

Problem Set 7

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): Political Science

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.

```
Call:
glm(formula = PAN.visits.06 ~ competitive.district + marginality.06 +
PAN.governor.06, family = poisson, data = mexico_elections)
```

Deviance Residuals:

```
Min      1Q   Median      3Q      Max
-2.1441  -0.3596  -0.1742  -0.0783  15.2935
```

Coefficients:

```
Estimate Std. Error z value Pr(>|z|)
(Intercept)      -3.9304    0.1747 -22.503  <2e-16 ***
competitive.district -0.4594    0.3276  -1.402    0.161
marginality.06     -2.0981    0.1210 -17.343  <2e-16 ***
PAN.governor.06     -0.2073    0.1660  -1.249    0.212
```

```
Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
```

(Dispersion parameter for poisson family taken to be 1)

```
Null deviance: 1433.83  on 2392  degrees of freedom
Residual deviance: 963.57  on 2389  degrees of freedom
(4 observations deleted due to missingness)
AIC: 1255.9
```

Number of Fisher Scoring iterations: 7

There is no evidence that supports PAN presidential candidates visit swing districts more as the test statistic was came out to be -1.402 and the p value was 0.161 which is large.

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

```
(Intercept) competitive.district      marginality.06
-3.9304467      -0.4594186      -2.0981427
PAN.governor.06
-0.2073147
```

For states that are "safe seat" district without a PAN-affiliated governor, by increasing the measure of poverty by 1, this increases the odds of visits from the winning PAN presidential candidate by the multiplication factor of 0.12.

For states that are "safe seat" district with a measure of poverty of 0, by having a PAN-affiliated governor increases the odds of visits from the winning PAN presidential candidate by 0.81.

- (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`).

The estimated mean number of visits from the winning PAN presidential candidate for a hypothetical district is 0.01.

Question 2 (50 points): Biology

We'll be using data from a longitudinal sleep study of under 20 undergraduate students ($n=18$), which took place over the course of 10 days to see if sleep deprivation has any effect on participants' reaction time. Load the data through the `lmer` package.

1. Create a "pooled" linear model where you regress `Days` on the outcome `Reaction`. Make sure to run regression diagnostics to check if the variance around the regression line is equal for every year.
2. Fit an "un-pooled" regression model with varying intercepts for patient (include an additive factor for patient) and save the fitted values.
3. Fit a "un-pooled" regression model with varying slopes of time (days) for patients (include only the interaction `Days:Subject`) and save the fitted values.
4. Fit an "un-pooled" regression model with varying intercepts for patients with varying slopes of time (days) by patient (include the interaction and constituent terms of `Days` and `Subject`, `Days + Subject + Days:Subject`) and save the fitted values.
5. Fit a "semi-pooled" multi-level model with varying-intercept for subject and varying-slope of day by subject. Is it worthwhile for us to run a multi-level model with varying effects of time by subject? Why? Compare your model from part 5 to the other completely "pooled" or "un-pooled models".