

Problem Set 3

Applied Stats/Quant Methods 1

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`.

```
1 fit_reg1=lm(voteshare ~ difflog , data=data_ps3)
2 summary(fit_reg1)
3 ####
4 ##### Answer/ R_code output ##
5 ####
6 #Call:
7 #lm(formula = voteshare ~ difflog , data = data_ps3)
8
9 #Residuals:
10 #   Min       1Q   Median       3Q      Max
11 #-0.26832 -0.05345 -0.00377  0.04780  0.32749
```

```

12
13 #Coefficients:
14 #   Estimate Std. Error t value Pr(>|t|)
15 #(Intercept) 0.579031    0.002251  257.19   <2e-16 ***
16 #   difflog    0.041666    0.000968   43.04   <2e-16 ***
17 #   -----
18 #   Signif. codes:  0      ***      0.001      **      0.01      *      0.05      .
19 #   0.1      1
20 #Residual standard error: 0.07867 on 3191 degrees of freedom
21 #Multiple R-squared:  0.3673, Adjusted R-squared:  0.3671
22 #F-statistic: 1853 on 1 and 3191 DF,  p-value: < 2.2e-16

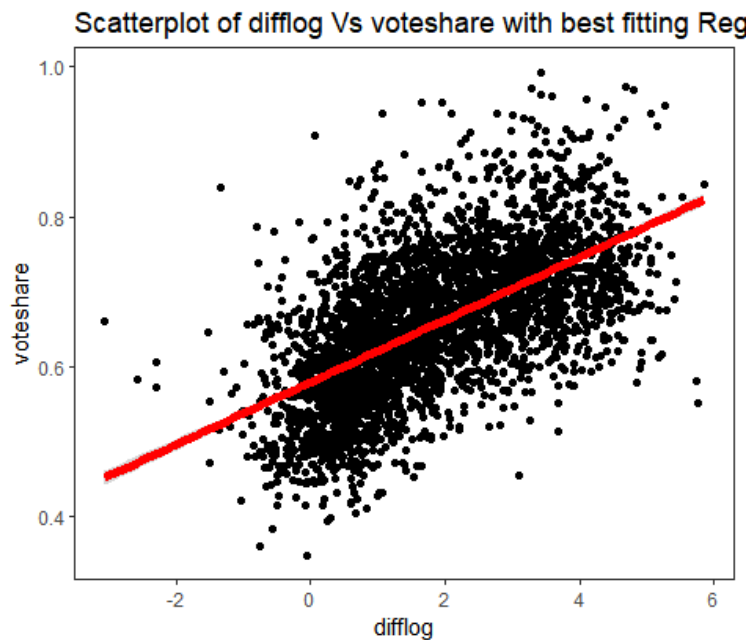
```

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

```

1 # Method: Using ggplot2
2
3 ggplot(data_ps3, aes(x=difflog, y=voteshare)) +
4   geom_point() +
5   geom_jitter() +
6   geom_smooth(method="lm", formula=y ~ x, se=T, color="red", lwd=2) +
7   theme_bw() + theme(panel.grid=element_blank()) +
8   ggtitle("Scatterplot of difflog Vs voteshare with best fitting
9           Regression Line")

```



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

```

1 # 1.3 Save the residuals of the model in a separated object:
2
3 my_residual_q1 <-residuals(fit_reg1)
4 head(my_residual_q1) # The first 6 observations of residuals:
5 # Answer:
6 #      1      2      3      4      5
7 # -0.0004227622 -0.0316840149 -0.0045514943  0.0386688767  0.0355287965
8 # 0.0322832521
9 ##
10 tail(my_residual_q1) # The last 6 observations of residuals:
11 ##### Answer:
12 #      3188      3189      3190      3191      3192
13 # 0.018604721  0.048283877  0.023159323 -0.040639860 -0.065834625
14 # 0.007829042

```

4. Write the prediction equation.

```

1 ##### 1.4 Write the Predicted Equation:
2
3 # y_hat= Beta_0_hat + Beta_1_hat* difflog imply that voteshare^ =
4 # 0.579031 + 0.041666*difflog_i
5 ##
6 # Here We can see clearly that the estimated slope Beta1_hat=0.041666,
7 # means that for every one unit increase in the
8 # difference in campaign spending in favor of the incumbent, the
9 # estimated vote share for the incumbent is expected
10 # to increase on average by 0.041666.
11 ##
12 #The estimated y-intercept or beta0_hat=0.579031 means when the spending
13 # difference is zero (difflog=0) the estimated
14 #vote share(vote_share) for the incumbent is equal to 0.579031.
15 #####
16 # Note: This implies that on average, an increase in campaign spending by
17 # the incumbent compared to the challenger is
18 # associated with an increase in the incumbent's vote share.

```

Question 2

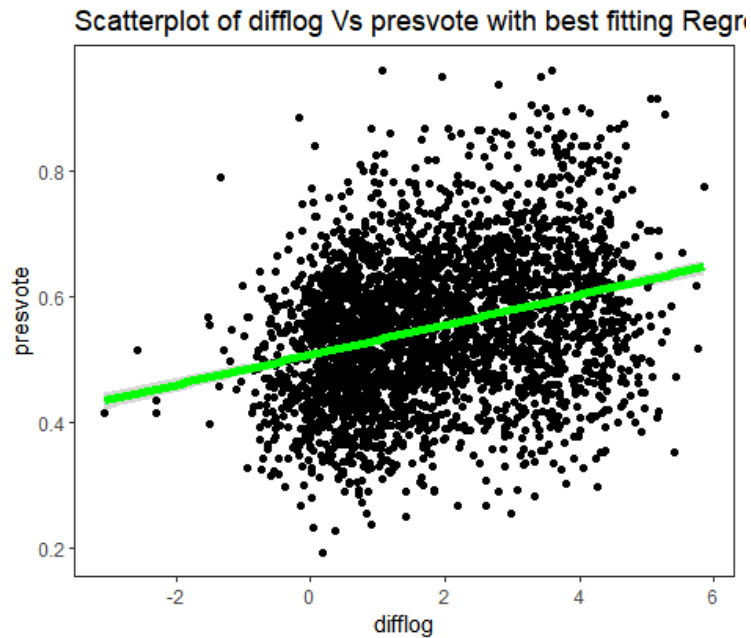
We are interested in knowing how the difference between incumbent and challenger's spending and the vote share of the presidential candidate of the incumbent's party are related.

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

```
1 # Question 2.1 Run a regression where the outcome variable is presvote
  and the explanatory variable is difflog.
2 #### Answer:
3 fit_reg2=lm(presvote ~ difflog , data=data_ps3)
4 summary(fit_reg2)
5 ##### Answer/ R-code output ##
6 #####
7 #Call:
8 #lm(formula = presvote ~ difflog , data = data_ps3)
9
10 #Residuals:
11 #   Min       1Q   Median       3Q      Max
12 # -0.32196 -0.07407 -0.00102  0.07151  0.42743
13
14 #Coefficients:
15 #   Estimate Std. Error t value Pr(>|t|)
16 # (Intercept)  0.507583    0.003161  160.60  <2e-16 ***
17 #   difflog      0.023837    0.001359   17.54  <2e-16 ***
18 #   ---
19 #   Signif. codes:  0   ***    0.001   **    0.01   *    0.05   .
20 #   0.1      1
21
22 #Residual standard error: 0.1104 on 3191 degrees of freedom
23 #Multiple R-squared:  0.08795, Adjusted R-squared:  0.08767
24 #F-statistic: 307.7 on 1 and 3191 DF, p-value: < 2.2e-16
```

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

```
1 #Question 2.2 Make a scatterplot of the two variables and add the
  regression line
2 #### Answer:
3 # Method: Using ggplot2
4
5 ggplot(data_ps3, aes(x=difflog , y=presvote)) +
6   geom_point() +
7   geom_jitter() + # I use the jitter() function to avoid overlapping the
  points
8   geom_smooth(method="lm" , formula=y ~ x, se=T, color="green" , lwd=2) +
9   theme_bw() + theme(panel.grid=element_blank()) +
10  ggtitle("Scatterplot of difflog Vs presvote with best fitting
  Regression Line")
```



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

```

1 # Question 2. 3 Save a Residuals of the model in a separate object.
2 ####
3 my_residual_q2 <- residuals(fit_reg2)
4 head(my_residual_q2) # The first 6 observations of residuals:
5 #### Answer/R-Output of Residuals
6
7 #      1      2      3      4      5
8 # 0.005605594 0.037578519 -0.053134788 -0.052993694
9 # -0.045842994 0.074339701
10 #####
11 # The Last six observations of the residuals models:
12 ####
13 tail(my_residual_q2)
14 ### Answer/ R output
15 ###

```

15	#	3188	3189	3190	3191	3192
		3193				
16	#	-0.017727276	-0.033198949	-0.002119851	0.032545042	
		0.036938994	0.035795200			

4. Write the prediction equation.

```

1 # Question 2. 4 Write the prediction equation.
2 ####
3 #### Answer:
4 ##
5 # presvote_hat=0.507583 + 0.023837*difflog
6 ##
7 # The estimated slope beta1_hat=0.023837 means that for every one unit
  increase in the difference in campaign spending
8 #in favor of the incumbent, the estimated vote share for the presidential
  candidate of the incumbent's party is
9 # expected to increase on average by 0.023837.
10 #####
11 #The estimated y-intercept=beta0_hat=0.507583 means when the spending
  difference is zero (i.e., difflog=0), the
12 # estimated vote share for the presidential candidate of the incumbent's
  party (presvote = 0.507583).
13 ####
14 # Note: This implies that, on average, an increase in campaign spending
  by the incumbent compared to the challenger
15 # is associated with an increase in the vote share of the incumbent's
  party's presidential candidate.

```

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.

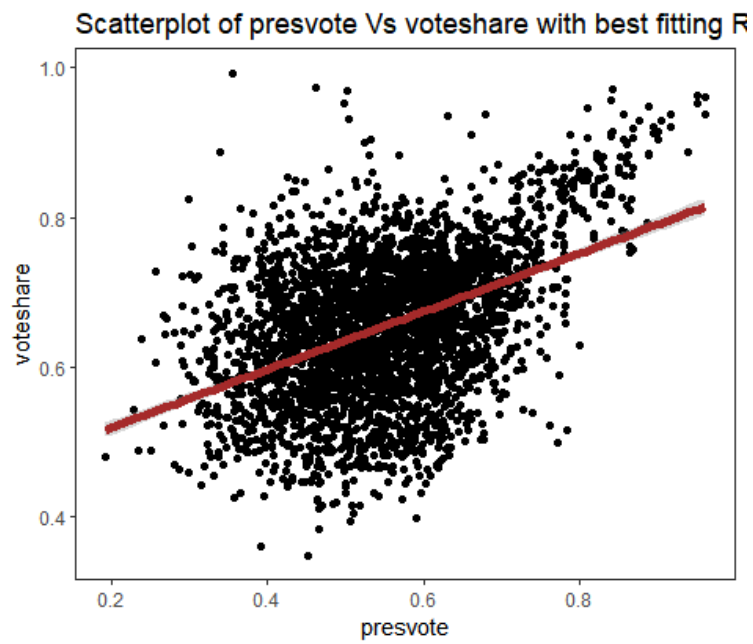
1. Run a regression where the outcome variable is `voteshare` and the explanatory variable is `presvote`.

```
1 # Question 3. 1 Run a regression where the outcome variable is voteshare
  and the explanatory variable is presvote.
2 ###
3 fit_reg3=lm(voteshare ~ presvote, data=data_ps3)
4 summary(fit_reg3)
5 # Answer / R Outcome ###
6 ##
7 #Call:
8 #lm(formula = voteshare ~ presvote, data = data_ps3)
9
10 #Residuals:
11 #   Min       1Q   Median       3Q      Max
12 # -0.27330 -0.05888  0.00394  0.06148  0.41365
13
14 #Coefficients:
15 #   Estimate Std. Error t value Pr(>|t|)
16 # (Intercept)  0.441330   0.007599   58.08  <2e-16 ***
17 #   presvote    0.388018   0.013493   28.76  <2e-16 ***
18 #   ---
19 #   Signif. codes:  0   ***    0.001   **    0.01   *    0.05   .
20 #   0.1      1
21
22 #Residual standard error: 0.08815 on 3191 degrees of freedom
23 #Multiple R-squared:  0.2058, Adjusted R-squared:  0.2056
24 #F-statistic: 827 on 1 and 3191 DF, p-value: < 2.2e-16
```

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

```
1 # Question 3.2 Make a scatterplot of the two variables and add the
  regression line.
2 ###
3 ## Answer / R -Output.
4 #
5 ### Answer:
6 # Method: Using ggplot2
7
8 ggplot(data_ps3, aes(x=presvote, y=voteshare)) +
9   geom_point() +
10  geom_jitter() +
11  geom_smooth(method="lm", formula=y ~ x, se=T, color="brown", lwd=2) +
12  theme_bw() + theme(panel.grid=element_blank()) +
```

```
13 ggtitle("Scatterplot of presvote Vs voteshare with best fitting  
14 Regression Line")  
14 #####
```



3. Write the prediction equation.

```
1 ##### Question 3.3 Write the prediction equation
2 #####
3 # Answer:
4
5 # voteshare_hat= 0.441330 + 0.388018*presvote
6 #####
7 # The estimated slope=beta1_hat=0.388018 means for every one unit
   increase in the vote share of the presidential
8 #candidate of the incumbent's party, the estimated vote share for the
   incumbent is expected to increase on average
9 # by 0.388018
10 #####
11 # The y-intercept=beta0_hat=0.441330 means when the vote share of the
   presidential candidate is zero (presvote=0),
12 # the estimate vote share for the incumbent is equal to 0.441330.
```

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.

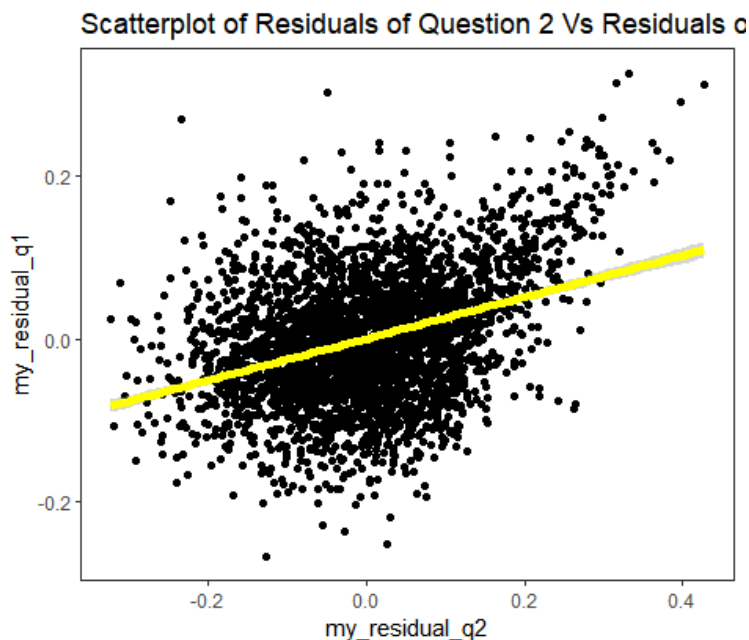
```
1 #### Answer:
2 #####
3 # recall both residuals:
4 my_residual_q1 # Residual from Question 1 (Outcome variable)
5 my_residual_q2 # Residual from Question 2 (Explanatory variable)
6 #####
7 fit_reg4 = lm(my_residual_q1 ~ my_residual_q2, data= data_ps3 )
8 summary(fit_reg4)
9 ##### Answer/R Output
10 ###
11 # Call:
12 #lm(formula = my_residual_q1 ~ my_residual_q2, data = data_ps3)
13
14 #Residuals:
15 #   Min       1Q   Median       3Q      Max
16 # -0.25928 -0.04737 -0.00121  0.04618  0.33126
17
18 #Coefficients:
19 #   Estimate Std. Error t value Pr(>|t|)
20 # (Intercept)  -5.934e-18  1.299e-03    0.00      1
21 # my_residual_q2  2.569e-01  1.176e-02   21.84 <2e-16 ***
22 # ---
23 #   Signif. codes:  0   ***    0.001   **    0.01   *    0.05   .
24 #   0.1      1
25
26 #Residual standard error: 0.07338 on 3191 degrees of freedom
27 #Multiple R-squared:  0.13, Adjusted R-squared:  0.1298
28 #F-statistic:  477 on 1 and 3191 DF, p-value: < 2.2e-16
```

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

```

1 #####Question 4.2 Make a scatterplot of the two residuals and add the
  regression line:
2 #####
3 # Method: ggplot2
4 #
5 # Method: Using ggplot2
6
7 ggplot(data_ps3, aes(x=my_residual_q2, y=my_residual_q1)) +
8   geom_point() +
9   geom_jitter() +
10  geom_smooth(method="lm", formula=y ~ x, se=T, color="yellow", lwd=2) +
11  theme_bw() + theme(panel.grid=element_blank()) +
12  ggtitle("Scatterplot of Residuals of Question 2 Vs Residuals of
  Question 1 with best fitting Regression Line")

```



3. Write the prediction equation.

```

1 # Question 4. 3 Write the prediction equation:
2 ##### Answer:
3 ###
4 # my_residual_q1_hat = -5.934e*(-18) + 2.569e*(-1)*my_residual_q2
5 #####
6 # The estimated slope=beta1_hat=2.569e*(-01) means that for each
  increase of one unit in residual of question 2
7 # (my_residual_q2), the value of residual of question 1 (my_residual_q1)
  is expected to increase by an average of
8 # 2.569e*(-1) units.
9 #####

```

```
10 # The y-intercept=beta0_hat=-5.934e**(-18) represent the value when
    residual of question 2 (i.e., my_residual_q2=0).
```

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 fit_reg5 = lm(voteshare ~ difflog + presvote, data=data_ps3)
2 summary(fit_reg5)
3 #####
4 # Answer/ RStudio Output
5 #####
6 #Call:
7 #lm(formula = voteshare ~ difflog + presvote, data = data_ps3)
8
9 #Residuals:
10 #   Min       1Q   Median       3Q      Max
11 #-0.25928 -0.04737 -0.00121  0.04618  0.33126
12
13 #Coefficients:
14 #   Estimate Std. Error t value Pr(>|t|)
15 #(Intercept)  0.4486442   0.0063297   70.88  <2e-16 ***
16 #   difflog     0.0355431   0.0009455   37.59  <2e-16 ***
17 #   presvote     0.2568770   0.0117637   21.84  <2e-16 ***
18 #   ---
19 #   Signif. codes:  0   ***    0.001   **    0.01   *    0.05   .
20 #   0.1      1
21
22 #Residual standard error: 0.07339 on 3190 degrees of freedom
23 #Multiple R-squared:  0.4496, Adjusted R-squared:  0.4493
24 #F-statistic: 1303 on 2 and 3190 DF, p-value: < 2.2e-16
```

2. Write the prediction equation.

```
1 # voteshare_hat = 0.4486442 + 0.0355431*difflog + 0.2568770*presvote
2 #
3 #The y-intercept=beta0_hat=0.4486442 represent the estimate vote share
  for the incumbent when both variables difflog
4 # and presvote are zero (i.e., difflog=presvote=0).In this context it
  represents the expected baseline vote share for
5 # the incumbent##
6 ##
7 #The coefficient for difflog (beta1_hat=0.0355431) means that for every
  one unit increase in the difference in
8 # spending in favor of the of the incumbent the estimated vote for the
  incumbent is expected to increase on average
9 # by 0.0355431 assuming the president's popularity remainn constant.
```

```

10 #####
11 # The coefficient for presvote (beta2_hat=0.2568770) means for every one
    unit increase in the president's popularity,
12 #the estimate vote share for the incumbent is expected to increase by
    0.2568770, assuming the difference in spending
13 # remain constant.

```

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

```

1 # Here the output of the both residuals are equal:
2 ###Residuals of question 4:
3 #   Min      1Q   Median      3Q      Max
4 # -0.25928 -0.04737 -0.00121  0.04618  0.33126
5 ###
6 ###Residuals of Question 5:
7 #   Min      1Q   Median      3Q      Max
8 # -0.25928 -0.04737 -0.00121  0.04618  0.33126
9 ##
10 #One of the possible explanation is related to perfect multicollinearity
    between the independent variables presvote
11 # and difflog that are perfectly correlated. On the other hand the
    difflog variable it's appeared in both models,
12 # this could lead to identical residuals.

```