# 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`.

```r
1  # This analyse will be about a dataset called incumbents subset.csv
2  inc.sub <- read.csv("https://raw.githubusercontent.com/ASDS-TCD/StatsI_
      Fall2023/main/datasets/incumbents_subset.csv")
3
4  # Check out the head and tail.
5  head(inc.sub)
6  tail(inc.sub)
7
8  # We will select only the columns "voteshare" and "difflog" to create a
      dataframe with the relevant variables
9  selected_data <- inc.sub[, c("voteshare", "difflog")]
10
11 # Simple linear regression
12 model <- lm(voteshare ~ difflog, data = selected_data)
13
14 # Summary of the regression model
15 summary(model)
```

```
Call:
lm(formula = voteshare ~ difflog, data = selected_data)

Residuals:
     Min       1Q   Median       3Q      Max
-0.26832 -0.05345 -0.00377  0.04780  0.32749

Coefficients:
             Estimate Std. Error t value Pr(>|t|)
(Intercept) 0.579031   0.002251  257.19   <2e-16 ***
difflog     0.041666   0.000968   43.04   <2e-16 ***
---
Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

Residual standard error: 0.07867 on 3191 degrees of freedom
Multiple R-squared:  0.3673,     Adjusted R-squared:  0.3671
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.

```r
1    plot(inc.sub$difflog, inc.sub$voteshare)
2  abline(model, col = "red")
```

```
Call:
lm(formula = voteshare ~ difflog + presvote, data = selected_data)

Residuals:
     Min       1Q   Median       3Q      Max
-0.25928 -0.04737 -0.00121  0.04618  0.33126

Coefficients:
             Estimate Std. Error t value Pr(>|t|)
(Intercept) 0.4486442  0.0063297   70.88   <2e-16 ***
difflog     0.0355431  0.0009455   37.59   <2e-16 ***
presvote    0.2568770  0.0117637   21.84   <2e-16 ***
---
Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

Residual standard error: 0.07339 on 3190 degrees of freedom
Multiple R-squared:  0.4496,    Adjusted R-squared:  0.4493
F-statistic:  1303 on 2 and 3190 DF,  p-value: < 2.2e-16
```

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

```
residuals <- resid(model)

# Summary of the regression model
summary(model)

# Residuals
head(residuals)
```

4. Write the prediction equation.

# 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`.

```r
# We will select only the variables "presvote" and "difflog"
selected_data <- inc.sub [, c("presvote", "difflog")]
# Linear regression model
model_presvote <- lm(presvote ~ difflog, data = inc.sub)
# Summary
summary(model_presvote)
```

```
Call:
lm(formula = presvote ~ difflog, data = inc.sub)

Residuals:
     Min       1Q   Median       3Q      Max
-0.32196 -0.07407 -0.00102  0.07151  0.42743

Coefficients:
             Estimate Std. Error t value Pr(>|t|)
(Intercept) 0.507583   0.003161  160.60   <2e-16 ***
difflog     0.023837   0.001359   17.54   <2e-16 ***
---
Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

Residual standard error: 0.1104 on 3191 degrees of freedom
Multiple R-squared:  0.08795,   Adjusted R-squared:  0.08767
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.

```r
plot(inc.sub$difflog, inc.sub$presvote)
# Regression
abline(model_presvote, col = "red")

```

```
Call:
lm(formula = voteshare ~ difflog + presvote, data = selected_data)

Residuals:
     Min       1Q    Median       3Q      Max
-0.25928 -0.04737 -0.00121  0.04618  0.33126

Coefficients:
             Estimate Std. Error t value Pr(>|t|)
(Intercept) 0.4486442  0.0063297   70.88   <2e-16 ***
difflog     0.0355431  0.0009455   37.59   <2e-16 ***
presvote    0.2568770  0.0117637   21.84   <2e-16 ***
---
Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

Residual standard error: 0.07339 on 3190 degrees of freedom
Multiple R-squared:  0.4496,    Adjusted R-squared:  0.4493
F-statistic:  1303 on 2 and 3190 DF,  p-value: < 2.2e-16
```

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

```
1 # Residuals
2 residuals_presvote <- resid(model_presvote)
3
4 # Summary of the regression model
5 summary(model_presvote)
6
7 # Residuals
8 head(residuals_presvote)
```

4. Write the prediction equation.

# 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`.

```r
# We will select the columns 'voteshare' and 'presvote'
selected_data <- inc.sub[, c("voteshare", "presvote")]

# Linear regression model
model_voteshare <- lm(voteshare ~ presvote, data = selected_data)

# Summary of the regression model
summary(model_voteshare)
```

```
Call:
lm(formula = voteshare ~ presvote, data = selected_data)

Residuals:
     Min       1Q   Median       3Q      Max
-0.27330 -0.05888  0.00394  0.06148  0.41365

Coefficients:
            Estimate Std. Error t value Pr(>|t|)
(Intercept) 0.441330   0.007599   58.08   <2e-16 ***
presvote    0.388018   0.013493   28.76   <2e-16 ***
---
Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

Residual standard error: 0.08815 on 3191 degrees of freedom
Multiple R-squared:  0.2058,     Adjusted R-squared:  0.2056
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.

```r
# Scatterplot
plot(selected_data$presvote, selected_data$voteshare)

# Regression
abline(model_voteshare, col = "red")
```

```
Call:
lm(formula = voteshare ~ difflog + presvote, data = selected_data)

Residuals:
     Min       1Q   Median       3Q      Max
-0.25928 -0.04737 -0.00121  0.04618  0.33126

Coefficients:
             Estimate Std. Error t value Pr(>|t|)
(Intercept) 0.4486442  0.0063297   70.88   <2e-16 ***
difflog     0.0355431  0.0009455   37.59   <2e-16 ***
presvote    0.2568770  0.0117637   21.84   <2e-16 ***
---
Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

Residual standard error: 0.07339 on 3190 degrees of freedom
Multiple R-squared:  0.4496,    Adjusted R-squared:  0.4493
F-statistic:  1303 on 2 and 3190 DF,  p-value: < 2.2e-16
```

3. Write the prediction equation.

# 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.

```
# 4 The residuals from Question 1 are stored in residuals and residuals
    from Question 2 are stored in residuals_presvote
residuals_model <- lm(residuals ~ residuals_presvote)

# Summary of the regression model
summary(residuals_model)
```

```
Call:
lm(formula = residuals ~ residuals_presvote)

Residuals:
     Min       1Q   Median       3Q      Max
-0.25928 -0.04737 -0.00121  0.04618  0.33126

Coefficients:
                     Estimate Std. Error t value Pr(>|t|)
(Intercept)        -5.934e-18  1.299e-03    0.00        1
residuals_presvote  2.569e-01  1.176e-02   21.84   <2e-16 ***
---
Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

Residual standard error: 0.07338 on 3191 degrees of freedom
Multiple R-squared:  0.13,     Adjusted R-squared:  0.1298
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.

```
# Scatterplot with regression line
plot(residuals_presvote, residuals)
abline(residuals_model, col = "red")
```

```
Call:
lm(formula = voteshare ~ difflog + presvote, data = selected_data)

Residuals:
     Min       1Q    Median       3Q      Max
-0.25928 -0.04737 -0.00121  0.04618  0.33126

Coefficients:
             Estimate Std. Error t value Pr(>|t|)
(Intercept) 0.4486442  0.0063297   70.88   <2e-16 ***
difflog     0.0355431  0.0009455   37.59   <2e-16 ***
presvote    0.2568770  0.0117637   21.84   <2e-16 ***
---
Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

Residual standard error: 0.07339 on 3190 degrees of freedom
Multiple R-squared:  0.4496,    Adjusted R-squared:  0.4493
F-statistic:  1303 on 2 and 3190 DF,  p-value: < 2.2e-16
```

3. Write 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 # We will select the columns 'voteshare', 'difflog', and 'presvote'
2 selected_data <- inc.sub[, c("voteshare", "difflog", "presvote")]
3
4 # Linear regression model
5 model_combined <- lm(voteshare ~ difflog + presvote, data = selected_data
    )
6
7 # Summary of the regression model
8 summary(model_combined)
```

```
Call:
lm(formula = voteshare ~ difflog + presvote, data = selected_data)

Residuals:
     Min       1Q   Median       3Q      Max
-0.25928 -0.04737 -0.00121  0.04618  0.33126

Coefficients:
             Estimate Std. Error t value Pr(>|t|)
(Intercept) 0.4486442  0.0063297   70.88   <2e-16 ***
difflog     0.0355431  0.0009455   37.59   <2e-16 ***
presvote    0.2568770  0.0117637   21.84   <2e-16 ***
---
Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

Residual standard error: 0.07339 on 3190 degrees of freedom
Multiple R-squared:  0.4496,    Adjusted R-squared:  0.4493
F-statistic:  1303 on 2 and 3190 DF,  p-value: < 2.2e-16
```

2. Write the prediction equation.

Write the prediction equation.

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