# Question 1 (40 points): Political Science

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

```
1  getwd()
2  setwd("//Users/lucykinnear/Desktop/ASDS/Statistics")
3  mydata = read.csv("Users/lucykinnear/Desktop/ASDS/Statistics\\02_Stats_
       Data_Orig.csv")
4
5  #Assign Observed Frequencies
6  Fo_1 <- 14
7  Fo_2 <- 6
8  Fo_3 <- 7
9  Fo_4 <- 7
10 Fo_5 <- 7
11 Fo_6 <- 1
12
13 #Assign Expected Frequencies
14 Fe_1 <- ((27/42) * 21)
15 Fe_2 <- ((27/42) * 13)
16 Fe_3 <- ((27/42) * 8)
17 Fe_4 <- ((15/42) * 21)
18 Fe_5 <- ((15/42) * 13)
19 Fe_6 <- ((15/42) * 8)
20
21 #Calculate chisq
22 chisq <- (
23   ((Fo_1 - Fe_1)^2/Fe_1) + ((Fo_2 - Fe_2)^2/Fe_2) +
24     ((Fo_3 - Fe_3)^2/Fe_3) + ((Fo_4 - Fe_4)^2/Fe_4) +
25     ((Fo_5 - Fe_5)^2/Fe_5) + ((Fo_6 - Fe_6)^2/Fe_6)
26 )
27 chisq
```

$\chi^2 = 3.7911...$

(b) Now calculate the p-value from the test statistic you just created (in R). What do you conclude if $\alpha = .1$?

```
1  # df = (rows -1)(columns -1) = 1(2) = 2
2  pchisq <- pchisq(3.791168, df = 2, lower.tail = FALSE)
3  pchisq
```

pchisq = 0.15023... The p-value is bigger than alpha, therefore we cannot reject the null hypothesis. We should not be surprised to observe these results. Class did not meaningfully impact whether participants were asked for bribes.

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

```r
table <- matrix(c(14, 6, 7, 7, 7, 1), ncol=3, byrow=TRUE)
colnames(table) <- c('Not Stopped', 'Bribed', 'Warned')
rownames(table) <- c('Upper Class', 'Lower Class')
table <- as.table(table)
chisq_test <- chisq.test(table)
chisq_test
chisq_test$stdres

             Not Stopped      Bribed      Warned
Upper Class    0.3220306  -1.6419565   1.5230259
Lower Class   -0.3220306   1.6419565  -1.5230
```

(d) How might the standardized residuals help you interpret the results?
We can compare the standardized residuals in the table to see which category of variables have the largest difference between the expected and the actual counts relative to size. Standardized residuals are the raw residuals divided by the square root of the expected counts. Positive standardized residuals indicate that there were more occurrences of this outcome than expected. The negative standardized residuals indicate that there were less occurrences of this outcome than expected. Small standardised residuals tell us that the prediction line is a good fit for the data.

# Question 2 (20 points): Economics

(a) State a null and alternative (two-tailed) hypothesis.
Null hypothesis: whether or not there is a reservation policy does not impact the number of new or repaired drinking water facilities in the villages. Alternative hypothesis: whether or not there is a reservation policy does impact the number of new or repaired drinking water facilities in the villages. Two-Tailed Tests:

```
1  z  /2 =  qnorm (1        /2)
```

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

```
1  Q2_data <- read.csv(file="~/Desktop/women.csv")
2  Q2_data_of_interest = lm(Q2_data$water ~ Q2_data$reserved)
3  Q2_data_of_interest
4  summary(Q2_data_of_interest)
5  Call:
6  lm(formula = Q2_data$water ~ Q2_data$reserved)
7
8  Coefficients:
9       (Intercept)   Q2_data$reserved
10           14.738             9.252
11
12  > summary(Q2_data_of_interest)
13
14  Call:
15  lm(formula = Q2_data$water ~ Q2_data$reserved)
16
17  Residuals:
18      Min       1Q   Median       3Q      Max
19  -23.991  -14.738   -7.865    2.262  316.009
20
21  Coefficients:
22                   Estimate  Std. Error  t value  Pr(>|t|)
23  (Intercept)        14.738       2.286    6.446  4.22e-10 ***
24  Q2_data$reserved    9.252       3.948    2.344    0.0197 *
25  ---
26  Signif. codes:  0    ***    0.001    **    0.01    *    0.05    .    0.1
                 1
27
28  Residual standard error: 33.45 on 320 degrees of freedom
29  Multiple R-squared:  0.01688, Adjusted R-squared:  0.0138
30  F-statistic: 5.493 on 1 and 320 DF,  p-value: 0.0197
```

Our p-value is 0.0197. We can reject the null, if alpha is 0.05, but not if alpha is 0.01 ! This means, we have evidence to support the alternative hypothesis. Although, the determinant coefficient (multiple r2) is quite low.

(c) Interpret the coefficient estimate for reservation policy.
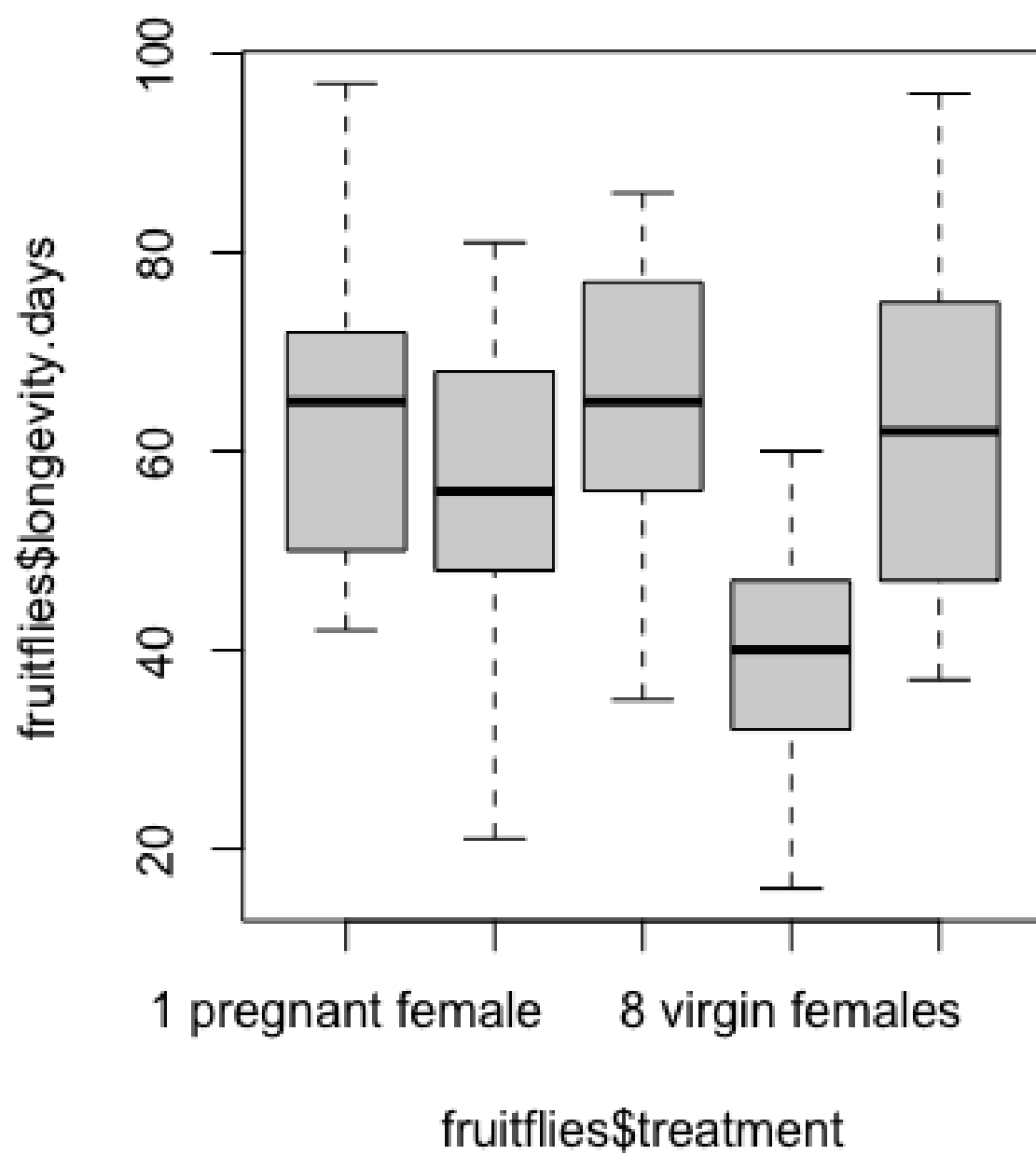
The coefficient of is 9.252 for reservation policy. Changes in the independent variable are associated with changes in the dependent variable at the population level. This variable is statistically significant and probably a worthwhile addition to your regression model. A positive coefficient indicates that as the value of the independent variable increases, the mean of the dependent variable also tends to increase. The coefficient value signifies how much the mean of the dependent variable changes given a one-unit shift in the independent variable while holding other variables in the model constant. For our data, if women are reserved places on the council, water outcomes improve.

# Question 3 (40 points): Biology

1. Import the data set and obtain summary statistiscs and examine the distribution of the overall lifespan of the fruitflies.
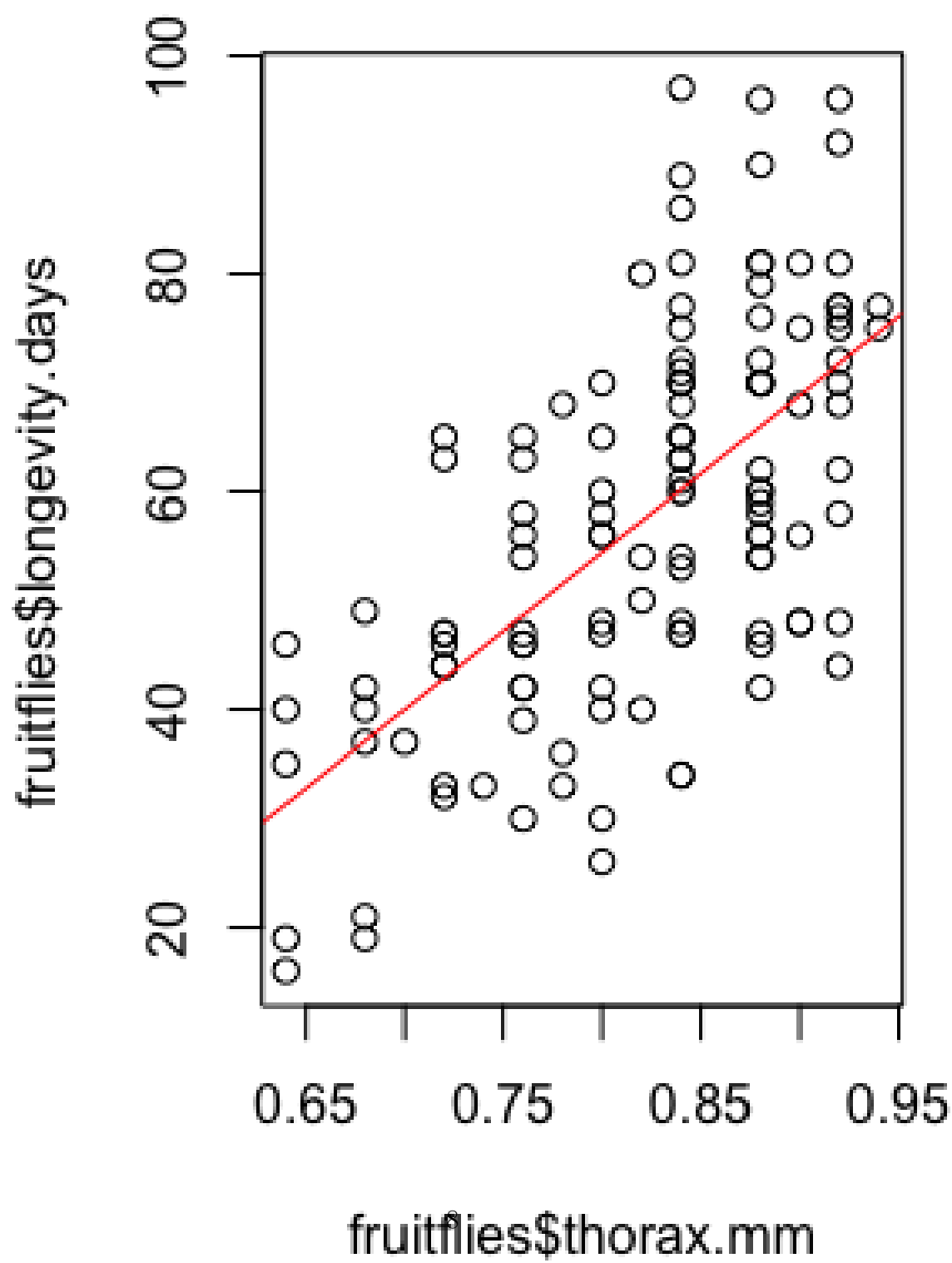
```
1 fruitflies <- read.csv(file="~/Desktop/fruitflies.csv")
2 summary(fruitflies)
3 hist(fruitflies$longevity.days)
4 means
5  1 pregnant female    1 virgin female 8 pregnant females    8 virgin
     females    no females added
6              64.80                  56.76                  63.36
     38.72                63.56
7
8 boxplot(fruitflies$longevity.days ~ fruitflies$treatment)
```

Fruitflies lived a minimum of 16 days and a maximum of 97 days, with a median life span of 58 days. Noteably, this is close to the average lifespan of a fruitfly, 57.4 days - potentially indicating a normal distribution. Groups with 8 virgin females had shortest average lifespan, groups with 1 pregnant female the longest.

02_fruitflies_his.png

2. Plot `lifespan` vs `thorax`. Does it look like there is a linear relationship? Provide the plot. What is the correlation coefficient between these two variables?
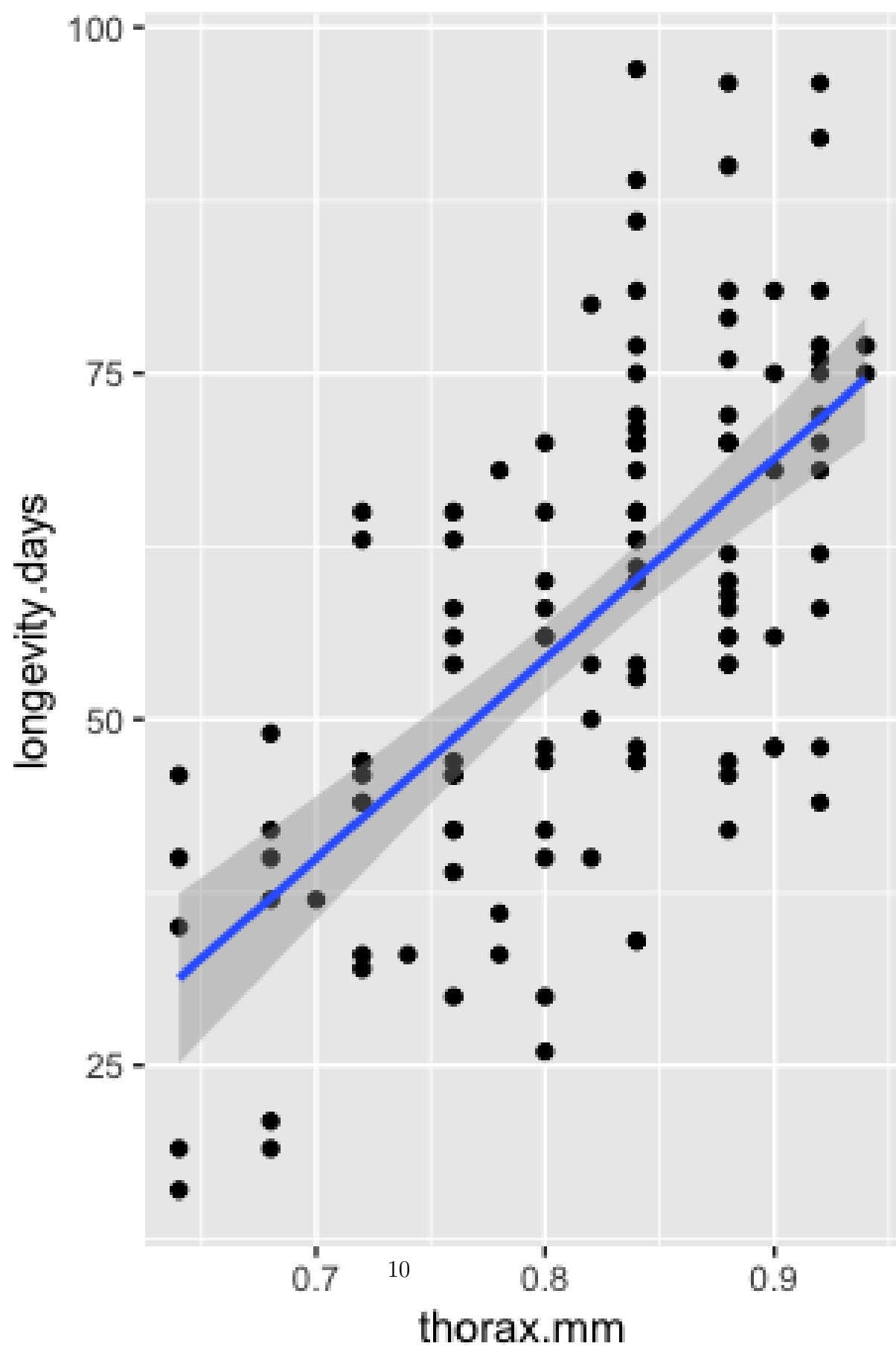
There does appear to be a positive linear relationship.

```
1  cor ( f r u i t f l i e s $ l o n g e v i t y . days , f r u i t f l i e s $ t h o r a x .mm)
```

The correlation coefficient is 0.63640...

3. Regress `lifespan` on `thorax`. Interpret the slope of the fitted model.

```
1  q3_lifespan_thorax = lm( fruitflies$longevity . days ~ fruitflies$thorax .mm)
2  summary ( q3_lifespan_thorax )
3  plot ( fruitflies$longevity . days ~ fruitflies$thorax .mm)
4  abline (lm( fruitflies$longevity . days ~ fruitflies$thorax .mm) , col = "red")
5  library ( ggplot2 )
6  ggplot ( aes ( longevity . days , thorax .mm) , data = fruitflies ) +
7     geom_point ( position = "jitter")
8  abline (lm( fruitflies$longevity . days ~ fruitflies$thorax .mm) , col = "red")
9
10 ggplot ( aes ( thorax .mm, longevity . days ) , data = fruitflies ) +
11    geom_point () +
12    geom_smooth ( method = "lm", formula = y ~ x)
```

The thorax.mm coefficient is 144.33 with tiny p-value.

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

   Extract p-values Null: slope of regression is 0 From summary function:

```
1    summary(q3_lifespan_thorax)
2
3  Call:
4  lm(formula = fruitflies$longevity.days ~ fruitflies$thorax.mm)
5
6  Residuals:
7      Min      1Q   Median      3Q     Max
8  -28.415   -9.961    1.132    9.265   36.812
9
10  Coefficients:
11                         Estimate  Std. Error  t value  Pr(>|t|)
12  (Intercept)              -61.05       13.00   -4.695   7.0e-06 ***
13  fruitflies$thorax.mm     144.33       15.77    9.152   1.5e-15 ***
14  ---
15  Signif. codes:  0    ***    0.001    **    0.01    *    0.05    .    0.1
                1
16
17  Residual standard error: 13.6 on 123 degrees of freedom
18  Multiple R-squared:  0.4051,   Adjusted R-squared:  0.4003
19  F-statistic: 83.76 on 1 and 123 DF,   p-value: 1.497e-15
20
21
```

   p-value way below alpha of 0.05 or 0.01, can reject the null - there is a significant linear relationship here.

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

   Using the formula of confidence interval did not allow me to find the confidence interval because argument was not numeric or logical. Using the function `confint()` in R was

   fruitful.

```
1      confint(q3_lifespan_thorax, level = 0.9)
2                              5 %       95 %
3  (Intercept)          -82.60361   -39.4998
4  fruitflies$thorax.mm  118.19616  170.4700
5
```

   The 90% confidence interval is 118.2 - 170.5.

6. 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. predict

```
1    new.data = data.frame(thorax.mm = 0.8)
2    predict.lifespan = predict(lm(fruitflies$longevity.days ~
       fruitflies$thorax.mm), newdata = new.data)
3    predict.lifespan
4
```

   Lifespan prediction is 54.41.

   2. fitted

```
1    fitted = fitted(lm(fruitflies$longevity.days ~ fruitflies$thorax.mm))
2    data = cbind(fruitflies$thorax.mm, unname(fitted))
3    colnames(data) = c("Thorax in mm", "FittedModel")
4    data = as.data.frame(data)
5    data.08 = data[which(data$thorax.mm == 0.8]
6
7    mean(data.08$FittedModel)
8
9
```

   The mean is 54.41 - i.e. the same as the predicted from part 1.

```
1    predict(lm(fruitflies$longevity.days ~ fruitflies$thorax.mm), newdata =
       new.data, interval = 'confidence')
2
3    predict(lm(fruitflies$longevity.days ~ fruitflies$thorax.mm), newdata =
       new.data, interval = 'prediction')
4
```

   fit = 54.41, lower = 59.92, upper = 56.91
   fit = 54.41, lower = 27.38, upper = 81.45
   Prediction interval is wider because it accounts for variability and error of the estimates. Average lifespan is the same for both, and expected lifespan in prediction is wider, between 27 and 8 days.

7. For a sequence of `thorax` values, draw a plot with their fitted values for `lifespan`, as well as the prediction intervals and confidence intervals.
   I have no clue how to do this !