

Problem Set 3 Answers

QTM 200: Applied Regression Analysis

Vanessa Wong

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 the course GitHub page in .pdf form.
- This problem set is due at the beginning of class on Monday, February 17, 2020. No late assignments will be accepted.
- Total available points for this homework is 100.

In this problem set, you will run several regressions and create an add variable plot (see the lecture slides) in R using the `incumbents_subset.csv` dataset. Include all of your code.

Question 1 (20 points)

We are interested in knowing how the difference in campaign spending between incumbent and challenger affects the incumbent's vote share.

1. Run a regression where the outcome variable is `voteshare` and the explanatory variable is `difflog`.

Linear model: $y = 0.5790 + 0.0417x$

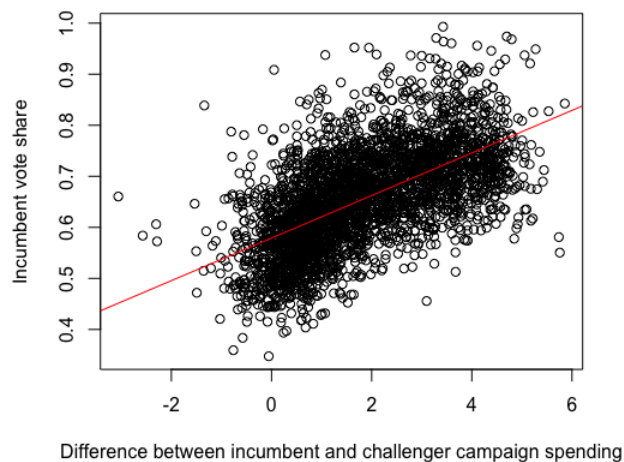
```
1 # x = difflog (explanatory)
2 # y = voteshare (outcome)
3
4 # run a regression where the outcome variable is voteshare and the
  explanatory variable is difflog.
5 ymean <- mean(incumbents$voteshare)
```

```

6 xmean <- mean(incumbents$difflog)
7 ysum <- sum(incumbents$voteshare)
8 xsum <- sum(incumbents$difflog)
9 yy <- (incumbents$voteshare) - (ymean)
10 xx <- (incumbents$difflog) - (xmean)
11 yyxxsum <- sum(yy*xx)
12 xxsq <- (xx)^2
13 sumxxsq <- sum(xxsq)
14 betaincumbs <- yyxxsum/sumxxsq
15 betaincumbs
16 # beta = 0.0417
17 alphaincumbs <- ymean - (betaincumbs*xmean)
18 alphaincumbs
19 # alpha = 0.5790
20 # linear model: y = 0.5790 + 0.0417x
21 # check work
22 incumbreg <- lm(incumbents$voteshare ~ incumbents$difflog)
23 incumbreg

```

2. Make a scatterplot of the two variables and add the regression line.



```

1 plot(incumbents$difflog, incumbents$voteshare,
2       xlab="Difference between incumbent and challenger campaign spending",
3       ylab="Incumbent vote share")
4 abline(a=0.5790, b=0.0417, col="red")

```

3. Save the residuals of the model in a separate object.

```
1 residscheck <- residuals(incumbreg)
2 residscheck
```

4. Write the prediction equation.

$$\hat{y} = 0.5790 + 0.0417x$$

Question 2 (20 points)

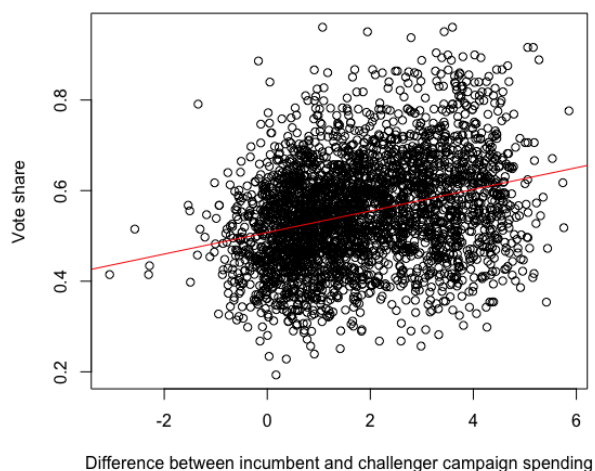
We are interested in knowing how the difference between incumbent and challenger's spending and the vote share of the presidential candidate of the incumbent's party are related.

1. Run a regression where the outcome variable is `presvote` and the explanatory variable is `difflog`.

Linear model: $y = 0.5076 + 0.0238x$

```
1 y2mean <- mean(incumbents$presvote)
2 y2sum <- sum(incumbents$presvote)
3 yy2 <- (incumbents$presvote) - y2mean
4 yyxxsum2 <- sum(xx*yy2)
5 betaincumbs2 <- yyxxsum2/sumxxsq
6 betaincumbs2
7 # beta = 0.0238
8 alphaincumbs2 <- y2mean - (betaincumbs2*xmean)
9 alphaincumbs2
10 # alpha = 0.5076
11 # linear model: y = 0.5076 + 0.0238x
12 # check work
13 incumbreg2 <- lm(incumbents$presvote ~ incumbents$difflog)
14 incumbreg2
```

2. Make a scatterplot of the two variables and add the regression line.



```
1 plot(incumbents$difflog, incumbents$presvote,
```

```

2      xlab="Difference between incumbent and challenger campaign spending"
      , ylab="Vote share")
3 abline(a=0.5076, b=0.0238, col="red")

```

3. Save the residuals of the model in a separate object.

```

1 resids2 <- residuals(incumbreg2)
2 resids2

```

4. Write the prediction equation.

$$\hat{y} = 0.5076 + 0.0238x$$

Question 3 (20 points)

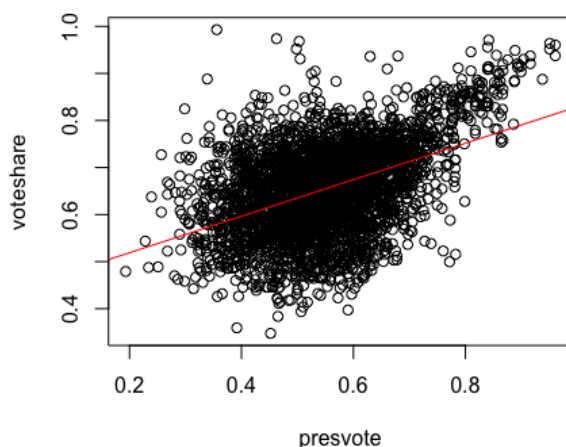
We are interested in knowing how the vote share of the presidential candidate of the incumbent's party is associated with the incumbent's electoral success.

1. Run a regression where the outcome variable is **voteshare** and the explanatory variable is **presvote**.

Linear model: $y = 0.4413 + 0.3880x$

```
1 xmean3 <- mean(incumbents$presvote)
2 ymean3 <- mean(incumbents$voteshare)
3 xsum3 <- sum(incumbents$presvote)
4 ysum3 <- sum(incumbents$voteshare)
5 yy3 <- (incumbents$voteshare) - (ymean3)
6 xx3 <- (incumbents$presvote) - (xmean3)
7 yyxxsum3 <- sum(yy3*xx3)
8 xxsq3 <- (xx3)^2
9 sumxxsq3 <- sum(xxsq3)
10 beta3 <- yyxxsum3/sumxxsq3
11 beta3
12 # beta = 0.3880
13 alpha3 <- ymean3 - (beta3*xmean3)
14 alpha3
15 # alpha = 0.4413
16 # linear model: y= 0.4413 + 0.3880x
17 # check work
18 reg3 <- lm(incumbents$voteshare~incumbents$presvote)
19 reg3
```

2. Make a scatterplot of the two variables and add the regression line.



```
1 plot(incumbents$presvote, incumbents$voteshare,
2       xlab="presvote", ylab="voteshare")
3 abline(a=0.4413, b=0.3880, col="red")
```

3. Write the prediction equation.

$$\hat{y} = 0.4413 + 0.3880x$$

Question 4 (20 points)

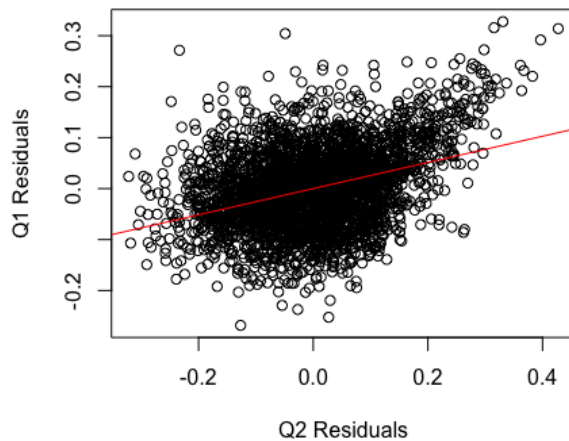
The residuals from part (a) tell us how much of the variation in **voteshare** is *not* explained by the difference in spending between incumbent and challenger. The residuals in part (b) tell us how much of the variation in **presvote** is *not* explained by the difference in spending between incumbent and challenger in the district.

1. Run a regression where the outcome variable is the residuals from Question 1 and the explanatory variable is the residuals from Question 2.

Linear model: $y = -4.860 \times 10^{-18} + 0.2569x$

```
1 xmean4 <- mean(resids2)
2 ymean4 <- mean(residscheck)
3 xsum4 <- sum(resids2)
4 ysum4 <- sum(residscheck)
5 yy4 <- residscheck - (ymean4)
6 xx4 <- resids2 - (xmean4)
7 yyxxsum4 <- sum(yy4*xx4)
8 xxsq4 <- (xx4)^2
9 sumxxsq4 <- sum(xxsq4)
10 beta4 <- yyxxsum4/sumxxsq4
11 beta4
12 # beta = 0.2569
13 alpha4 <- ymean4 - (beta4*xmean4)
14 alpha4
15 # alpha = -4.860 x 10^(-18)
16 # linear model: y = -4.860 x 10^(-18) + 0.2569x
17 # check work
18 reg4 <- lm(residscheck ~ resids2)
19 reg4
20 summary(reg4)
```

2. Make a scatterplot of the two residuals and add the regression line.



```
1 plot(resids2, residscheck,
2      xlab="Q2 Residuals", ylab="Q1 Residuals")
3 abline(reg4, col="red")
```

3. Write the prediction equation.

$$\hat{y} = -4.860 \times 10^{-18} + 0.2569x$$

Question 5 (20 points)

What if the incumbent's vote share is affected by both the president's popularity and the difference in spending between incumbent and challenger?

1. Run a regression where the outcome variable is the incumbent's `voteshare` and the explanatory variables are `difflog` and `presvote`.

```
1 # [,2] = presvote
2 Y3 <- matrix(incumbents$voteshare, nrow=3193)
3 Y3
4 XtX <- t(X3) %*% X3
5 Xty <- t(X3) %*% Y3
6 Xty
7 XtX.inv <- solve(XtX)
8 XtX
9 b <- XtX.inv %*% Xty
10 b
11 dim(incumbents)[1]
12
13 # check work
14 lm(incumbents$voteshare ~ incumbents$difflog + incumbents$presvote)$
   coefficients
```

```

15 multireg2 <- lm(incumbents$voteshare ~ incumbents$difflog + incumbents$
    presvote)
16 summary(multireg2)
17 # linear model: y = 0.4486 + 0.0355xdifflog + 0.2569xpresvote

```

2. Write the prediction equation.

$$\hat{y} = 0.4486 + 0.0355x_{difflog} + 0.2569x_{presvote}$$

3. What is it in this output that is identical to the output in Question 4? Why do you think this is the case?

Answered in tex file

NOTE: The final R file with answers is called PS3_answers4. I had to resave the file multiple times because I was making changes to the R file as I was compiling.

```

1 #####
2 # load libraries
3 # set wd
4 # clear global .envir
5 #####
6
7 # remove objects
8 rm(list=ls())
9 # detach all libraries
10 detachAllPackages <- function() {
11   basic.packages <- c("package:stats", "package:graphics", "package:grDevices",
12     "package:utils", "package:datasets", "package:methods", "package:base")
13   package.list <- search()[ifelse(unlist(gregexpr("package:", search()))==1,
14     TRUE, FALSE)]
15   package.list <- setdiff(package.list, basic.packages)
16   if (length(package.list)>0) for (package in package.list) detach(package,
17     character.only=TRUE)
18 }
19 detachAllPackages()
20
21 # load libraries
22 pkgTest <- function(pkg){
23   new.pkg <- pkg[!(pkg %in% installed.packages()[, "Package"])]
24   if (length(new.pkg))
25     install.packages(new.pkg, dependencies = TRUE)
26   sapply(pkg, require, character.only = TRUE)
27 }
28
29 # here is where you load any necessary packages
30 # ex: stringr
31 # lapply(c("stringr"), pkgTest)

```

```

29
30 lapply(c("stringr"), pkgTest)
31
32 # set working directory, import datasets
33 setwd("~/GitHub/QTM200Spring2020/problem_sets/PS3")
34 incumbents <- read.csv("incumbents_subset.csv")
35
36 #####
37 # Problem 1:
38 #####
39 # x = difflog (explanatory)
40 # y = voteshare (outcome)
41
42 # run a regression where the outcome variable is voteshare and the explanatory
    variable is difflog.
43 ymean <- mean(incumbents$voteshare)
44 xmean <- mean(incumbents$difflog)
45 ysum <- sum(incumbents$voteshare)
46 xsum <- sum(incumbents$difflog)
47 yy <- (incumbents$voteshare) - (ymean)
48 xx <- (incumbents$difflog) - (xmean)
49 yyxxsum <- sum(yy*xx)
50 xxsq <- (xx)^2
51 sumxxsq <- sum(xxsq)
52 betaincumbs <- yyxxsum/sumxxsq
53 betaincumbs
54 # beta = 0.0417
55 alphaincumbs <- ymean - (betaincumbs*xmean)
56 alphaincumbs
57 # alpha = 0.5790
58 # linear model:  $y = 0.5790 + 0.0417x$ 
59 # check work
60 incumbreg <- lm(incumbents$voteshare ~ incumbents$difflog)
61 incumbreg
62
63 # make a scatterplot of the two variables and add the regression line
64 plot(incumbents$difflog, incumbents$voteshare,
65      xlab="Difference between incumbent and challenger campaign spending",
66      ylab="Incumbent vote share")
67 abline(a=0.5790, b=0.0417, col="red")
68
69 # save the residuals of the model in a separate object
70 residscheck <- residuals(incumbreg)
71 residscheck
72
73 # write the prediction equation
74 #yhat =  $0.5790 + 0.0417x$ 
75 #####
76 # Problem 2:
77 #####

```

```

78 # x = difflog (explanatory)
79 # y = presvote (outcome)
80
81 # run a regression where the outcome variable is presvote and the explanatory
    variable is difflog
82 y2mean <- mean(incumbents$presvote)
83 y2sum <- sum(incumbents$presvote)
84 yy2 <- (incumbents$presvote) - y2mean
85 yyxxsum2 <- sum(xx*yy2)
86 betaincumbs2 <- yyxxsum2/sumxxsq
87 betaincumbs2
88 # beta = 0.0238
89 alphaincumbs2 <- y2mean - (betaincumbs2*xmean)
90 alphaincumbs2
91 # alpha = 0.5076
92 # linear model:  $y = 0.5076 + 0.0238x$ 
93 # check work
94 incumbreg2 <- lm(incumbents$presvote ~ incumbents$difflog)
95 incumbreg2
96
97 # make a scatterplot of the two variables and add the regression line
98 plot(incumbents$difflog, incumbents$presvote,
99       xlab="Difference between incumbent and challenger campaign spending",
       ylab="Vote share")
100 abline(a=0.5076, b=0.0238, col="red")
101
102 #save the residuals of the model in a separate object
103 resids2 <- residuals(incumbreg2)
104 resids2
105
106 # write the prediction equation
107 #  $\hat{y} = 0.5076 + 0.0238x$ 
108
109 #####
110 # Problem 3:
111 #####
112 # x = presvote (explanatory)
113 # y = voteshare (outcome)
114
115 # run a regression where the outcome variable is voteshare and the explanatory
    variable is presvote
116 xmean3 <- mean(incumbents$presvote)
117 ymean3 <- mean(incumbents$voteshare)
118 xsum3 <- sum(incumbents$presvote)
119 ysum3 <- sum(incumbents$voteshare)
120 yy3 <- (incumbents$voteshare) - (ymean3)
121 xx3 <- (incumbents$presvote) - (xmean3)
122 yyxxsum3 <- sum(yy3*xx3)
123 xxsq3 <- (xx3)^2
124 sumxxsq3 <- sum(xxsq3)
125 beta3 <- yyxxsum3/sumxxsq3

```

```

126 beta3
127 # beta = 0.3880
128 alpha3 <- ymean3 - (beta3*xmean3)
129 alpha3
130 # alpha = 0.4413
131 # linear model: y= 0.4413 + 0.3880x
132 # check work
133 reg3 <- lm(incumbents$voteshare~incumbents$presvote)
134 reg3
135
136 # make a scatterplot of the two variables and add the regression line
137 plot(incumbents$presvote, incumbents$voteshare,
138       xlab="presvote", ylab="voteshare")
139 abline(a=0.4413, b=0.3880, col="red")
140
141 # write the prediction equation
142 # yhat = 0.4413 + 0.3880x
143
144 #####
145 # Problem 4:
146 #####
147 # x = resids2 (explanatory)
148 # y = residscheck (outcome)
149
150 # run a regression where the outcome variable is the residuals from Q1 and the
    explanatory variable is the residuals from Q2
151 xmean4 <- mean(resids2)
152 ymean4 <- mean(residscheck)
153 xsum4 <- sum(resids2)
154 ysum4 <- sum(residscheck)
155 yy4 <- residscheck - (ymean4)
156 xx4 <- resids2 - (xmean4)
157 yyxxsum4 <- sum(yy4*xx4)
158 xxsq4 <- (xx4)^2
159 sumxxsq4 <- sum(xxsq4)
160 beta4 <- yyxxsum4/sumxxsq4
161 beta4
162 # beta = 0.2569
163 alpha4 <- ymean4 - (beta4*xmean4)
164 alpha4
165 # alpha = -4.860 x 10^(-18)
166 # linear model: y = -4.860 x 10^(-18) + 0.2569x
167 # check work
168 reg4 <- lm(residscheck ~ resids2)
169 reg4
170 summary(reg4)
171
172 # make a scatterplot of the two residuals and add the regression line
173 plot(resids2, residscheck,
174       xlab="Q2 Residuals", ylab="Q1 Residuals")
175 abline(reg4, col="red")

```

```

176
177 # write the prediction equation
178 # yhat = -4.860 x 10-18 + 0.2569x
179
180 #####
181 # Problem 5:
182 #####
183 # x = difflog , presvote (explanatory)
184 # y = voteshare (outcome)
185
186 # run a regression where the outcome variable is the incumbent's voteshare and
187 # the explanatory variables are difflog and presvote
188 X3 <- cbind(incumbents$difflog, incumbents$presvote)
189 X3
190 # [,1] = difflog
191 # [,2] = presvote
192 dim(incumbents)[1]
193 Y3 <- matrix(incumbents$voteshare, nrow=3193)
194 Y3
195 XtX <- t(X3) %*% X3
196 Xty <- t(X3) %*% Y3
197 Xty
198 XtX.inv <- solve(XtX)
199 XtX
200 b <- XtX.inv %*% Xty
201 b
202 # check work
203 lm(incumbents$voteshare ~ incumbents$difflog + incumbents$presvote)
204 $coefficients
205 multireg2 <- lm(incumbents$voteshare ~ incumbents$difflog +
206 incumbents$presvote)
207 summary(multireg2)
208 # linear model: y = 0.4486 + 0.0355xdifflog + 0.2569xpresvote
209
210
211 # write the prediction equation
212 # yhat = 0.4486 + 0.0355xdifflog + 0.2569xpresvote
213
214 # what is it in this output that is identical to the output in Q4? why do you
215 # think this is the case?
216 # the slope of the regression of voteshare on presvote is the same as the
217 # slope of of the regression of Q1 residuals on Q2 residuals.
218 # i think this is the case because the regression of Q1 residuals on Q2
219 # residuals tells us how much variation in presvote
220 # *and* voteshare is not explained by difflog. in both cases, the slopes
221 # describe the relationship between voteshare and presvote, while
222 # holding difflog constant. As such, they have the same slope.

```