

# Problem Set 2

## Applied Stats/Quant Methods 1

Due: October 16, 2022

### 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.
- This problem set is due before 23:59 on Sunday October 16, 2022. No late assignments will be accepted.
- Total available points for this homework is 80.

### 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/manually (even better if you can do "by hand" in R).

**Answer:**

The  $\chi^2$  test statistic was calculated manually in R, using the code below. The value for the  $\chi^2$  test statistic was found to be 3.79.

```

1 # Input observed frequencies
2 o_uppernotstopped <- 14
3 o_upperbribe <- 6
4 o_upperstopped <- 7
5 o_lowernotstopped <- 7
6 o_lowerbribe <- 7
7 o_lowerstopped <- 1
8
9 # Calculate row totals for the explanatory variable (class)
10 total_upper <- sum(14,6,7)
11 total_lower <- sum(7,7,1)
12
13 # Calculate column totals for the response variable (bribe outcome)
14 total_notstopped <- sum(14,7)
15 total_bribe <- sum(6,7)
16 total_stopped <- sum(7,1)
17
18 # Calculate overall sample size
19 total_sample <- sum(total_upper, total_lower)
20
21 # Calculate expected frequencies that would satisfy a null hypothesis
22 # of independence
23 e_uppernotstopped <- (total_upper*total_notstopped/total_sample)
24 e_upperbribe <- (total_upper*total_bribe/total_sample)
25 e_upperstopped <- (total_upper*total_stopped/total_sample)
26 e_lowernotstopped <- (total_lower*total_notstopped/total_sample)
27 e_lowerbribe <- (total_lower*total_bribe/total_sample)
28 e_lowerstopped <- (total_lower*total_stopped/total_sample)
29
30 # Calculate chi-squared test statistic by 1) for each cell, squaring the
31 # differences between the observed and expected frequencies and then
32 # dividing
33 # that square by the expected frequency, and 2) summing these values.
34 chi_sqrd <-
  (((o_uppernotstopped - e_uppernotstopped)^2)/e_uppernotstopped) +

```

```

35  (((o_upperbribe - e_upperbribe)^2)/e_upperbribe) +
36  (((o_upperstopped - e_upperstopped)^2)/e_upperstopped) +
37  (((o_lowernotstopped - e_lowernotstopped)^2)/e_lowernotstopped) +
38  (((o_lowerbribe - e_lowerbribe)^2)/e_lowerbribe) +
39  (((o_lowerstopped - e_lowerstopped)^2)/e_lowerstopped)
40
41  chi_sqrd

```

- (b) Now calculate the p-value from the test statistic you just created (in R).<sup>2</sup> What do you conclude if  $\alpha = 0.1$ ?

**Answer:**

The degrees of freedom is first calculated by multiplying the number of rows minus one by the number of columns minus one, and is found to equal 2. The P-value is then calculated using the `pchisq` function. Its value is found to be approximately 0.15.

The code below illustrates this process:

```

1  # Calculate degrees of freedom (df)
2  df <- (2-1)*(3-1)
3  df
4
5  # Calculate P-value for the test statistic
6  p_value <- pchisq(chi_sqrd, df, lower.tail = FALSE)
7  p_value

```

As our p-value ( $\approx 0.15$ ) is higher than our  $\alpha = 0.1$  threshold, we do not find sufficient evidence to reject the null hypothesis that there is no relationship between a driver's class and the likelihood of an officer soliciting a bribe after committing a minor traffic offence.

---

<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.

	Not Stopped	Bribe requested	Stopped/given warning
Upper class	0.32	-1.64	1.52
Lower class	-0.32	1.64	-1.52

The code below was used to calculate the standardised residuals in R:

```

1 # Calculate the standardised residuals for each cell
2 sr_uppernotstopped <- (o_uppernotstopped - e_uppernotstopped) /
3   sqrt(e_uppernotstopped*(1-(total_upper/total_sample))*(1-(total_
4     notstopped/total_sample)))
5 sr_upperbribe <- (o_upperbribe - e_upperbribe) /
6   sqrt(e_upperbribe*(1-(total_upper/total_sample))*(1-(total_bribe/total_
7     sample)))
8 sr_upperstopped <- (o_upperstopped - e_upperstopped) /
9   sqrt(e_upperstopped*(1-(total_upper/total_sample))*(1-(total_stopped/
10     total_sample)))
11 sr_lowernotstopped <- (o_lowernotstopped - e_lowernotstopped) /
12   sqrt(e_lowernotstopped*(1-(total_lower/total_sample))*(1-(total_
13     notstopped/total_sample)))
14 sr_lowerbribe <- (o_lowerbribe - e_lowerbribe) /
15   sqrt(e_lowerbribe*(1-(total_lower/total_sample))*(1-(total_bribe/total_
16     sample)))
17 sr_lowerstopped <- (o_lowerstopped - e_lowerstopped) /
18   sqrt(e_lowerstopped*(1-(total_lower/total_sample))*(1-(total_stopped/
19     total_sample)))

```

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

## Question 2 (40 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).

---

<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.

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

(c) Interpret the coefficient estimate for reservation policy.