

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 for (i in 1:nrows) {
17   for (j in 1:ncols) {
18     expected[i,j] <- row_tots[i] * col_tots[j] / totals
19   }
20 }
21
22 # calculate difference between observed and expected
```

```

23 o_e <- observed
24 for (i in 1:nrows) {
25   for (j in 1:ncols) {
26     o_e[i,j] <- (observed[i,j] - expected[i,j])^2 / expected[i,j]
27   }
28 }
29
30 #calculate chi-squared value & degrees of freedom
31 chi_sq_val <- sum(o_e)
32 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.

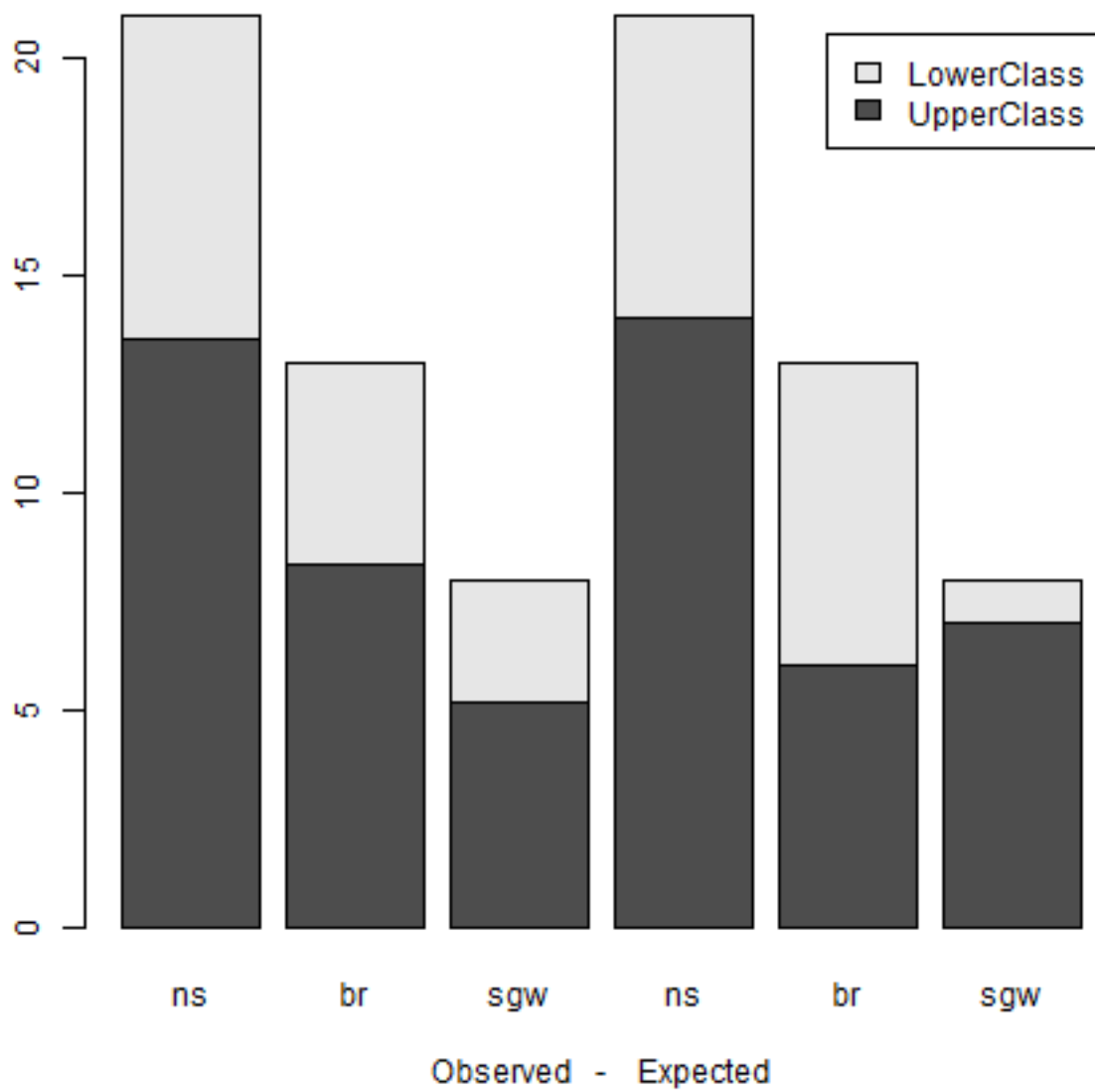


Figure 1: Observed vs Expected values for traffic stop. ns = Not Stopped; sgw = Stopped Given Warning; br = Bribe Requested

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

Table 1: Standardised Residuals

	NotStopped	BribeRequested	StoppedGivenWarning
UpperClass	0.32	-1.64	1.52
LowerClass	-0.32	1.64	-1.52

(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.

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")
```

Use 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 lm(female ~ reserved, data=policy)
```

This results in the following output:

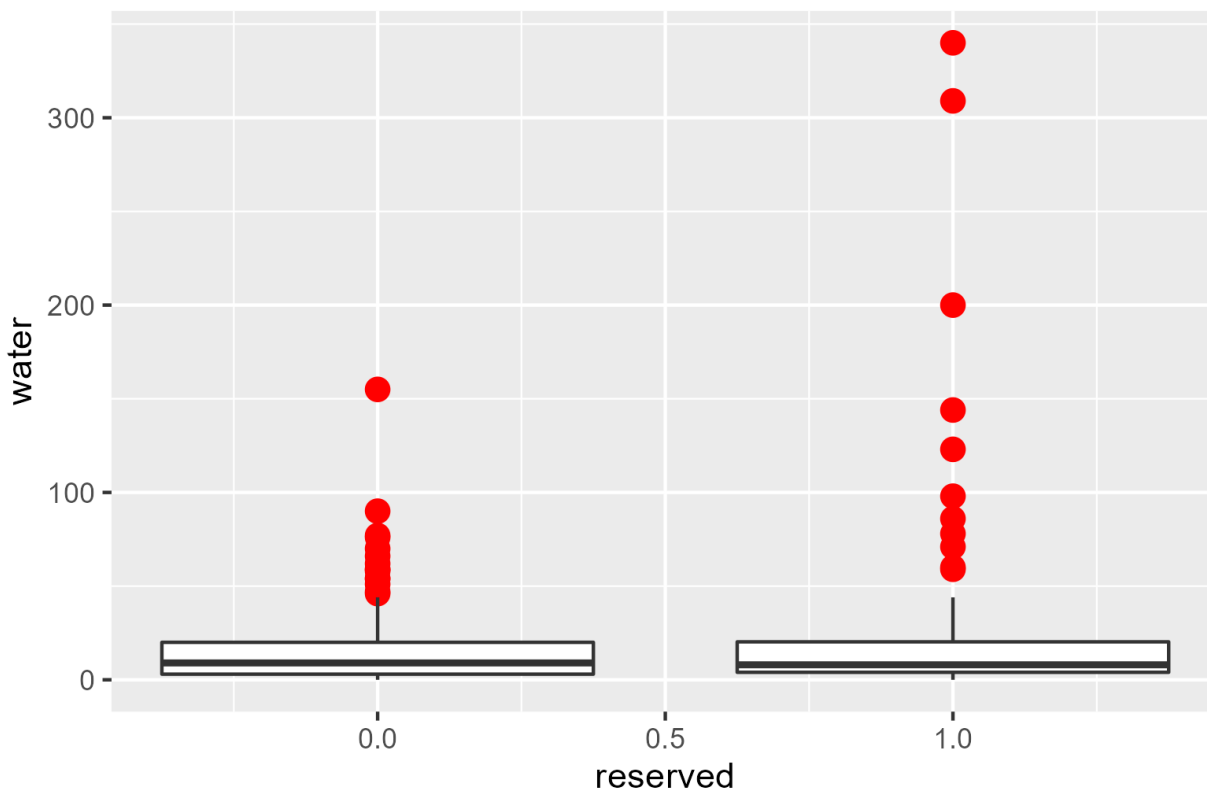
Table 2: Pearson Linear Regression - Water Reserved

	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	

(c) Interpret the coefficient estimate for reservation policy.

The

Figure 3: Boxplot of number of drinking water projects, grouped by reserved



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_df <- data.frame(observed, row.names = rows)
68 names(observed_df) <- cols
69 print(observed_df)
70
71 pairs(observed)
72
73 #\item [(a)]
74 #Calculate the  $\chi^2$  test statistic by hand/manually\\
75
76 ###----- 0 start listing of code from here
77 ncols <- length(observed[1,])
78 nrows <- length(observed[,1])
79
80 # get totals
81 row_tots <- vector("double", nrows)
82 col_tots <- vector("double", ncols)
83
84 totals <- sum(observed) # total number of observations
85
86 # calculate row and column totals, e.g, total for NotStopped, UpperClass, etc
87 for (i in 1:nrows) {row_tots[i] <- sum(observed[i, ])}
88 for (i in 1:ncols) {col_tots[i] <- sum(observed[, i])}
89
90 #get expected = row total * column total / total observations
91 expected <- observed
92 for (i in 1:nrows) {

```

```

93   for (j in 1:ncols) {
94     expected[i,j] <- row_tots[i] * col_tots[j] / totals
95   }
96 }
97
98 # calculate difference between observed and expected
99 o_e <- observed
100 for (i in 1:nrows) {
101   for (j in 1:ncols) {
102     o_e[i,j] <- (observed[i,j] - expected[i,j])^2 / expected[i,j]
103   }
104 }
105
106 #calculate chi-squared value & degrees of freedom
107 chi_sq_val <- sum(o_e)
108 df = (nrows-1) * (ncols-1)
109
110 cat(str_glue("The chi-squared statistic is {round(chi_sq_val,3)}"))
111 cat(str_glue("The chi-squared degrees of freedom is {df}"))
112
113 # plot of observed and expected values
114 png("obs_exp.png")
115 barplot(cbind(expected, observed), legend.text = rows,
116         names.arg = c("ns", "br", "sgw", "ns", "br", "sgw"),
117         xlab = "Observed - Expected")
118 dev.off()
119
120 #\item [(b)]
121 #Now calculate the p-value from the test statistic you just created R
122 # .\footnote{Remember frequency should be  $> 5$  for all cells, but let's
123 # calculate
124 # the p-value here anyway.} What do you conclude if  $\alpha = 0.1$ ?\\
125
126 p_value <- pchisq(chi_sq_val, df=df, lower.tail=FALSE)
127 alpha <- 0.1
128
129 # p > alpha, can't reject null
130 if (p_value > alpha) txt <- "cannot " else txt <- ""
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

```

```

140 for (i in 1:nrows) {
141   row_prop<- (1 - (row_tots [i] / totals))
142   for (j in 1:ncols) {
143     col_prop<- (1- (col_tots[j] / totals))
144     z[i,j] <- (observed[i,j] - expected[i,j]) /sqrt (expected[i,j]* row_prop
      * col_prop)
145   }
146 }
147 z_df <- data.frame(z, 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                 digits=2,
157                 summary = FALSE,
158                 style = "apsr",
159                 table.placement = "h",
160                 label = "StandardisedResiduals",
161                 rownames = TRUE
162                 )
163
164
165 # check result
166 chisq.test(observed)
167
168 #https://www.rdocumentation.org/packages/stargazer/versions/5.2.3/topics/
   stargazer
169
170 # \item [(d)] How might the standardized residuals help you interpret the
   results?
171
172 # fewer upper class individuals asked for bribes and more given warnings;
173 # the contribution from lower class drivers expected to give bribes is nearly
174 # equivalent to the contribution from upper class drivers getting warnings
175
176 #
   #####
177 # Problem 2
178 #####
179
180 #Question 2 (40 points): Economics
181 #Chattopadhyay and Duflo were interested in whether women promote different
   policies
182 # than men.
183 # Answering this question with observational data is pretty difficult due to
   potential

```

```

184 # confounding problems (e.g. the districts that choose female politicians are
185 # likely to systematically differ in other aspects too). Hence, they exploit a
186 # randomized policy experiment in India, where since the mid-1990s, 1/3 of
187 # village council heads have been randomly reserved for women. A subset of the
    data
188 # from West Bengal can be found at the following link:
189 #   \url{https://raw.githubusercontent.com/kosukeimai/qss/master/PREDICTION/
    women.csv}
190
191 # Each observation in the data set represents a village and there are two
    villages
192 # associated with one GP (i.e. a level of government is called "GP").
193 # Figure~\ref{fig:women_desc} below shows the names and descriptions of the
    variables
194 # in the dataset. The authors hypothesize that female politicians are more
    likely to
195 # support policies female voters want. Researchers found that more women
    complain about
196 # the quality of drinking water than men. You need to estimate the effect of
    the
197 # reservation policy on the number of new or repaired drinking water
    facilities
198 #in the villages.
199 # Names and description of variables from Chattopadhyay and Duflo (2004)
200 # 1 'GP' Identifier for the Gram Panchayat &nbsp;&nbsp; 
201 # 2 'village' identifier for each village
202 # 3 'reserved' binary variable indicating whether the GP was reserved for
    women leaders or not
203 # 4 'female' binary variable indicating whether the GP had a female leader or
    not
204 # 5 'irrigation' variable measuring the number of new or repaired irrigation
    facilities in the village since the reserve policy started
205 # 6 'water' variable measuring the number of new or repaired drinking-water
    facilities in the village since the reserve policy started
206
207 #\item [(a)] State a null and alternative (two-tailed) hypothesis.
208 # null: no diff in incidence of new or repaired drinking-water facilities
209 # in the village since the reserve policy started
210 # ie 'water' is independent of 'reserved'
211 # alternate: the incidence of new or repaired drinking-water facilities is
212 # correlated to the reservation policy
213
214
215 policy <- read.csv("https://raw.githubusercontent.com/kosukeimai/qss/master/
    PREDICTION/women.csv")
216 #write.csv(policy,"Data/policy.csv")
217 policy<-read.csv("Data/policy.csv")
218
219 summary(policy)
220 pairs(policy[4:7])
221

```



```

222 lm(female ~ reserved , data=policy)
223
224 sum(policy$reserved)
225 sum(policy$female)
226
227
228  $\item [(b)]$  Run a bivariate regression to test this hypothesis in  $\texttt{\textbackslash texttt{R}}$ 
    (include your code!).
229
230 water <- lm(water ~ reserved , data = policy)
231 summary(water)
232
233 p<- ggplot(policy , aes(reserved , water , colour=female))
234 p + geom_jitter()
235
236
237 p<- ggplot(policy , aes(reserved , water , group_by(reserved)))
238 p + geom_boxplot(outlier.colour = 'red' , outlier.size = 3 , aes(group=reserved)
    )
239 ggsave("water_boxplot.png")
240
241
242 output_stargazer(water , outputFile="water_model.tex" , type = "latex" ,
243                 appendVal=FALSE,
244                 title="Pearson Linear Regression - Water ~ Reserved" ,
245                 style = "apsr" ,
246                 label = "water_reserved"
247 )
248
249
250
251 output_stargazer(z_df , outputFile="std_residuals.tex" , type = "latex" ,
252                 appendVal=FALSE,
253                 title="Standardised Residuals" ,
254                 digits=2,
255                 summary = FALSE,
256                 style = "apsr" ,
257                 table.placement = "h" ,
258                 label = "StandardisedResiduals" ,
259                 rownames = TRUE
260 )
261
262
263
264
265 ###
266
267 water_female <- lm(water ~ female , data = policy)
268
269 summary(water_female)
270

```

```

271 plot(water)
272 #\item [(c)] Interpret the coefficient estimate for reservation policy.
273
274 with(policy, plot(water, reserved))
275 p<- ggplot(policy)
276 p+ geom_jitter(aes(reserved, water, colour=female))
277
278 p<- ggplot(policy, aes(reserved, water), colour=female) + geom_jitter()
279 p+ facet_wrap(vars(female))

```