

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.¹ 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.

¹Fried, 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 χ^2 test statistic by hand (even better if you can do "by hand" in R).

```

1 discrimination <- matrix(c(14, 6, 7, 7, 7, 1), ncol=3, byrow=TRUE)
2 colnames(discrimination) <- c("Not Stopped", "Bribe Requested", "Stopped -
  Given Warning")
3 rownames(discrimination) <- c("Upper Class", "Lower Class")
4 discrimination <- as.table(discrimination)
5
6
7 totals <- sum(discrimination)
8 #Sum Not Stopped
9 not_stopped <- sum(14 + 7)
10 #Sum Bribe Requested
11 bribe_requested <- sum(6 + 7)
12 #Sum Warning
13 warning <- (7 + 1)
14 #Sum Upper Class
15 upper <- sum(14 + 6 + 7)
16 #Sum Lower Class
17 lower <- sum(7 + 7 + 1)
18
19 #Expected Not Stopped, Upper Class
20 fe1 <- (upper/totals)*not_stopped
21 #Expected Bribe requested, Upper Class
22 fe2 <- (upper/totals)*bribe_requested
23 #Expected Warning Given, Upper Class
24 fe3 <- (upper/totals)*warning
25 #Expected Not Stopped, Lower Class
26 fe4 <- (lower/totals)*not_stopped
27 #Expected Bribe requested, Lower Class
28 fe5 <- (lower/totals)*bribe_requested
29 #Expected Warning Given, Lower Class
30 fe6 <- (lower/totals)*warning
31
32 ##Observed-Expected squared divided by expected
33 #Not Stopped, Upper Class
34 chisquare1 <- (14-fe1)^2/fe1
35 #Expected Bribe requested, Upper Class
36 chisquare2 <- (6-fe2)^2/fe2
37 #Expected Warning Given, Upper Class
38 chisquare3 <- (7-fe3)^2/fe3
39 #Expected Not Stopped, Lower Class
40 chisquare4 <- (7-fe4)^2/fe4
41 #Expected Bribe requested, Lower Class

```

```

42 chisquare5 <- (7-fe5)^2/fe5
43 #Expected Warning Given, Lower Class
44 chisquare6 <- (1-fe6)^2/fe6
45
46 #Test Statistic
47 chi_square_stat <- sum(chisquare1 ,chisquare2 ,chisquare3 ,chisquare4 ,
    chisquare5 ,chisquare6 )
48 #3.791168

```

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

```

1 # df = (rows-1)columns-1)
2 df_pchisq <- (2-1)*(3-1)
3 pchisq(chi_square_stat , df = df_pchisq , lower.tail = FALSE)
4 ### With a significance of alpha = 0.1 and a p-value of 0.150, I fail to
    reject the null hypothesis and conclude that
5 ### the class and bribe solicitation are statistically 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.

	Not Stopped	Bribe requested	Stopped/given warning
Upper class	0.1651982	-1.093382	1.182516
Lower class	-0.2014441	1.260318	-1.303106

```

1 z1 <- (14-fe1)/sqrt(fe1*(1-(upper/totals)*(1-(not_stopped/totals))))
2 z2 <- (6-fe2)/sqrt(fe2*(1-(upper/totals)*(1-(bribe_requested/totals))))
3 z3 <- (7-fe3)/sqrt(fe3*(1-(upper/totals)*(1-(warning/totals))))
4 z4 <- (7-fe4)/sqrt(fe4*(1-(lower/totals)*(1-(not_stopped/totals))))
5 z5 <- (7-fe5)/sqrt(fe5*(1-(lower/totals)*(1-(bribe_requested/totals))))
6 z6 <- (1-fe6)/sqrt(fe6*(1-(lower/totals)*(1-(warning/totals))))
7 matrix(c(z1, z2, z3, z4, z5, z6), ncol = 3, byrow = TRUE)

```

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

```

1 #Standardized residuals show us how far away each observed value is from
  the "expectation"
2 #Thus, standardized residuals can tell us the accuracy of the expected
  values and the significance of each cell to the chi squared statistic

```

Question 2 (20 points): Economics

Chattopadhyay and Duflo were interested in whether women promote different policies than men.³ 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).

³Chattopadhyay 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.

```
1 download.file("https://raw.githubusercontent.com/kosukeimai/qss/master/
  PREDICTION/women.csv", "women.csv")
2 women <- read.csv("women.csv")
3 head(women)
4 #2a
5 #Null Hypothesis: The Reservation policy has no effect on the number of
  new or repaired drinking water facilities in the villages.
6 #Alternative Hypothesis: The Reservation policy has an effect 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 #Standardize variation
2 mean_water_y <- mean(women$water)
3 mean_reserved_x <- mean(women$reserved)
4 sum_y <- sum(women$water)
5 sum_x <- sum(women$reserved)
6
7 b_hat <- sum((women$water - mean(women$water)) * (women$reserved - mean(
  women$reserved))) /
8   sum((women$reserved - (mean_reserved_x))^2)
9
10 a_hat <- mean_water_y - (b_hat*mean_reserved_x)
11 #14.738
12
13 #Check
14 lm(women$water ~ women$reserved) #14.738
```

(c) Interpret the coefficient estimate for reservation policy.

¹ #For every increase in reservation for women leaders , there was a 9.252 increase in new or repaired drinking water facilities .

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

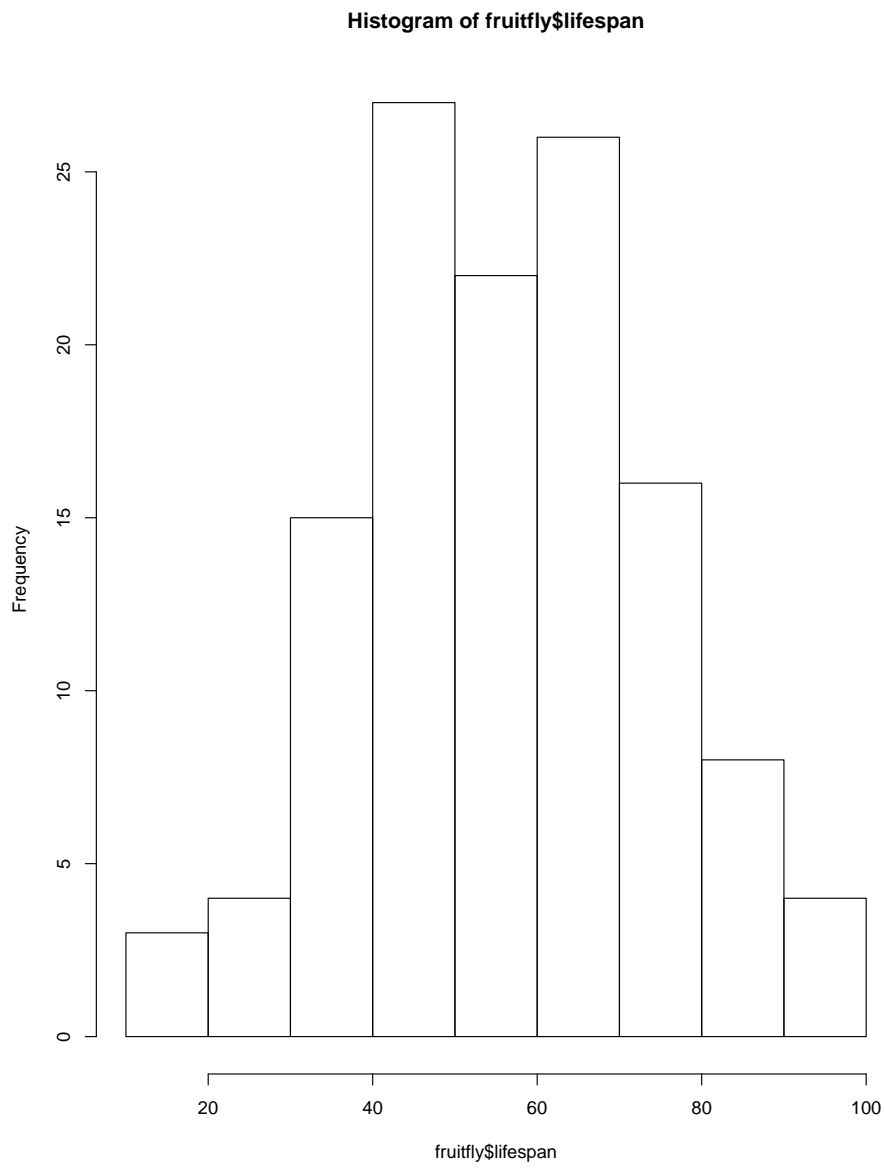
<code>No</code>	serial number (1-25) within each group of 25
<code>type</code>	Type of experimental assignment
	1 = no females
	2 = 1 newly pregnant female
	3 = 8 newly pregnant females
	4 = 1 virgin female
	5 = 8 virgin females
<code>lifespan</code>	lifespan (days)
<code>thorax</code>	length of thorax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

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

```

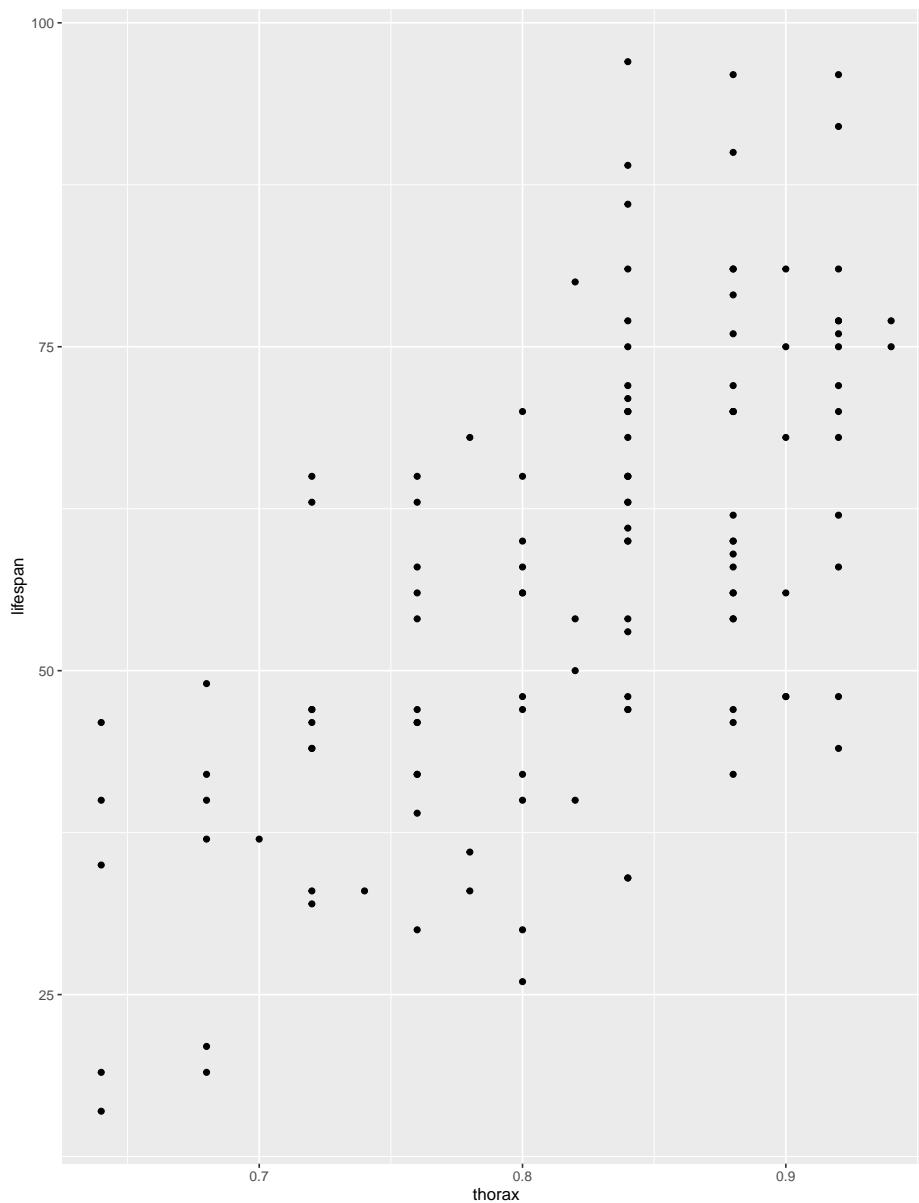
1 fruitfly <- read.csv("fruitfly.csv")
2
3 #3.1
4 summary(fruitfly)
5 # No      type      lifespan      thorax      sleep
6 # Min.    : 1      Min.    :1      Min.    :16.00      Min.    :0.640      Min.    : 1.00
7 # 1st Qu.: 7      1st Qu.:2      1st Qu.:46.00      1st Qu.:0.760      1st Qu.:13.00
8 # Median :13      Median :3      Median :58.00      Median :0.840      Median :20.00
9 # Mean   :13      Mean   :3      Mean   :57.44      Mean   :0.821      Mean   :23.46
10 # 3rd Qu.:19      3rd Qu.:4      3rd Qu.:70.00      3rd Qu.:0.880      3rd Qu.:29.00
11 # Max.   :25      Max.   :5      Max.   :97.00      Max.   :0.940      Max.   :83.00
12 hist(fruitfly$lifespan)
```

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



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?

```
1 library(tidyverse)
2 qplot(x = thorax, y = lifespan, data = fruitfly)
3 #Yes, there seems to be a linear correlation
4
5 cor(fruitfly$thorax, fruitfly$lifespan, method="pearson") #0.6364835
6 #Check if null p = 0
7 cor.test(fruitfly$thorax, fruitfly$lifespan)
8 #0.6364835
```



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

```
1 lm(fruitfly$lifespan ~ fruitfly$thorax)
2 #For every 1 mm increase in thorax size, lifespan increases by 144.33
  days
```

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

```
1 cor.test(fruitfly$lifespan, fruitfly$thorax) #P-value = 1.497e-15
2
3 #Pearson's product-moment correlation
4
5 #data: fruitfly$lifespan and fruitfly$thorax
6 #t = 9.1521, df = 123, p-value = 1.497e-15
7 #alternative hypothesis: true correlation is not equal to 0
8 #95 percent confidence interval:
9 #0.5188709 0.7304479
10 #sample estimates:
11 #cor
12 #0.6364835
13
14 ## With a p-value of 1.497e-15, I reject the null hypothesis that there
   is no correlation in the population between thorax and lifespan.
```

5. 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 ##p-value = 1.497e-15 (see 3.4)
2 t.test(lm(fruitfly$lifespan~fruitfly$thorax), level = 0.90)
3 #Formula based on t values:
4 confint(lm(fruitfly$lifespan~fruitfly$thorax), level = 0.90)
5 #[118.19616,170.4700]
```

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 new_fruitfly <- fruitfly
2 new_fruitfly$thorax <- 0.8
3 #Confidence intervals around expected values of lifespan
4 pred_interval <- predict(lm(fruitfly$lifespan~fruitfly$thorax), newdata=
  new_fruitfly , interval = "prediction", level = 0.95)
5 conf_interval <- predict(lm(fruitfly$lifespan~fruitfly$thorax), newdata=
  new_fruitfly , interval = "confidence", level = 0.95)
6
7 #3.6.2
8 #Expected Values of Lifespan:
9 predict(lm(fruitfly$lifespan~fruitfly$thorax), newdata=new_fruitfly , se.
  fit = TRUE)
```

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.

```
1 fitted_lifespan <- predict(lm(fruitfly$lifespan ~ fruitfly$thorax), newdata
   =new_fruitfly, se.fit = TRUE)
2 reg <- lm(fruitfly$lifespan ~ fruitfly$thorax)
3
4 new_df <- cbind(fruitfly, pred_interval, conf_interval)
5 names(new_df)[9] <- "fit_conf"
6 names(new_df)[10] <- "lwr_conf"
7 names(new_df)[11] <- "upr_conf"
8
9 ggplot(new_df, aes(x=thorax, y=lifespan))+
10   geom_point()+
11   geom_smooth(method=lm, se=TRUE)+
12   geom_line(aes(y=lwr), color = "red", linetype = "dashed")+
13   geom_line(aes(y=upr), color = "red", linetype = "dashed")
```

