

# Problem Set 4

QTM 200: Applied Regression Analysis

Due: February 24, 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 24, 2020. No late assignments will be accepted.
- Total available points for this homework is 100.

## Question 1 (50 points): Economics

In this question, use the **prestige** dataset in the **car** library. First, run the following commands:

```
install.packages(car)
library(car)
data(Prestige)
help(Prestige)
```

We would like to study whether individuals with higher levels of income have more prestigious jobs. Moreover, we would like to study whether professionals have more prestigious jobs than blue and white collar workers.

- (a) Create a new variable **professional** by recoding the variable **type** so that professionals are coded as 1, and blue and white collar workers are coded as 0 (Hint: **ifelse**.)

```
1 #Q1
2 P<-Prestige#rename dataset
3 #Create a new variable
4 professional <- ifelse(P$type == "prof", "1", "0")
```

- (b) Run a linear model with prestige as an outcome and **income**, **professional**, and the interaction of the two as predictors (Note: this is a continuous  $\times$  dummy interaction.)

```
1 #Run a linear model
2 lm1 <- lm(prestige ~ income + professional + income:professional, data=P)
3 summary(lm1)
4 #Coefficients:
5 #Estimate Std. Error t value Pr(>|t|)
6 #(Intercept)      21.1422589    2.8044261     7.539 2.93e-11 ***
7 # income           0.0031709    0.0004993     6.351 7.55e-09 ***
8 # professional1    37.7812800    4.2482744     8.893 4.14e-14 ***
9 # income:professional1 -0.0023257    0.0005675    -4.098 8.83e-05 ***
10 #Residual standard error: 8.012 on 94 degrees of freedom
11 #(4 observations deleted due to missingness)
12 #Multiple R-squared:  0.7872, Adjusted R-squared:  0.7804
13 #F-statistic: 115.9 on 3 and 94 DF, p-value: < 2.2e-16
```

- (c) Write the prediction equation based on the result.

```
1 #y = 21.14 + 0.003Xi + 37.78Di - 0.002XiDi + epsiloni (with sigma square
  as 0,8.012^2)
```

(d) Interpret the coefficient for **income**.

```
1 #For professional , the regression line would result in a 37.78 unit  
   upward shift , regardless of income. No effect for white and blue  
   collar workers
```

(e) Interpret the coefficient for **professional**.

```
1 #For professional , the regression line would result in a 37.78 unit  
   upward shift , regardless of income. No effect for white and blue  
   collar workers
```

- (f) What is the effect of a \$1,000 increase in income on prestige score for professional occupations? In other words, we are interested in the marginal effect of income when the variable **professional** takes the value of 1. Calculate the change in  $\hat{y}$  associated with a \$1,000 increase in income based on your answer for (c).

```
1 #marginal difference for $1,000 increase in income
2 #y = 21.14 + 0.003Xi + 37.78Di - 0.002XiDi + epsiloni
3 #when professional=Di=1
4 #y1 = 21.14 + 0.003(Xi+1000) + 37.78 - 0.002(Xi+1000) + epsiloni
5 #y^=0.003*1000-0.002*1000=1
```

- (g) What is the effect of changing one's occupations from non-professional to professional when her income is \$6,000? We are interested in the marginal effect of professional jobs when the variable **income** takes the value of 6,000. Calculate the change in  $\hat{y}$  based on your answer for (c).

```
1 #marginal difference for income of $6,000
2 #y = 21.14 + 0.003Xi + 37.78Di - 0.002XiDi + epsiloni
3 #when income=Xi=6000
4 #y2 = 21.14 + 0.003*6000 + 37.78Di - 0.002*6000Di + epsiloni
5 #y^=37.781-0.002*6000=25.781
```

## Question 2 (50 points): Political Science

Researchers are interested in learning the effect of all of those yard signs on voting preferences.<sup>1</sup> Working with a campaign in Fairfax County, Virginia, 131 precincts were randomly divided into a treatment and control group. In 30 precincts, signs were posted around the precinct that read, “For Sale: Terry McAuliffe. Don’t Sellout Virginia on November 5.”

Below is the result of a regression with two variables and a constant. The dependent variable is the proportion of the vote that went to McAuliffe’s opponent Ken Cuccinelli. The first variable indicates whether a precinct was randomly assigned to have the sign against McAuliffe posted. The second variable indicates a precinct that was adjacent to a precinct in the treatment group (since people in those precincts might be exposed to the signs).

Impact of lawn signs on vote share	
Precinct assigned lawn signs (n=30)	0.042 (0.016)
Precinct adjacent to lawn signs (n=76)	0.042 (0.013)
Constant	0.302 (0.011)

Notes:  $R^2=0.094$ ,  $N=131$

- (a) Use the results to determine whether having these yard signs in a precinct affects vote share (e.g., conduct a hypothesis test with  $\alpha = .05$ ).

```
1 #Q2
2 #determine having yard signs
3 #for a two-tailed hypothesis with df=n-k
4 #n=131 and k=3
5 (0.042 - 0) / 0.016 = 2.625
6 #pvalue for 2.625 is 0.00972 < 0.05, so it affects the vote share
```

---

<sup>1</sup>Donald P. Green, Jonathan S. Krasno, Alexander Coppock, Benjamin D. Farrer, Brandon Lenoir, Joshua N. Zingher. 2016. “The effects of lawn signs on vote outcomes: Results from four randomized field experiments.” *Electoral Studies* 41: 143-150.

- (b) Use the results to determine whether being next to precincts with these yard signs affects vote share (e.g., conduct a hypothesis test with  $\alpha = .05$ ).

```
1 #determine having yard signs next to precincts
2 #for a two-tailed hypothesis with df=n-k
3 #n=131 and k=3
4 (0.042-0)/0.013=3.231
5 #pvalue for 3.231 is 0.001568 < 0.05, so it affects the vote share
```

- (c) Interpret the coefficient for the constant term substantively.

```
1 #Since Precinct assigned lawn signs and Precinct adjacent to lawn signs
   are both 0 in this case, the proportion of the vote that went to Ken
   Cuccinelli, the constant is 0.302
```

- (d) Evaluate the model fit for this regression. What does this tell us about the importance of yard signs versus other factors that are not modeled?

```
1 #Evaluate the model fit for regression
2 #use F-test for the model
3 n=131
4 k=3
5 F.test<-((0.094/k)/((1-0.094)/(n-k-1)))
6 df1<-k
7 df2<-n-k-1
8 pvalue<-df(F.test, df1, df2)
9 #pvalue=0.007141957 < 0.05
10 #Thus we reject the null hypothesis and conclude that at least one
    coefficient is significant
```