

# Problem Set 2/

Quant Methods 1/Due: October 14, 2024

Zengyuan Zhao / zhaoze@tcd.ie

## 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 Monday October 14, 2024. 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 #Q1.1
2 #The X-squared statistic is calculated by subtracting the square of the
  expected value from the observed value and dividing by the expected
  value.
3 #The question gives a 2*3 contingency table.
4 resulting_data <- matrix(c(14, 6, 7,
5                             7, 7, 1),
6                             nrow = 2, byrow = TRUE,
7                             dimnames = list(c("Upper class", "Lower class"),
8                                             c("Not Stopped", "Bribe requested",
9                                             "Stopped/given warning")))
9 # Calculate summary rows
10 row_sums <- rowSums(resulting_data)
11 col_sums <- colSums(resulting_data)
12 total_sum <- sum(resulting_data)
13 # Create a new matrix including raw data and summary rows
14 resulting_table <- rbind(resulting_data, c(col_sums, total_sum))
15 resulting_table <- cbind(resulting_table, c(row_sums, total_sum))
16 rownames(resulting_table) <- c("Upper class", "Lower class", "Total")
17 colnames(resulting_table) <- c("Not Stopped", "Bribe requested", "Stopped
  /given warning", "Total")
18 resulting_table
19 #Calculate expected frequency
20 Ex_Upper_NotStopped <- (27*21)/42
21 Ex_Upper_Bribe <- (13*27)/42
22 Ex_Upper_Stopped <- (8*27)/42
23 Ex_Lower_NotStopped <- (21*15)/42
24 Ex_Lower_Bribe <- (13*15)/42
25 Ex_Lower_Stopped <- (8*15)/42
26 #In order to be more intuitive, make an expected frequency table
27 Expected_table <- matrix(c(Ex_Upper_NotStopped, Ex_Upper_Bribe, Ex_Upper_
  Stopped,
28                             Ex_Lower_NotStopped, Ex_Lower_Bribe, Ex_Lower_
  Stopped),
29                             nrow = 2, byrow = TRUE,
30                             dimnames = list(c("Upper class", "Lower class"),
31                                             c("Not Stopped", "Bribe
  requested", "Stopped/given warning")))
32 Expected_table
33 #Calculate X^2
34 X_squared_statistic <- (14-13.5)^2/13.5 + (6-8.357143)^2/8.357143 +

```

```

35 (7-5.142857)^2/5.142857 + (7-7.5)^2/7.5 +
36 (7-4.642857)^2/4.642857 + (1-2.857143)^2/2.857143
37 X_squared_statistic

```

First, we need to calculate the row total, column total, and grand total, as shown in Figure 1

Figure 1: Summary table

```
> resulting_table
```

	Not Stopped	Bribe requested	Stopped/given warning	Total
Upper class	14	6	7	27
Lower class	7	7	1	15
Total	21	13	8	42

Next, we need to calculate the expected frequency, as shown in Figure 2

Figure 2: Expected Frequency table

```
> Expected_table
```

	Not Stopped	Bribe requested	Stopped/given warning
Upper class	13.5	8.357143	5.142857
Lower class	7.5	4.642857	2.857143

Finally, according to the X-squared statistical calculation formula, we get:

Figure 3: X-squared statistical calculation

```

> X_squared_statistic
[1] 3.791169

```

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

---

<sup>2</sup>Remember frequency should be  $> 5$  for all cells, but let's calculate the p-value here anyway.

```

1 #Q1.2
2 #Calculate degrees of freedom
3 df <- (2-1)*(3-1)
4 df
5 # Calculate p_value
6 p_value <- pchisq(X_squared_statistic, df, lower.tail = FALSE)
7 p_value

```

After calculation, the p value is 0.1502305

Figure 4: p value

```

> p_value
[1] 0.1502305

```

If  $\alpha$  is 0.1, then the confidence interval is 90, and the p-value is greater than 0.1, so the null hypothesis cannot be rejected. Therefore, there is no significant correlation between the driver's class status and whether he is asked for a bribe.

- (c) Calculate the standardized residuals for each cell and put them in the table below.

	Not Stopped	Bribe requested	Stopped/given warning
Upper class	0.3220306	-1.641957	1.523026
Lower class	-0.3220306	1.641957	-1.523026

```

1 #1.3
2 #Formula for standardized residuals:(observed-expected)/sqrt(expected(1-
  row prop.)*(1-column prop.))
3 #row prop = row total/grand total
4 #column prop = column/grand total
5 sr_Upper_NotStopped <- (14-13.5)/sqrt(13.5*(1-27/42)*(1-21/42))
6 sr_Upper_Bribe <- (6-8.357143)/sqrt(8.357143*(1-27/42)*(1-13/42))
7 sr_Upper_Stopped <- (7-5.142857)/sqrt(5.142857*(1-27/42)*(1-8/42))
8 sr_Lower_NotStopped <- (7-7.5)/sqrt(7.5*(1-15/42)*(1-21/42))
9 sr_Lower_Bribe <- (7-4.642857)/sqrt(4.642857*(1-15/42)*(1-13/42))
10 sr_Lower_Stopped <- (1-2.857143)/sqrt(2.857143*(1-15/42)*(1-8/42))

```

```

11 sr_table <- matrix(c(sr_Upper_NotStopped, sr_Upper_Bribe, sr_Upper_
12   Stopped,
13   sr_Lower_NotStopped, sr_Lower_Bribe, sr_Lower_
14   Stopped),
15   nrow = 2, byrow = TRUE,
16   dimnames = list(c("Upper class", "Lower class"),
17   c("Not Stopped", "Bribe
18   requested", "Stopped/given warning")))
19 sr_table
20 sr_Upper_NotStopped
21 sr_Upper_Bribe
22 sr_Upper_Stopped
23 sr_Lower_NotStopped
24 sr_Lower_Bribe
25 sr_Lower_Stopped

```

Figure 5: standardized residuals

```

> sr_table
      Not Stopped Bribe requested Stopped/given warning
Upper class  0.3220306      -1.641957      1.523026
Lower class -0.3220306       1.641957     -1.523026
> sr_Upper_NotStopped
[1] 0.3220306
> sr_Upper_Bribe
[1] -1.641957
> sr_Upper_Stopped
[1] 1.523026
> sr_Lower_NotStopped
[1] -0.3220306
> sr_Lower_Bribe
[1] 1.641957
> sr_Lower_Stopped
[1] -1.523026

```

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

(1)The value of the standardized residual for the upper class who has not been stopped is positive and relatively small (0.3220306), less than 1 and close to 0, indicating that this observation is slightly higher than the predicted value of the model, but the degree of deviation is not large.

(2)The standardized residual value for the upper class who has been asked for a bribe is negative and relatively large (-1.641957), with an absolute value greater than 1 and less than 2, indicating that this observation is significantly lower than the predicted value of the model, and the degree of deviation is large, but it is still within a reasonable range.

(3)The value of the upper class who has been stopped or warned is positive and relatively large (1.523026), greater than 1 but less than 2, indicating that this observation is significantly higher than the predicted value of the model, although the degree of deviation is large, but it is still within a reasonable range.

Lower Class

(1)The value of the standardized residual for the lower class who has not been stopped is negative and relatively small (-0.3220306), with an absolute value less than 1 and close to 0, indicating that this observation is slightly lower than the predicted value of the model, but the degree of deviation is not large.

(2)The standardized residual value of the lower class being asked for bribes is positive and relatively large (1.641957), greater than 1 and less than 2, indicating that this observation is significantly higher than the model's predicted value, and the degree of deviation is large, but it is still within a reasonable range.

(3)The standardized residual value of the lower class being stopped or warned is negative and relatively large (-1.523026), and the absolute value is greater than 1 but less than 2, indicating that this observation is significantly lower than the model's predicted value. Although the degree of deviation is large, it is still within a reasonable range.

Conclusion

The standardized residual values of the upper and lower classes being treated differently are all between  $[-2, 2]$ , all within the normal range, and there are no abnormal values. The results passed the standardized residual test, and the hypothesis test model met the requirements.

## 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 6 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 6: 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 reservation policy for women as village government leaders has no effect on the number of new or repaired drinking water facilities in the villages.

Alternative Hypothesis: The reservation policy for women as village government leaders has positive/negative 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 #read data
2 women_data <- read.csv("https://raw.githubusercontent.com/kosukeimai/qss/
  master/PREDICTION/women.csv")
3 women_data
4 #inspect data
5 inspect_null <- sum(is.na(women_data))
6 inspect
7 str(women_data)
8 class(women_data)
9 #bivariate regression
10 model1 <- lm(women_data$irrigation ~ women_data$reserved, data = women_
  data)
11 summary(model1)
12 model2 <- lm(women_data$water ~ women_data$reserved, data = women_data)
13 summary(model2)
```

By observing the data and combining it with the variable explanation table, we can see that the values of the "women" variable and the "reserved" variable are the same. Therefore, I decided to use the "reserved" variable as the independent variable and "irrigation" and "water" as the dependent variables to explore their regression relationship.



Figure 7: The binary regression between irrigation and reserved

Call:

```
lm(formula = women_data$irrigation ~ women_data$reserved, data = women_data)
```

Residuals:

Min	1Q	Median	3Q	Max
-3.388	-3.388	-3.019	-1.019	86.612

Coefficients:

	Estimate	Std. Error	t value	Pr(> t )
(Intercept)	3.3879	0.6498	5.214	3.33e-07 ***
women_data\$reserved	-0.3693	1.1220	-0.329	0.742

---

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

Residual standard error: 9.506 on 320 degrees of freedom

Multiple R-squared: 0.0003385, Adjusted R-squared: -0.002785

F-statistic: 0.1084 on 1 and 320 DF, p-value: 0.7422

From the above figure, we can see that the p-value of the regression analysis of "irrigation" and "reserved" is 0.7422, which is much larger than 0.1, and the null hypothesis cannot be rejected. Therefore, there is no significant correlation between the implementation of the reservation policy and the number of irrigation facilities.

Next, I will calculate the regression relationship between the variables water and reserved.

Figure 8: The binary regression between water and reserved

```
Call:
lm(formula = women_data$water ~ women_data$reserved, data = women_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 ***
women_data$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
```

From the above figure, we can see that the p-value of the regression analysis of "water" and "reserved" is 0.0197, which is less than 0.05, and the null hypothesis is rejected. Therefore, there is a significant correlation between the implementation of the reservation policy and the number of drinking water facilities.

(c) Interpret the coefficient estimate for reservation policy.

The coefficient estimate of the regression model of water and reserved is 9.252, which is a positive number and shows a positive correlation. This usually means that if the independent variable reserved increases by one unit, the dependent variable water may increase by 9.252 units. Under this premise, it is possible that for every leadership position reserved for a woman, 9.252 drinking water facilities will be added.

Assuming that the variables are linear and considering that the intercept is 14.738, we can get the relationship between reserved and water:

$$Y(\text{water}) = 9.252 * X(\text{reserved}) + 14.738$$