# Problem Set 3

## Applied Stats II

## Due: March 26, 2023

## 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 23:59 on Sunday March 26, 2023. No late assignments will be accepted.

## 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.

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

   MY ANSWER FOR QUESTION 1

```
1
2  ############################################################
3  # Question 1
4  ############################################################
5
6  # Part 1
7  # Required to construct and interpret an unordered multinomial model
8  # Given — the response variable GDPWdiff is difference between year t and
        t−1
9  # possible catagories "pos" and "neg"
10 # and "no change" for the reference category
11
12 # expanding on this and to make the variables categorical let:
13 # no change = 0
14 # pos = 1
15 # neg = 2
16
17 # note that within() evaluates the expression and creates a copy of the
        original data frame ps3Data
18 ps3Data <- within(ps3Data, {
19   # this initializes the variable
20   GDPWdiff.cat <- NA
21   # for the neg value
22   GDPWdiff.cat[GDPWdiff < 0] <- "neg"
23   # for the no change value
24   GDPWdiff.cat[GDPWdiff == 0] <- "no change"
25   # for the pos value
26   GDPWdiff.cat[GDPWdiff > 0] <- "pos"}
27 )
28
29 # check that everything looks okay
30 print(head(ps3Data,n=5))
31
32 # X COUNTRY CTYNAME YEAR GDPW OIL REG    EDT GDPWlag GDPWdiff GDPWdifflag
        GDPWdifflag2 GDPWdiff.cat
33 # 1 1        1 Algeria 1965 6620   1   0  1.45    6502      118
        419            1071             pos
34 # 2 2        1 Algeria 1966 6612   1   0  1.56    6620       −8
        118             419            neg
35 # 3 3        1 Algeria 1967 6982   1   0 1.675    6612      370
        −8             118             pos
```

```
36 # 4 4        1 Algeria 1968 7848   1   0 1.805     6982       866
      370              −8          pos
37 # 5 5        1 Algeria 1969 8378   1   0  1.95     7848       530
      866             370          pos

39 # the next step is to turn the data into factors or convert numeric to a
      factor
40 # can use either as.factor() or factor() which is wrapper for factor but
      allows quick return if the input factor
41 # is already a factor
42 ps3Data$GDPWdiff.cat <− factor(ps3Data$GDPWdiff.cat, levels = c("no
      change", "pos", "neg"))
43 # ps3Data$GDPWdiff.cat <− as.factor(ps3Data$GDPWdiff.cat)

45 # check that everything looks okay
46 print(summary(ps3Data$GDPWdiff.cat))

48 # no change      pos         neg
49 #        16      2600        1105

51 # REG and OIL should also be factors
52 ps3Data$REG <− as.factor(ps3Data$REG)
53 ps3Data$OIL <− as.factor(ps3Data$OIL)

55 # the next step is to address the reference category and create the p
      regression model
56 # the relevel() function doesn't afffect the original dataset
57 # lm(x ~ y + relevel(b, ref = "3")) ... just working out how to implement
      it
58 ps3Data$GDPWdiff2 <− relevel(ps3Data$GDPWdiff.cat, ref = "no change")

60 # returns the following
61 # weights:  12 (6 variable)
62 # initial   value 4087.936326
63 # iter   10 value 2340.076844
64 # final   value 2339.385155
65 # converged

67 ps3Multinomial <− multinom(ps3Data$GDPWdiff2 ~ REG + OIL, data = ps3Data)

69 # check that everything looks okay
70 print(summary(ps3Multinomial))

72 # returns the following
73 # Call:
74 #  multinom(formula = ps3Data$GDPWdiff2 ~ REG + OIL, data = ps3Data)

76 # Coefficients:
77 #        (Intercept)      REG1      OIL1
78 # pos     4.533759 1.769007 4.576321
79 # neg     3.805370 1.379282 4.783968
```

```r
80
81 # Std. Errors:
82 #         (Intercept)       REG1        OIL1
83 # pos    0.2692006 0.7670366 6.885097
84 # neg    0.2706832 0.7686958 6.885366
85
86 # Residual Deviance: 4678.77
87 # AIC: 4690.77
88
89 ## Add this stargazer table to the TeXShop / latex document file
90 stargazer(ps3Multinomial, title = "unordered multinominal logit")
91
92
93 ############# now we can visualize the data with a jitter plot
94 ggplot(ps3Data, aes(x = OIL , y = GDPWdiff.cat)) +
95   geom_jitter(alpha = .5, color="purple") +
96   theme(axis.text.x = element_text(angle = 45, hjust = 1, vjust = 1)) +
97   theme(legend.position="bottom")
98
99
100 ggplot(ps3Data, aes(x = REG, y = GDPWdiff.cat)) +
101   geom_jitter(alpha = .5, color="purple") +
102   theme(axis.text.x = element_text(angle = 45, hjust = 1, vjust = 1)) +
103   theme(legend.position="bottom")
104
105
106 ### get the P−values and the estimated cutoff points and the coefficients
107
108 # Interpretation
109 ## I got a little bit confused here. Referencing the lecture slides (Week
        8: Multinomial Logit Regression)
110 ## and Slide 11. replicating exp(coef(multinom_model)[,c(1:5)])
111 ## to exp(coef(ps3Multinomial)[ ,c(1:3)]) this was giving some unusual
       values
112 ## removing the exp or exponent and the indexing was not returning unusal
        values
113 # coef(ps3Multinomial) ??
114
115 ps3table <- coef(summary(ps3Multinomial))
116 ps3table
117
118 # this returns
119 #             (Intercept)      REG1        OIL1
120 # pos           4.533759 1.769007 4.576321
121 # neg           3.805370 1.379282 4.783968
122
123
124 #### Interpretations continued
125 #### REG1  'pos'
126 # for REG1 'pos' = for a unit change in REG going from 0 to 1, non
      democracy to a democracy, the log−odds are that there will
```

```r
127 # be a pos change in the GDP from one year to the next increase by
        (1.769007) when all other variables in the multinom_model are
128 # held constant and the ref category is "no change"
129
130 #### REG1 'neg'
131 # for REG1 'neg' = for a unit change in REG1 going from 0 to 1, non
        democracy to a democracy, the log-odds are that there
132 # will be a neg change in the GDP from one year to the next increase by
        (1.379282) when all other variables in the multinom_model
133 # are held constant and the ref category is "no change"
134
135 #### OIL1 'pos'
136 # OIL1 'pos' when there is a unit change in the OIL variable, where it
        increases from 0 to 1
137 # this indicates that the average ratio of fuel exports to total exports
        in 1984-86 exceeded by 50%
138 # the log-odds here means  that there will be a pos difference in the
        total GDP in a COUNTRY from one year to the next year
139 # and ths results in an increase of (4.576321) and the other variables in
         the model are held constant
140
141
142 ### OIL 'neg'
143 # for OIL 'neg' when there is a unit change in the OIL variable from
        going from 0 to 1, then the average ratio of fuel exports
144 # to total exports in 1984-86 exceeded by 50%, while the log-odds results
         in a neg difference in the total GDP of a country
145 # from one year to the next and increases by (4.783968)  while all all
        other variables are held constant
146
147
148 ########## Ordered multinominal logit - Q1 - Part 2
149 # running the ordered logit
150 # referencing lecture slide 45
151 # Hess=TRUE to have the model return the observed information matrix from
        optimization (called the Hessian)
152 # which is used to get standard errors.
153 # ref. https://stats.oarc.ucla.edu/r/dae/ordinal-logistic-regression/
154
155 ps3_ordered_multi <- polr(GDPWdiff2 ~ REG + OIL, data = ps3Data, Hess = T
       )
156
157 ## referencing lecture slide 46
158 test <- exp(cbind(OR = coef(ps3_ordered_multi), confint(ps3_ordered_multi
       ))) # odds ratio
159 test
160 ### This returns
161 # OR              2.5 %        97.5 %
162 # REG1 0.7000737  0.6042257    0.8102918
163 # OIL1 1.2593051  1.0029960    1.5754005
164
```

```
165
166
167  summary(ps3_ordered_multi)
168
169  ## this returns
170
171  #Coefficients:
172  #      Value Std.      Error     t value
173  #REG1  −0.3566     0.07485    −4.764
174  #OIL1   0.2306     0.11510     2.003
175
176  #Intercepts:
177  #                Value    Std. Error  t value
178  #no change|pos  −5.5846    0.2534     −22.0376
179  #pos|neg         0.7491    0.0479      15.6475
180
181  #Residual Deviance: 4692.109
182  #AIC:  4700.109
```

# 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.

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

(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`).

MY ANSWER FOR QUESTION 2

```
1
2
3
4  ############################################################
5  # Question 2
```

```r
6  ###################################################################

7
8  # import the data MexicaMuniData.csv
9  mexData <- read.csv("/Users/marklikeman/desktop/ASDS-2023/applied-stats
       -2-2023/problemset03/MexicoMuniData.csv",
10                      stringsAsFactors = FALSE)

11
12 # inspecting the data
13 head(mexData, n=5)
14 tail(mexData)
15 str(mexData)
16 summary(mexData)
17 stargazer(mexData, title = "MexicoMuniData")

18
19 # a quick look at the loaded data
20 # MunicipCode pan.vote.09 marginality.06 PAN.governor.06 PAN.visits.06
       competitive.district
21 # 1    1001        0.283         -1.831              0               5
                      1
22 # 2    1002        0.352         -0.620              0               0
                      1
23 # 3    1003        0.359         -0.875              0               0
                      1
24 # 4    1004        0.238         -0.747              0               0
                      1
25 # 5    1005        0.378         -1.234              0               0
                      1

26
27 #### looking at the question 2 details the variables outcome Pan.visits
       .06 looks to be of interest here
28 #### the main predictor of interest is whether the district was highly
       contested, plus measure of property
29 #### and PAN.governor.06  ...
30 #### PAN.governor.06 and competitive.district are the binary variables,
       (1 = close/swing district, 0 = 'safe seat')

31
32
33 # note that within() evaluates the expression and creates a copy of the
       original dataset
34 mexData <- within(mexData, {
35   PAN.governor.06 <- as.logical(PAN.governor.06) # binary
36   competitive.district <- as.logical(competitive.district)} # binary
37 )

38
39 ############################################### Part a
40 ## running the poisson regression
41 pm <- glm(PAN.visits.06 ~ competitive.district + marginality.06 + PAN.
       governor.06, data = mexData, family = poisson)
42 summary(pm)

43
44 ## this returns
```

```r
45
46 ################################################################
47 # Deviance  Residuals:
48 # Min          1Q     Median           3Q         Max
49 # −2.2309    −0.3748    −0.1804    −0.0804    15.2669
50
51 # Coefficients:
52 #                                  Estimate  Std.  Error  z  value  Pr(>|z|)
53 # (Intercept)                      −3.81023      0.22209  −17.156    <2e−16 ***
54 #   competitive.districtTRUE  −0.08135      0.17069   −0.477    0.6336
55 # marginality.06                   −2.08014      0.11734  −17.728    <2e−16 ***
56 # PAN.governor.06TRUE          −0.31158      0.16673   −1.869    0.0617 .
57
58 #   Signif.  codes:  0     ***     0.001     **     0.01     *     0.05     .
        0.1            1
59
60 # (Dispersion  parameter  for  poisson  family  taken  to  be  1)
61
62 # Null  deviance:  1473.87   on  2406   degrees  of  freedom
63 # Residual  deviance:   991.25   on  2403   degrees  of  freedom
64 # AIC:  1299.2
65
66 # Number  of  Fisher  Scoring  iterations:  7
67 ################################################################
68
69
70 # looking  at  the  poisson  model  it  appears  to  be  the  case  that  when
        changing  from  a  safe  swing  seat,  this  decreases
71 # the  log−odds  that  there  will  be  a PAN  presidential  candidates  visit
        while  holding  all  the  other  visits  contant
72
73
74 ## for  the  TeXShop  template
75 stargazer(pm,  title = "Poisson  Model")
76
77
78 #### visualizing  the  data
79 ggplot(data = NULL,  aes(x = pm$fitted.values,  y = mexData$PAN.visits.06))
        +
80   geom_jitter(alpha = 0.5) +
81   geom_abline(color = "purple") +
82   theme(legend.position="bottom")
83
84
85
86
87 ########################################### Part  b
88
89 ### the  coef  marginality.06  is  −> marginality.06      −2.08014
90 ### this  indicates  that  for  a  unit  one  increase  in  a  measure  of  property,
        the  log−odds  of  a PAN  presidential  candidates  visit
```

```r
### will decrease by factor 2.080 which indicates that poorer districts
     have a low probility of getting a visit from
### a AN presidential candidate


############################################### Part c


#### referring to the ouput from poisson regression above

# Coefficients:
#                              Estimate Std. Error z value Pr(>|z|)
# (Intercept)                  -3.81023    0.22209 -17.156   <2e-16 ***
#  competitive.districtTRUE    -0.08135    0.17069  -0.477   0.6336
# marginality.06               -2.08014    0.11734 -17.728   <2e-16 ***
# PAN.governor.06TRUE          -0.31158    0.16673  -1.869   0.0617 .

### this used to estimate the mean

# (intercept*1)  + (competitive.districtTRUE*1) + (marginality.06*0) + (
     PAN.governor.06TRUE*1)
part_c <- exp((-3.81023*1) + (-0.08135*1) + (2.08014*0) + (-0.31158*1))
part_c
## this gives
## 0.01494827

## interpretation
# the estimated mean for the amount of times that a PAN presidential
     candidate winning in 2006 is 0.01494827
```