

Naive Aggregate Model Creation

Theodore Dounias

10/10/2018

Introduction

For the rest of this write-up, assume the following:

$$y_i \sim \text{Bernoulli}(p_vote)$$

Where $y_i \in \{0, 1\}$ is the probability that the i -th ballot was completed. Indexes and variables coded according to the following table:

Purpose/Unit	Variable	Index
Ballot	u	i
Individual	z	j
County	x	k
Election	w	l
Number of Vars	n	–
General Index	–	i’

Additionally, I do veer away a bit from the structure we talked about on Monday; I do not linearly add to the model until it reaches a final stage. The reasoning here is that there is no exact linear path to follow; there is an overarching unit of observation—the ballot—and all the rest are dependent between each other. For instance, adding a variable for Party at the ballot level would not significantly change the way I later add percentage of white residents at the county level. Therefore, the way I proceed is the following: I “build” the models step by step and separately for each group of variables (grouping by unit of observation). Then I present one example of what a model using two of these initial “building blocks” would look like. Since this is fairly generalizable, I then proceed directly to the full model which includes all different variables.

No Data Guess

If receiving a ballot with no information, I would predict that the probability that an additional ballot was a vote in favor would be equal to turnout, as calculated through all other ballots. Therefore:

$$p_vote_i = \frac{\#votes\ cast}{\#ballots}$$

Estimation with only one type of data

There are four levels of data I will go through here: County, Election, Person, and Ballot.

County Level

Assume that the ballot I am trying to assess completion for has the name of the county it is from written on it. There are two ways I can think of for predicting p_vote . First, assume that each different county has a different, independent p_vote . Therefore, in model-lingo this would look like:

$$p_vote_i \sim \text{logit}^{-1}\left(\sum_{k=1}^{64} x_{k,i}\beta_k\right)$$

Where k counts over the 64 counties of Colorado, and x_k is an indicator variable for each county. If I, quite reasonably, throw away the assumption of independence—these counties are, after all, in the same state and the same country—I could also fit a mixed effects model as such:

$$p_vote \sim \text{logit}^{-1}(a_{k[i]}),$$

$$a_k \sim N(\gamma_0, \sigma_\alpha^2)$$

Where $\alpha_{k[i]}$ varies by county, constrained by its standard deviation and γ_0 , an intercept coefficient. Let's say now that along with the one ballot, I was given a short list of $n^{\text{county vars}}$ other county-level variables, be they discrete, continuous, or indicators. The two models would then look like:

$$p_vote_i \sim \text{logit}^{-1}\left(\sum_{k=1}^{64} x_k\beta_k + \sum_{i'=1}^{n^{\text{county vars}}} x_{k[i],i'}\beta_{i'+64}\right)$$

Where $x_{k[i],i'}$ is the k -th value of the i' -th variable. If, as before, I do not assume independence, the model can be written as:

$$p_vote \sim \text{logit}^{-1}(a_{k[i]}),$$

$$a_k \sim N\left(\gamma_0 + \sum_{i'=1}^{n^{\text{county vars}}} x_{k[i],i'}\gamma_{i'}, \sigma_\alpha^2\right)$$

In the case of my specific data, for the time being I have county-level data for white population and urban population, so $n^{\text{county vars}} = 2$.

Individual Level

Assuming that I know the voter ID of the individual that cast their ballot, I can treat this piece of information in about the same way that I did for county as described above. This means that the following is mostly an exercise in maintaining notation constant. For these purposes, let n^{ID} be the number of total unique voter IDs—individuals—that I have data on, and j an indice that sums over all individuals. Also let z_j be an indicator variable for each individual. Then:

$$p_vote_i \sim \text{logit}^{-1}\left(\sum_{j=1}^{n^{ID}} z_j\beta_j\right)$$

And the second model, not assuming independence, would be:

$$p_vote \sim \text{logit}^{-1}(\delta_{j[i]}),$$

$$\delta_j \sim N(\zeta_0, \sigma_\delta^2)$$

Again, in a similar way to county level data, there are variables at an individual level, thus making it relatively easy to build further models. Let's say now that along with the one ballot, I was given a short

list of $n^{\text{indiv vars}}$ other individual-level variables, be they discrete, continuous, or indicators. The two models would then look like:

$$p_vote_i \sim \text{logit}^{-1} \left(\sum_{j=1}^{n^{ID}} z_j \beta_j + \sum_{i'=1}^{n^{\text{indiv vars}}} z_{j[i],i'} \beta_{i'+n^{ID}} \right)$$

Where $z_{j[i],l}$ is the j -th value of the i '-th variable. If, as before, I do not assume independence, the model can be written as:

$$p_vote \sim \text{logit}^{-1}(\delta_{j[i]}),$$

$$\delta_j \sim N(\zeta_0 + \sum_{i'=1}^{n^{\text{indiv vars}}} z_{j[i],i'} \delta_{i'}, \sigma_\delta^2)$$

In the case of my specific data, for the time being I have individual-level data for gender, so $n^{\text{indiv vars}} = 1$.

Election Level

Again as previously, four models come from including election level data. The first two are assuming I only knew what specific election the ballot comes from. Let $w_{i'}$ be an indicator variable for each election and n^{elect} the number of elections. The model assuming independence, with $w_{i'}$ being indicator variables for each election, is:

$$p_vote_i \sim \text{logit}^{-1} \left(\sum_{l=1}^{n^{\text{elect}}} w_l \beta_l \right)$$

Again, as previously, it would be safe to assume that each election is not held in a vacuum. Adding mixed effects this model would be:

$$p_vote \sim \text{logit}^{-1}(\eta_{l[i]}),$$

$$\eta_l \sim N(\nu_0, \sigma_\nu^2)$$

Again, in a similar way to county and individual level data, I add in variables at an election level. Let's say now that along with the one ballot, I was given a short list of $n^{\text{election vars}}$ other election-level variables, be they discrete, continuous, or indicators. The two models would then look like:

$$p_vote_i \sim \text{logit}^{-1} \left(\sum_{l=1}^{n^{\text{elect}}} w_l \beta_l + \sum_{i'=1}^{n^{\text{election vars}}} w_{l[i],i'} \beta_{i'+n^{\text{elect}}} \right)$$

Where $w_{l[i],i'}$ is the l -th value of the i '-th variable.

Assuming independence:

$$p_vote \sim \text{logit}^{-1}(\eta_{l[i]}),$$

$$\eta_l \sim N(\nu_0 + \sum_{i'=1}^{n^{\text{election vars}}} w_{l[i],i'} \nu_{i'}, \sigma_\nu^2)$$

For the time being I have two different variables that describe individual elections: date and type. Note that the above models may not be the best way to describe dates! An alternative could be fitting a glm, with some smoothing spline function for year. As for type, this would include four distinct indicators; one for each election type.

Ballot Level

In this section I assume that the ballot has some key features written on it, like the voting method, age, or party registration of the person that filled it out. A mixed effects model here would make no sense, since all the data is at the same unit of observation. Therefore, when adding ballot level variables, the model would look like:

$$p_vote_i \sim \text{logit}^{-1}(\beta_0 + \sum_{i'=1}^{n^{\text{ballot vars}}} u_{i,i'} \beta_{i'})$$

Where $u_{i,i'}$ is the i -th value of the i' -th variable, and $n^{\text{ballot vars}}$ is the number of ballot level variables. For now, I have data on voting method, age, and party. Voting method is coded as a binary variable with value one if the method was a Mail Vote. Party includes four distinct indicators for REP, DEM, Other, and Unaffiliated. Age is tricky; for now the options would be: straight up inclusion as an integer, inclusion as a cubic polynomial, inclusion as a 2nd degree polynomial, inclusion in some form of spline function.

Estimation with two types of data

After the work of setting up the four models at four different levels of observation, combining them in twos should be fairly straightforward. To avoid being needlessly cumulative, I will pursue this combination for County and Individual level only—instead of the six different possible combinations.

With the assumption that both counties and individuals are independent of one another, I proceed to the first type of model:

$$p_vote_i \sim \text{logit}^{-1}(\sum_{k=1}^{64} x_k \beta_k + \sum_{i'=1}^{n^{\text{county vars}}} x_{k[i],i'} \beta_{i'+64} + \sum_{j=1}^{n^{ID}} z_j \beta_{j+n^{\text{county vars}}+64} + \sum_{i'=1}^{n^{\text{indiv vars}}} z_{j[i],i'} \beta_{i'+n^{ID}+n^{\text{county vars}}+64})$$

This is large and clunky. It includes variables as described above: indicators for each county and individual, and all individual or county-level variables. For the corresponding mixed-effects model, I assume the tree-like structure we discussed on Monday. The hierarchy has two “levels”, with the second level consisting of two different regressions:

$$\begin{aligned} p_vote &\sim \text{logit}^{-1}(\delta_{j[i]} + a_{k[i]}), \\ a_k &\sim N(\gamma_0 + \sum_{i'=1}^{n^{\text{county vars}}} x_{k[i],i'} \gamma_{i'}, \sigma_\alpha^2) \\ \delta_j &\sim N(\zeta_0 + \sum_{i'=1}^{n^{\text{indiv vars}}} z_{j[i],i'} \delta_{i'}, \sigma_\delta^2) \end{aligned}$$

Estimation with the full dataset

I now proceed to include variables from all units of observation into one model. The first model, assuming independence, is:

$$p_vote_i \sim \text{logit}^{-1} \left(\sum_{k=1}^{64} x_k \beta_* + \sum_{i'=1}^{n^{\text{county vars}}} x_{k[i],i'} \beta_* + \sum_{j=1}^{n^{ID}} z_j \beta_* + \sum_{i'=1}^{n^{\text{indiv vars}}} z_{j[i],i'} \beta_* + \right. \\ \left. \sum_{l=1}^{n^{\text{elect}}} w_l \beta_* + \sum_{i'=1}^{n^{\text{election vars}}} w_{l[i],i'} \beta_* + \sum_{i'=1}^{n^{\text{ballot vars}}} u_{i,i'} \beta_* \right)$$

You will notice that I have omitted the subscript for all beta coefficients. This is because after two or three parameters, this becomes very, very large. I think it's reasonable to assume increasing indexes for different beta coefficients from left to right in this expression.

The mixed effects model will again operate on two “levels” of hierarchy, but the second level will now include three distinct regressions. Caveats for variables like age and date should be noted from previous sections.

$$p_vote \sim \text{logit}^{-1} \left(\sum_{i'=1}^{n^{\text{ballot vars}}} u_{i,i'} \beta_{i'} + \delta_{j[i]} + a_{k[i]} + \eta_{l[i]} \right),$$

$$a_k \sim N(\gamma_0 + \sum_{i'=1}^{n^{\text{county vars}}} x_{k[i],i'} \gamma_{i'}, \sigma_\alpha^2)$$

$$\delta_j \sim N(\zeta_0 + \sum_{i'=1}^{n^{\text{indiv vars}}} z_{j[i],i'} \delta_{i'}, \sigma_\delta^2)$$

$$\eta_l \sim N(\nu_0 + \sum_{i'=1}^{n^{\text{election vars}}} w_{l[i],i'} \nu_{i'}, \sigma_\nu^2)$$