# Problem Set 2

## Applied Stats II

## Due: February 18, 2024

## 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 February 18, 2024. No late assignments will be accepted.

We're interested in what types of international environmental agreements or policies people support (Bechtel and Scheve 2013). So, we asked 8,500 individuals whether they support a given policy, and for each participant, we vary the (1) number of countries that participate in the international agreement and (2) sanctions for not following the agreement.

Load in the data labeled `climateSupport.RData` on GitHub, which contains an observational study of 8,500 observations.

- Response variable:

  - `choice`: 1 if the individual agreed with the policy; 0 if the individual did not support the policy

- Explanatory variables:

  - `countries`: Number of participating countries [20 of 192; 80 of 192; 160 of 192]
  - `sanctions`: Sanctions for missing emission reduction targets [None, 5%, 15%, and 20% of the monthly household costs given 2% GDP growth]

Please answer the following questions:

1. Remember, we are interested in predicting the likelihood of an individual supporting a policy based on the number of countries participating and the possible sanctions for non-compliance.

   Fit an additive model. Provide the summary output, the global null hypothesis, and $p$-value. Please describe the results and provide a conclusion.
   Interpretation:

   Since the dependent variable is binary, the entire question uses a logistic regression model. And use the Likelihood ratio test to test the validity of the coefficients. Global null hypothesis : all coefficients (coefficients of independent variables) are zero, and there is no significant relationship between independent variables and dependent variable , that is, the independent variables in the model have no predictive power on the dependent variable.

|  | Model |
| --- | --- |
| (Intercept) | $-0.27^{***}$ |
|  | (0.05) |
| countries80 of 192 | $0.34^{***}$ |
|  | (0.05) |
| countries160 of 192 | $0.65^{***}$ |
|  | (0.05) |
| sanctions5% | $0.19^{**}$ |
|  | (0.06) |
| sanctions15% | $-0.13^{*}$ |
|  | (0.06) |
| sanctions20% | $-0.30^{***}$ |
|  | (0.06) |
| AIC | 11580.26 |
| BIC | 11622.55 |
| Log Likelihood | $-5784.13$ |
| Deviance | 11568.26 |
| Num. obs. | 8500 |

$^{***}p < 0.001$; $^{**}p < 0.01$; $^{*}p < 0.05$

Table 1: Model

```
1  install.packages("texreg")
2  library(texreg)
3
4  install.packages("xtable")
5  library(xtable)
```

```r
6
7  # load data
8  load(url("https://github.com/ASDS-TCD/StatsII_Spring2024/blob/main/
        datasets/climateSupport.RData?raw=true"))
9  load("/Users/daisy/Downloads/climateSupport.RData")
10 ls()
11
12 # Preview data
13 summary(climateSupport)  # View summary statistics of the data
14 head(climateSupport)     # View the first few rows of the data
15
16 # Convert data format of DV into binary outcome
17 climateSupport$choice <- ifelse(climateSupport$choice == "Supported",
        1, 0)
18
19 # Change data format of IVs to unordered factors
20 climateSupport$countries <- factor(climateSupport$countries, ordered
        = FALSE)
21 climateSupport$sanctions <- factor(climateSupport$sanctions, ordered
        = FALSE)
22
23 # Use the relevel function to set the baseline level and fit logistic
         regression model
24 climateSupport$countries <- relevel(climateSupport$countries, ref = "
        20 of 192")
25 climateSupport$sanctions <- relevel(climateSupport$sanctions, ref = "
        None")
26 model <- glm(choice ~ ., family = binomial(link = "logit"), data =
        climateSupport)
27
28 # Output the results of the model
29 summary(model)
30
31 # Create LaTeX table for the model
32 texreg(model)
33
34 ## Likelihood ratio test
35
36 # H0: beta_1 = beta_2 = ...= beta_5 = 0
37 # H1: at least one slope is not equal to 0
38 #   Create a null model
39 null_model <- glm(choice ~ 1,
40                   data = climateSupport,
41                   family = "binomial")
42
43 #  Run an anova test on the model compared to the null model
44 anova_result_1 <- anova(null_model, model, test = "LRT")
45
46 # print results
47 print(anova_result_1)
48
```

```
49  # Create LaTeX table for the model
50
51  # Create a data frame containing ANOVA results
52  result_table_1 <- xtable(anova_result_1)
53
54  # Convert to xtable format
55  print(result_table_1, include.rownames = FALSE)
```

| Resid. Df | Resid. Dev | Df | Deviance | Pr(>Chi) |
|-----------|-----------|-----|----------|----------|
| 8499 | 11783.41 | | | |
| 8494 | 11568.26 | 5 | 215.15 | 0.0000 |

From the table, the difference in degrees of freedom is 5, while the difference in bias is 215.15, and the p-value is extremely small (2.2e-16), indicating that the difference in bias between the two models is significant.This suggests that a model that includes the countries and sanctions variables has an advantage in explaining the data relative to the null model.

2. If any of the explanatory variables are significant in this model, then:

   (a) For the policy in which nearly all countries participate [160 of 192], how does increasing sanctions from 5% to 15% change the odds that an individual will support the policy? (Interpretation of a coefficient)
   Interpretation:
   Calculate the logarithm of the odds in two different situations, then make the difference, and then exponentialize the difference to get the change in odds ratio.

```
1   # (a)
2
3   # Fitted logistic regression model is:
4   # logit^(p/(1-p)) = beta0 + beta_1*countries_80 + beta_2*countries_
        160 +
5   # beta_3*sanctions_5% + beta_4*sanctions_15% + beta_5*sanctions_20%
6
7   # After substituting the coefficients, get:
8   # logit^(p/(1-p)) = -0.27266 + 0.33636*countries_80 + 0.64835*
        countries_160 +
9   # 0.19186*sanctions_5% - 0.13325*sanctions_15% - 0.30356*sanctions_
        20%
10
11  # according to: logit^(p_5%/(1-p_5%)) =   -0.27266 + 0.64835*countries
        _160 + 0.19186*sanctions_5%
12  log_odd_p_5 <- -0.27266 + 0.64835*1 + 0.19186*1
13
14  # According to: logit^(p_15%/(1-p_15%)) =   -0.27266 + 0.64835*
        countries_160 - 0.13325*sanctions_15%
15  log_odd_p_15 <- -0.27266 + 0.64835*1 - 0.13325*1
```

```
16
17 # Subtract two log odds, get the difference of log odds
18 delta = log_p_15 - log_p_5
19
20 # Indexed results according to the formula: OR = exp^(delta)
21 exp(delta)
22
23 # Print results:0.7224479
```

For the policy in which nearly all countries participate [160 of 192], increasing sanctions from 5% to 15% will increase the odds that an individual will support the policy by 0.7224479.

(b) What is the estimated probability that an individual will support a policy if there are 80 of 192 countries participating with no sanctions?

```
1 # (b)
2
3 # After substituting the coefficients, get:
4 # logit^(p_80/(1-p_80)) =   -0.27266 + 0.33636*countries_80
5 log_odd_p_80 <- -0.27266 + 0.33636*1
6
7 # According to the formula: P = exp(logit^(P/(1-P))/(1+exp(logit^(P/
      (1-P))))
8 p_80 = exp(log_odd_p_80) / (1 + exp(log_odd_p_80))
9
10 # Print the results
11 print(p_80)
12
13 # Results : 0.5159196
```

Interpretation:
The estimated probability that an individual will support a policy if there are 80 of 192 countries participating with no sanctions is about 51.6%.

(c) 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.
  Interpretation:

  Use the chi-square test. First fit a logistic regression model containing interaction terms, then use anova() to compare with the original model, and finally analyze the p value in the anova result to see the significance of the coefficient.

  ```
  1 # (c)
  2
  3 # Fit a logistic regression model that includes interaction terms
  4 interaction_model <- glm(choice ~ countries * sanctions, family =
        binomial(link = "logit"), data = climateSupport)
  ```

```r
5
6 # Use anova() to compare with the addictive model
7 anova_result_2 <- anova(interaction_model, model, test = "Chi")
8
9 # Print results
10 print(anova_result_2)
11
12 # Create LaTeX table for the model
13
14 # Create a data frame containing ANOVA results
15 result_table_2 <- xtable(anova_result_2)
16
17 # Convert to xtable format
18 print(result_table_2, include.rownames = FALSE)
```

| Resid. Df | Resid. Dev | Df | Deviance | Pr(>Chi) |
|---|---|---|---|---|
| 8488 | 11561.97 | | | |
| 8494 | 11568.26 | -6 | -6.29 | 0.3912 |

Since the p value is 0.3912 which is bigger than 0.05, it indicates that the deviation difference between the two models is not significant, that is, the interaction term has no significant impact on the explanatory power of the model.The answers to 2a and 2b would not potentially change if we included the interaction term in this model.