

# Problem Set 2

Applied Stats/Quant Methods 1  
Zhuo Zhang/Student ID: 23346227

Due: October 15, 2023

## 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 15, 2023. No late assignments will be accepted.

## Question 1: 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).

```

1 # Define the observed frequency matrix
2 observed_data <- matrix(c(14, 6, 7, 7, 7, 1), nrow = 2, byrow = TRUE)
3
4 # Perform a chi-squared test
5 test_result <- chisq.test(observed_data)
6
7 # Calculate row sums for the contingency table
8 row_sums <- rowSums(observed_data)
9
10 # Calculate column sums for the table
11 col_sums <- colSums(observed_data)
12
13 # Compute the overall total count
14 total_count <- sum(observed_data)
15
16 # Calculate the expected frequency table
17 expected_data <- outer(row_sums, col_sums) / total_count
18
19 # Calculate the chi-squared values
20 chi_squared_values <- (observed_data - expected_data)^2 / expected_data
21
22 # Calculate the chi-squared test statistic
23 chi_squared_statistic <- sum(chi_squared_values)
24 cat("Chi-Squared Test Statistic: ", chi_squared_statistic, "\n")

```

**Result:**

Chi-Squared Test Statistic: 3.79116

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

```
1 # Degrees of freedom
2 # Calculate the degrees of freedom for the chi-squared test
3 df <- (nrow(observed_data) - 1) * (ncol(observed_data) - 1)
4
5 # Calculate the p-value for the chi-squared test
6 p_value <- 1 - pchisq(chi_squared_statistic, df)
7
8 # Set the significance level (alpha)
9 alpha <- 0.1
10
11 # Print the calculated p-value
12 cat("P-Value: ", p_value, "\n")
13
14 # Perform a hypothesis test based on the p-value and significance level
15 if (p_value <= alpha) {
16   cat("Reject the null hypothesis. There is evidence of an association
17     between variables.\n")
18 } else {
19   cat("Fail to reject the null hypothesis. There is no strong evidence of
20     an association between variables.\n")
21 }
```

**Result:**

Fail to reject the null hypothesis. There is no strong evidence of an association between variables.

---

<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			
Lower class			

```

1 # Calculate the standardized residuals for each cell
2 standardized_residuals <- (observed_data - expected_data) / sqrt(expected
  _data)
3
4 # Create an empty result table
5 result_table <- matrix(NA, nrow = 2, ncol = 3)
6 colnames(result_table) <- c("Not Stopped", "Bribe requested", "Stopped/
  given warning")
7 rownames(result_table) <- c("Upper class", "Lower class")
8
9 # Format the standardized residuals with two decimal places
10 result_table <- format(standardized_residuals, digits = 2)
11
12 # Set row names for the result table
13 rownames(result_table) <- c("Upper class", "Lower class")
14
15 print(result_table)

```

### result

	Not Stopped	Bribe requested	Stopped/given warning
Upper class	0.14	-0.82	0.82
Lower class	-0.18	1.09	-1.10

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

- Looking at the cell where “Upper class” individuals are “Not Stopped,” the standardized residual is 0.14. Since it is close to zero, it suggests that the observed and expected frequencies for this cell are reasonably close, and there may not be a strong association between being in the upper class and not getting stopped.
- For the “Upper class” group and “Bribe requested” cell, the standardized residual is -0.82. This negative value signals that the observed frequency of upper-class individuals requesting a bribe is lower than what would be expected. It suggests that being in the upper class might be associated with a lower likelihood of requesting a bribe.
- For the cell corresponding to “Upper class” and “Stopped/given warning,” the standardized residual is 0.82, a positive value. This suggests that the observed frequency of upper-class individuals being stopped or given a warning is higher than expected. It implies that being in the upper class might be associated with a higher likelihood of being stopped or given a warning.
- For the cell corresponding to “Lower class” and “Not Stopped,” the standardized residual is -0.18, which is close to zero. It indicates that there may not be a strong association between being in the lower class and not getting stopped.
- For the cell corresponding to “Lower class” and “Bribe requested,” the standardized residual is 1.09, indicating that the observed frequency of lower-class individuals requesting a bribe is higher than expected. It suggests that being in the lower class might be associated with a higher likelihood of requesting a bribe.
- For the “Lower class” group and “Stopped/given warning,” the standardized residual is -1.10, a negative value. It suggests that the observed frequency of lower-class individuals being stopped or given a warning is lower than expected. This implies that being in the lower class might be associated with a lower likelihood of experiencing these interactions.

## Question 2: 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).

Name	Description
<b>GP</b>	An identifier for the Gram Panchayat (GP)
<b>village</b>	identifier for each village
<b>reserved</b>	binary variable indicating whether the GP was reserved for women leaders or not
<b>female</b>	binary variable indicating whether the GP had a female leader or not
<b>irrigation</b>	variable measuring the number of new or repaired irrigation facilities in the village since the reserve policy started
<b>water</b>	variable measuring the number of new or repaired drinking-water facilities in the village since the reserve policy started

---

<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.

- Null Hypothesis : The retention policy for women among village committee leaders has no impact on the number of new or repaired drinking water facilities in the village. Mathematically speaking, this can be expressed as: where is the coefficient of the “reserved” variable in the regression model.
- Alternative Hypothesis: The retention policy of village committee leaders towards women has an impact on the number of new or repaired drinking water facilities in the village. Mathematically speaking, this can be expressed as: this is a double tailed substitution hypothesis, as we are testing whether the coefficients differ significantly from zero in any direction.

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

```
1 data <- read.csv("https://raw.githubusercontent.com/kosukeimai/qss/master/PREDICTION/women.csv")
2 #Run the bivariate regression
3 model <- lm(water ~ reserved, data = data)
4 ## Summarize the regression results
5 summary(model)
```

### Result:

Call:

```
lm(formula = water ~ reserved, data = data)
```

Residuals:

	Min	1Q	Median	3Q	Max
	-23.991	-14.738	-7.865	2.262	316.009

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.

- The coefficient estimate for the “reserved” variable stands at 9.252, bearing a standard error of 3.948, along with an associated p-value of 0.0197. This coefficient signifies the alteration in the count of newly established or refurbished drinking water facilities within villages that correspond to a one-unit shift in the “reserved” variable. This variable determines whether the village council head’s position is set aside for women.
  - Coefficient Value: The positive coefficient estimate of 9.252 implies that, on average, when the village council head’s position is reserved for women (i.e., “reserved” = 1), there is an approximate increment of 9.252 in the number of fresh or fixed drinking water facilities in contrast to situations with no reservation policy (i.e., “reserved” = 0).
  - Statistical Significance: The coefficient bears statistical significance, as indicated by a p-value of 0.0197, which is below the conventional significance threshold of 0.05. This suggests that the implementation of a reservation policy for women in village council head positions has a noteworthy impact on the count of new or repaired drinking water facilities in villages.