

# Problem Set 6

QTM 200: Applied Regression Analysis

Due: May 6, 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 before midnight on Wednesday, May 6, 2020. No late assignments will be accepted.
- Total available points for this homework is 100.

## Question 1 (50 points): Biology

Load in the data labelled `cholesterol.csv` on GitHub, which contains an observational study of 315 observations.

- Response variable:
  - **cholCat**: 1 if the individual has high cholesterol; 0 if the individual does not have high cholesterol
- Explanatory variables:
  - **sex**: 1 Male; 0 Female
  - **fat**: grams of fat consumed per day

Please answer the following questions:

1. We are interested in predicting the cholesterol category based on sex and fat intake.
  - (a) Fit an additive model. Provide the summary output, the global null hypothesis, and  $p$ -value. Please describe the results and provide a conclusion.

```
1 #1)
2 model_additive <- lm(cholCat ~ sex + fat, data=cholesterol)
3 summary(model_additive)
4
5 #H0: Neither sex nor fat has an effect on the cholesterol.
6 #H1: At least one variable has a relationship with the cholesterol.
```

Since the  $p$ -values for both variables (sex and fat) are smaller than 0.05, we have enough evidence to reject the null and conclude that the two variables have relationship with the cholesterol.

2. If explanatory variables are significant in this model, then
  - (a) For women, how does increasing their fat intake by 1 gram per day change their odds on being in the high cholesterol group? (Interpretation of a coefficient)

```
1 cholesterol = -0.1303597 + 0.1894 * sex + 0.0082466 * fat
```

For females,  $\text{sex} = 0$ , the probability of having high cholesterol increases by 0.008 if their fat increases by 1g per day.

- (b) For men, how does increasing their fat intake by 1 gram per day change their odds on being in the high cholesterol group? (Interpretation of a coefficient)

```
1 cholesterol = -0.1303597 + 0.1894 * sex + 0.0082466 * fat
```

For males,  $\text{sex} = 1$ , the probability of having high cholesterol increases by 0.189 if their fat increases by 1g per day.

- (c) What is the estimated probability of a woman with a fat intake of 100 grams per day being in the high cholesterol group?

```
1 #sex = 0, fat = 100,
2 cholesterol = -0.1303597 + 0.1894 * 0 + 0.0082466 * 100 = 0.6943003
```

The estimated probability of a woman with a fat intake of 100 grams per day being in the high cholesterol group is 0.694.

(d) Would the answers to 2a and 2b potentially change if we included the interaction term in this model? Why?

- Perform a test to see if including an interaction is appropriate.

```
1 model_int <- lm(cholCat ~ sex + fat + sex*fat, data=cholesterol)
2 summary(model_int)
```

```
Residual standard error: 0.4022 on 311 degrees of freedom
Multiple R-squared:  0.3588, Adjusted R-squared:  0.3526
F-statistic: 58.01 on 3 and 311 DF,  p-value: < 2.2e-16
```

No, the answer would not change. After we add in the interaction term, we do not see a significant increase in the explanatory power.

## Question 2 (50 points): Political Economy

We are interested in how governments' management of public resources impacts economic prosperity. Our data come from Alvarez, Cheibub, Limongi, and Przeworski (1996) and is labelled `gdpChange.csv` on GitHub. The dataset covers 135 countries observed between 1950 or the year of independence or the first year for which data on economic growth are available ("entry year"), and 1990 or the last year for which data on economic growth are available ("exit year"). The unit of analysis is a particular country during a particular year, for a total  $> 3,500$  observations.

- Response variable:
  - `GDPWdiff`: Difference in GDP between year  $t$  and  $t-1$ . Possible categories include: "positive", "negative", or "no change"
- Explanatory variables:
  - `REG`: 1=Democracy; 0=Non-Democracy
  - `OIL`: 1=if the average ratio of fuel exports to total exports in 1984-86 exceeded 50%; 0= otherwise

Please answer the following questions:

1. Construct and interpret an unordered multinomial logit with `GDPWdiff` as the output and "no change" as the reference category, including the estimated cutoff points and coefficients.

```
1 gdp_output <- relevel(gdp$GDPWdiff, ref="no change")
2 model_unorder <- multinom(out ~ REG + OIL, data=gdp, Hess = T)
3 summary(model_unorder)
```

Coefficients:

(Intercept)	REG	OIL
negative	3.805370	1.379282
positive	4.533759	1.769007

2. Construct and interpret an ordered multinomial logit with `GDPWdiff` as the outcome variable, including the estimated cutoff points and coefficients.

```
1 gdp_out <- ordered(gdp_output, levels=c("negative", "no change", "positive"))
2 polr(gdp_out ~ REG + OIL, data=gdp, Hess=T)
```