

Problem Set 3

Tolga Bag - 23371290

Due: November 19, 2022

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 GitHub.
- This problem set is due before 23:59 on Sunday November 19, 2023. No late assignments will be accepted.

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

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`. I am checking the table and running the regression

```
1 head(inc.sub)
2 summary(inc.sub)
3 str(inc.sub) #I inspected the dataset.
4 regmod_a <- lm(voteshare ~ difflog, data = inc.sub) #I run the regression
  model
```

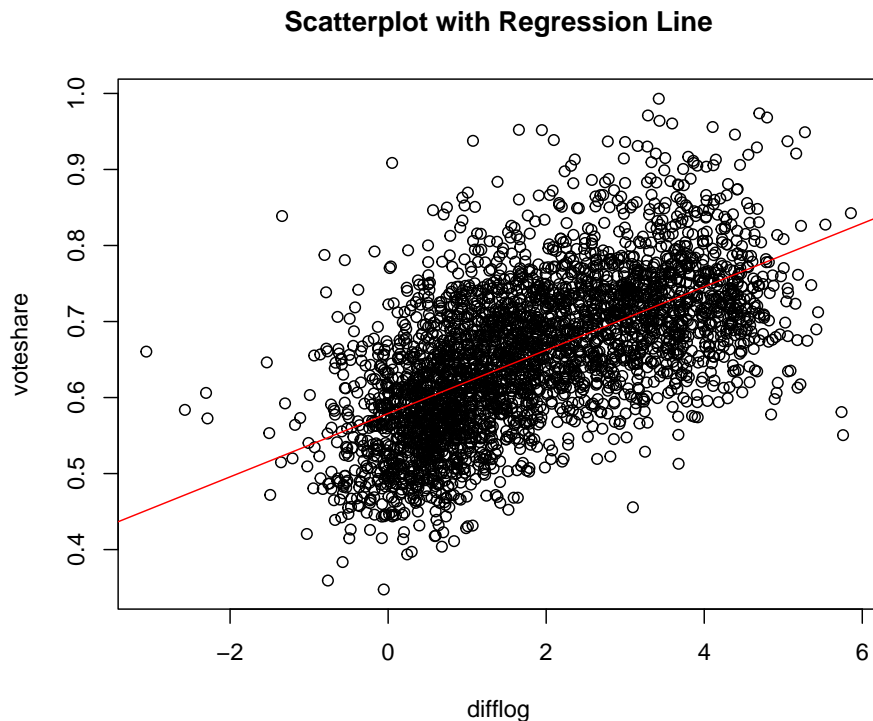
`voteshare` is the dependent or outcome variable and `difflog` is the independent or explanatory variable. My null hypothesis is there is no relationship between `voteshare` and `difflog`.

```
1 summary(regmod_a) #p value is 0.0000000000000002, so I can reject the
   null
```

p value is 0.0000000000000002, so I can reject the null hypothesis. Slope is 0.041666. In other words, for each unit change in difflog, voteshare is expected to increase by that amount. And it is a statistically significant result per the p value.

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

```
1 #I prepare a scatter plot with voteshare as the dependent variable
   against
2 #difflog as the independent variable.
3 plot(inc.sub$difflog, inc.sub$voteshare,
4       xlab = "difflog",
5       ylab = "voteshare",
6       main = "Scatterplot with Regression Line")
7 #I use the abline function to add the regression line per https://www.
   geeksforgeeks.org/adding-straight-lines-to-a-plot-in-r-programming-
   abline-function/
8 abline(regmod_a, col = "red") #I made it red to make it clear.
9 #as many observations are close to the regression line, it looks like a
   strong
10 #relationship and strong linearity.
```



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

```
1 residuals_a <- residuals(regmod_a)
2 summary(residuals_a) #my residuals are on this object. The symmetry
3 #between the minimum and maximum suggest a normally distributed model.
```

4. Write the prediction equation.

```
1 #I need to obtain the coefficients for the prediction equation.
2 coef_a <- coef(regmod_a)
3 #I need the intercept and slope of the coefficients to make the equation.
4 intercept_a <- coef_a[1]
5 slope_a <- coef_a[2]
6 #I make some research to find the code for the prediction. This provides
  detailed
7 #information: https://www.dataquest.io/blog/statistical-learning-for-
  predictive-modeling-r/
8 #I also used chatgpt to fox the errors I came across.
9 prediction_equation <- function(difflog) {
10   voteshare_prediction = intercept_a + slope_a * difflog
11   return(voteshare_prediction)
12 }
13 #I use the equation to check the voteshare value for a difflog value of
  1.42
14 prediction_equation(1.42)
15 #it predicts a voteshare value of 0.6381969. The prediction equation
  works.
```

Question 2

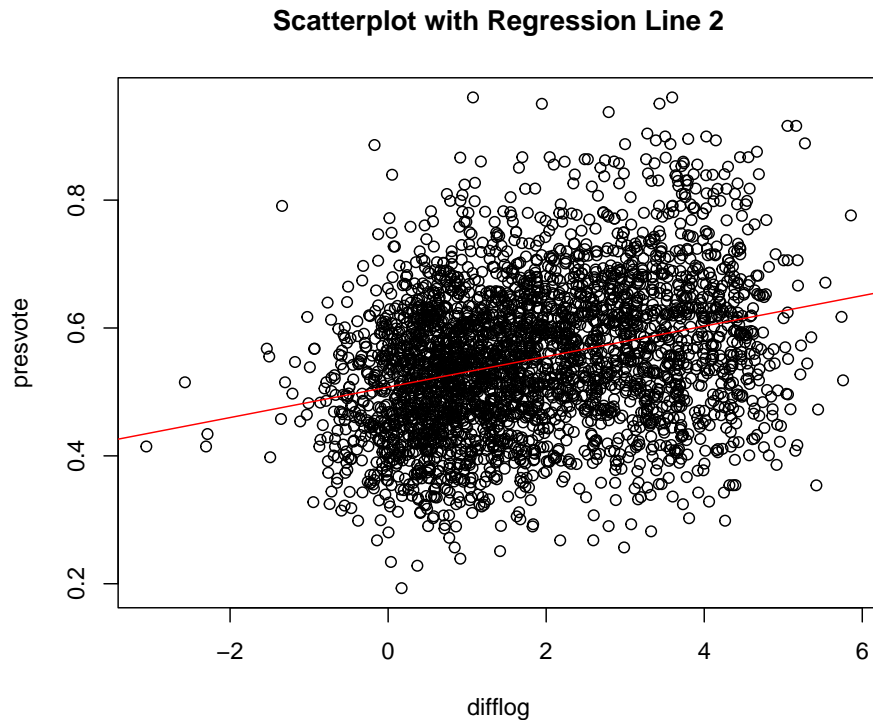
We are interested in knowing how the difference between incumbent and challenger's spending and the vote share of the presidential candidate of the incumbent's party are related.

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

```
1 regmod_b <- lm(presvote ~ difflog, data = inc.sub) #I run the regression
  model
2 #presvote is the dependent or outcome variable and difflog is the
  independent or
3 #explanatory variable.
4 summary(regmod_b) #p value is 0.0000000000000002, so I can reject the
  null
5 #hypothesis. Slope is 0.023837. In other words, for each unit change in
  difflog,
6 #voteshare is expected to increase by that amount. And it is a
  statistically
7 #significant result per the p value.
```

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

```
1 plot(inc.sub$difflog, inc.sub$presvote,
2       xlab = "difflog",
3       ylab = "presvote",
4       main = "Scatterplot with Regression Line 2")
5 #I make it just as the first one
6 abline(regmod_b, col = "red") #Again, I made it red to make it clear.
```



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

```
1 residuals_b <- residuals(regmod_b)
2 summary(residuals_b) #my residuals are on this object.
```

4. Write the prediction equation.

```
1 #I need to obtain the coefficients for the prediction equation.
2 coef_b <- coef(regmod_b)
3 #I need the intercept and slope of the coefficients to make the equation.
4 intercept_b <- coef_b[1]
5 slope_b <- coef_b[2]
6 prediction_equation_b <- function(difflog) {
7   voteshare_prediction_b = intercept_b + slope_b * difflog
8   return(voteshare_prediction_b)
9 }
```

```

10 #I use the equation to check the voteshare value for a difflog value of
    1.42
11 prediction_equation_b(1.42)
12 #it predicts a voteshare value of 0.5414322. The prediction equation
    works.

```

Question 3

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 regmod_c <- lm(voteshare ~ presvote, data = inc.sub)
2 summary(regmod_c)

```

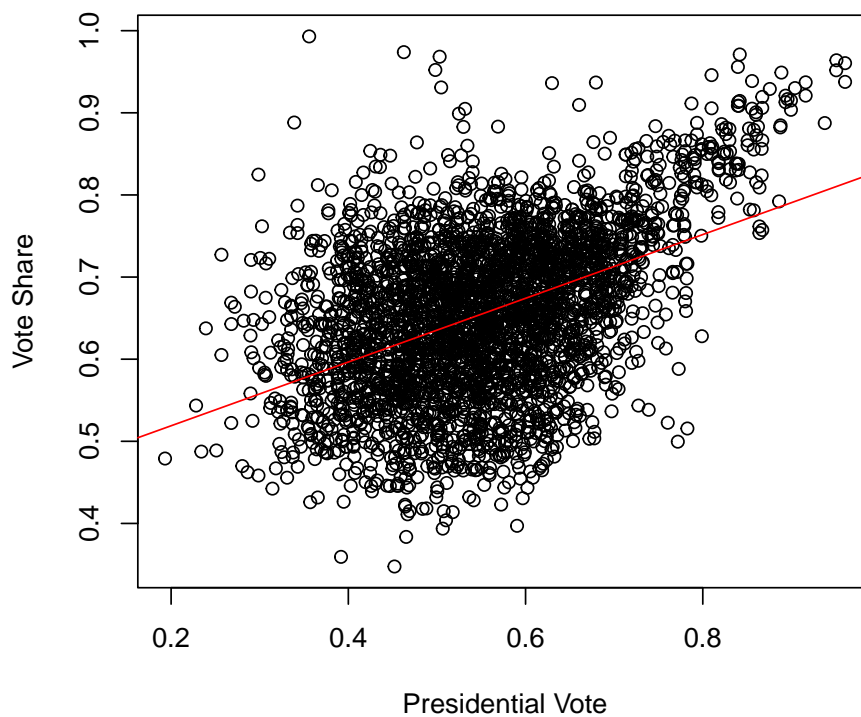
2. Make a scatterplot of the two variables and add the regression line.

```

1 plot(inc.sub$presvote, inc.sub$voteshare,
2      xlab = "Presidential Vote",
3      ylab = "Vote Share",
4      main = "Scatterplot of Vote Share vs. Presidential Vote")
5 abline(regmod_c, col = "red")

```

Scatterplot of Vote Share vs. Presidential Vote



3. Write the prediction equation.

```
1 #Again, #I need to obtain the coefficients for the prediction equation.
2 coef_c <- coef(regmod_c)
3 #I need the intercept and slope of the coefficients to make the equation.
4 intercept_c <- coef_c[1]
5 slope_c <- coef_c[2]
6 prediction_equation_c <- function(difflog) {
7   voteshare_prediction_c = intercept_c + slope_c * difflog
8   return(voteshare_prediction_c)
9 }
10 #I use the equation to check the voteshare value for a difflog value of
11   1.42
12 prediction_equation_c(1.42) #it returns 0.99 and it works.
```

Question 4

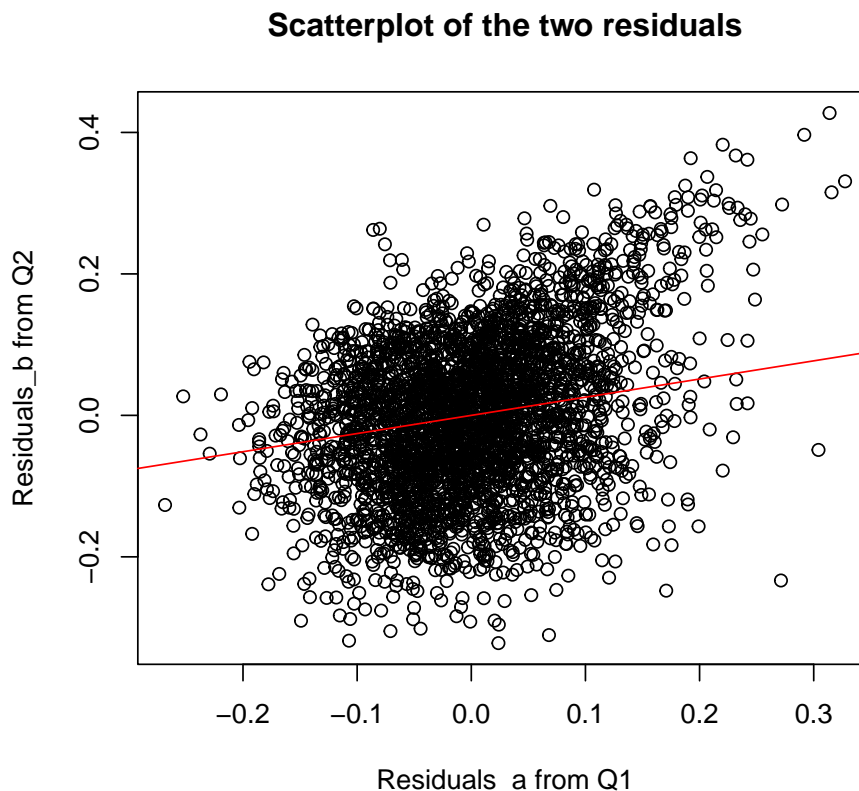
The residuals from part (a) tell us how much of the variation in **voteshare** is *not* explained by the difference in spending between incumbent and challenger. The residuals in part (b) tell us how much of the variation in **presvote** is *not* explained by the difference in spending between incumbent 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 regmod_d <- lm(residuals_a ~ residuals_b, data = inc.sub)
2 summary(regmod_d)
```

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

```
1 plot(residuals_a, residuals_b,
2       xlab = "Residuals_a from Q1",
3       ylab = "Residuals_b from Q2",
4       main = "Scatterplot of the two residuals")
5 abline(regmod_d, col = "red")
```



3. Write the prediction equation.

```
1 #Again, #I need to obtain the coefficients for the prediction equation.
```

Question 5

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 regmod_e <- lm(voteshare ~ difflog + presvote, data = inc.sub) #this is a
  multi
2 #variate regression, so I edit the code accordingly.
3 summary(regmod_e) #per p values this is a good fit of a model.
```

2. Write the prediction equation.

```
1 #again, I start with the coefficients.
2 coef_e <- coef(regmod_e)
3 #I need to assign the intercept and slopes to variables:
4 intercept_e <- coef_e[1]
5 slopee_difflog <- coef_e[2]
6 slopee_presvote <- coef_e[3]
7
8 prediction_equation_e <- function(difflog, presvote) {
9   voteshare_prediction_e = intercept_e + slopee_difflog * difflog +
10     slopee_presvote * presvote
11   return(voteshare_prediction_e)
12 }
13 #I test my model. The first input is difflog and the second is presvote
14 #to predict
15 #voteshare
16 prediction_equation_e(1.42, 0.44)
17 #it gives me 0.61. It works.
```

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 residuals are exactly the same. I think it is because both regressions are checking the effect of `difflog`. The one in Question 4 does it indirectly through the residuals of models that already accounted `difflog`, while the multi variate model does it directly.