

Answer Key: Problem Set 2

QTM 200: Applied Regression Analysis

Jeffrey Ziegler

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.¹ As part of the study, confederate made illegal left turns across traffic to draw the attention of the police officers. Two of the confederates were upper class drivers and two were lower class drivers. 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.

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

¹Fried, Brian J, Paul Lagunes, and Atheendar Venkataramani. 2010. "Corruption and Inequality at the Crossroad: A Multimethod Study of Bribery and Discrimination in Latin America". *Latin American Research Review*. 45 (1): 76-97.

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

$$\text{Expected} = \frac{\sum_{\text{Row}} * \sum_{\text{Column}}}{\sum_N}$$

$$\chi^2 = \sum_N \frac{\text{Observed}_i - \text{Expected}_i}{\text{Expected}_i}$$

Let's first try by ourselves:

```

1 # create matrix to conduct chi-square test
2 trafficViolations <- matrix(c(14, 6, 7, 7, 7, 1), byrow=T, nrow=2)
3 rownames(trafficViolations) <- c("Upper class", "Lower class")
4 colnames(trafficViolations) <- c("Not stopped", "Bribe", "Stopped/warned")
5 # by hand approach
6 # create function from chi-square test github.io
7 byHandChiSquare <- function(table){
8   # turn into table
9   observedValues <- as.table(table)
10  # create sums (row, column, and total)
11  grandSum <- sum(observedValues)
12  sumRow <- rowSums(observedValues)
13  sumCol <- colSums(observedValues)
14  # calculate expected values for each observation
15  # check "?outer" to see that this takes the outer product
16  # of the row and col sum divided by the total sum
17  expectedValues <- outer(sumRow, sumCol, "*") / grandSum
18  v <- function(r, c, n) c * r * (n - r) * (n - c) / n^3
19  V <- outer(sumRow, sumCol, v, grandSum)
20
21  dimnames(expectedValues) <- dimnames(observedValues)
22  # create function that calculates each cell residual variance
23  # essentially formula on p. 225 in Agresti and Finlay(2009)
24  test_statistic <- sum((abs(table - expectedValues))^2 / expectedValues)
25  df <- (nrow(observedValues) - 1L) * (ncol(observedValues) - 1L)
26  p_value <- pchisq(test_statistic, df, lower.tail = FALSE)
27  adjusted_residuals <- (observedValues - expectedValues) / sqrt(
28    expectedValues * (1 - sumRow / grandSum) * (1 - sumCol / grandSum))
29  standardized_residuals <- (observedValues - expectedValues) / sqrt(V)
30  # return values
31  return(list(statistic = test_statistic,
32             df = df,
33             p.value = p_value,
34             observed = observedValues,
35             expected = expectedValues,
36             adj_res = adjusted_residuals,
37             std_res = standardized_residuals))
38 }
39 byHandChiSquare(table=trafficViolations)

```

```

$statistic
[1] 3.791168

$df
[1] 2

$p.value
[1] 0.1502306

$observed
      Not_Stopped Bribe_Requested Stopped_Given_Warning
Upper_Class      14              6                    7
Lower_Class       7              7                    1

$expected
      Not_Stopped Bribe_Requested Stopped_Given_Warning
Upper_Class     13.5           8.357143             5.142857
Lower_Class      7.5           4.642857             2.857143

$adj_res
      Not_Stopped Bribe_Requested Stopped_Given_Warning
Upper_Class  0.3220306        -1.5164259             1.6491029
Lower_Class -0.2740361         1.9295276            -1.5230259

$std_res
      Not_Stopped Bribe_Requested Stopped_Given_Warning
Upper_Class  0.3220306        -1.6419565             1.5230259
Lower_Class -0.3220306         1.6419565            -1.5230259

```

Now we can check to make sure:

```

1 # run chi square test with built in function
2 chisq.test(trafficViolations)

```

Pearson's Chi-squared test

```

data:trafficViolations
X-squared=3.7912, df=2, p-value=0.1502

```

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

```
pchisq(3.79, df = (2-1)*(3-1), lower.tail = FALSE) = 0.1502306
```

P-value checks out to our "hand" calculation and the built in function. Cannot reject the null that the two variables of interest are independent.

²Remember 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.*

We can do this by hand (see above function), or the standardized residuals are stored in the `chisq.test` object. We're reporting the standardized residuals, $(\text{observed} - \text{expected}) / \sqrt{V}$, where V is the residual cell variance (Agresti, 2007, section 2.4.5 for the case where x is a matrix, $n * p * (1 - p)$ otherwise).

```
1 # use function to extract standardized residuals
2 chisq.test(trafficViolations)$stdres
```

	Not Stopped	Bribe requested	Stopped/given warning
Upper class	0.322	-1.642	1.523
Lower class	-0.322	1.642	-1.523

(d) *How might the standardized residuals help you interpret the results?*

From the frequency table, it is already clear that there is no obvious pattern for a relationship between rows and columns. Further, the standardized residuals turn out to be quite small, which only supports us to be more confident about the lack of the dependency relationship. None of the standardized residuals indicate any of the cells are more or less than we would expect if the two variables were independent. Nevertheless, they do not tell us much, we need the chi-squared test to make conclusions in either case, i.e. whether variables are dependent or not.

Question 2 (20 points): Economics

Chattopadhyay and Duflo were interested in studying the causal effect of having female politicians on policy outcomes.³ Do women promote different policies than men? 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 "women.csv" 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 will be asked 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
GP	An identifier for the Gram Panchayat (GP)
village	identifier for each village
reserved	binary variable indicating whether the GP was reserved for women leaders or not
female	binary variable indicating whether the GP had a female leader or not
irrigation	variable measuring the number of new or repaired irrigation facilities in the village since the reserve policy started
water	variable measuring the number of new or repaired drinking-water facilities in the village since the reserve policy started

³Raghabendra Chattopadhyay and Esther Duflo. (2004). "Women as Policy Makers: Evidence from a Randomized Policy Experiment in India. *Econometrica*, Vol. 72, No. 5, pp. 1409-1443.

- (a) *State a null and alternative (two-tailed) hypothesis.*

Null: Having reserved seats for female politicians does not change the number drinking water facilities in the villages.

Alternative: The reservation policy has an effect on policy outcomes.

$$H_0 : \beta = 0$$

$$H_A : \beta \neq 0$$

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

After we load our dataset into our working environment, we execute our regression model in which the number of new or repaired water facilities is explained by whether there are reserved seats for female leaders. We then investigate the estimated coefficients of the model using `summary()`.

```
1 # read in women data from online .csv
2 women <- read.csv("https://raw.githubusercontent.com/kosukeimai/qss/
  master/PREDICTION/women.csv")
3 # run regression model with water regressed on whether there are reserved
  seats for women
4 regression_model_problem2 <- lm(water ~ reserved, data=women)
5 # get summary of model with coefficient estimates
6 summary(regression_model_problem2)
```

Coefficients:

	Estimate	Std. Error	t value	Pr(> t)
(Intercept)	14.738	2.286	6.446	4.22e-10 ***
reserved	9.252	3.948	2.344	0.0197 *

Signif. codes: 0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

Residual standard error: 33.45 on 320 degrees of freedom

Multiple R-squared: 0.01688, Adjusted R-squared: 0.0138

F-statistic: 5.493 on 1 and 320 DF, p-value: 0.0197

- (c) *Interpret the coefficient estimate for reservation policy.*

Having reserved seats for female politicians increase the number drinking water facilities in the villages, by 9.2 units. The estimated coefficient is statistically differentiable from zero at the $\alpha = 0.05$ level because the p-value < 0.05 (≈ 0.02).

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*.⁴

No	serial number (1-25) within each group of 25
type	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
lifespan	lifespan (days)
thorax	length of thorax (mm)
sleep	percentage of each day spent sleeping

- (a) Import the data set and obtain summary statistics and examine the distribution of the overall lifespan of the fruitflies.

First, let's load in our data and investigate it by using the `summary()` function, as well as plotting the distribution of lifespan.

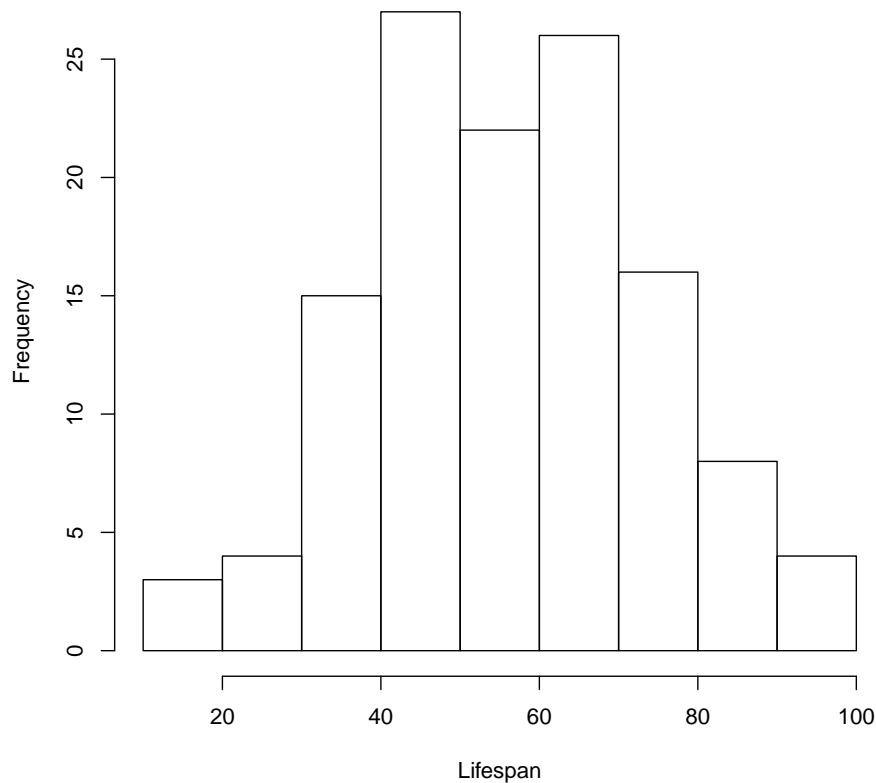
```
1 # (a) import the data set
2 fruitfly <- read.csv("https://raw.githubusercontent.com/jeffreyziegler/
   QTM200Spring2020/master/problem_sets/PS2/fruitfly.csv")
3 # summarize the statistics in the data set
4 summary(fruitfly)
5 # show the distribution of the overall lifespan of the fruitflies
6 pdf("plot3_a.pdf")
7 hist(fruitfly$lifespan, main="", xlab="Lifespan")
8 dev.off()
```

No	type	lifespan	thorax	sleep
Min. : 1	Min. :1	Min. :16.00	Min. :0.640	Min. : 1.00
1st Qu.: 7	1st Qu.:2	1st Qu.:46.00	1st Qu.:0.760	1st Qu.:13.00

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

Median :13	Median :3	Median :58.00	Median :0.840	Median :20.00
Mean :13	Mean :3	Mean :57.44	Mean :0.821	Mean :23.46
3rd Qu.:19	3rd Qu.:4	3rd Qu.:70.00	3rd Qu.:0.880	3rd Qu.:29.00
Max. :25	Max. :5	Max. :97.00	Max. :0.940	Max. :83.00

Figure 2: Histogram of `lifespan`.



- (b) *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?*

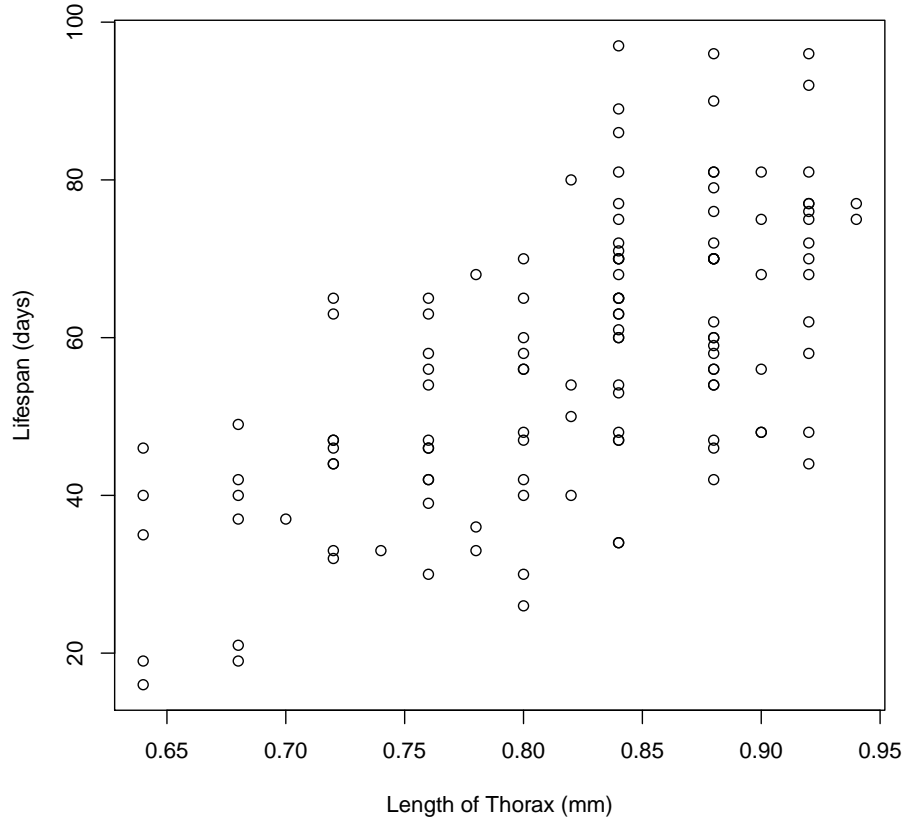
Let's create a scatter plot of the relationship between `thorax` and `lifespan`. We can see in Figure 3 that there appears to be a positive relationship between the two variables, which is confirmed if we investigate the correlation coefficient (using `cor()`), which is ≈ 0.64 .

```

1 pdf("plot3_b.pdf")
2 plot(fruitfly$thorax, fruitfly$lifespan,
3       xlab = "Length of Thorax (mm)", ylab = "Lifespan (days)")
4 dev.off()

```


Figure 3: Scatter plot of `thorax` and `lifespan`.



```
1 # calculate correlation coefficient between lifespan and thorax
2 cor(fruitfly$thorax, fruitfly$lifespan)
```

```
[1] 0.6364835
```

(c) Regress *lifespan* on *thorax*. Interpret the slope of the fitted model.

```
1 # (c) # Run the regression of lifespan on thorax
2 regression_model_problem3 <- lm(lifespan ~ thorax, data=fruitfly)
3 # get summary statistics for linear regression model
4 summary(regression_model_problem3)
```

Using the estimated coefficients and the summary statistics of from the linear regression model, we can calculate the fitted model as $\hat{y} = -61.05 + 144.33x$. The slope of the fitted model is 144.33 and we can interpret the slope as such: when the length of `thorax` increases by 1 mm, the average lifespan a fruitfly increases by 144.33 days.

- (d) *Test for a significant linear relationship between **lifespan** and **thorax**. Provide and interpret your results of your test.*

Table 1: Estimated regression coefficient from model executed in part (c).

	Estimate	Std. Error	t value	Pr(> t)
(Intercept)	-61.05	13.00	-4.69	0.00
thorax	144.33	15.77	9.15	0.00

For our hypothesis test that the slope $\beta_1 = 0$ in Table 1, the associated test statistic $t = 9.152$ and the p -value $= 1.5e - 10$. The p-value is much less than 0.05 (since level of significance $\alpha = 0.05$), so we can reject the null hypothesis ($\beta_1 = 0$). In other words, there is a statistically differentiable linear relationship from zero between the lifespan of the fruitflies and the length of thorax of the fruitflies.

- (e) *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.*

Let's use the formula to calculate the confidence interval where $\hat{\beta}_1 = 144.33$, t -score $= 1.65$ and $se = 15.77$.

```

1 # (e) calculate the confidence interval by formula
2 pointEst <- 144.33
3 se <- 15.77
4
5 # get the t-score
6 t <- qt(0.95, 25*5-2)
7
8 # create the upper and lower bounds
9 lower_CI <- pointEst - t*se
10 upper_CI <- pointEst + t*se

```

The resulting interval is [118.19, 170.47] around $\hat{\beta} = 144.33$. The confidence interval does not include zero, which is consistent with the hypothesis test we did in the previous question.

- *Now, try using the function `confint()` in R.*

Surprise, surprise, we get the same answer!

```

1 # now try confint
2 confint(regression_model_problem3, "thorax", level = 0.9)

```

- (f) 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 # (f) prediction
2 # store the two variables in x and y
3 x <- fruitfly$thorax
4 y <- fruitfly$lifespan
5
6 # predict() function to predict an individual fruitfly lifespan
7 predict(lm(y~x), newdata = data.frame(x=0.8), interval = "prediction",
8         level = 0.90)
9
10 # predict() function to predict the average lifespan of fruitflies
11 predict(lm(y~x), newdata = data.frame(x=0.8), interval = "confidence",
12         level = 0.90)

```

First, let's calculate the predicted individual fruitfly lifespan with thorax value of `0.8mm`. The resulting estimated individual value of lifespan for a fruitfly with `0.8mm` length of thorax is 54.41 days and the corresponding 90% prediction interval is [31.78, 77.05]. The resulting estimated average lifespan for fruitflies with `0.8mm` length of thorax is 54.41 days and the corresponding 90% confidence interval is [52.33, 56.50].

Notice that the prediction interval for an estimated individual lifespan is much wider than the confidence interval for an estimated average lifespan. This is just the same as what we would expect since there is more variability in individual responses than in average responses.

- (g) 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 # (g) create plot of confidence and prediction intervals
2 # find fitted values and prediction and confidence intervals
3 # going from .64 to .94 since these are the max and min values of thorax
4 prep.a <- predict(lm(y~x), newdata = seq(min(x), max(x), 0.003),
5         interval = "confidence")
6
7 prep.i <- predict(lm(y~x), newdata = seq(min(x), max(x), 0.003),
8         interval = "prediction")
9
10 # open a plot to show the fitted values of lifespan, the prediction and
11 # confidence intervals
12 pdf("plot3_g.pdf")
13 matplot(newSeq$x, cbind(prepare.a, prepare.i[, -1]), lty = c(1,2,2,3,3), type =
14         "l",
15         col = c("black", "blue", "green", "red", "purple"),

```

```

11     xlab = "Length of Thorax (mm)", ylab = "Fitted Value of Lifespan
      (days)")
12
13 # Add a legend to the plot to show which line represents which value
14 legend("topleft", legend = c("Fitted Value", "Lower Bound for Confidence
      Interval",
15     "Upper Bound for Confidence Interval", "Lower Bound for Prediction
      Interval",
16     "Upper Bound for Prediction Interval"), lty = c(1,2,2,3,3),
17     col = c("black", "blue", "green", "red", "purple"), cex = 0.85)
18 dev.off()

```

Figure 4: Plot of fitted values, confidence and prediction intervals.

