

Problem Set 2

Applied Stats/Quant Methods 1

Due: October 16, 2022

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) The χ^2 test statistic is calculated as follows:

Read in the data as a matrix.

```
1 observed <- matrix( c (14, 6, 7, 7, 7, 1), nrow = 2, byrow = TRUE)
```

Calculate the expected values, then calculate the difference between the observed and expected values for each sub-category. Calculate the contribution to the χ^2 statistic.

expected number in class * number of outcomes / total number

difference observed - expected

contribution difference²/expected)

For example, for the sub-category ‘Upper Class’ and ‘Not Stopped’:

Upper Class, Not Stopped	
observed	14
expected	13.5 = (27 * 21 / 42)
difference	0.5 = (14 - 13.5)
chi sq contribution	0.0185 = (0.5) ² / 13.5

```
1 ncols <- length(observed[1,])
2 nrows <- length(observed[,1])
3
4 # get totals
5 row_tots <- vector("double", nrows)
6 col_tots <- vector("double", ncols)
7
8 totals <- sum(observed) # total number of observations
9
10 # calculate row and column totals, e.g, total for NotStopped, UpperClass,
    etc
11 for (i in 1:nrows) {row_tots[i] <- sum(observed[i,])}
12 for (i in 1:ncols) {col_tots[i] <- sum(observed[,i])}
13
14 #get expected = row total * column total / total observations
15 expected <- observed
16
17 for (i in 1:nrows) {
18   for (j in 1:ncols) {
19     expected[i,j] <- row_tots[i] * col_tots[j] / totals
```

```

20 }
21 }
22
23 # calculate difference between observed and expected
24 o_e <- observed
25 o_e <- (o_e - expected)^2 / expected
26
27 #calculate chi-squared value & degrees of freedom
28 chi_sq_val <- sum(o_e)
29 df = (nrows-1) * (ncols-1)

```

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

```

1 p_value <- pchisq(chi_sq_val, df=df, lower.tail=FALSE)
2 alpha <- 0.1

```

The p-value is 15.02%, alpha is 10%

We cannot reject the null hypothesis that the two sets are from the same population

1 observed cell(s) with less than 5 values

The observed and expected values are shown in Figure 1

The results of the builtin R `chisq.test` function are as follows:

Pearson's Chi-squared test

data: observed

X-squared = 3.7912, df = 2, p-value = 0.1502

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

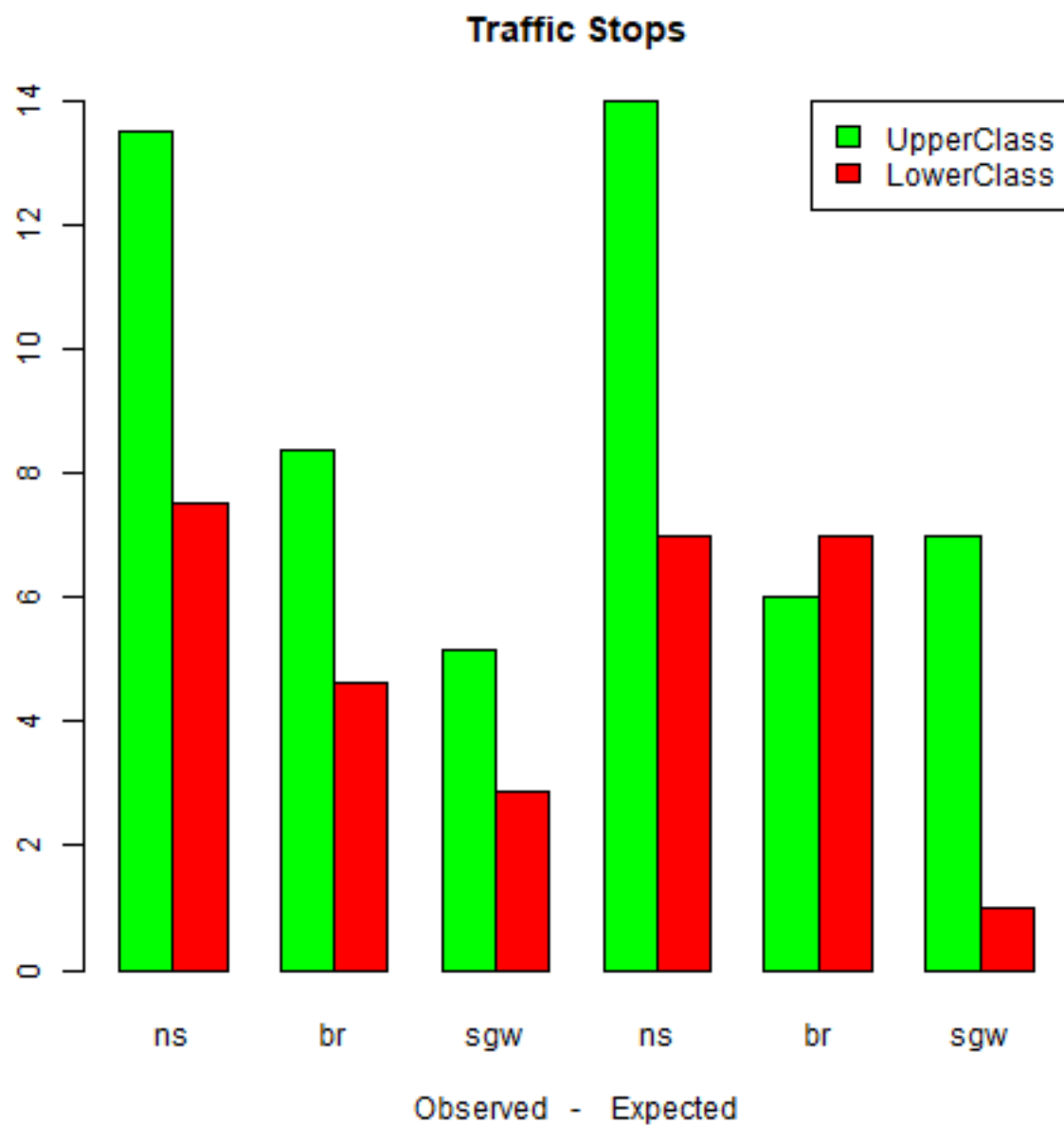


Figure 1: Observed vs Expected values for traffic stop. ns = Not Stopped; br = Bribe Requested; sgw = Stopped Given Warning

(c) The standardized residuals are set out in the table below:

Table 1: Standardised Residuals

	NotStopped	BribeRequested	StoppedGivenWarning
UpperClass	0.322	-1.642	1.523
LowerClass	-0.322	1.642	-1.523

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

The biggest contribution to the residuals was from the 'Bribe Requested' variable - fewer upper class individuals were expected to hand over bribes. The difference between the two groups appears to be a combination of fewer upper class drivers being expected to hand over bribes and more of them being given a warning instead the opposite outcome occurring for lower class drivers.

We are not rejecting the null hypothesis, so we are concluding that there may not be any significant relationship between class and the outcomes experienced during traffic stops. The combined effect from the different experiences of the two groups was not enough to convince us that class predicts whether or not a driver is asked for a bribe.

Question 2 (40 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 2 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 2: Names and description of variables from Chattopadhyay and Duflo (2004).

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

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.

Null The reservation policy has no effect on the number of new or repaired drinking water facilities in the villages.

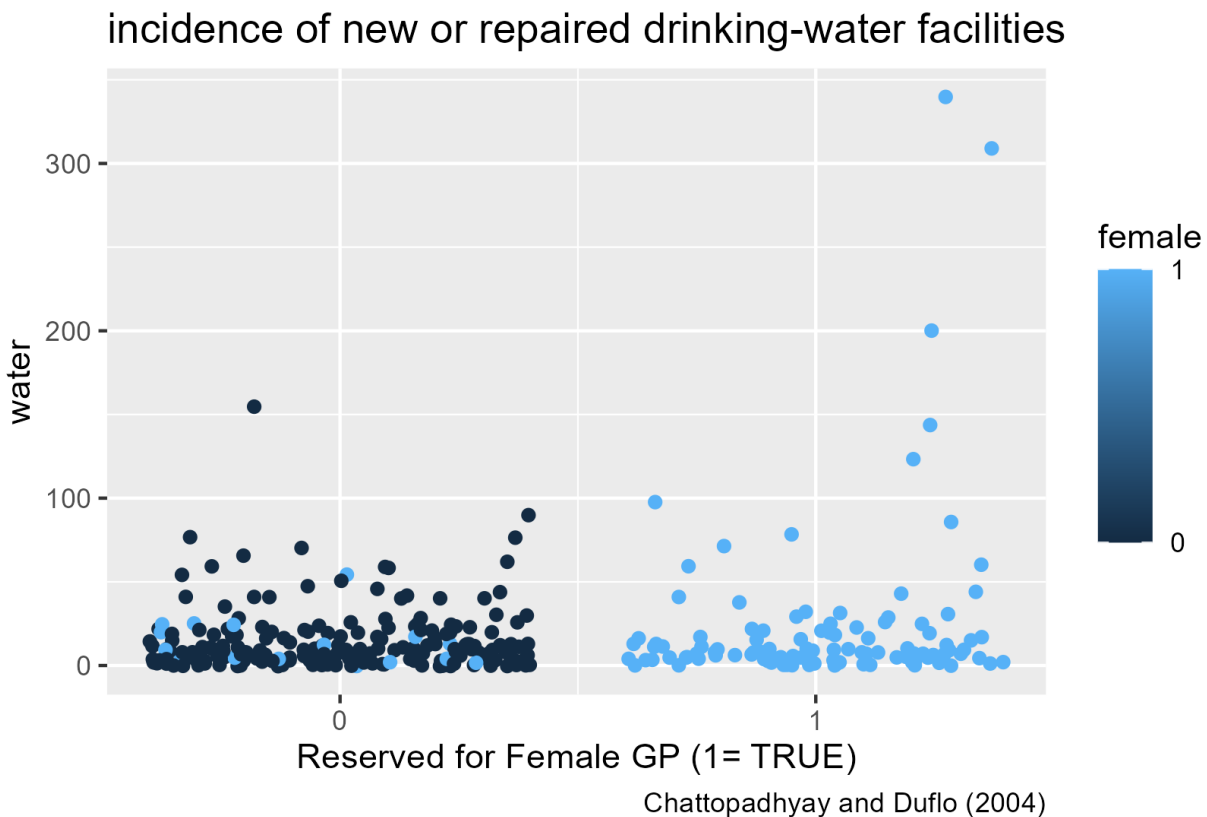
Alternate The reservation policy does have an effect on the number of new or repaired drinking water facilities in the villages.

(b) Bivariate regression to test this hypothesis:.

Import the data.

```
1 policy <- read_csv("https://raw.githubusercontent.com/kosukeimai/qss/master/PREDICTION/women.csv")
```

Figure 3: Drinking water projects, grouped by $reserved = [1, 0]$



The analysis used the builtin R function `lm` to investigate the relationship between the number of new or repaired drinking water facilities in the villages and the binary variable indicating whether the GP was reserved for women leaders or not.

```
1 water <- lm(water ~ reserved , data = policy)
```

This gives the following model results:

Table 2: Pearson Linear Regression - Water Reserved

	water
reserved	9.252** (3.948)
Constant	14.738*** (2.286)
N	322
R ²	0.017
Adjusted R ²	0.014
Residual Std. Error	33.446 (df = 320)
F Statistic	5.493** (df = 1; 320)
*p < .1; **p < .05; ***p < .01	

The estimate for β_0 is 14.738; the estimate for β_1 is 9.252, where $y = \beta_0 + \beta_1 * x$; the response variable (y) is the incidence of investment in drinking water projects; the explanatory variable (x) is 1 if the GP position is reserved for a woman, 0 otherwise. The pvalue is 0.0197, so at a confidence level of 5%, we reject the null hypothesis that the two variables are independent.

(c) Interpret the coefficient estimate for reservation policy.

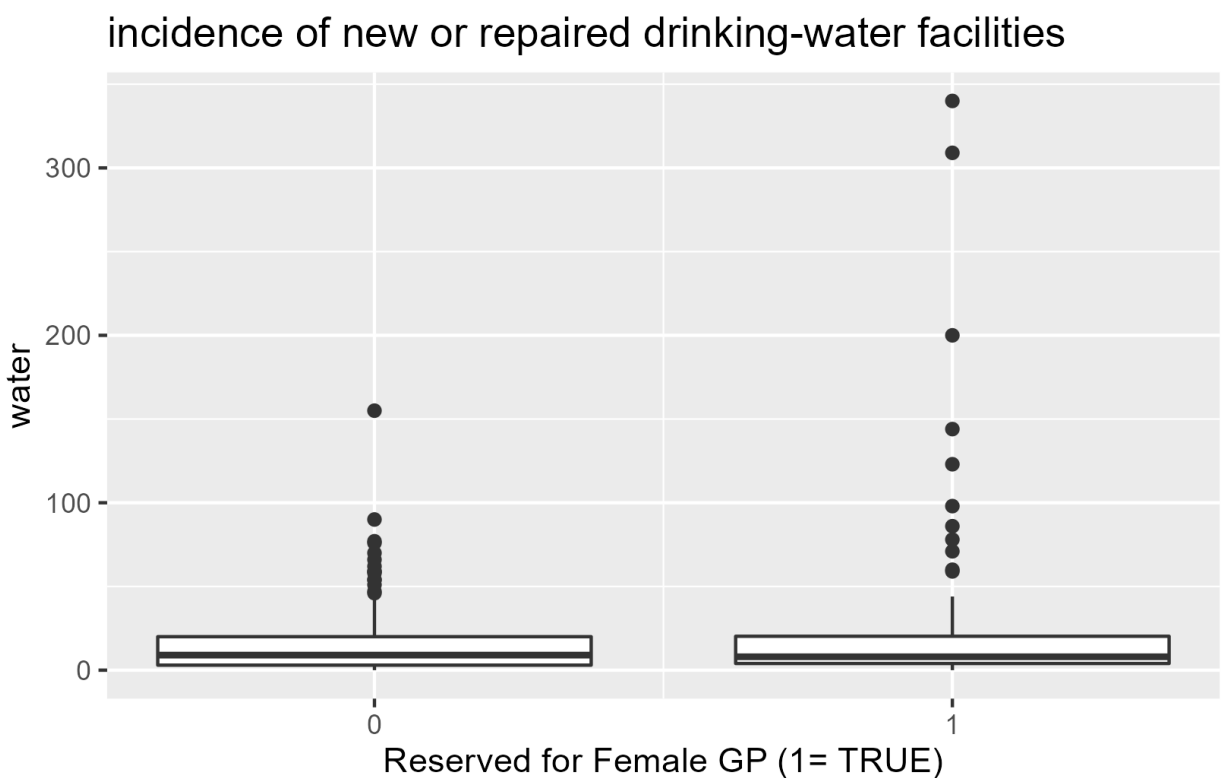
We expect that where the GP position is not reserved for a female, the average number of drinking water projects will be 14.738 and that this will increase by 9.252 if the position is reserved.

Caveats

Outliers

On inspection, it is clear that the data, and the model, are significantly affected by outliers (see Figure 4 and Table 3).

Figure 4: Boxplot of number of drinking water projects, grouped by reserved



Chattopadhyay and Duflo (2004)

Table 3: Outliers in water incidence

reserved	mean_water	count_water	q3	iqr	outlier_limit
0	68.267	15	c('75%' = 20)	17	c('75%' = 45.5)
1	142.545	11	c('75%' = 20.25)	16.25	c('75%' = 44.625)

The data was modelled with outliers excluded and the results were as in Table 4

The estimate for β_0 is 10.7035; the estimate for β_1 is -0.1571 (p-value = 0.9015). Using this data, we cannot reject the hypothesis that water projects and reserved status are independent. The expected number of drinking water projects decreases by 0.1571 if the village is reserved for a female GP.

However, we have no data to support the idea that the outliers are bad data. We are more likely to conclude that the data is heavily skewed.

Table 4: Pearson Linear Regression - Water Reserved - excluding outliers

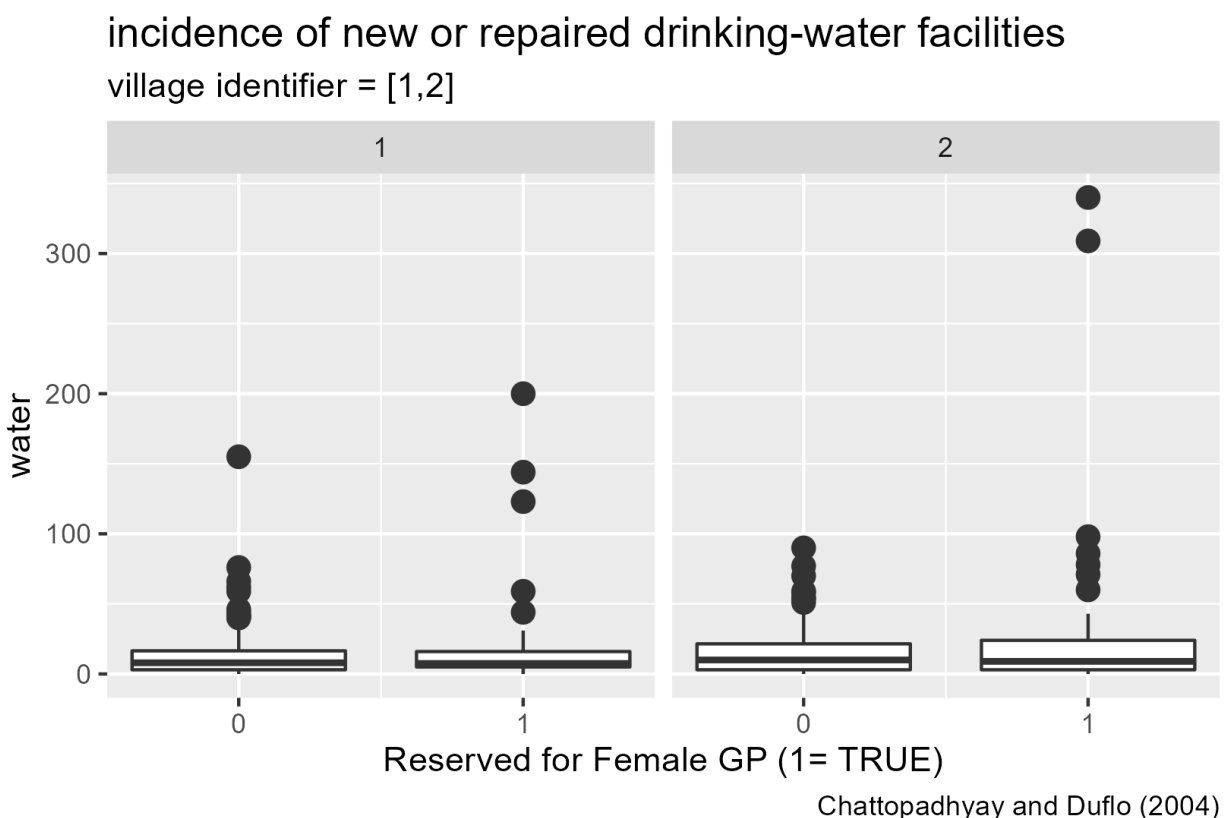
	water
reserved	-0.157 (1.268)
Constant	10.704*** (0.726)
N	296
R ²	0.0001
Adjusted R ²	-0.003
Residual Std. Error	10.243 (df = 294)
F Statistic	0.015 (df = 1; 294)
*p < .1; **p < .05; ***p < .01	

Villages

The assumption in using a linear regression model is that each village is a separate case and each case is independent. However, in this study each GP is associated with two villages, so there is a risk that the values for each village are not independent.

As seen in Figure 5, the profile for the two sets of data has some differences, mainly the extra high values of the outliers in the *village == 2* dataset.

Figure 5: Drinking water projects, grouped by Village



A χ^2 test was run on binned values, and this did not reject the hypothesis that the two samples were from the same population.

Pearson's Chi-squared test

```
data: one_counts and two_counts
X-squared = 12, df = 9, p-value = 0.2133
```

The linear model with the two villages combined (so our units are now GPs, not villages), gives the same expected values, but with lower confidence as we now have fewer data points.

When the two sets of villages are considered separately the estimate for β_0 is 13.907 for $village == 1$ (p-value = 0.2506) and 13.374 for $village == 2$ (pvalue = 0.04172)

This suggests that splitting or combining our data by village does not add greatly to our information about whether *reserved* is a predictor for *water*.

Table 5: Pearson Linear Regression - Water Reserved - Village = 1

	water
reserved	5.130 (4.450)
Constant	13.907*** (2.577)
N	161
R ²	0.008
Adjusted R ²	0.002
Residual Std. Error	26.656 (df = 159)
F Statistic	1.329 (df = 1; 159)
*p < .1; **p < .05; ***p < .01	

Table 6: Pearson Linear Regression - Water Reserved - Village = 2

	water
reserved	13.374** (6.515)
Constant	15.570*** (3.773)
N	161
R ²	0.026
Adjusted R ²	0.020
Residual Std. Error	39.028 (df = 159)
F Statistic	4.215** (df = 1; 159)
*p < .1; **p < .05; ***p < .01	

Appendix - Code

```
1 #####
2 # Imelda Finn, 22334657
3 # POP77003 – Stats I
4 # clear global .envir, load libraries, set wd
5 #####
6
7 # remove objects
8 rm(list=ls())
9
10 # detach all libraries
11 detachAllPackages <- function() {
12   basic.packages <- c("package:stats", "package:graphics", "package:grDevices"
13     , "package:utils", "package:datasets", "package:methods", "package:base")
14   package.list <- search()[ifelse(unlist(gregexpr("package:", search()))==1,
15     TRUE, FALSE)]
16   package.list <- setdiff(package.list, basic.packages)
17   if (length(package.list)>0) for (package in package.list) detach(package,
18     character.only=TRUE)
19 }
20 detachAllPackages()
21
22 # load libraries
23 pkgTest <- function(pkg){
24   new.pkg <- pkg[!(pkg %in% installed.packages()[, "Package"])]
25   if (length(new.pkg))
26     install.packages(new.pkg, dependencies = TRUE)
27   supply(pkg, require, character.only = TRUE)
28 }
29
30 # load necessary packages
31 lapply(c("ggplot2", "stargazer", "tidyverse", "stringr"), pkgTest)
32
33 # function to save output to a file that you can read in later to your docs
34 output_stargazer <- function(outputFile, appendVal=TRUE, ...) {
35   output <- capture.output(stargazer(...))
36   cat(paste(output, collapse = "\n"), "\n", file=outputFile, append=appendVal)
37 }
38
39
40 # set working directory to current parent folder
41 setwd(dirname(rstudioapi::getActiveDocumentContext()$path))
42
43 #####
44 # Problem 1
45 #####
46
47 #Question 1 (40 points): Political Science
48
49 #The following table was created using the data from a study run in a major
50 # Latin American city.
```

```

48 # As part of the experimental treatment in the study, one employee of the
    research
49 # team was chosen to make illegal left turns across traffic to draw the
    attention
50 # of the police officers on shift. Two employee drivers were upper class, two
    were
51 # lower class drivers, and the identity of the driver was randomly assigned
    per
52 # encounter. The researchers were interested in whether officers were more or
    less
53 # likely to solicit a bribe from drivers depending on their class (officers
    use
54 # phrases like, ‘‘We can solve this the easy way’’ to draw a bribe).
55 # The table below shows the resulting data.
56
57
58 #& Not Stopped & Bribe requested & Stopped/given warning \\
59 #Upper class & 14 & 6 & 7 \\
60 #Lower class & 7 & 7 & 1 \\
61 observed <- matrix( c (14, 6, 7, 7, 7, 1), nrow = 2, byrow = TRUE)
62
63 # create data structure with named dimensions
64 cols <- c("NotStopped", "BribeRequested", "StoppedGivenWarning")
65 rows <- c("UpperClass", "LowerClass")
66
67 observed <- matrix( c (14, 6, 7, 7, 7, 1), nrow = 2, byrow = TRUE)
68
69 #plot(jitter(observed))
70
71 #\item [(a)]
72 #Calculate the  $\chi^2$  test statistic by hand/manually\\
73
74 ###----- 0 start listing of code from here
75 ncols <- length(observed[1,])
76 nrows <- length(observed[,1])
77
78 # get totals
79 row_tots <- vector("double", nrows)
80 col_tots <- vector("double", ncols)
81
82 totals <- sum(observed) # total number of observations
83
84 # calculate row and column totals, e.g, total for NotStopped, UpperClass, etc
85 for (i in 1:nrows) {row_tots[i] <- sum(observed[i, ])}
86 for (i in 1:ncols) {col_tots[i] <- sum(observed[, i])}
87
88 #get expected = row total * column total / total observations
89 expected <- observed
90
91 for (i in 1:nrows) {
92   for (j in 1:ncols) {

```

```

93     expected[i,j] <- row_tots[i] * col_tots[j] / totals
94   }
95 }
96
97 # calculate difference between observed and expected
98 o_e <- observed
99 o_e <- (o_e - expected)^2 / expected
100
101 #calculate chi-squared value & degrees of freedom
102 chi_sq_val <- sum(o_e)
103 df = (nrows-1) * (ncols-1)
104
105 cat(str_glue("The chi-squared statistic is {round(chi_sq_val,3)}"))
106 cat(str_glue("The chi-squared degrees of freedom is {df}"))
107
108 # plot of observed and expected values
109 png("graphics/obs-exp.png")
110 barplot(cbind(expected, observed), legend.text = rows,
111         names.arg = c("ns", "br", "sgw", "ns", "br", "sgw"),
112         args.legend = list(x = "topright"),
113         main = "Traffic Stops", beside = TRUE, col = c("green", "red"),
114         xlab = "Observed - Expected")
115 dev.off()
116
117
118 #\item [(b)]
119 #Now calculate the p-value from the test statistic you just created R
120 # .\footnote{Remember frequency should be  $\geq 5$  for all cells, but let's
121 # calculate
122 # the p-value here anyway.} What do you conclude if  $\alpha = 0.1$ ?\\
123
124 p_value <- pchisq(chi_sq_val, df=df, lower.tail=FALSE)
125 alpha <- 0.1
126
127 # p > alpha, can't reject null
128 if (p_value > alpha) txt <- "cannot " else txt <- ""
129
130 # should have min of 5 values in each observed cell
131 cells_under <- length(observed[observed < 5])
132
133 cat(str_glue("The p-value is {round(p_value*100,2)}%, alpha is {alpha*100}%."))
134 cat(str_glue("We {txt}reject the null hypothesis that the two sets are from
135 the\n same population."))
136 cat(str_glue("note: {cells_under} observed cell(s) with less than 5 values."))
137
138 # \item [(c)] Calculate the standardized residuals for each cell and put them
139 # in the table below.
140
141 z <- observed
142 for (i in 1:nrows) {

```

```

140 row_prop<- (1 - (row_tots [i] / totals))
141 for (j in 1:ncols) {
142   col_prop<- (1- (col_tots[j] / totals))
143   z[i,j] <- (observed[i,j] - expected[i,j]) /sqrt (expected[i,j]* row_prop
      * col_prop)
144 }
145 }
146
147 z_df <- data.frame(round(z,3), row.names = rows)
148 names(z_df) <- cols
149
150 print(z_df)
151
152 # output results for Zij values to .tex file
153 output_stargazer(z_df, outputFile="std_residuals.tex", type = "latex",
154                 appendVal=FALSE,
155                 title="Standardised Residuals",
156                 summary = FALSE,
157                 style = "apsr",
158                 table.placement = "htb",
159                 label = "StandardisedResiduals",
160                 rownames = TRUE
161                 )
162
163
164 # check result
165 chisq.test(observed)
166 # Pearson's Chi-squared test
167
168 #data:  observed
169 #X-squared = 3.7912, df = 2, p-value = 0.1502
170
171 # \item [(d)] How might the standardized residuals help you interpret the
      results?
172
173 # fewer upper class individuals asked for bribes and more given warnings;
174 # the contribution from lower class drivers expected to give bribes is nearly
175 # equivalent to the contribution from upper class drivers getting warnings
176
177 #
      #####

178 # Problem 2
179 #####
180
181 #Question 2 (40 points): Economics
182 #Chattopadhyay and Duflo were interested in whether women promote different
      policies
183 # than men.
184 # Answering this question with observational data is pretty difficult due to
      potential

```

```

185 # confounding problems (e.g. the districts that choose female politicians are
186 # likely to systematically differ in other aspects too). Hence, they exploit a
187 # randomized policy experiment in India, where since the mid-1990s, 1/3 of
188 # village council heads have been randomly reserved for women. A subset of the
    data
189 # from West Bengal can be found at the following link:
190 #   \url{https://raw.githubusercontent.com/kosukeimai/qss/master/PREDICTION/
    women.csv}
191
192 # Each observation in the data set represents a village and there are two
    villages
193 # associated with one GP (i.e. a level of government is called "GP").
194 # Figure~\ref{fig:women_desc} below shows the names and descriptions of the
    variables
195 # in the dataset. The authors hypothesize that female politicians are more
    likely to
196 # support policies female voters want. Researchers found that more women
    complain about
197 # the quality of drinking water than men. You need to estimate the effect of
    the
198 # reservation policy on the number of new or repaired drinking water
    facilities
199 #in the villages.
200 # Names and description of variables from Chattopadhyay and Duflo (2004)
201 # 1 'GP' Identifier for the Gram Panchayat &nbsp;&nbsp; 
202 # 2 'village' identifier for each village
203 # 3 'reserved' binary variable indicating whether the GP was reserved for
    women leaders or not
204 # 4 'female' binary variable indicating whether the GP had a female leader or
    not
205 # 5 'irrigation' variable measuring the number of new or repaired irrigation
    facilities in the village since the reserve policy started
206 # 6 'water' variable measuring the number of new or repaired drinking-water
    facilities in the village since the reserve policy started
207
208 #\item [(a)] State a null and alternative (two-tailed) hypothesis.
209 # null: no diff in incidence of new or repaired drinking-water facilities
210 # in the village since the reserve policy started
211 # ie 'water' is independent of 'reserved'
212 # alternate: the incidence of new or repaired drinking-water facilities is
213 # correlated to the reservation policy
214
215
216 policy <- read_csv("https://raw.githubusercontent.com/kosukeimai/qss/master/
    PREDICTION/women.csv")
217 #write.csv(policy,"Data/policy.csv")
218 policy<-read_csv("Data/policy.csv")
219
220 summary(policy)
221
222 plot(policy$water)

```

```

223 boxplot(policy$water)
224 # lots of outliers, distribution is skewed right (mean > median)
225
226 plot(policy$water, policy$irrigation)
227
228 pairs(policy[4:7])
229
230 sum(policy$reserved)    # 108 of 322 villages have reserved GP (54 GPs)
231 sum(policy$female)     # 124 of 322 villages have female GP (62 GPs)
232
233 #\item [(b)] Run a bivariate regression to test this hypothesis
234
235 water <- lm(water ~ reserved , data = policy)
236 summary(water)
237
238 output_stargazer(water, outputFile="water_model.tex", type = "latex",
239                  appendVal=FALSE,
240                  title="Pearson Linear Regression – Water ~ Reserved",
241                  style = "apsr",
242                  table.placement = "htb",
243                  label = "model:water_reserved"
244 )
245
246 #-----
247
248 cor(policy$water, policy$reserved) # .1299
249 # water increases with increase in reserved (ie reserved = TRUE), not strong
250
251 p<- ggplot(policy, aes(reserved, water, colour=female, group_by(female)))
252 p + geom_jitter() +
253   scale_x_continuous(breaks = seq(0, 1, by = 1)) +
254   scale_color_continuous(breaks = seq(0, 1, by = 1)) +
255   labs(title = "incidence of new or repaired drinking–water facilities",
256        x = "Reserved for Female GP (1= TRUE)",
257        caption = "Chattopadhyay and Duflo (2004)",
258        alt = "Boxplot of incidence of new or repaired drinking–water
259              facilities, by reserved [1,0]",
260 )
261 ggsave("graphics/resrwd_water.png")
262
263 ##=====
264 # consider outliers
265
266 p<- ggplot(policy, aes(reserved, water, group_by(reserved)))
267 p + geom_boxplot(aes(group=reserved)) +
268   scale_x_continuous(breaks = seq(0, 1, by = 1)) +
269   labs(title = "incidence of new or repaired drinking–water facilities",
270        x = "Reserved for Female GP (1= TRUE)",
271        caption = "Chattopadhyay and Duflo (2004)",
272        alt = "Boxplot of incidence of new or repaired drinking–water facilities,
273              by reserved [1,0]",

```

```

272 )
273 ggsave("graphics/resrwd_water_boxplot.png")
274
275 outliers_tbl <- policy %>%
276   group_by(reserved) %>%
277   mutate(iqr = IQR(water), q3 = quantile(water, .75), outlier_limit = q3 + iqr
278     * 1.5 ) %>%
279   filter(water > outlier_limit ) %>%
280   mutate(mean_water = round(mean(water),3), count_water = n()) %>%
281   select(reserved, mean_water, count_water, q3, iqr, outlier_limit) %>%
282   unique()
283
284 #reserved mean_water count_water q3 iqr outlier_limit
285 #<dbl> <dbl> <int> <dbl> <dbl> <dbl>
286 # 1 0 68.3 15 20 17 45.5
287 # 2 1 143. 11 20.2 16.2 44.6
288
289 output_stargazer(outliers_tbl, outputFile="water_outliers.tex", type = "latex"
290   ,
291     appendVal=FALSE,
292     title="Outliers in water incidence",
293     summary = FALSE,
294     style = "apsr",
295     digits= 3,
296     table.placement = "htb",
297     label = "tab:wateroutliers",
298     rownames = FALSE
299 )
300
301 # there are fewer outliers in the reserved=1 cohort, but their average
302 # value is significantly higher
303
304 no_outlier_water <- policy %>%
305   group_by(reserved) %>%
306   mutate(outlier_limit = quantile(water, .75) + IQR(water) * 1.5) %>%
307   ungroup() %>%
308   filter(water <= outlier_limit)
309
310 outlier_model <- lm(water ~ reserved , data = no_outlier_water)
311
312 output_stargazer(outlier_model, outputFile="outlier_model.tex", type = "latex"
313   ,
314     appendVal=FALSE,
315     title="Pearson Linear Regression – Water ~ Reserved –
316     excluding outliers",
317     style = "apsr",
318     table.placement = "htb",
319     label = "tab:noOutliers"
320 )

```



```

319
320 summary(outlier_model)
321
322
323 # coefficient for beta0 goes to -0.1571 - with no significance
324 # (p-value is 0.9015, df= 294)
325 # same result if exclude sample outliers (ie not by reserved)
326 ###=====
327 # assumption is that each village is a separate case and each case is
    independent
328 # but, each GP relates to 2 villages - need to check for impact of combining
    villages
329
330 # inspect data
331 p<- ggplot(policy, aes(reserved, water, group_by(reserved)))
332 p + geom_boxplot(outlier.size = 3, aes(group=reserved)) +
333   scale_x_continuous(breaks = seq(0, 1, by = 1)) +
334   labs(title = "incidence of new or repaired drinking-water facilities",
335        subtitle = "village identifier = [1,2]",
336        x = "Reserved for Female GP (1= TRUE)",
337        caption = "Chattopadhyay and Duflo (2004)",
338        alt = "Boxplot of incidence of new or repaired drinking-water
    facilities, by reserved [1,0]",
339   ) +
340   facet_wrap(policy$village)
341 ggsave("graphics/village_water_boxplot.png")
342
343 reserved_water_tab <- policy %>%
344   group_by(reserved) %>%
345   summarise(n = n(), sum_water = sum(water)) %>%
346   mutate(prop_reserved = round(n / sum(n), 4), sum_water, prop_water_reserved
    =
347     round(sum_water / sum(sum_water), 4)) %>% # mutate after our
    summarise to find the proportion
348   arrange(desc(prop_reserved))
349
350 str(reserved_water_tab)
351 reserved_water_tab
352
353 sum(policy$water)
354
355 # see if villages are from same population
356 one_village_policy <- policy %>%
357   group_by(GP) %>%
358   filter(village ==1)
359
360 two_village_policy <- policy %>%
361   group_by(GP) %>%
362   filter(village ==2)
363
364

```

```

365 hist(one_village_policy$water)$counts
366 #[1] 130 16 8 3 0 0 1 2 0 1
367 hist(two_village_policy$water)$counts
368 #[1] 146 13 0 0 0 0 2
369 hist(one_village_policy$water)$breaks
370 #[1] 0 20 40 60 80 100 120 140 160 180 200
371
372 # coerce counts of water variable into suitably sized bins
373 one_counts <- hist(one_village_policy$water, breaks = c(0, 20, 40, 60, 350))$
  counts
374 two_counts <- hist(two_village_policy$water, breaks = c(0, 20, 40, 60, 350))$
  counts
375 # run chisq test - null: both from same population
376
377 chi_village <- chisq.test(one_counts, two_counts)
378
379
380 villagetab <- matrix(c(one_counts, two_counts), nrow = 2, byrow = TRUE)
381 chi_village
382
383 output_stargazer(tibble(villagetab), outputFile="village_bins.tex", type = "
  latex",
384
  appendVal=FALSE,
385 title="Binned data for village dataset comparison",
386 summary = FALSE,
387 style = "apsr",
388 table.placement = "htb",
389 label = "tab:villageBins",
390 rownames = TRUE
391 )
392
393
394 # Pearson's Chi-squared test
395
396 #data: one_counts and two_counts
397 #X-squared = 12, df = 9, p-value = 0.2133
398
399
400 #tibble('village1' = one_counts, 'village2' = two_counts )
401
402 # run regression model on each set of villages
403 one_model <- lm(water ~ reserved, data = one_village_policy)
404 two_model <- lm(water ~ reserved, data = two_village_policy)
405
406 summary(one_model)
407 summary(two_model)
408
409 output_stargazer(one_model, outputFile="village_model.tex", type = "latex",
410
  appendVal=FALSE,
411 title="Pearson Linear Regression - Water ~ Reserved - Village
  = 1",

```

```

412         style = "apsr",
413         table.placement = "htb",
414         label = "tab:village1"
415     )
416
417 output_stargazer(two_model, outputFile="village_model.tex", type = "latex",
418                 appendVal=TRUE,
419                 title="Pearson Linear Regression – Water ~ Reserved – Village
420                     = 2",
421                 style = "apsr",
422                 table.placement = "htb",
423                 label = "tab:village2"
424             )
425
426 # or combine the villages
427
428 combined_village_policy <- policy %>%
429   group_by(GP) %>%
430   mutate (sum_water = sum(water), sum_irrigation = sum(irrigation)) %>%
431   select(GP, reserved, female, sum_water, sum_irrigation) %>%
432   unique()
433
434 # run model – scaled by 1/2 to get equivalent values to 1 village coefficients
435 cvp <- lm(sum_water/2 ~ reserved, data = combined_village_policy)
436
437 summary(cvp)
438
439 output_stargazer(cvp, outputFile="villages_combined.tex", type = "latex",
440                 appendVal=FALSE,
441                 title="Pearson Linear Regression – Water ~ Reserved –
442                     Villages combined",
443                 style = "apsr",
444                 table.placement = "htb",
445                 label = "tab:combinedVillages"
446             )

```