

# Problem Set 3

## QTM 200: Applied Regression Analysis

Due: February 17, 2020

### 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`.

```
1 #1a. Run a regression where the outcome variable is voteshare and the
   #   explanatory variable is difflog.
2 #calculate sums and means
3 mean.difflog <- mean(incumbents_subset$difflog)
4 mean.vote <- mean(incumbents_subset$voteshare)
5 sum(incumbents_subset$difflog)
```

```

6 sum(incumbents_subset$voteshare)
7 numerator <- sum((incumbents_subset$voteshare - mean(incumbents_subset$
  voteshare))*(incumbents_subset$difflog - mean(incumbents_subset$
  difflog)))
8 denominator <- sum((incumbents_subset$difflog - mean(incumbents_subset$
  difflog))^2)
9 #calculate regression coefficient
10 beta.hat <- numerator/denominator
11 beta.hat
12 alpha.hat <- mean.vote - (beta.hat*mean.difflog)
13 alpha.hat
14 #calculate p-value
15 sd.vote <- sd(incumbents_subset$voteshare)
16 se.vote <- sd.vote/sqrt(sum((incumbents_subset$voteshare - mean.difflog)
  ^2))
17 TS <- (beta.hat - 0) / se.vote
18 TS
19 p <- 2*pt(abs(TS), df=(length(incumbents_subset$voteshare)), lower.tail =
  F)
20 p
21 #check model
22 lm1 <- lm(voteshare ~ difflog, data=incumbents_subset)
23 summary(lm1)

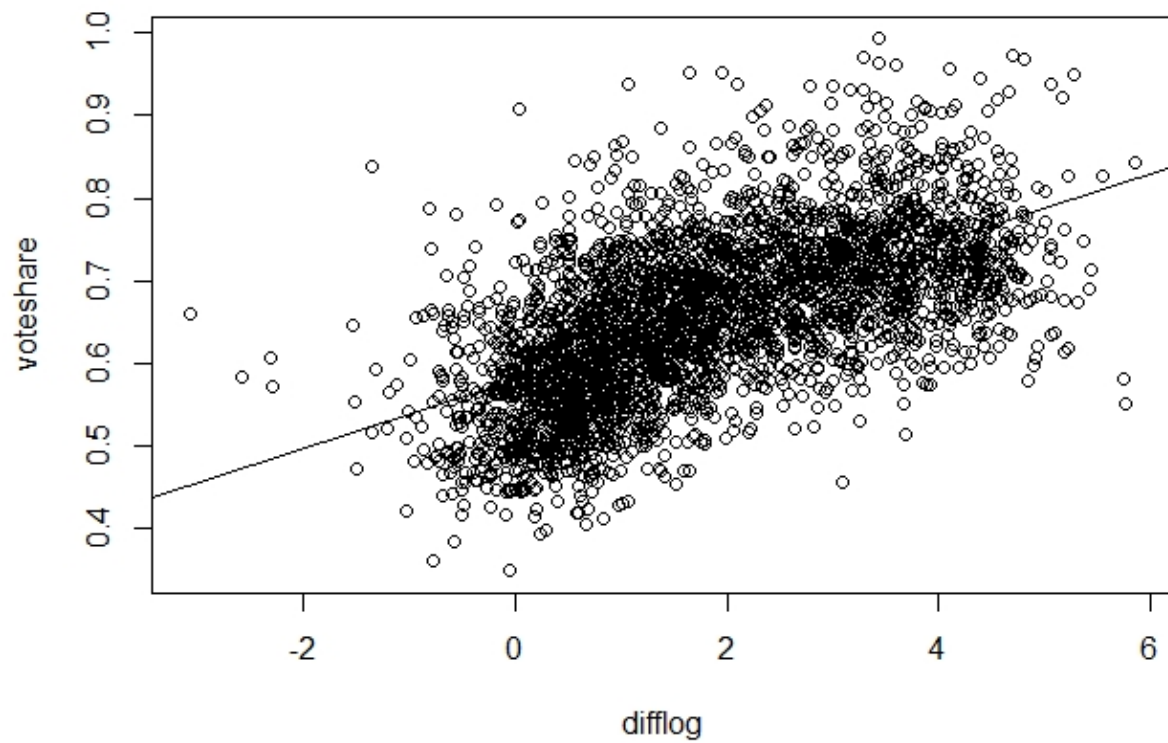
```

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

```

1 #1b. Make a scatterplot of the two variables and add the regression line.
2 library(ggplot2)
3 plot(voteshare ~ difflog, data=incumbents_subset)
4 abline(lm(voteshare ~ difflog, data=incumbents_subset))

```



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

```
1 residual.lm1 <- incumbents_subset$voteshare - 0.579031 - (0.041666*  
  incumbents_subset$difflog)
```

4. Write the prediction equation.  $Y_i = 0.579031 + 0.041666X_i + \text{error}_i$

## Question 2 (20 points)

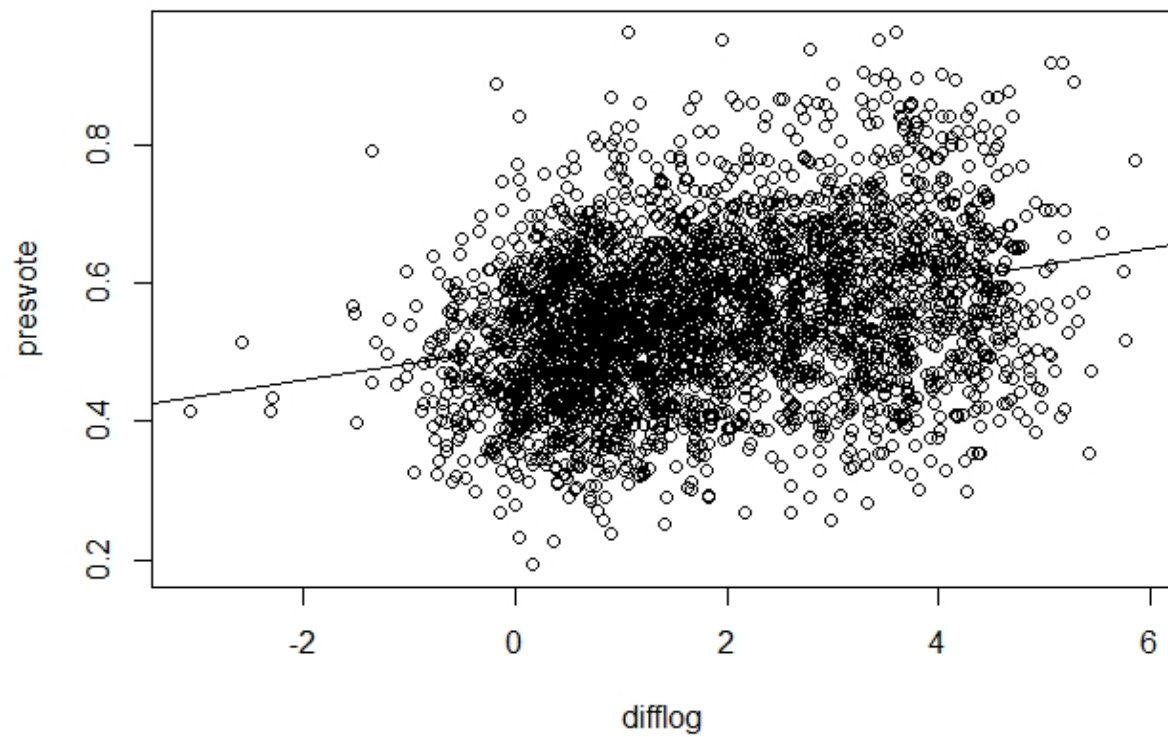
We are interested in knowing how the difference between incumbent and challenger's spending and the vote share of the presidential candidate of the incumbent's party are related.

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

```
1 #question 2a. Run a regression where the outcome variable is presvote and
  the explanatory variable is difflog.
2 #calculate sums and means
3 mean.difflog <- mean(incumbents_subset$difflog)
4 mean.presvote <- mean(incumbents_subset$presvote)
5 sum(incumbents_subset$difflog)
6 sum(incumbents_subset$presvote)
7 numerator2 <- sum((incumbents_subset$presvote - mean(incumbents_subset$
  presvote))*(incumbents_subset$difflog - mean(incumbents_subset$difflog
  )))
8 denominator2 <- sum((incumbents_subset$difflog - mean(incumbents_subset$
  difflog))^2)
9 #calculate regression coefficient
10 beta.hat2 <- numerator2/denominator2
11 beta.hat2
12 alpha.hat2 <- mean.presvote - (beta.hat2*mean.difflog)
13 alpha.hat2
14 #calculate p-value
15 sd.presvote <- sd(incumbents_subset$presvote)
16 se.presvote <- sd.presvote/sqrt(sum((incumbents_subset$presvote - mean.
  difflog)^2))
17 TS2 <- (beta.hat2 - 0) / se.presvote
18 TS2
19 p2 <- 2*pt(abs(TS2), df=(length(incumbents_subset$presvote)), lower.tail
  = F)
20 p2
21 #check model
22 lm2 <- lm(presvote ~ difflog, data=incumbents_subset)
23 summary(lm2)
```

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

```
1 plot(presvote ~ difflog, data=incumbents_subset)
2 abline(lm(presvote ~ difflog, data=incumbents_subset))
```



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

```
1 residual.lm2 <- incumbents_subset$presvote - 0.507583 - (0.023837 *  
  incumbents_subset$difflog)
```

4. Write the prediction equation.  $Y_i = 0.507583 + 0.023837X_i + \text{error}_i$

## Question 3 (20 points)

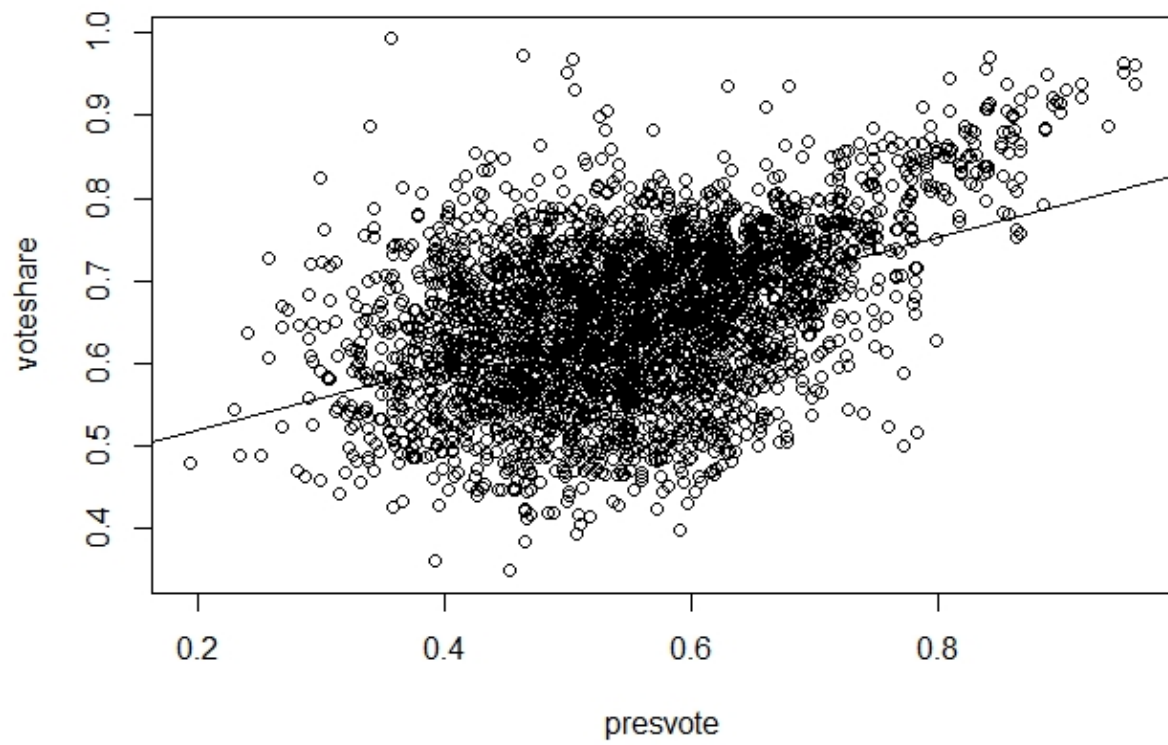
We are interested in knowing how the vote share of the presidential candidate of the incumbent's party is associated with the incumbent's electoral success.

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

```
1 #question 3a. Run a regression where the outcome variable is voteshare
  and the explanatory variable is presvote.
2 #calculate sums and means
3 mean.presvote
4 mean.vote
5 sum(incumbents_subset$presvote)
6 sum(incumbents_subset$voteshare)
7 numerator3 <- sum((incumbents_subset$voteshare - mean(incumbents_subset$
  voteshare))*(incumbents_subset$presvote - mean(incumbents_subset$
  presvote)))
8 denominator3 <- sum((incumbents_subset$presvote - mean(incumbents_subset$
  presvote))^2)
9 #calculate regression coefficient
10 beta.hat3 <- numerator3/denominator3
11 beta.hat3
12 alpha.hat3 <- mean.vote - (beta.hat3*mean.presvote)
13 alpha.hat3
14 #calculate p-value
15 sd.vote
16 se.vote2 <- sd.vote/sqrt(sum((incumbents_subset$voteshare - mean.
  presvote)^2))
17 TS3 <- (beta.hat3 - 0) / se.vote2
18 TS3
19 p3 <- 2*pt(abs(TS3), df=length(incumbents_subset$voteshare), lower.tail
  = F)
20 p3
21 #check model
22 lm3 <- lm(voteshare ~ presvote, data=incumbents_subset)
23 summary(lm3)
```

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

```
1 plot(voteshare ~ presvote, data=incumbents_subset)
2 abline(lm(voteshare ~ presvote, data=incumbents_subset))
```



3. Write the prediction equation.  $Y_i = 0.441330 + 0.388018X_i + \text{error}_i$

## 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.

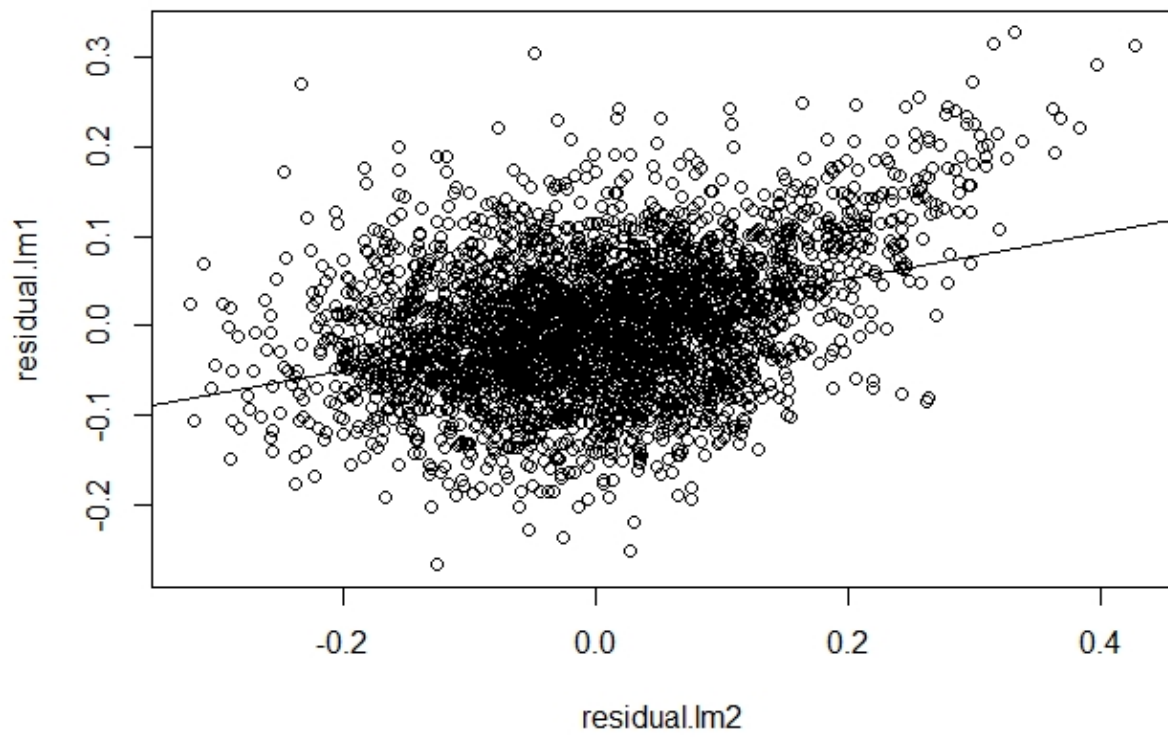
1. Run a regression where the outcome variable is the residuals from Question 1 and the explanatory variable is the residuals from Question 2.

```
1 mean.res1 <- mean(residual.lm1)
2 mean.res2 <- mean(residual.lm2)
3 sum(residual.lm1)
4 sum(residual.lm2)
5 numerator4 <- sum((residual.lm1 - mean(residual.lm1))*(residual.lm2 -
  mean(residual.lm2)))
6 denominator4 <- sum((residual.lm2 - mean(residual.lm2))^2)
7 #calculate regression coefficient
8 beta.hat4 <- numerator4/denominator4
9 beta.hat4
10 alpha.hat4 <- mean.res1 - (beta.hat4*mean.res2)
11 alpha.hat4
12 #calculate p-value
13 sd.r <- sd(residual.lm1)
14 se4 <- sd.r/sqrt(sum((residual.lm1 - mean.res2)^2))
15 TS4 <- (beta.hat4 - 0) / se4
16 TS4
17 p4 <- 2*pt(abs(TS4), df=length(residual.lm1), lower.tail = F)
18 p4
19 #check
20 lm4 <- lm(residual.lm1~residual.lm2)
21 summary(lm4)
```

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

```
1 plot(residual.lm1~residual.lm2)
2 abline(lm(residual.lm1 ~ residual.lm2))
```





3. Write the prediction equation.  $\hat{Y} = 1.086e-07 + 0.256877X + \text{error}$

## 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 #question 5a. Run a regression where the outcome variable is the
  incumbent's voteshare and the explanatory variables are difflog and
  presvote.
2 #subset data
3 sub.incumb <- incumbents_subset[,c("voteshare", "difflog", "presvote")]
4 #estimate regression by hand
5 lm_by_hand <- function(sub.incumb, difflog, presvote, voteshare)
6 X <- as.matrix(cbind(rep(1, dim(sub.incumb)[2]), sub.incumb[, "difflog", "
  presvote"])))
7 Y <- sub.incumb[, "voteshare"]
8 #calculate betas
9 betas <- solve((t(X)%*%X)) %*% (t(X)%*%Y)
10 #estimate sigma-squared
11 sigma_squared <- sum((Y-X%*%betas)^2)/nrow(X)-ncol(X)
12 #create-variance-covariance matrix for betas
13 var_covar_mat <-sigma_squared*solve(t(X)%*%X)
14 #SEs for coefficient estimates
15 SEs <- sqrt(diag(var_covar_mat))
16
17 #check
18 lm5 <- lm(voteshare~difflog+presvote, data=incumbents_subset)
19 summary(lm5)
```

2. Write the prediction equation.  $Y_i = 0.4486442 + 0.0355431X_1 + 0.2568770X_2 + \text{error}$

(X1 is difflog, X2 is 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? The min, 1Q, median, 3Q and max for the residuals are identical.

I think this is the case because the residual model from question 4 takes into account the same three variables as the multiple regression, so they have the same residual values.