

# 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.<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.

	Not Stopped	Bribe requested	Stopped/given warning
Upper class	14	6	7
Lower class	7	7	1

(a) The  $\chi^2$  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  $\chi^2$  statistic.

**expected** number in class \* number of outcomes / total number

**difference** observed - expected

---

<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.

**contribution**  $\text{difference}^2/\text{expected}$ )

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) The p-value from the test statistic is calculated as follows. If  $\alpha = 0.1$ , we cannot reject the null hypothesis that both subgroups are from the same population (i.e. the difference in experience is not statistically significant).

```
1 p_value <- pchisq(chi_sq_val, df=df, lower.tail=FALSE)
```

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

(c) The standardized residuals are set out in Table 1:

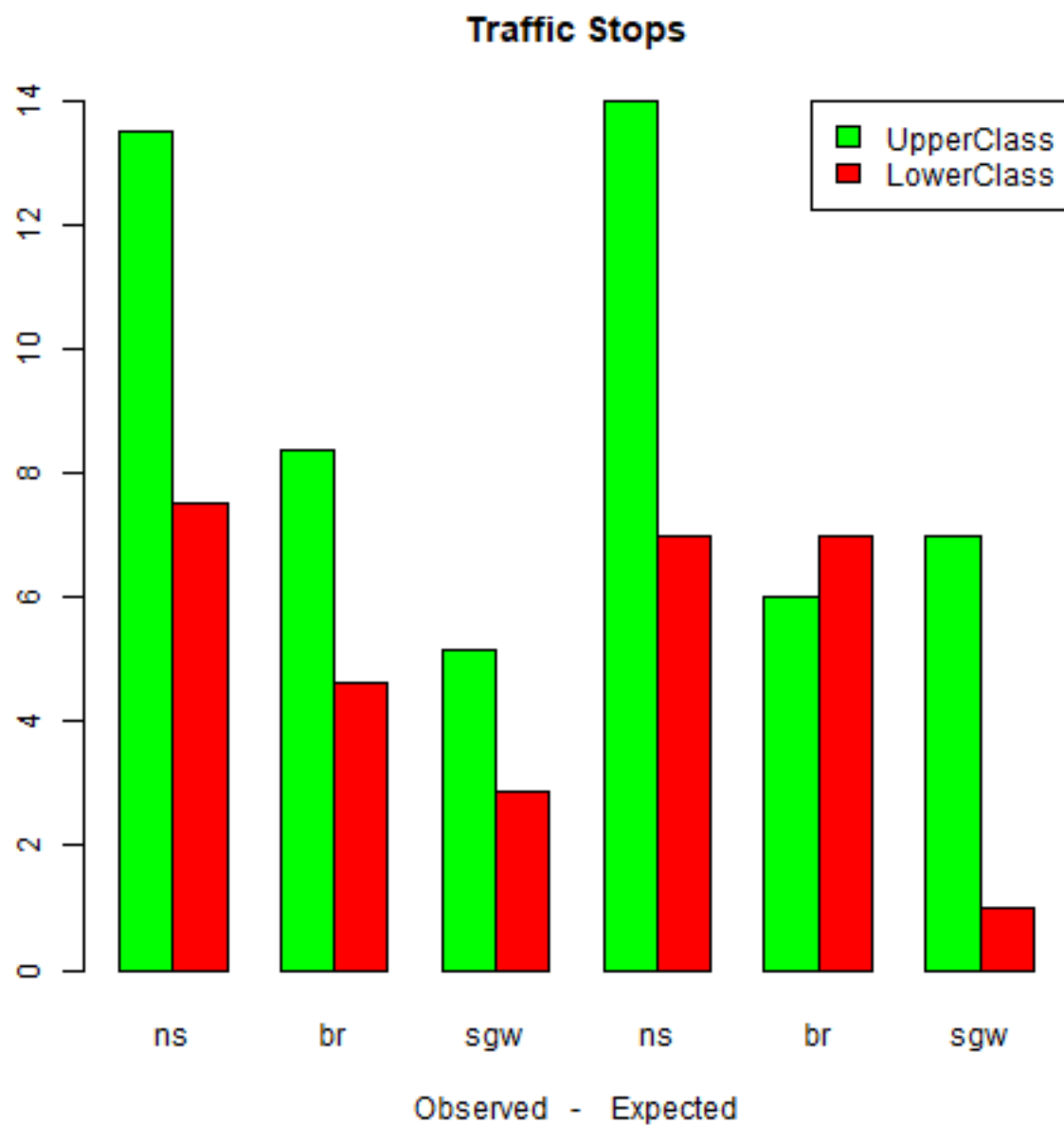


Figure 1: Observed vs Expected values for traffic stop. ns = Not Stopped; br = Bribe Requested; sgw = Stopped Given Warning

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.<sup>2</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 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
<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>2</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** The reservation policy has no effect on the number of new or repaired drinking water facilities in the villages.

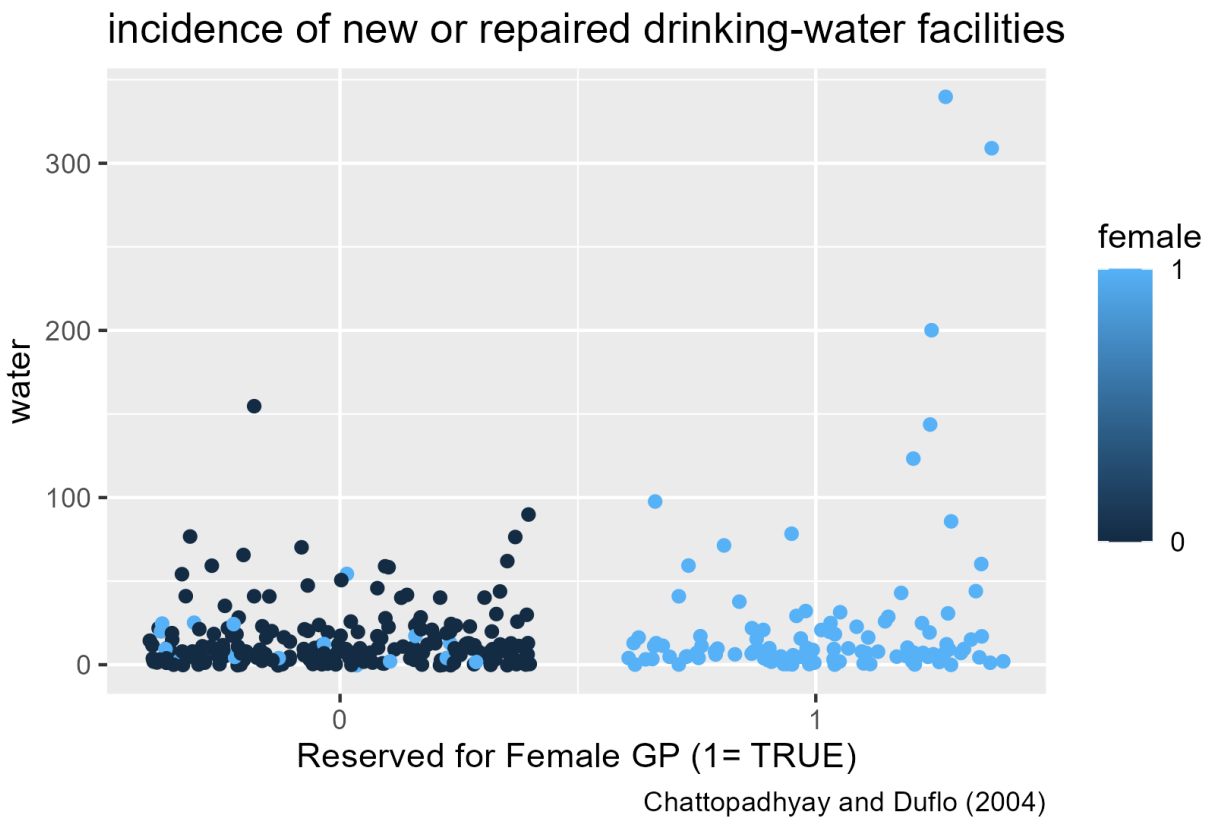
**Alternate** The reservation policy does have an effect on the number of new or repaired drinking water facilities in the villages.

(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 relationship between the

$x = \text{reserved}$  binary variable indicating whether the GP was reserved for women leaders or not



**y = water** numeric variable denoting the number of new or repaired drinking water facilities in the villages

was modelled using the Pearson model for linear regression. The code is as follows:

```
1 #\item [(b)] Run a bivariate regression to test this hypothesis
2 water <- policy$water           # ie y = response var
3 reserved <- policy$reserved      # ie x = explanatory var
4
5 mean_water <- mean(water)
6 mean_reserved <- mean(reserved)
7
8 n <- length(water)
9
10 # calculate sum of squares for reserved and water
11 ssxx <- sum((reserved - mean_reserved)^2)
12 ssyy <- sum((water - mean_water)^2)
13 ssxy <- sum((reserved - mean_reserved)*(water - mean_water) )
14 # calculate covariance
15 covxy <- ssxy / n
16 # check result
17 cov(x = reserved, y = water, method = "pearson")
18 #calculate correlation coefficient
19 corxy <- ssxy / sqrt(ssxx * ssyy)
20
21 #calculate estimates for coefficients
22 beta1 <- ssxy / ssxx
23 beta0 <- mean_water - mean_reserved*beta1
24
25 # calculate standard error values
26 sse <- sum((water-(beta0 + beta1*reserved))^2)
27 se <- sqrt(sse / (n-2))
28
29 # calculate standard errors for coefficients
30 s_beta1 <- se * sqrt(1/ssxx)
31 s_beta0 <- se * sqrt((1/n + mean_reserved ^2 / ssxx))
32
33 # calculate the t-test statistics for coefficients
34 t_beta1 <- beta1 / s_beta1
35 t_beta0 <- beta0 / s_beta0
36
37 # calculate r^2 and p values
38 r2 <- 1 - (sse / ssyy)
39 p_beta1 <- 2*pt(t_beta1, df=n-2, lower.tail = FALSE)
40 p_beta0 <- 2*pt(t_beta0, df=n-2, lower.tail = FALSE)
```

This gives the following results:

Table 2: coefficients for linear regression model water - reserved

	estimate	Std Error	t value	pr(> t )
intercept	14.738	2.286	6.446	$4.216474e - 10$
reserved	9.252	3.948	2.344	$1.970398e - 02$

Table 3: results for linear regression model water - reserved

residual error	degrees of freedom	$R^2$	covariance	correlation
33.4457	320	0.0169	2.0624	0.1299

The estimate for  $\beta_0$  is 14.738; the estimate for  $\beta_1$  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. The  $R^2$  value suggests that our model accounts for less than 2% of the variance in our water values.

- (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

Code in Appendix.

## Outliers

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

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

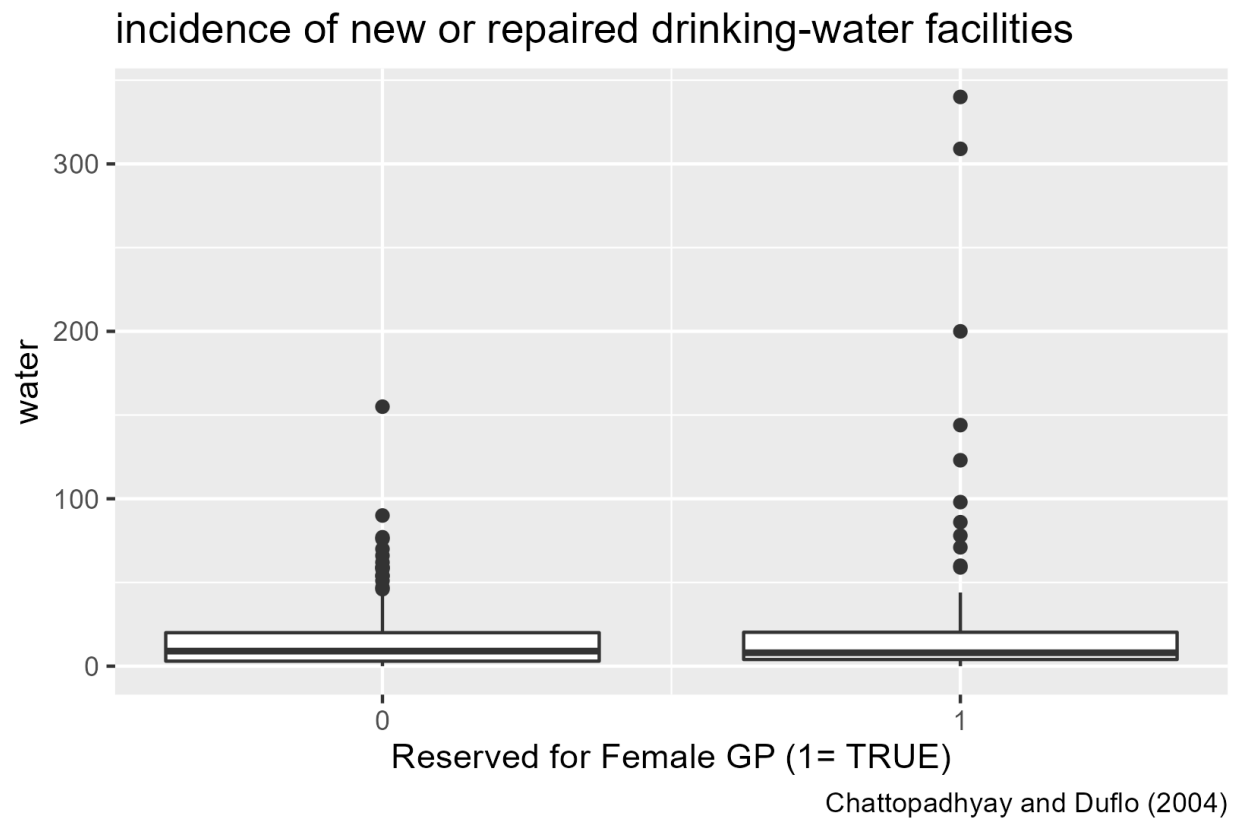


Table 4: Outliers in water incidence

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 5

The estimate for  $\beta_0$  is 10.7035; the estimate for  $\beta_1$  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 5: Pearson Linear Regression - Water Reserved - excluding outliers

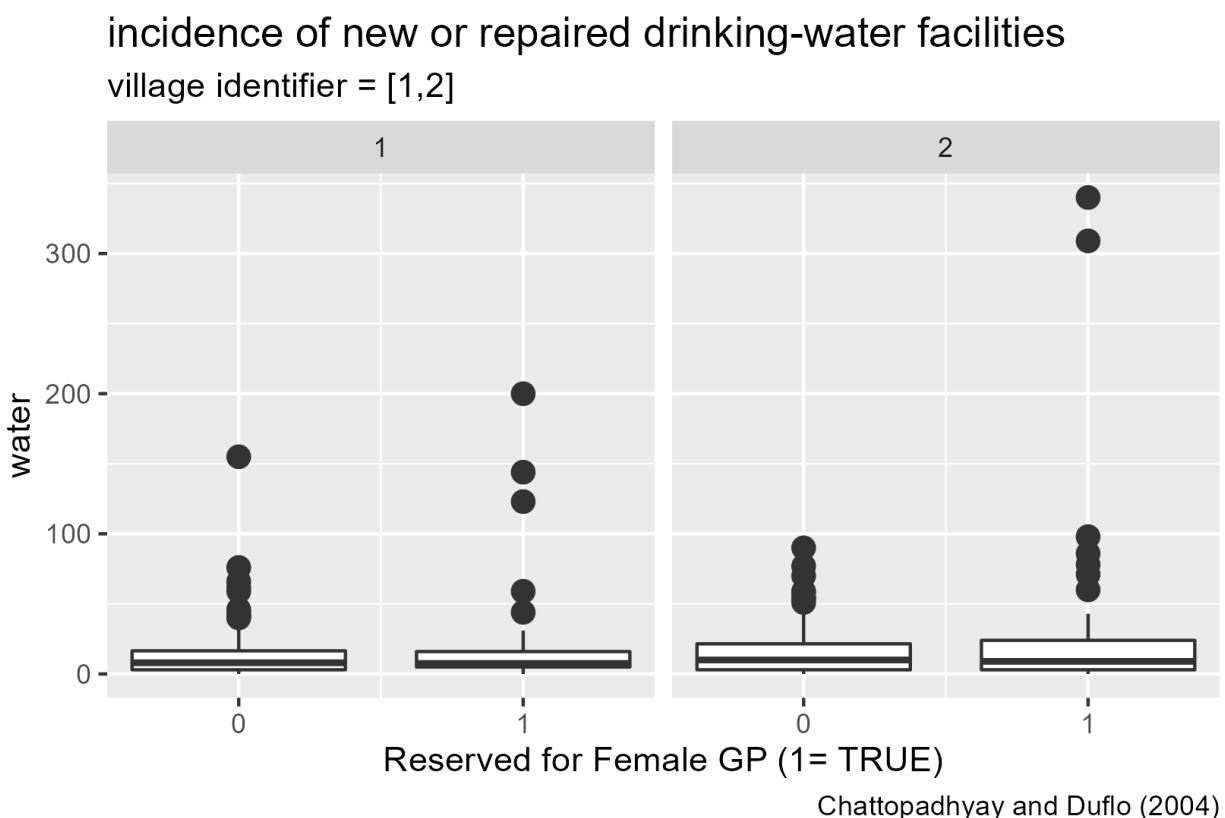
	water
reserved	-0.157 (1.268)
Constant	10.704*** (0.726)
N	296
R <sup>2</sup>	0.0001
Adjusted R <sup>2</sup>	-0.003
Residual Std. Error	10.243 (df = 294)
F Statistic	0.015 (df = 1; 294)
*p < .1; **p < .05; ***p < .01	

## 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  $\chi^2$  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  $\beta_1$  is 5.130 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 6: Pearson Linear Regression - Water Reserved - Village = 1

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

Table 7: Pearson Linear Regression - Water Reserved - Village = 2

	water
reserved	13.374** (6.515)
Constant	15.570*** (3.773)
N	161
R <sup>2</sup>	0.026
Adjusted R <sup>2</sup>	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 #\item [(a)]
68 #Calculate the  $\chi^2$  test statistic by hand/manually\\
69
70 ###----- 0 start listing of code from here
71 ncols <- length(observed[1,])
72 nrows <- length(observed[,1])
73
74 # get totals
75 row_tots <- vector("double", nrows)
76 col_tots <- vector("double", ncols)
77
78 totals <- sum(observed) # total number of observations
79
80 # calculate row and column totals, e.g, total for NotStopped, UpperClass, etc
81 for (i in 1:nrows) {row_tots[i] <- sum(observed[i, ])}
82 for (i in 1:ncols) {col_tots[i] <- sum(observed[, i])}
83
84 #get expected = row total * column total / total observations
85 expected <- observed
86
87 for (i in 1:nrows) {
88   for (j in 1:ncols) {
89     expected[i,j] <- row_tots[i] * col_tots[j] / totals
90   }
91 }
92

```

```

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

```

```

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

174 # Problem 2
175 #####
176
177 #Question 2 (40 points): Economics
178 #Chattopadhyay and Duflo were interested in whether women promote different
    policies
179 # than men.
180 # Answering this question with observational data is pretty difficult due to
    potential
181 # confounding problems (e.g. the districts that choose female politicians are
182 # likely to systematically differ in other aspects too). Hence, they exploit a
183 # randomized policy experiment in India, where since the mid-1990s, 1/3 of
184 # village council heads have been randomly reserved for women. A subset of the

```

```

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

```

```

223
224 pairs(policy[4:7])
225
226 sum(policy$reserved)    # 108 of 322 villages have reserved GP (54 GPs)
227 sum(policy$female)     # 124 of 322 villages have female GP (62 GPs)
228
229 #\item [(b)] Run a bivariate regression to test this hypothesis
230 water <- policy$water           # ie y = response var
231 reserved <- policy$reserved     # ie x = explanatory var
232
233 mean_water <- mean(water)
234 mean_reserved <- mean(reserved)
235
236 n <- length(water)
237
238 # calculate sum of squares for reserved and water
239 sxxx <- sum((reserved - mean_reserved)^2)
240 ssyy <- sum((water - mean_water)^2)
241 ssxy <- sum((reserved - mean_reserved)*(water - mean_water) )
242 # calculate covariance
243 covxy <- ssxy / n
244 # check result
245 cov(x = reserved, y = water, method = "pearson")
246 #calculate correlation coefficient
247 corxy <- ssxy / sqrt(sxxx * ssyy)
248
249 #calculate estimates for coefficients
250 betal <- ssxy / sxxx
251 beta0 <- mean_water - mean_reserved*betal
252
253 # calculate standard error values
254 sse <- sum((water-(beta0 + betal*reserved))^2)
255 se <- sqrt(sse / (n-2))
256
257 # calculate standard errors for coefficients
258 s_betal <- se * sqrt(1/sxxx)
259 s_beta0 <- se * sqrt((1/n + mean_reserved ^2 / sxxx))
260
261 # calculate the t-test statistics for coefficients
262 t_betal <- betal / s_betal
263 t_beta0 <- beta0 / s_beta0
264
265 # calculate r^2 and p values
266 r2 <- 1 - (sse / ssyy)
267 p_betal <- 2*pt(t_betal, df=n-2, lower.tail = FALSE)
268 p_beta0 <- 2*pt(t_beta0, df=n-2, lower.tail = FALSE)
269
270 # output results as two tables
271 cols <- c("estimate ", "Std Error", "t value", "pr(>|t|)")
272 rows <- c("intercept", "reserved")
273

```



```

274 beta_vals <- data.frame(matrix(c(round(beta0, 3), round(s_beta0, 3),
275                                round(t_beta0, 3), p_beta0,
276                                round(beta1, 3), round(s_beta1, 3), round(t_beta1, 3), p-
                                beta1),
                                nrow = 2, byrow = TRUE), row.names = rows)
277
278
279 names(beta_vals) <- cols
280
281 print(beta_vals)
282
283 # output results for beta values to .tex file
284 output_stargazer(beta_vals, outputFile="policy_model.tex", type = "latex",
285                  appendVal=FALSE,
286                  title="coefficients for linear regression model water -
                reserved ",
287                  summary = FALSE,
288                  style = "apsr",
289                  table.placement = "htbp!",
290                  label = "tab:coefficients",
291                  rownames = TRUE
292 )
293
294 result_vals <- tibble('residual error' = round(se, 4), 'degrees of freedom' = n
                -2,
295                      'R^2' = round(r2, 4), 'covariance' = round(covxy, 4),
296                      'correlation' = round(corxy, 4))
297
298
299
300 output_stargazer(result_vals, outputFile="policy_model.tex", type = "latex",
301                  appendVal=TRUE,
302                  title="results for linear regression model water - reserved "
                ,
303                  summary = FALSE,
304                  style = "apsr",
305                  label = "tab:results",
306                  rownames = FALSE
307 )
308
309 result_cols <- tibble(round(se, 4), n-2, round(r2, 4),
310                      round(covxy, 4), round(corxy, 4))
311 names(result_cols) <- c("residual error", "degrees of freedom", "R^2",
312                      "covariance", "correlation")
313
314 #check r^2
315 r<-cov(reserved, water) / (sd(reserved)* sd(water))
316
317 #check correlation coefficient
318 cor(policy$water, policy$reserved) # .1299
319 # water increases with increase in reserved (ie reserved = TRUE), not strong
320

```

```

321
322 output_stargazer(water_model, outputFile="water_model.tex", type = "latex",
323                   appendVal=FALSE,
324                   title="Pearson Linear Regression – Water ~ Reserved",
325                   style = "apsr",
326                   table.placement = "htbp!",
327                   label = "model:water_reserved"
328 )
329
330
331
332 water_model <- lm(water ~ reserved , data = policy)
333 summary(water_model)
334 #Residuals:
335 #Min      1Q  Median      3Q      Max
336 #-23.991 -14.738  -7.865   2.262 316.009
337
338 #Coefficients:
339 #   Estimate Std. Error t value Pr(>|t|)
340 #(Intercept)  14.738      2.286   6.446 4.22e-10 ***
341 # reserved      9.252      3.948   2.344 0.0197 *
342 # ———
343 # Signif. codes:  0      ***      0.001      **      0.01      *      0.05      .      0.1
344                    1
345
346 #Residual standard error: 33.45 on 320 degrees of freedom
347 #Multiple R-squared:  0.01688, Adjusted R-squared:  0.0138
348 #F-statistic: 5.493 on 1 and 320 DF, p-value: 0.0197
349
350
351 output_stargazer(water_model, outputFile="water_lm.tex", type = "latex",
352                   appendVal=FALSE,
353                   title="Pearson Linear Regression – Water ~ Reserved",
354                   style = "apsr",
355                   table.placement = "htbp!",
356                   label = "model:water_reserved"
357 )
358
359 #—————
360
361
362
363
364 p<- ggplot(policy , aes(reserved , water , colour=female , group_by(female)))
365 p + geom_jitter() +
366   scale_x_continuous(breaks = seq(0, 1, by = 1)) +
367   scale_color_continuous(breaks = seq(0, 1, by = 1)) +
368   labs(title = "incidence of new or repaired drinking–water facilities",
369        x = "Reserved for Female GP (1= TRUE)",
370        caption = "Chattopadhyay and Duflo (2004)",

```

```

371     alt = "Boxplot of incidence of new or repaired drinking-water
372     facilities , by reserved [1,0]",
373   )
374   ggsave("graphics/resrwd-water.png")
375   ##
376   # consider outliers
377
378   p<- ggplot(policy , aes(reserved , water , group_by(reserved)))
379   p + geom_boxplot( aes(group=reserved)) +
380     scale_x_continuous(breaks = seq(0, 1, by = 1)) +
381     labs(title = "incidence of new or repaired drinking-water facilities",
382          x = "Reserved for Female GP (1= TRUE)",
383          caption = "Chattopadhyay and Duflo (2004)",
384          alt = "Boxplot of incidence of new or repaired drinking-water facilities ,
385          by reserved [1,0]",
386        )
387   ggsave("graphics/resrwd-water-boxplot.png")
388
389   outliers_tbl <- policy %>%
390     group_by(reserved) %>%
391     mutate(iqr = IQR(water), q3 = quantile(water, .75), outlier_limit = q3 + iqr
392           * 1.5 ) %>%
393     filter(water > outlier_limit ) %>%
394     mutate(mean_water = round(mean(water),3), count_water = n()) %>%
395     select(reserved , mean_water , count_water , q3 , iqr , outlier_limit) %>%
396     unique()
397
398   #reserved mean_water count_water q3 iqr outlier_limit
399   #<dbl> <dbl> <int> <dbl> <dbl> <dbl>
400   # 1 0 68.3 15 20 17 45.5
401   # 2 1 143. 11 20.2 16.2 44.6
402
403   output_stargazer(outliers_tbl , outputFile="water_outliers.tex" , type = "latex"
404     ,
405     appendVal=FALSE,
406     title="Outliers in water incidence",
407     summary = FALSE,
408     style = "apsr",
409     digits= 3,
410     table.placement = "htbp!",
411     label = "tab:wateroutliers",
412     rownames = FALSE
413   )
414
415   # there are fewer outliers in the reserved=1 cohort , but their average
416   # value is significantly higher
417
418   no_outlier_water <- policy %>%

```

```

418 group_by(reserved) %>%
419 mutate(outlier_limit = quantile(water, .75) + IQR(water) * 1.5) %>%
420 ungroup() %>%
421 filter(water <= outlier_limit)
422
423 outlier_model <- lm(water ~ reserved, data = no_outlier_water)
424
425 output_stargazer(outlier_model, outputFile="outlier_model.tex", type = "latex"
426 ,
427               appendVal=FALSE,
428               title="Pearson Linear Regression - Water ~ Reserved -
429               excluding outliers",
430               style = "apsr",
431               table.placement = "htbp!",
432               label = "tab:noOutliers"
433 )
434
435 summary(outlier_model)
436
437 # coefficient for beta0 goes to -0.1571 - with no significance
438 # (p-value is 0.9015, df= 294)
439 # same result if exclude sample outliers (ie not by reserved)
440 # =====
441 # assumption is that each village is a separate case and each case is
442 # independent
443 # but, each GP relates to 2 villages - need to check for impact of combining
444 # villages
445
446 # inspect data
447 p <- ggplot(policy, aes(reserved, water, group_by(reserved)))
448 p + geom_boxplot(outlier.size = 3, aes(group=reserved)) +
449   scale_x_continuous(breaks = seq(0, 1, by = 1)) +
450   labs(title = "incidence of new or repaired drinking-water facilities",
451         subtitle = "village identifier = [1,2]",
452         x = "Reserved for Female GP (1= TRUE)",
453         caption = "Chattopadhyay and Duflo (2004)",
454         alt = "Boxplot of incidence of new or repaired drinking-water
455               facilities, by reserved [1,0]",
456         ) +
457   facet_wrap(policy$village)
458 ggsave("graphics/village_water_boxplot.png")
459
460 reserved_water_tab <- policy %>%
461   group_by(reserved) %>%
462   summarise(n = n(), sum_water = sum(water)) %>%
463   mutate(prop_reserved = round(n / sum(n), 4), sum_water, prop_water_reserved
464         =
465         round(sum_water / sum(sum_water), 4)) %>% # mutate after our
466   summarise to find the proportion
467   arrange(desc(prop_reserved))

```

```

462
463 str(reserved_water_tab)
464 reserved_water_tab
465
466 sum(policy$water)
467
468 # see if villages are from same population
469 one_village_policy <- policy %>%
470   group_by(GP) %>%
471   filter(village ==1)
472
473 two_village_policy <- policy %>%
474   group_by(GP) %>%
475   filter(village ==2)
476
477
478 hist(one_village_policy$water)$counts
479 #[1] 130 16 8 3 0 0 1 2 0 1
480 hist(two_village_policy$water)$counts
481 #[1] 146 13 0 0 0 0 2
482 hist(one_village_policy$water)$breaks
483 #[1] 0 20 40 60 80 100 120 140 160 180 200
484
485 # coerce counts of water variable into suitably sized bins
486 one_counts <- hist(one_village_policy$water, breaks = c(0, 20, 40, 60, 350))$
  counts
487 two_counts <- hist(two_village_policy$water, breaks = c(0, 20, 40, 60, 350))$
  counts
488 # run chisq test - null: both from same population
489
490 chi_village <- chisq.test(one_counts, two_counts)
491
492
493 villagetab <- matrix(c(one_counts, two_counts), nrow = 2, byrow = TRUE)
494 chi_village
495
496 output_stargazer(tibble(villagetab), outputFile="village_bins.tex", type = "
  latex",
497
  appendVal=FALSE,
498   title="Binned data for village dataset comparison",
499   summary = FALSE,
500   style = "apsr",
501   table.placement = "htbp!",
502   label = "tab:villageBins",
503   rownames = TRUE
504 )
505
506
507 # Pearson's Chi-squared test
508
509 #data: one_counts and two_counts

```

```

510 #X-squared = 12, df = 9, p-value = 0.2133
511
512
513 #tibble('village1' = one_counts, 'village2' = two_counts )
514
515 # run regression model on each set of villages
516 one_model <- lm(water ~ reserved, data = one_village_policy)
517 two_model <- lm(water ~ reserved, data = two_village_policy)
518
519 summary(one_model)
520 summary(two_model)
521
522 output_stargazer(one_model, outputFile="village_model.tex", type = "latex",
523                  appendVal=FALSE,
524                  title="Pearson Linear Regression - Water ~ Reserved - Village
525                      = 1",
526                  style = "apsr",
527                  table.placement = "htb",
528                  label = "tab:village1"
529 )
530 output_stargazer(two_model, outputFile="village_model.tex", type = "latex",
531                  appendVal=TRUE,
532                  title="Pearson Linear Regression - Water ~ Reserved - Village
533                      = 2",
534                  style = "apsr",
535                  table.placement = "htb",
536                  label = "tab:village2"
537 )
538
539 # or combine the villages
540
541 combined_village_policy <- policy %>%
542   group_by(GP) %>%
543   mutate (sum_water = sum(water), sum_irrigation = sum(irrigation)) %>%
544   select(GP, reserved, female, sum_water, sum_irrigation) %>%
545   unique()
546
547 # run model - scaled by 1/2 to get equivalent values to 1 village coefficients
548 cvp <- lm(sum_water/2 ~ reserved, data = combined_village_policy)
549
550 summary(cvp)
551
552 output_stargazer(cvp, outputFile="villages_combined.tex", type = "latex",
553                  appendVal=FALSE,
554                  title="Pearson Linear Regression - Water ~ Reserved -
555                      Villages combined",
556                  style = "apsr",
557                  table.placement = "htb",
558                  label = "tab:combinedVillages"

```

```
558 )  
559  
560  
561 # refs  
562  
563 # Foundations of Statistics for Data Scientists; with R and Python  
564 # https://en.wikipedia.org/wiki/Least\_squares  
565 #
```