

Problem Set 3

Applied Stats/Quant Methods 1

Due: November 13, 2025

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 Thursday November 13, 2025. 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 # read in data
2 inc.sub <- read.csv("https://raw.githubusercontent.com/ASDS-TCD/StatsI_
   2025/main/datasets/incumbents_subset.csv")
3
4 View(inc.sub)
5
6 #I first create a model which includes both variables and then summarise
   it.
7
```

```

8 model_q1 <- lm(voteshare ~ difflog, data = inc.sub)
9 summary(model_q1)

```

Table 1:

	<i>Dependent variable:</i>
	voteshare
difflog	0.042*** (0.001)
Constant	0.579*** (0.002)
Observations	3,193
R ²	0.367
Adjusted R ²	0.367
Residual Std. Error	0.079 (df = 3191)
F Statistic	1,852.791*** (df = 1; 3191)
<i>Note:</i>	*p<0.1; **p<0.05; ***p<0.01

The summary of this regression shows that the difference in campaign spending between the incumbent and the challenger (difflog) has, on average, a positive effect of approximately 0.042 on the incumbent's vote share (voteshare). This relationship is statistically significant.

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

```

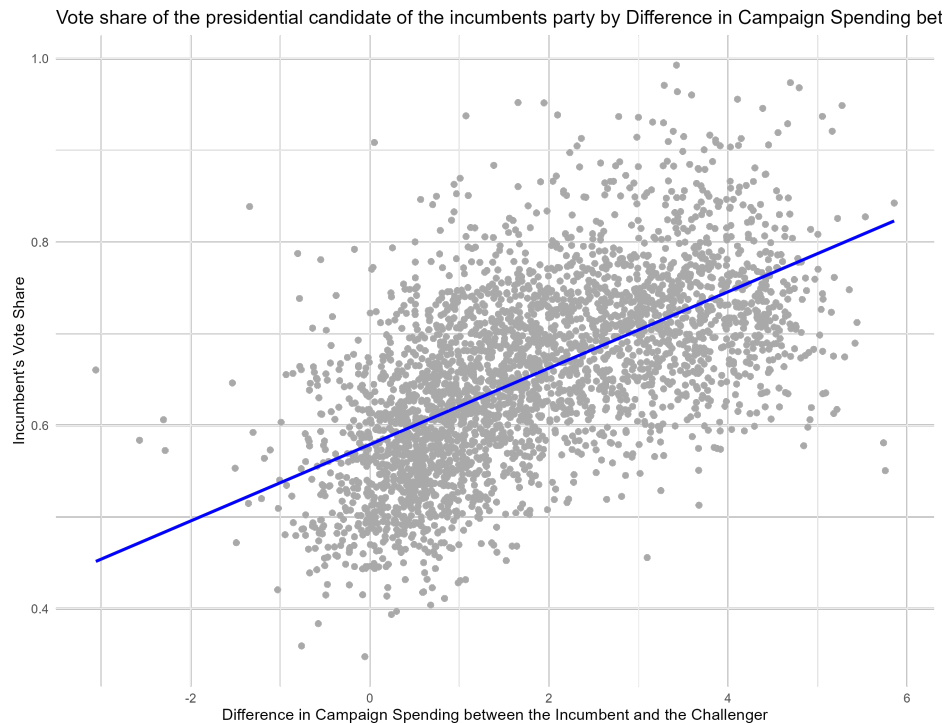
1 #To create this scatterplot of the two variables, I will load the system
2 #"ggplot2", which is know for having good visuals.
3
4 library(ggplot2)
5 library(stargazer)
6
7 #Then, I will create the scatterplot using, again, "difflog" as the
8   explanatory
9   #variable, and "voteshare" as the outcome variable.
10
11 ggplot1 <- ggplot(inc.sub, aes(x = difflog, y = voteshare)) +
12   geom_point(colour = "darkgrey") +
13   geom_smooth(method = "lm", colour = "blue", se = FALSE) +

```

```

13 labs(
14   title = "Vote share of the presidential candidate of the incumbents
15     party by Difference in Campaign Spending between the Incumbent and the
16     Challenge",
17   x = "Difference in Campaign Spending between the Incumbent and the
18     Challenger",
19   y = "Incumbent's Vote Share"
20 ) +
21 theme_minimal()

```



This scatterplot shows the positive correlation that we had identified with the regression analysis.

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

```

1 #I will create a new object using the function "residuals()" and I will
2   name it
3   # "residuals_model_q1"
4 residuals_model_q1 <- residuals(model_q1)
5 mean(residuals_model_q1)
6 residuals_model_q1

```

4. Write the prediction equation.

The general formula is $Y = a + bX$.

In this case, it is:

Incumbent vote share = Intercept + slope * difflog

According to the regression's results:

$$Y = 0.579031 + 0.041666 * X$$

$$\text{voteshare} = 0.579031 + 0.041666 * \text{difflog}$$

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 model_q2 <- lm(presvote ~ difflog, data = inc.sub)
2 summary(model_q2)
```

Table 2:

<i>Dependent variable:</i>	
	presvote
difflog	0.024*** (0.001)
Constant	0.508*** (0.003)
Observations	3,193
R ²	0.088
Adjusted R ²	0.088
Residual Std. Error	0.110 (df = 3191)
F Statistic	307.715*** (df = 1; 3191)
<i>Note:</i> *p<0.1; **p<0.05; ***p<0.01	

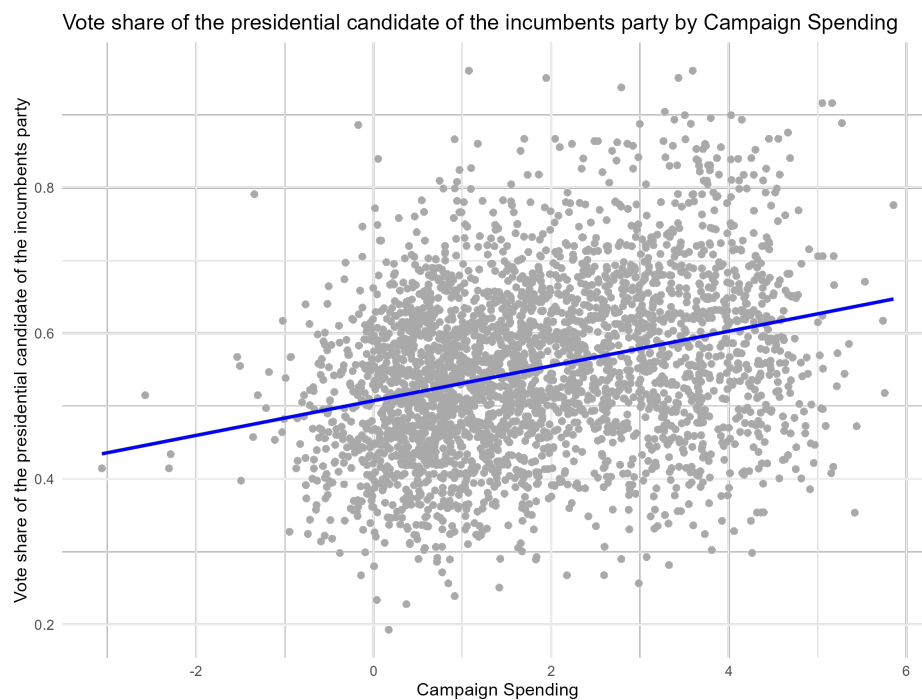
The summary of this regression shows that the difference in campaign spending between the incumbent and the challenger (`difflog`) has, on average, a positive effect of approximately 0.024 on the vote share of the presidential candidate of the incumbent's party (`presvote`). This relationship is statistically significant.

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

```

1 stargazer(model_q2)
2
3 ggplot2 <- ggplot(inc.sub, aes(x = difflog, y = presvote)) +
4   geom_point(colour = "darkgrey") +
5   geom_smooth(method = "lm", colour = "blue", se = FALSE) +
6   labs(
7     title = "Vote share of the presidential candidate of the incumbents
8     party by Campaign Spending",
9     x = "Campaign Spending",
10    y = "Vote share of the presidential candidate of the incumbents party
11    "
12  ) +
13  theme_minimal()

```



This scatterplot shows the positive correlation that we had identified with the regression analysis.

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

```

1 residuals_model_q2 <- residuals(model_q2)
2 mean(residuals_model_q2)
3 residuals_model_q2

```

4. Write the prediction equation.

The general formula is $Y = a + bX$.

In this case, it is:

Incumbent vote share = Intercept + slope * Difference in spending

According to the regression's results:

$$Y = 0.507583 + 0.023837 * X$$

$$\text{presvote} = 0.507583 + 0.023837 * \text{difflog}$$

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 model_q3 <- lm(voteshare ~ presvote, data = inc.sub)
2 summary(model_q3)
```

Table 3:

<i>Dependent variable:</i>	
	voteshare
presvote	0.388*** (0.013)
Constant	0.441*** (0.008)
Observations	3,193
R ²	0.206
Adjusted R ²	0.206
Residual Std. Error	0.088 (df = 3191)
F Statistic	826.950*** (df = 1; 3191)
<i>Note:</i> *p<0.1; **p<0.05; ***p<0.01	

The summary of this regression shows that the vote share of the presidential candidate of the incumbent's part (presvote) has, on average, a positive effect of approximately 0.388 on the incumbent's electoral success (voteshare). This relationship is statistically significant.

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

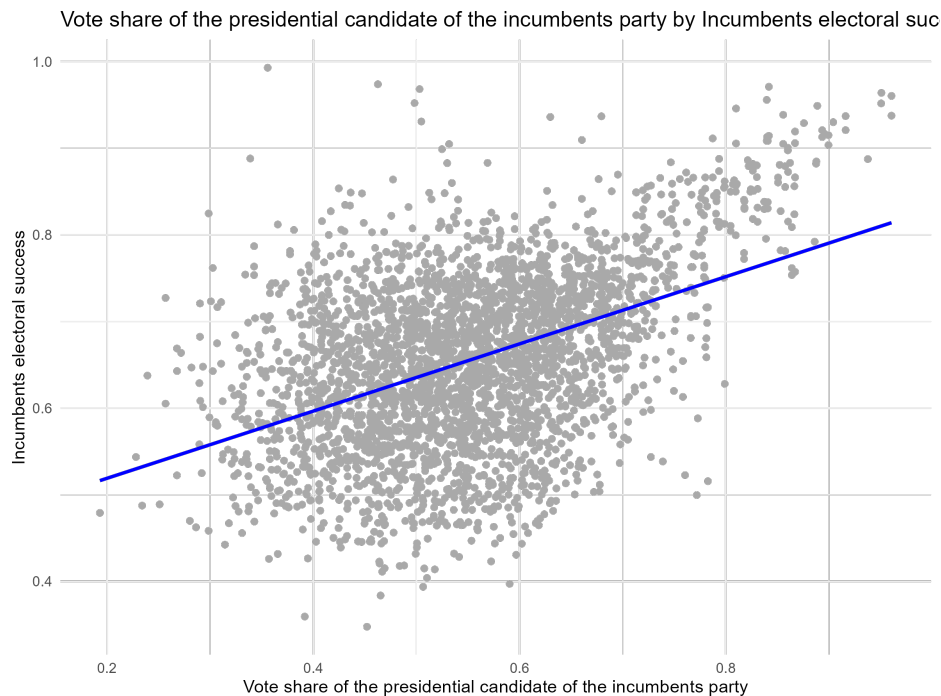
```
1 ggplot3 <- ggplot(inc.sub, aes(x = presvote, y = voteshare)) +
2   geom_point(colour = "darkgrey") +
3   geom_smooth(method = "lm", colour = "blue", se = FALSE) +
4   labs(
```



```

5   title = "Vote share of the presidential candidate of the incumbents
6   x = "Vote share of the presidential candidate of the incumbents party
7   y = "Incumbents electoral success"
8   ) +
9   theme_minimal()
10  ggsave("scatterplot3_PS03.png", plot = ggplot3, width = 8, height = 6,
        dpi = 300)

```



This scatterplot shows the positive correlation that we had identified with the regression analysis.

3. Write the prediction equation.

The general formula is $Y = a + bX$.

In this case, it is:

In this case, it is Incumbents electoral success = Intercept + slope * vote share of the presidential candidate

According to the regression's results:

$$Y = 0.441330 + 0.388018 * X$$

$$\text{voteshare} = 0.441330 + 0.388018 * \text{presvote}$$

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 model_residuais <- lm(residuals_model_q1 ~ residuals_model_q2, data = inc
  . sub)
2 summary(model_residuais)
```

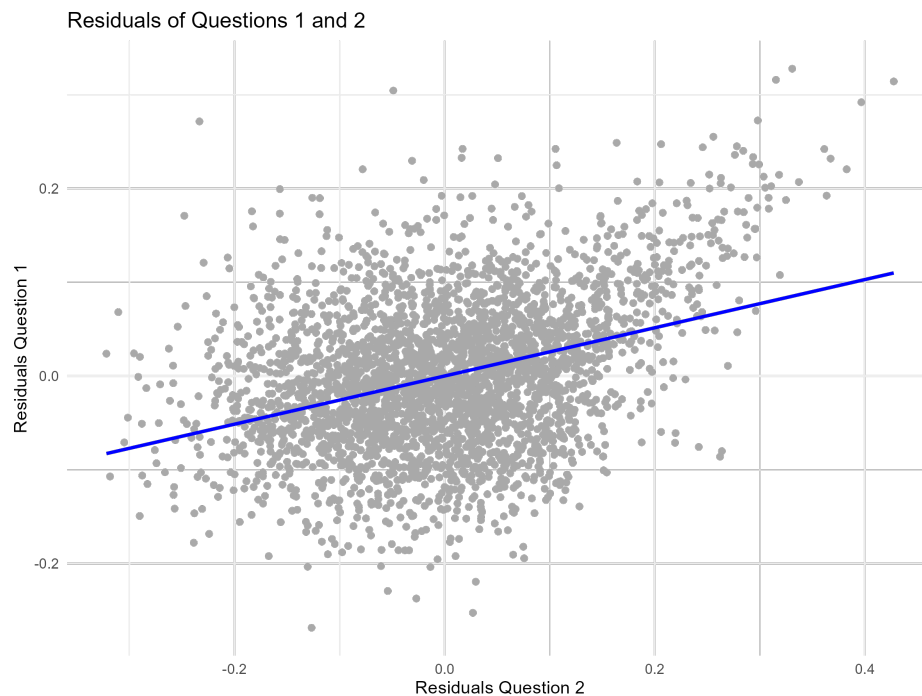
Table 4:

<i>Dependent variable:</i>	
	residuals_model_q1
residuals_model_q2	0.257*** (0.012)
Constant	−0.000 (0.001)
Observations	3,193
R ²	0.130
Adjusted R ²	0.130
Residual Std. Error	0.073 (df = 3191)
F Statistic	476.975*** (df = 1; 3191)
<i>Note:</i>	*p<0.1; **p<0.05; ***p<0.01

The summary of this regression shows that the residual variation in the incumbent's vote share **voteshare** that is not explained by campaign spending (**difflog**) is positively associated with the residual variation in the presidential candidate's vote share that is also (**presvote**) not explained by campaign spending. This indicates that, once the effect of **difflog** is accounted for, there remains a significant positive relationship (with an increase of 0.257 in residuals-model-q1 per one-unit increase in residuals-model-q2) between the two variables.

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

```
1 ggplot4 <- ggplot(inc.sub, aes(x = residuals_model_q2, y = residuals_
  model_q1)) +
2   geom_point(colour = "darkgrey") +
3   geom_smooth(method = "lm", colour = "blue", se = FALSE) +
4   labs(
5     title = "Residuals of Questions 1 and 2",
6     x = "Residuals Question 2",
7     y = "Residuals Question 1"
8   ) +
9   theme_minimal()
```



This scatterplot shows the positive correlation that we had identified with the regression analysis.

3. Write the prediction equation.

The general formula is $Y = a + bX$.

In this case, it is:

In this case, it is $\text{residuals-model-q1} = \text{Intercept} + \text{slope} * \text{residuals-model-q2}$

According to the regression's results:

```
Y = -5.520e-18 + 2.569e-01 * X  
residuals-model-q1 = -5.520e-18 + 2.569e-01 * residuals-model-q2
```

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?

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

```
1 model_q5 <- lm(voteshare ~ difflog + presvote, data = inc.sub)
2 summary(model_q5)
```

Table 5:

<i>Dependent variable:</i>	
	voteshare
difflog	0.036*** (0.001)
presvote	0.257*** (0.012)
Constant	0.449*** (0.006)
Observations	3,193
R ²	0.450
Adjusted R ²	0.449
Residual Std. Error	0.073 (df = 3190)
F Statistic	1,302.947*** (df = 2; 3190)
<i>Note:</i>	*p<0.1; **p<0.05; ***p<0.01

The summary of this regression shows that both the difference in campaign spending between the incumbent and the challenger (`difflog`) and the presidential candidate's vote share (`presvote`) have positive and statistically significant effects on the incumbent's vote share (`voteshare`). Holding `presvote` constant, a one-unit increase in `difflog` is associated with an average increase of 0.036 in the incumbent's vote share. At the same time, holding `difflog` constant, a one-unit increase in `presvote` corresponds to an average increase of 0.257 in `voteshare`. Importantly, this model explains about 45% of the variation in the incumbent's vote share, according to the R^2 value.

- (b) Write the prediction equation.

The general formula is $Y = a + bX$.

In this case, it is:

Incumbent's vote share = Intercept + slope1 * president's popularity
+ slope2 * difference in spending between incumbent and challenger

According to the regression's results:

$Y = 0.4486442 + 0.0355431 * X1 + 0.2568770 * X2$

$\text{voteshare} = 0.4486442 + 0.0355431 * \text{difflog} + 0.2568770 * \text{presvote}$

- (c) What is it in this output that is identical to the output in Question 4? Why do you think this is the case?

What is identical in this output and in Question 4's regression output are the values for residuals-model-q2 (in Question 4's regression) and presvote (in Question 5's regression). This happens because both models are basically accounting for the same correlation, although in different fashions.

In question 5 we regress voteshare on difflog and presvote. Therefore, the coefficient on presvote shows how much this variable has an effect on voteshare once we have controlled for the effect of difflog.

Question 4's model is, however, more complex. In this case, we've regressed residuals-model-q1 on residuals-model-q2. This means that, first, we have removed the effect of difflog on both variables voteshare and presvote, as we have taken the residuals. Then, the remaining variation on voteshare has to be explained by the remaining variation on presvote.