# Problem Set 2

### Applied Stats/Quant Methods 1

### Due: October 15, 2021

## 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 GitHub in `.pdf` form.

- This problem set is due before class on Friday October 15, 2021. 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.

---

[1]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).

```r
# (a) Calculate the x^2 statistic

# create a table in R for the given data

data_table <- matrix(c(14, 6, 7, 7, 7, 1), ncol=3, byrow=TRUE)
colnames(data_table) <- c('NotStopped','BribeRequested','
    StoppedGivenWarning')
rownames(data_table) <- c('UpperClass','LowerClass')
data_table <- as.table(data_table)
data_table

barplot(height = data_table,
        beside = TRUE,
        legend.text = TRUE,
        args.legend = list(x = "topleft", cex = 0.4, box.col = "purple"))

# build the dataset for the table

df <- data.frame (first_column  = c("NotStopped"),
                  second_column = c("BribeRequested"),
                  third_column = c("StoppedGivenWarning"))

# Using the chisq.test to get the x^2 test statistic or the chi square
    statistic
# x^2 = 3.7912
# DF = 2 (or n-1)
# value of alpha = .1 (for part b)
```

(b) Now calculate the p-value from the test statistic you just created (in R).[2] What do you
conclude if $\alpha = .1$?

```
1
2  # (b) Now calculate thep-value from the test statistic you just created (
        in R).
3  # What do you conclude if alpha = .1
4
5  chisq <- chisq.test(data_table)
6  chisq
7  tail(df, 2)
8
9
10 # x^2 = 3.7912
11 # DF = 2
12 # alpha value = .01
13 # critical value = 1.885618
14 # x^2(3) = 3.7912, p < 0.1 (the probability was set at 0.1)
15
16 # Find the Critical Value
17 qt(p=0.1, df=2, lower.tail = FALSE)
18 # This gives a critical value of 1.885618
19
20 # The p-value is 0.1502 and the alpha value = 0.1
21 # This concludes that alpha value 0.1 is less than the p-value of 0.1502
22 # X^2 > critical value
23 # reject the Null hypothesis
```

[2]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.

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

```
# (c) Calculate the standardized residual for each cell and put them in the
# table below

chisq$residuals
# this gives
#              NotStopped  BribeRequested  StoppedGivenWarning
# UpperClass   0.1360828     -0.8153742            0.8189230
#LowerClass   -0.1825742      1.0939393           -1.0987005
```

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

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

# to help identify outliers
```
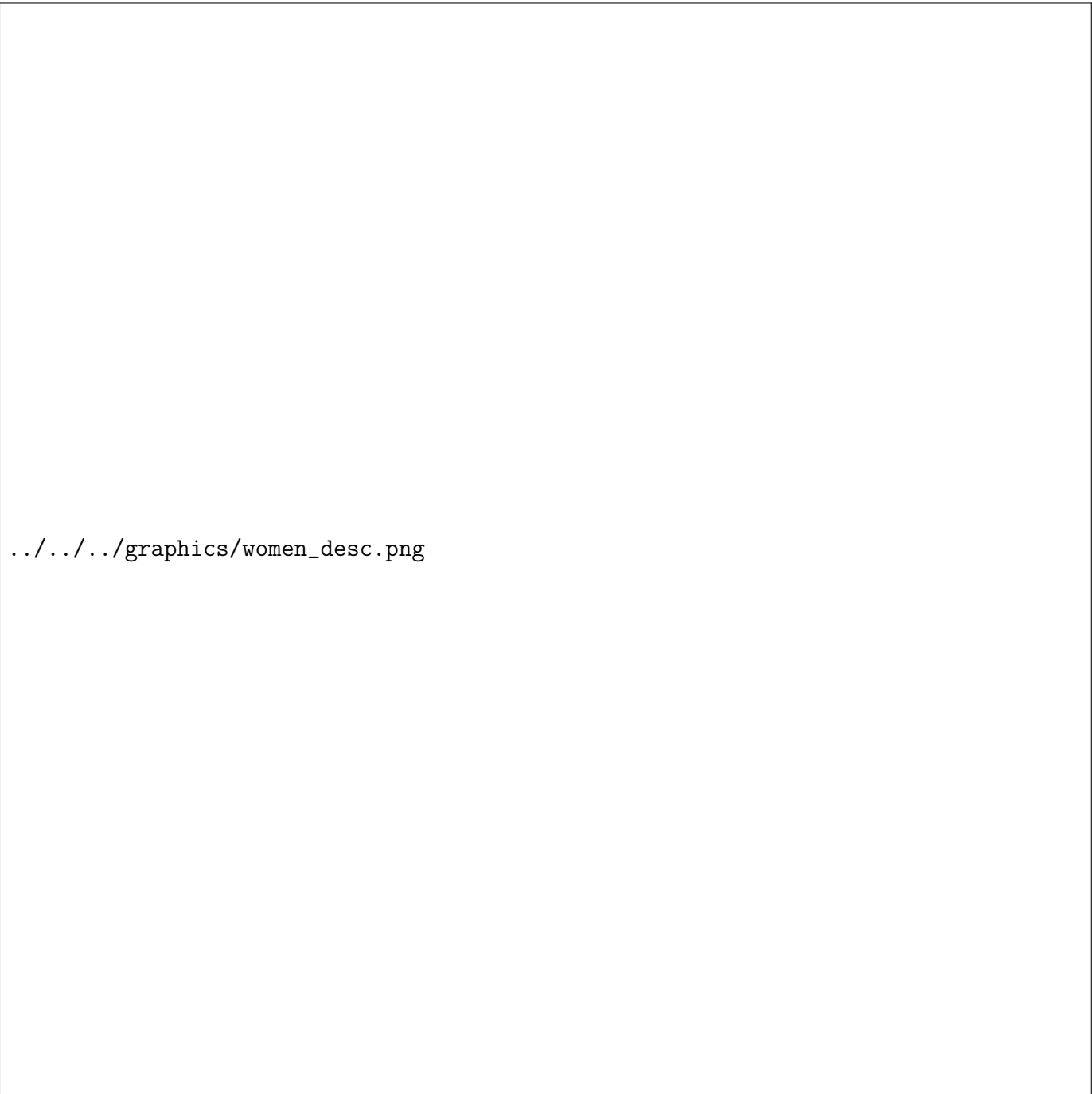
# Question 2 (20 points): Economics

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.

---

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

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

../../../graphics/women_desc.png

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

```
1
2 # (a) State a null and alternative (two-tailed) hypothesis
```

```
3
4 west_bengal <- read.csv("https://raw.githubusercontent.com/kosukeimai/qss
     /master/PREDICTION/women.csv")
5 attach(west_bengal)
6 west_bengal
7 summary(west_bengal)
8 str(west_bengal)
9
10 # H0 is the null hypothesis
11 # Ha is the alternative hypothesis
12
13 # H0: the reservation policy has no effect on the number of new or
     repaired drinking
14 # water facilities in the villages
15 # Ha: the reservation policy does have an effect on the number of new or
     repaired drinking
16 # water facilities in the villages
```

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

```
1
2 # (b) Run a bivariate regression to test this hypothesis
3 # prediction - for the data = women
4 # water depends/extends on reserved
5 res_policy <- lm(water ~ reserved, data = women)
6 summary(res_policy)
7
8 # this gives
9 #Coefficients:
10 #                Estimate Std. Error t value Pr(>|t|)
11 # (Intercept)    14.738      2.286    6.446 4.22e-10 ***
12 # reserved        9.252      3.948    2.344   0.0197 *
13
14 # and p-value of 0.0197
15 # F-statistic: 5.493 on 1 and 320 DF,  p-value: 0.0197
16
17 # taking an alpha of 5%
18 # alpha value = .05
19 # p-value = .0197
20 # the p-value is less than the alpha value that I am trying to test,
21 # so I can reject the null hypothesis
```

(c) Interpret the coefficient estimate for reservation policy.

```
1
2 # (c) Interpret the coefficient for reservation policy
3 # when measuring the correlation between two variables is trying to put a
4 # number on their association
```

# 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]

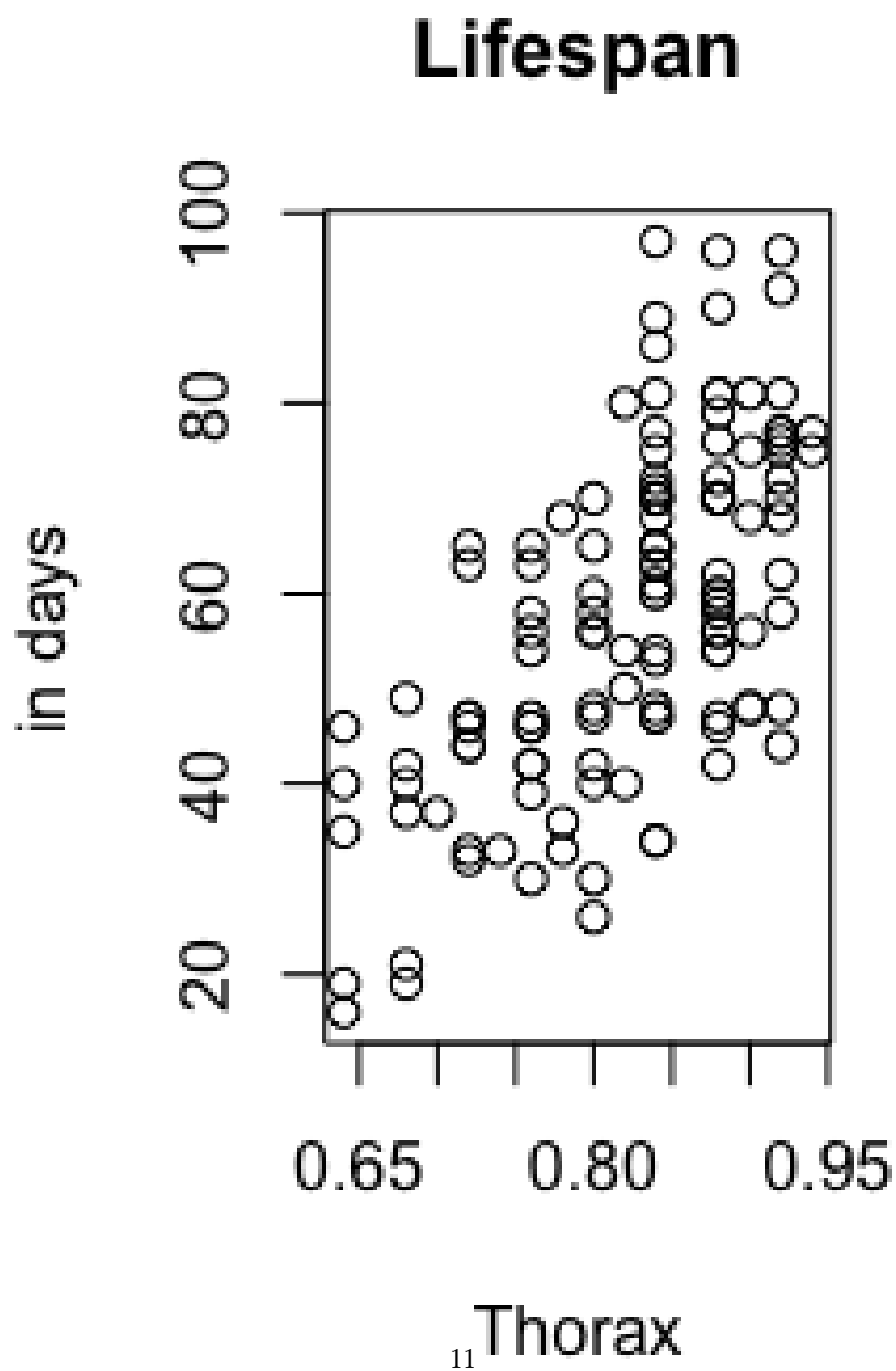| | |
|---:|---|
| No | serial number (1-25) within each group of 25 |
| type | Type of experimental assignment |
| | 1 = no females |
| | 2 = 1 newly pregnant female |
| | 3 = 8 newly pregnant females |
| | 4 = 1 virgin female |
| | 5 = 8 virgin females |
| lifespan | lifespan (days) |
| thorax | length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

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

```
# 1. Import the data set and obtain summary statistics and examine the
    distrubtion of
# the overall lifespan of the fruitflies.

# import data set
dat<-read.csv("http://stat2.org/datasets/FruitFlies.csv")
# or
dat<-read.table("http://stat2.org/datasets/FruitFlies.csv", sep=',',
    header=TRUE)
attach(dat)
dat
# can also get the descriptive statistics in R using Hmisec package
install.packages("Hmisc")
library(Hmisc)
describe(dat)
```
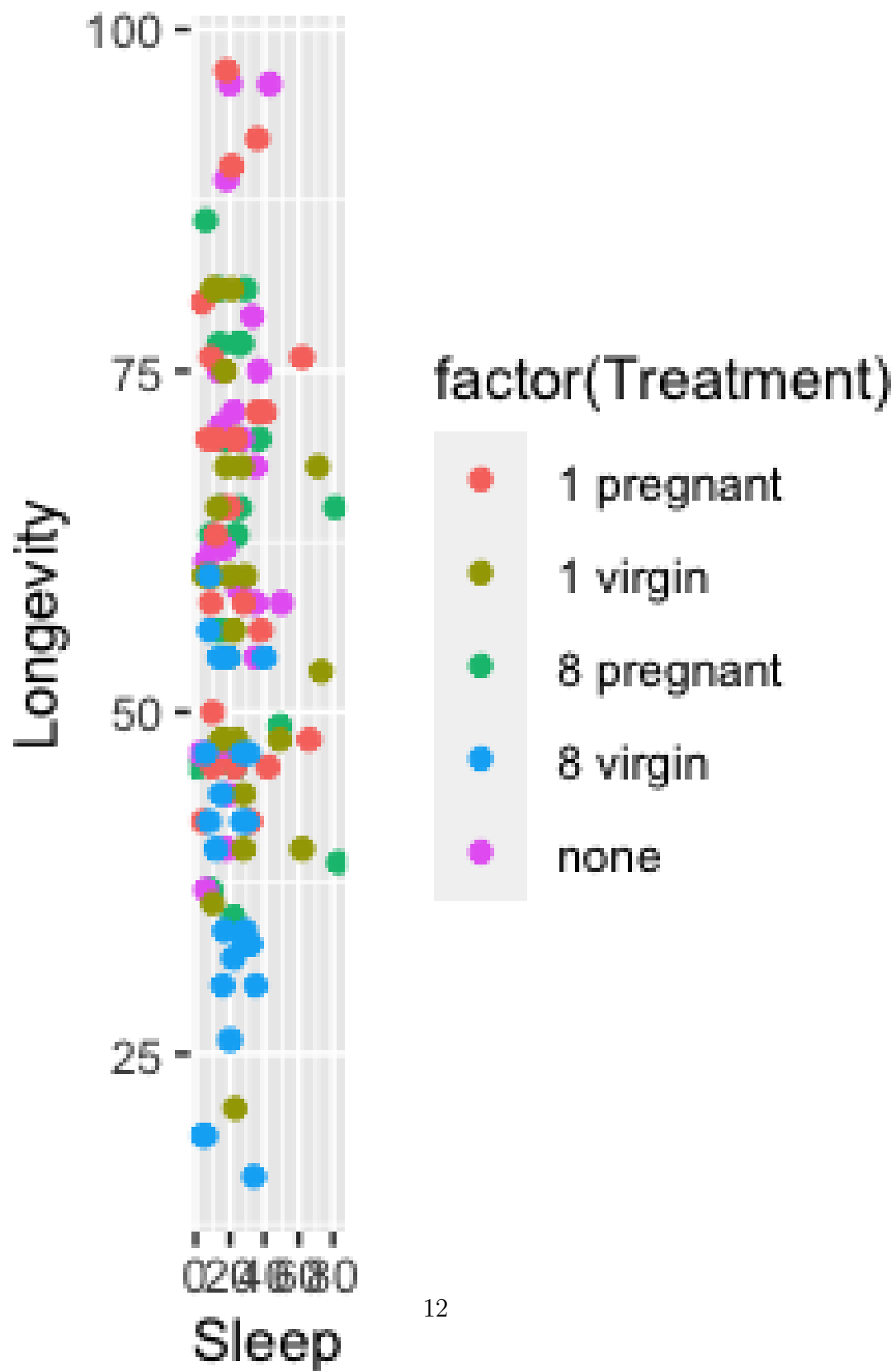
---

[4]Partridge 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?

```r
# 2. Plot lifespan vs thorax.
plot(Longevity~Thorax, xlab="Thorax",
     main="Lifespan", ylab= "in days")

boxplot(Longevity~Treatment, data=dat, frame =FALSE)
summary(dat)

qplot(Sleep, Longevity, data=dat, color=factor(Treatment))+
  labs(title="Longevity and sleep by treatment")

# The correlation coefficient between the two variables Longevity (or
#    lifespan) and thorax
# using the Pearson correlation
cor(dat$Longevity, dat$Thorax)
# the correlation deficit between the two variables lifespan vs thorax is
#    0.6364835

# this gives a summary of rhe statistics of the data
summary(dat)
# examines the distrubution of the overall lifespan of the fruitflies
summary(dat$Longevity)
```

Lifespan

in days

11 Thorax

Longevity and sleep by t

12

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

```
1
2  # 3. Regress lifespan on thorax. Interpret the slope of the fitted model
3
4
5  lm(dat$Longevity ~ dat$Thorax, data=dat)
6  # this gives:
7  # Call:
8  # lm(formula = dat$Longevity ~ dat$Thorax)
9
10 # Coefficients:
11 # (Intercept)     dat$Thorax
12 # -61.05          144.33
13
14 # 144.3 gives the slope of the line
```

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

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

- Use the formula of confidence interval.

- Use the function `confint()` in R .

```
1
2 # 5. Provide the 90 per cent confidence interval for the slope of the
      fitted model
3 # DF 123
4 qt(.05, 123, lower.tail = FALSE)
5
6 #Coefficients:
7 #   (Intercept)    dat$Thorax
8 #    -61.05         144.33
```

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?

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.