coefficient
11
12 (1/(125-1))*sum(round(standardized_fly_thorax, 2)) #the correlation
   coefficient is .636 which is a decently strong, and positive
   relationship that is explained between lifespan and thorax length.
```



The y-intercept is .66 mm of the thorax length when lifespan is theoretically 0 while the slope of the line is .002, as for every day gained in lifespan, the thorax length increases by .002 mm.

```
1 summary(lm(fruitfly$thorax ~ fruitfly$lifespan))
2
3 y <- .6597 + 0.0028*(fruitfly$lifespan) #creating the basic slope-
      intercept line with the gained intercept and slope from the summary.
4 plot(fruitfly$lifespan, y, type = "l") #creation of basic line graph to
      illustrate the intercept and slope of the relationship between
      lifespan and thorax length
```

Test for a significant linear relationship between `lifespan` and `thorax`. Provide and interpret your results of your test.

```
1 t.test(fruitfly$lifespan, fruitfly$thorax, alternative = "two.sided", var
      .equal = FALSE)
2 summary(lm(fruitfly$thorax ~ fruitfly$lifespan))
3
4 thorax_lifespan <- lm(fruitfly$thorax ~ fruitfly$lifespan) #creation of
      one variable under the line relationship of thorax and lifespan
5 sigma(thorax_lifespan) #finding sigma of this relationship to be 0.06
6
7 beta_se <- sigma(thorax_lifespan)/sqrt(sum((fruitfly$lifespan - mean(
      fruitfly$lifespan))^2))
8 beta_se #finding the standard error of the lifespan quantity to be .0003
9
10 2*pt((.0028 - 0)/beta_se, dim(fruitfly)[1]-2, lower.tail = F) #p-value for
      beta-1 is 1.69e-15, which is very similar to checking through summary
      function
11 summary(thorax_lifespan)
```



3. Provide the 90% confidence interval for the slope of the fitted model.

- Use the formula for typical confidence intervals to find the 90% confidence interval around the point estimate.
- Now, try using the function `confint()` in R.

```
1 summary(lm(fruitfly$thorax ~ fruitfly$lifespan))
2
3 b1 <- .0028 #identification of the slope beta1
4 se <- .0003 #standard error of the graph between the relationship of
   thorax and lifespan
5 n <- 125 #sample size
6
7 t_test <- abs(qt(.1/2, df= n-2)) #t test that is two tailed
8 t_test
9
10 left <- b1 - se*t_test #obtaining the values of left and right
11 left
12 right <- b1 + se*t_test
13 right
14
15 thorax_lifespan <- lm(fruitfly$thorax ~ fruitfly$lifespan)
16 confint(thorax_lifespan, level = .9) #double checking the confidence
   interval which we obtain to be .0023 to .0033
```

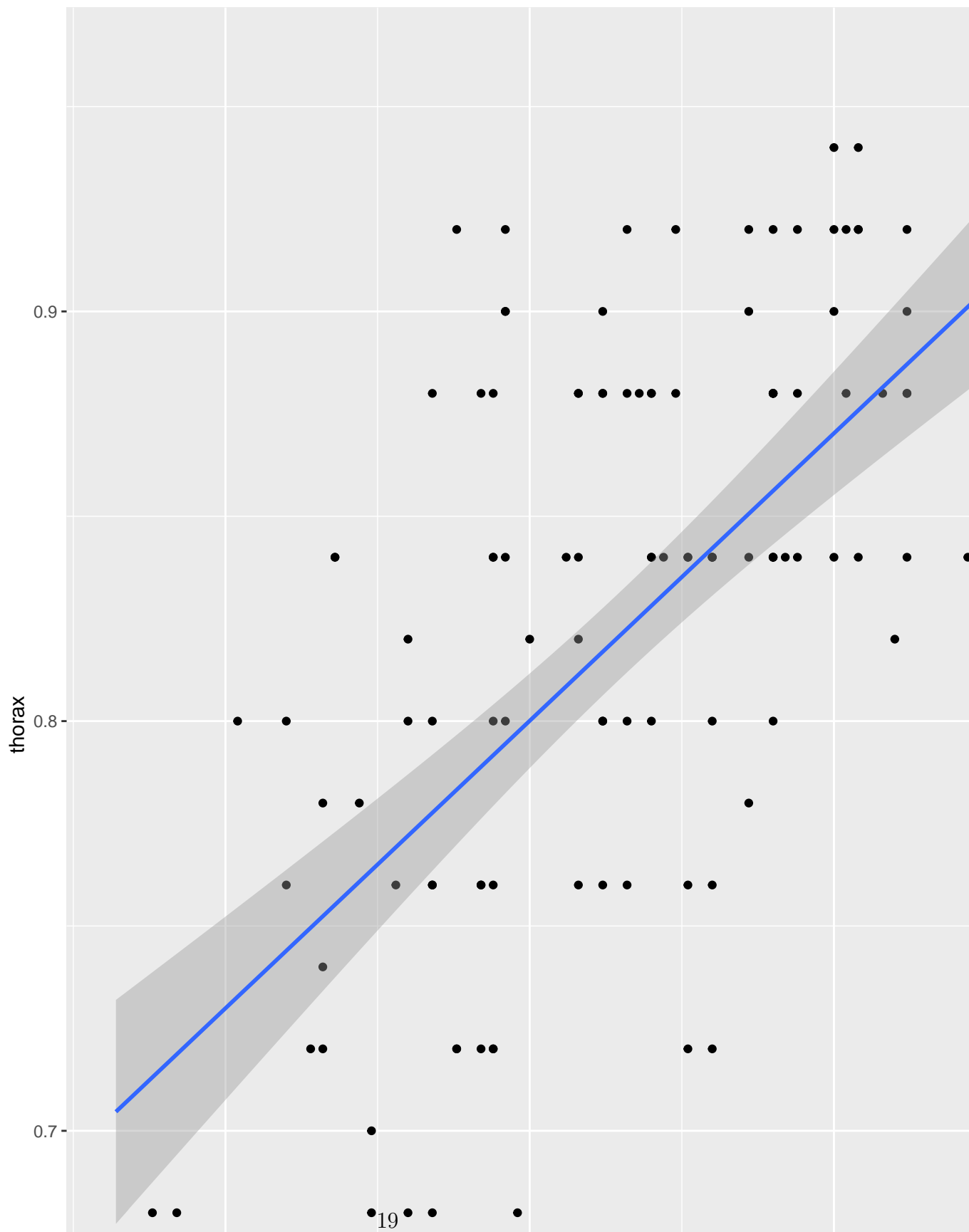
4. Use the `predict()` function in R to (1) predict an individual fruitfly's lifespan when `thorax=0.8` and (2) the average `lifespan` of fruitflies when `thorax=0.8` by the fitted model. This requires that you compute prediction and confidence intervals. What are the expected values of lifespan? What are the prediction and confidence intervals around the expected values?

```
1 summary(lm(fruitfly$thorax ~ fruitfly$lifespan))
2 ggplot(fruitfly, aes(lifespan, thorax)) +
```

```

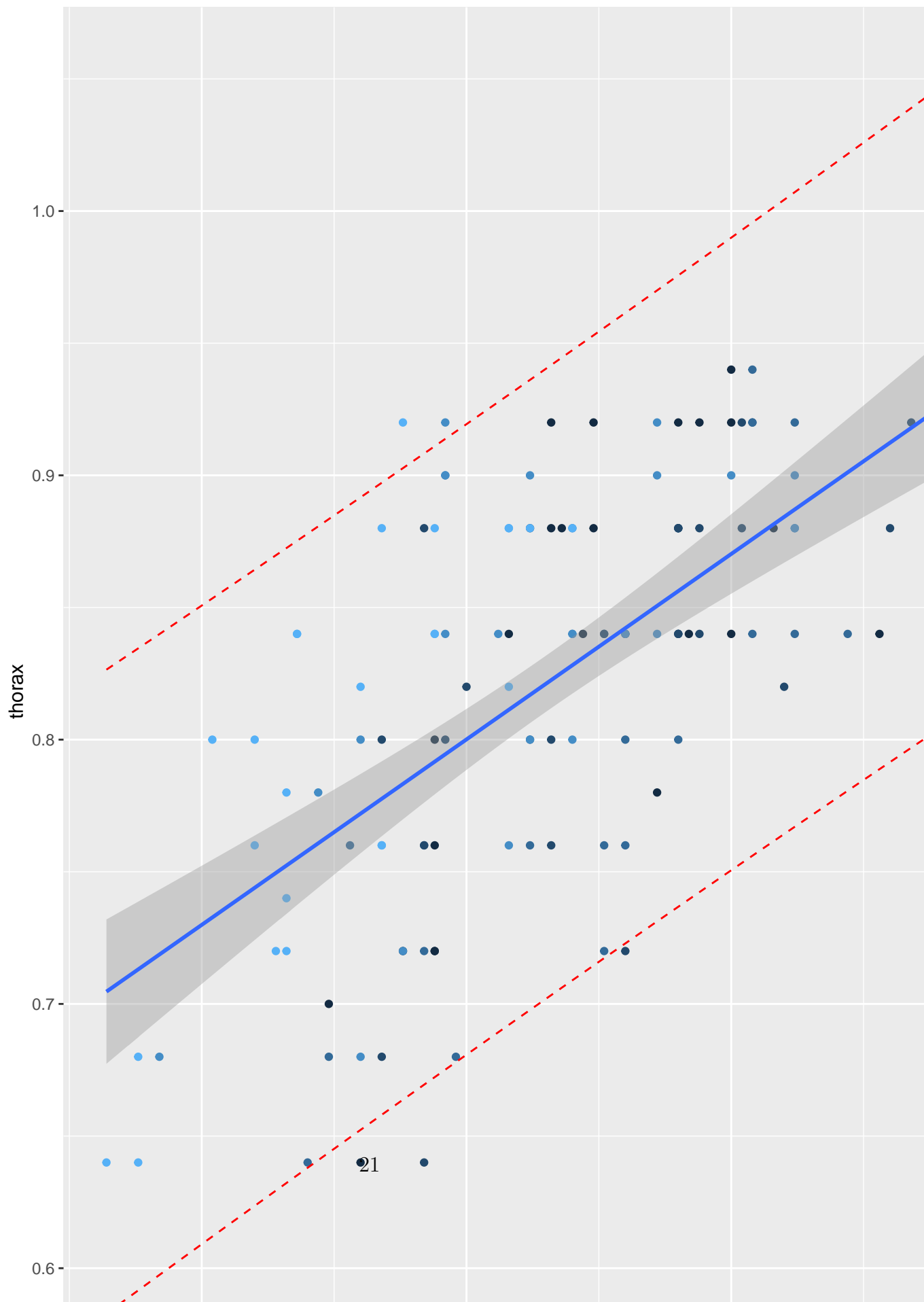
3 geom_point() +
4   stat_smooth(method = lm) #creating variable for graph with the grey
   color indicating the 95% confidence interval range from the regression
   line, whereas the red lines are the prediction lines, all under 95%.
5
6 new_fruitfly <- fruitfly
7 confidence <- as.data.frame(predict(thorax_lifespan, newdata = new_
   fruitfly, interval = "confidence")) #illustrates that on average per
   day in lifespan, the average thorax length ranges
8 mean(confidence$upr) #mean upper value in confidence interval is .8354
   for thorax length
9 mean(confidence$lwr) #mean lower value in confidence interval is .8065
   for thorax length
10
11 predict_data <- as.data.frame(predict(thorax_lifespan, newdata = new_
   fruitfly, interval = "prediction")) #illustrates that on average per
   day in lifespan, the average thorax length ranges
12 mean(predict_data$fit) #mean fitted value for predicted thorax length is
   .821

```



5. For a sequence of `thorax` values, draw a plot with their fitted values for `lifespan`, as well as the prediction intervals and confidence intervals.

```
1 new_graph <- ggplot(aes(lifespan, thorax), data = fruitfly) +  
2   geom_point(aes(color = type)) +  
3   geom_smooth(method = 'lm')  
4  
5 new_graph + geom_line(aes(y = predict_data$lwr), color = "red", linetype  
6   = "dashed") +  
7   geom_line(aes(y = predict_data$upr), color = "red", linetype = "dashed"  
8   )
```



Fitted values are existing on the blue straight regression line across the graph. The grey areas indicate the confidence interval of the fitted values along the regression line; the dashed, red lines, illustrate the prediction of 95