# Problem Set 3

## Applied Stats II

## Due: March 28, 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 in `.pdf` form.

- This problem set is due before class on Monday March 28, 2022. No late assignments will be accepted.

- Total available points for this homework is 80.

## Question 1

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 forwhich 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  data <- read.csv("C:/Users/crowl/Documents/Trinity/ASDS Course/POP77003_
      Applied_Statistical_Analysis_2/ProblemSet/PS3/template/GDPchange.csv")
2
3
4  library(foreign)
5  library(nnet)
6  library(stargazer)
7  library(jtools)
8  library(MASS)
9  library(nnet)
10  library(ggplot2)
11  install.packages("pscl")
12  library(pscl)
13
14  summary(data)
```

Firslty loading in my data and inspecting it.

```
1  data$GDPWdiff <- cut(x = data$GDPWdiff, breaks = c(-9257, -0.9, 0.1,
      7867)) #creating bounds of values for the levels, as they are real
      numbers
2  levels(data$GDPWdiff) <- c("negative", "no_change", "positive")
3
4  data$GDPWdiff
5
6
7  data$GDPWdiff <- factor(data$GDPWdiff , ordered = FALSE ) #how to make a
      table un-ordered
8
9  table(data$GDPWdiff)
10
11
12  ftable(xtabs(~ GDPWdiff + OIL + REG, data = data))
```

Firstly I inspect my data and see that I need to create a bounds of values for the 3 catagories of "negative", "no change", "positive". I then create a contingency table to help run a regression. Then I set the reference category to "no change" as specified

```
1  data$GDPWdiff = relevel(data$GDPWdiff, ref = "no_change") #setting the
      reference level to 0
2
3  # run model
4  mult.log <- multinom(GDPWdiff ~ OIL + REG, data = data)
5  summary(mult.log)
6  expc <- exp(coef(mult.log))
7
```

```
 8  stargazer(expc, type="latex")

 9
10  # get p values
11  z <- summary(mult.log)$coefficients/summary(mult.log)$standard.errors
12  z
13  p <- ((1 - pnorm(abs(z), 0, 1)) * 2) #2 tailed z test
14  p

15
16  stargazer(p, type="latex")

17
18  library(stargazer)
19  stargazer(mult.log, type="latex") #gives the P value and coefficents

20
21  head(fitted(mult.log)) # checking the fitted values

22

23
24  ###############
```

In this I set my reference catagory and run the regression, and obtain the P values and coefficents.

Table 1:

|          | (Intercept) | OIL     | REG   |
|----------|-------------|---------|-------|
| negative | 44.934      | 116.492 | 3.976 |
| positive | 93.118      | 95.344  | 5.867 |

These are the logit coefficients relative to the reference category of no change. For example, under 'REG', the 1.38 suggests that for one unit increase in 'REG' score, the logit coefficient for 'negative' relative to 'no change' will go increase by 1.38. In other words, if your REG score increases one unit, your chances of staying in the noc hange GDPWdiff category are higher compared to staying in low GDPWdiff. When the coefficients are expontiated keeping all other variables constant, if the oil score increases one unit, you are 116.492 times more likely to stay in the negative category as compared to the no change category. The coefficient, however, is not significant. keeping all other variables constant, if the oil score increases one unit, you are 95.344 times more likely to stay in the positive category as compared to the no change category. The coefficient, however, is not significant. The p values are not statistically significant and we fail to reject the null hypothesis of our model

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

```
1  ord.log <- polr(GDPWdiff ~ OIL + REG, data = data, Hess = TRUE)
2  summary(ord.log)
3
```

Table 2:

| | Dependent variable: | |
|---|---|---|
| | negative | positive |
| | (1) | (2) |
| OIL | 4.758 | 4.557 |
| | (6.823) | (6.823) |
| | | |
| REG | 1.380* | 1.769** |
| | (0.769) | (0.767) |
| | | |
| Constant | 3.805*** | 4.534*** |
| | (0.271) | (0.269) |
| | | |
| Akaike Inf. Crit. | 4,688.792 | 4,688.792 |
| Note: | *p<0.1; **p<0.05; ***p<0.01 | |

```
4  stargazer(ord.log, type="latex")
5
6  # Calculate a p value
7  ctable <- coef(summary(ord.log))
8  p <- pnorm(abs(ctable[, "t value"]), lower.tail = FALSE) * 2
9  ctable <- cbind(ctable, "p value" = p)
10 ctable
11
12 stargazer(ctable, type="latex")
13
14 # Calculate confidence intervals
15 ci <- confint(ord.log)
16 ci
17
18
19 # convert to odds ratio
20 exp(cbind(OR = coef(ord.log), ci))
```

For nations with REG(1) ie. Democracies, the odds of being more likely (i.e., very or somewhat likely versus unlikely) to have positive change in GDP is 105 percent more than undemocratic nations, holding constant all other variables. For nations with Oil(1) ie. 50 percent + exports, the odds of being more likely (i.e., very or somewhat likely versus unlikely) to have positive change in GDP is 84 percent more than non 50 percent oil exporting nations, holding constant all other variables.

Table 3:

| | Dependent variable: |
| --- | --- |
| | GDPWdiff |
| OIL | −0.172 |
| | (0.116) |
| REG | 0.409*** |
| | (0.075) |
| Observations | 3,720 |
| *Note:* | *p<0.1; **p<0.05; ***p<0.01 |

Table 4:

| | Value | Std. Error | t value | p value |
| --- | --- | --- | --- | --- |
| OIL | -0.172 | 0.116 | -1.483 | 0.138 |
| REG | 0.409 | 0.075 | 5.446 | 0.00000 |
| no_change|negative | -5.319 | 0.252 | -21.084 | 0 |
| negative|positive | -0.704 | 0.048 | -14.798 | 0 |

# Question 2

Consider the data set `MexicoMuniData.csv`, which includes municipal-level information from Mexico. The outcome of interest is the number of times the winning PAN presidential candidate in 2006 (`PAN.visits.06`) visited a district leading up to the 2009 federal elections, which is a count. Our main predictor of interest is whether the district was highly contested, or whether it was not (the PAN or their opponents have electoral security) in the previous federal elections during 2000 (`competitive.district`), which is binary (1=close/swing district, 0="safe seat"). We also include `marginality.06` (a measure of poverty) and `PAN.governor.06` (a dummy for whether the state has a PAN-affiliated governor) as additional control variables.

(a) Run a Poisson regression because the outcome is a count variable. Is there evidence that PAN presidential candidates visit swing districts more? Provide a test statistic and p-value.

```
1 Mex_data <- read.csv("C:/Users/crowl/Documents/Trinity/ASDS Course/
      POP77003_Applied_Statistical_Analysis_2/ProblemSet/PS3/template/
      MexicoMuniData.csv")
2
3 summary(Mex_data)
4 str(Mex_data)
5
6 with(Mex_data, #Testing to see if our response variable meets the
      assumptions for a Poisson test.
7    list(mean(PAN.visits.06), var(PAN.visits.06))) #The outcome would
      suggest it does as the expected mean and expected   are approx equal
```

loading in the data. As the mean is 0.09181554 and the variance is 0.6436861, which show they are approximately the same, we progress with running a Poisson regression.

```
1 Mex_data <- within(Mex_data, {
2    PAN.governor.06 <- as.logical(PAN.governor.06)
3    competitive.district <- as.logical(competitive.district)
4 })
5
6 mod.pos <- glm(PAN.visits.06 ~ ., data = Mex_data, family = poisson(link
      = "log"))
7 summary(mod.pos)
8
9
10 table(Mex_data$competitive.district) #this shows that FALSE is the
      reference category
11
12                                        #I want to change this to TRUE to
      easily compare my categories
13 Mex_data$competitive.district <- factor(Mex_data$competitive.district ,
      ordered = FALSE )
14
```

```
15 Mex_data$competitive.district2 <- relevel(Mex_data$competitive.district,
      ref= "TRUE")
16
17 table(Mex_data$competitive.district2) #shows the referennce has now
      changed
18
19 mod.pos1 <- glm(PAN.visits.06 ~ MunicipCode + pan.vote.09 + marginality
      .06 + PAN.governor.06 +competitive.district2, data = Mex_data, family
      = poisson)
20 summary(mod.pos1)
```

As competitive.district and PAN.governor.06 are binary (0/1) responses, I will turn them into logical (true/false) responses. I then run the model.

Table 5:

| | Dependent variable: |
|---|---|
| | PAN.visits.06 |
| MunicipCode | −0.00001 |
| | (0.00001) |
| pan.vote.09 | 0.232*** |
| | (0.052) |
| marginality.06 | −2.060*** |
| | (0.120) |
| PAN.governor.06 | −0.269 |
| | (0.167) |
| competitive.district2FALSE | −0.068 |
| | (0.182) |
| Constant | −3.761*** |
| | (0.244) |
| Observations | 2,407 |
| Log Likelihood | −639.535 |
| Akaike Inf. Crit. | 1,291.070 |
| *Note:* | *p<0.1; **p<0.05; ***p<0.01 |

In our model when the district is non competitive the intercept is -3.829, with all other factors held at zero. And -3.761 when the district is competitive, both of these

Table 6:

| | Dependent variable: |
|---|---|
| | PAN.visits.06 |
| MunicipCode | −0.00001 |
| | (0.00001) |
| pan.vote.09 | 0.232*** |
| | (0.052) |
| marginality.06 | −2.060*** |
| | (0.120) |
| PAN.governor.06 | −0.269 |
| | (0.167) |
| competitive.district | 0.068 |
| | (0.182) |
| Constant | −3.829*** |
| | (0.302) |
| Observations | 2,407 |
| Log Likelihood | −639.535 |
| Akaike Inf. Crit. | 1,291.070 |
| Note: | *p<0.1; **p<0.05; ***p<0.01 |

intercepts have p values that denote significance. When we change the model from non competitive to competitive we see a 0.64 difference in the coefficient intercept.

The coefficient for competitive.district is 0.6849. This means that the expected log count for a one-unit increase in competitive.district is 0.6849. However the p value of competitive.district (0.182) is not significant and therefore we must conclude that a districts competitiveness does not significantly influence the vistits of a presidential candidate.

(b) Interpret the `marginality.06` and `PAN.governor.06` coefficients.

```
1  cofs <- coef(mod.pos)
2  cofs
3
4  stargazer(cofs, type="latex")
5
6  confint(mod.pos) #confidence interval set at 95%
7
8  exp(coef(mod.pos)) #exponential model coefficients
9  exp(confint(mod.pos)) #CI for the exponential model coefficients
10
11
12 sjPlot::tab_model(mod.pos , show.intercept = TRUE, #shows the P value
13                   show.se = FALSE, dv.labels = "Predictors of presidental
      Visit", auto.label = TRUE, show.re.var = FALSE, show.icc =FALSE,
14                   show.r2 = FALSE, show.ngroups = FALSE, show.obs = FALSE
      )
```

marginality.06 has a coefficient of -2.06 and a p value of 0.12 which is statistically significant. PAN.governor.06 has a coefficient of -0.269 and a p value of 0.167 which is not statifcally significant. After this I then created a 95percent confidence interval and exponentiated the model's coefficients. This then produced a table showing each of the variables exponentiated coefficients.

there appears to be a significant association between marginality.06 and PAN.visits.06 IRR = 0.13 P value = ¡0.001 , 95percent CI = 0.10 − 0.16

there appears to be a non-significant association between PAN.governor.06 and PAN.visits.06 IRR = 0.76 P value = 0.107 , 95 percent CI = 0.55 − 1.05

Table 7:

| (Intercept) | MunicipCode | pan.vote.09 | marginality.06 | PAN.governor.06TRUE | competitive.dist |
|---|---|---|---|---|---|
| -3.829 | -0.00001 | 0.232 | -2.060 | -0.269 | 0.068 |

(c) Provide the estimated mean number of visits from the winning PAN presidential candidate for a hypothetical district that was competitive (`competitive.district`=1), had an average poverty level (`marginality.06` = 0), and a PAN governor (`PAN.governor.06`=1).

9

```
1  exp ( cofs [ 1 ]  +  cofs [ 2 ]  +  cofs [ 3 ]  +  cofs [ 4 ] *0  +  cofs [ 5 ] *1  +  cofs [ 6 ] *1 )
2
3  pred  <−  data . frame ( Mex_data $ MunicipCode ,
4                          Mex_data $ pan . vote .09 ,
5                           marginality .06  =  0 ,
6                          PAN . governor .06  =  TRUE,
7                          competitive . district  =  TRUE)
8
9  colnames ( pred )
10
11 names ( pred ) [ 1 ]  <−  ”MunicipCode” #had  to  rename  the  columns ,  or  they  would
        not  be  recognised
12 names ( pred ) [ 2 ]  <−  ”pan . vote .09 ”
13
14 colnames ( pred )
15
16
17 # check  with  predict ()  function
18  predict ( mod . pos ,  newdata  =  pred ,  type  =  ”response ”)
19
20 mod . zip  <−  zeroinfl (PAN. visits .06  ~  . ,  data  =  Mex_data ,  dist  =  ”poisson ”)
        #i  get  an  error  message  when  I  attempt  to  run  this ,  it  a  non  finite
21 summary ( mod . zip )
        #  vlaue  has  been  supplied ,  I  have  not  been  able  to  reslove  it .
```

When the model is exponentiated the expected intercept is 0.02172546 When the model
is exponentiated and the variables are selected the expected intercept is 0.02243164