

Analyzing Voter Registration in the 2016 Presidential Election

Andrew Amore

2022-11-21

1 Introduction

Political campaigns analyze historical election data to understand how demographic factors influence voter registration rates. This information can inform optimal advertising strategies which can drum up more votes and win elections, especially during high profile election years when voter turnout can be substantially higher. In this report 2016 election data from North Carolina counties was used to investigate how geography, gender, age and party affiliation influence registration behavior. The structure of the paper is as follows. . .

First, I provide some brief background on hierarchical modeling and motivate their usage under the voter registration context. Next, I review the data used in the analysis and discuss challenges arising from dated population estimates. Lastly, I specify several inference models and review their results for main questions of interest.

2 Background Information

A hierarchical model is a statistical inference tool that replicates the natural nesting structure of observed data. This structure can violate traditional modeling assumptions, like i.i.d. errors, that can lead to poor uncertainty estimation if ignored. To address these shortcomings, hierarchical models define parameters at different group levels and permits them to vary, which accounts for group heterogeneity and can improve uncertainty estimation. For this analysis, voter registration and population data are generally aggregated by geography (most naturally by county) during collection. Individuals may display different registration tendencies depending on the urban/rural environment of a particular geography, which lends itself well to hierarchical modeling, allowing parameters to vary by grouping structure.

Hierarchical models can be estimated using Bayesian or Frequentist methods. Under Bayesian estimation, parameters are treated as random quantities and derived from posterior distributions using Bayes rule. Posteriors provide more direct uncertainty quantification, via credible regions, that avoid the interpretation caveats of Frequentist based confidence intervals. However, this interpretability comes with an increased computational cost as Bayesian inference relies on expensive sampling methods compared to likelihood based optimization. This analysis predominantly uses Bayesian estimation methods using conjugate priors to avoid excessive computational costs. I'll now review the data used for the analysis.

3 Data Overview

To investigate main questions of interest data from the [2010 U.S. Census](#) was enhanced with North Carolina [voter registration](#) records from the 2016 Presidential election. To combine the data, demographic field values were standardized, as the coding structure varies slightly between State and Federal agencies. In total, there were 6,858 county level observations with unknown gender (~15%) which lack a population estimate from the Census data and these records were dropped. Irrelevant registration fields denoting precinct location were also removed. Metadata information for the combined dataset and a sample observation, can be viewed in Table 1.

Table 1: Metadata Information

Field Name	Description	Sample
Geography	County in North Carolina	BURKE
Age	Age Demographic Category	26–40
Gender	Gender Demographic Category	Female
Hispanic	Demographic Indicator of Hispanic Origin	Hispanic
Race	Race Demographic Category	WhiteAlone
VoterFreq	Number of Registered Voters in Specified Demographics	14
Freq	Total Population Count for Specified Demographics	128
TotalCountyPopulation	Total County Population	90912
PartyCd	Political Party Affiliation	DEM

3.1 Population Migration

The 2010 Census is assumed to represent the voter population during the 2016 election, however, **16.3%** of all observations have more registered voters than the demographic population estimates from the Census. Are we observing potential voter fraud or is historical Census information from six years ago too dated to reflect accurate estimates? To understand the scope of this issue, the difference between total registered voters and Census estimates were computed and aggregated by county. Figure 1 shows summaries for geographies with more than five observations for ease of viewing. Coincidentally, **16%** of counties have median population difference greater than **100 individuals**, however, **70%** of counties have median difference less than **50** and **35%** have medians less than **10**. The majority of moderately small differences assuages some concerns of major population shifts over the six year period. Population estimates for the “invalid” observations are set to the sum of registered voters, but could also be inflated with a correction factor learned from other observations.

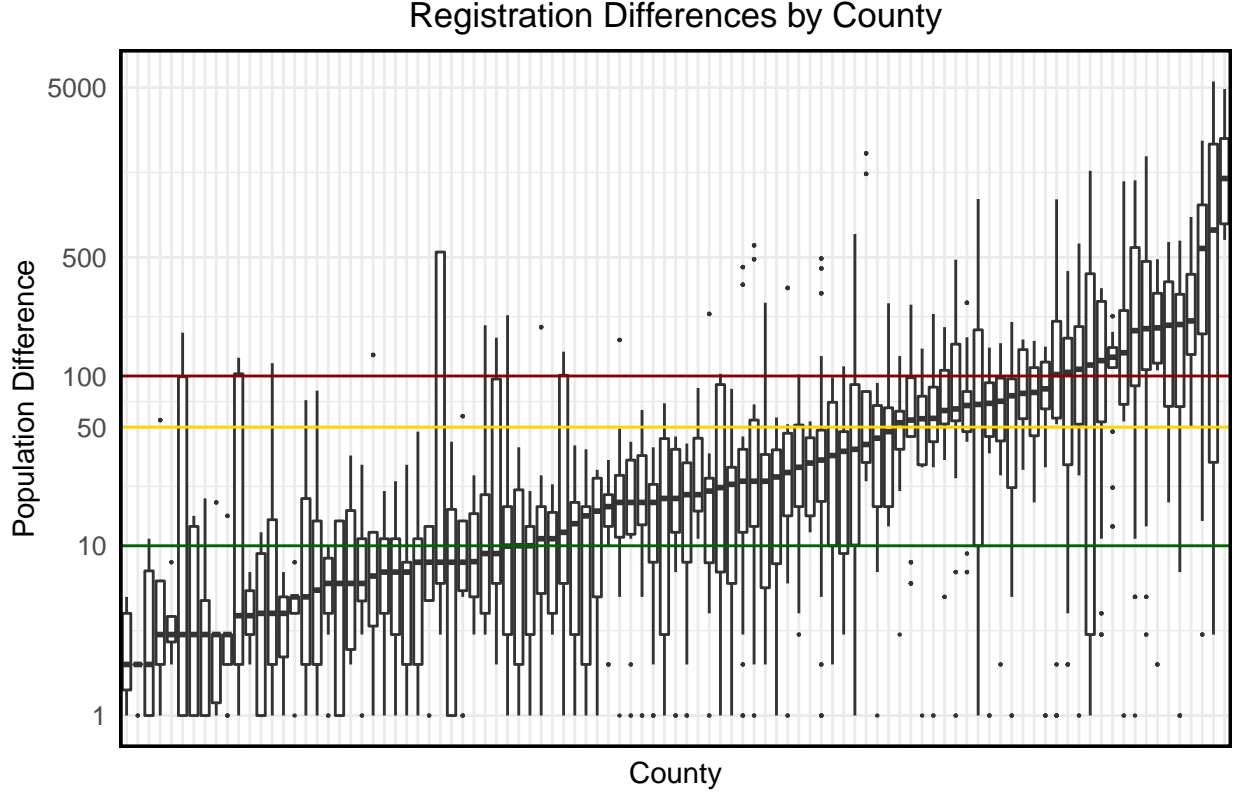


Figure 1: Registration Differences by County

To minimize error, 30 counties were randomly selected using sampling weights inversely proportional to the percentage of faulty observations (total registered voters > demographic population). This framework can be extended to bootstrap resample different county combinations for parameter estimation which reduces variation by steering our model towards more accurate data. An additional benefit of sub-sampling is a reduction in computation time for MCMC sampling for parameter estimation in Bayesian models.

3.2 Motivating a Multilevel Model

Voter data and population estimates are naturally grouped by geography, suggesting a multilevel model using the county of measurement. We are interested in understanding how covariates, like age and political affiliation, influence registration rates. A model, without any hierarchical components, can be formatted as a binomial regression, where registered voters represent the number of “successes” from the population (trials). To evaluate modeling decisions, like where to apply random effects, an exploratory data analysis was conducted. Figure 2 displays facets for different binomial regressions on registration probability for two demographic categories, color coded by county. (A) fits regressions by county/age with substantially different trends across/within facets. (B) fits regressions by county/race with lower facet variability. Additional plots can be found in the Appendix. Continuous covariates hindered posterior mixing with brms and were excluded.

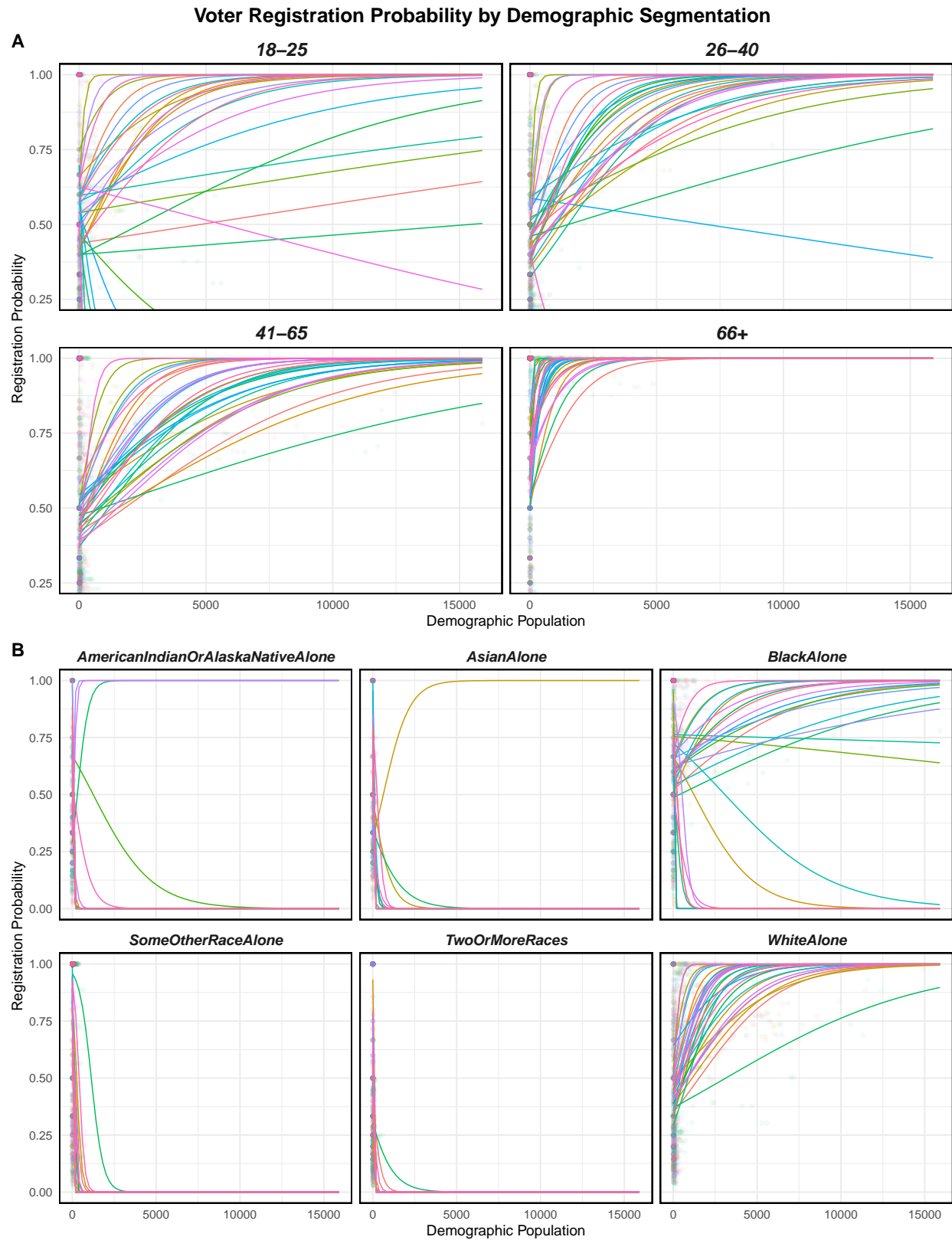


Figure 2: Voter Registration Behavior for (A) Age Category and (B) Race

Under the current model it is impossible to determine the effect of political party on registration probability as the Census doesn't collect political affiliation with population estimates and needs to be imputed. As an initial estimate, political affiliation can be assigned to unregistered voters from the observed political distribution of registered voters in each demographic group.

4 Multilevel Models

4.1 Model Specifications

The modeling framework can now be specified as follows. Let (y_i, n_i, p_i) denote the number of voter registrations, demographic population and registration probability respectively for observation i . y_i is assumed to be distributed as a binomial random variable.

$$y_i \sim \text{Binomial}(n_i, p_i)$$

Individuals can only register to vote once, but binomial trials sample independently with replacement. The independent draw assumption seems reasonable, as the probability of any one individual registering to vote is most likely independent of any other individual's decision (more on this later). Additionally, $\sim 70\%$ of demographic categories in our dataset have populations $n_i \geq 30$. However, the demographic samples can be considered population estimates and may warrant a hypergeometric distribution, but this induces a dependence between y_i 's which is unwarranted based on our independent registration assumption. Under the binomial model, several multilevel structures with varying complexity were evaluated. Table 3 displays the specifications.

Table 2: Random Effect Structures

	Fixed Effect Fields	Random Intercept Fields	Random Slope Fields
<i>Model I</i>	Gender, Hispanic, Race, PartyCd, Age	Geography	
<i>Model II</i>	Hispanic, Race, Age	PartyCD	Gender
<i>Model III</i>	Gender, Hispanic, Race	PartyCD	Age

4.1.1 Model I

Let μ denote a global intercept corresponding to a baseline demographic category with the following values: 'Female', 'Hispanic', 'AmericanIndianOrAlaskaNativeAlone', 'DEM' and '18-25'. Let θ_j denote a **county** random effect intercept for observations, i , in county j , X_i a vector of demographic indicator variables corresponding to column 1 in Table 2 and β the corresponding fixed effect estimates.

$$\text{Logit}(p_{ij}) = \mu + \theta_j + X_i\beta$$

Prior Specifications

$$\begin{aligned} \mu &\sim N(0, 1), \quad \beta \sim N(0, 1), \quad \theta_j \sim N(0, \sigma) \\ \sigma &\sim \text{HalfCauchy}\left(0, \frac{1}{2}\right) \end{aligned}$$

4.1.2 Model II

In accordance with the model specification in Table 2, let μ denote a global intercept corresponding to a baseline demographic category with the following values: ‘Hispanic’, ‘AmericanIndianOrAlaskaNativeAlone’ and ‘18-25’. Let Ω_p denote the **party affiliation** random effect for party p , let g_{ipk} denote the **gender** category, k , for observation i in political party p and Γ_{pk} the random effect. Finally, let X_i be a vector of demographic indicator variables corresponding to column 1 in Table 2 and β the corresponding fixed effect estimates.

$$\text{Logit}(p_{ipk}) = \mu + \Omega_p + g_{ipk}\Gamma_{pk} + X_i\beta$$

Prior Specifications

$$\Omega_p \sim N(0, \sigma), \Gamma_{pk} \sim N(0, \gamma), \sigma/\gamma \sim \text{HalfCauchy}\left(0, \frac{1}{2}\right)$$

$$\mu \sim N(0, 1), \beta \sim N(0, 1)$$

4.1.3 Model III

In accordance with the model specification in Table 2, let μ denote a global intercept corresponding to a baseline demographic category with the following values: ‘Female’, ‘Hispanic’ and ‘AmericanIndianOrAlaskaNativeAlone’. Let Ω_p denote the **party affiliation** random effect for party p , let a_{ipk} denote the **age** category, k for observation i in political party p and α_{pk} the random effect slope. Finally, let X_i be a vector of demographic indicator variables corresponding to column 1 in Table 2 and β the corresponding fixed effect estimates.

$$\text{Logit}(p_{ipk}) = \mu + \Omega_p + a_{ipk}\alpha_{pk} + X_i\beta$$

Prior Specifications

$$\Omega_p \sim N(0, \omega), \alpha_{pk} \sim N(0, \tau), \omega/\tau \sim \text{HalfCauchy}\left(0, \frac{1}{2}\right)$$

$$\mu \sim N(0, 1), \beta \sim N(0, 1)$$

5 Results

MCMC diagnostics, like trace plots, can be viewed in the Appendix for each model. The main questions of interest concern demographic factors influencing registration probabilities.

5.1 *How did different demographic subgroups register to vote?*

The fixed effect estimates from each model are plotted in Figure 3 with 95% credible regions (CR). Note, not all models include identical fixed effects, but the overlapping estimates tend to display more shrinkage with increased model complexity. Starting from the top of Figure 3, all Race indicators have significant impact on registration rate. Party affiliation displays no significant effect. Hispanic and Gender estimates are also significant. Older individuals are also more likely to register, but this can be attributed to having more opportunities/elections to do so.

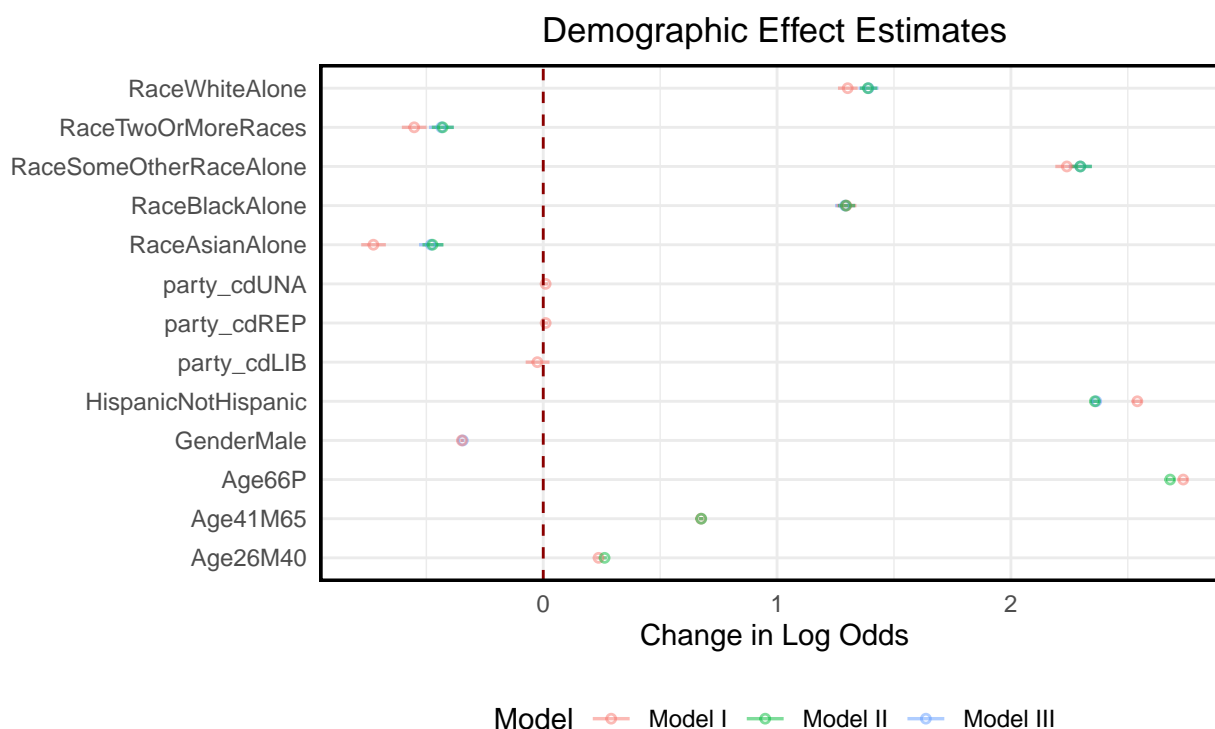


Figure 3: Fixed Effect Estimates From Each Model

5.2 *Did the overall odds of registering differ by county in 2016?*

Only 30 counties are included in this analysis and the list can be viewed in the Appendix. The random intercepts from Model I, with 95% CR are displayed in Figure 4. Tyrrell, the only county in the dataset without any “invalid” observations has the lowest registration effect, while **43%** of counties have no detectable effect.

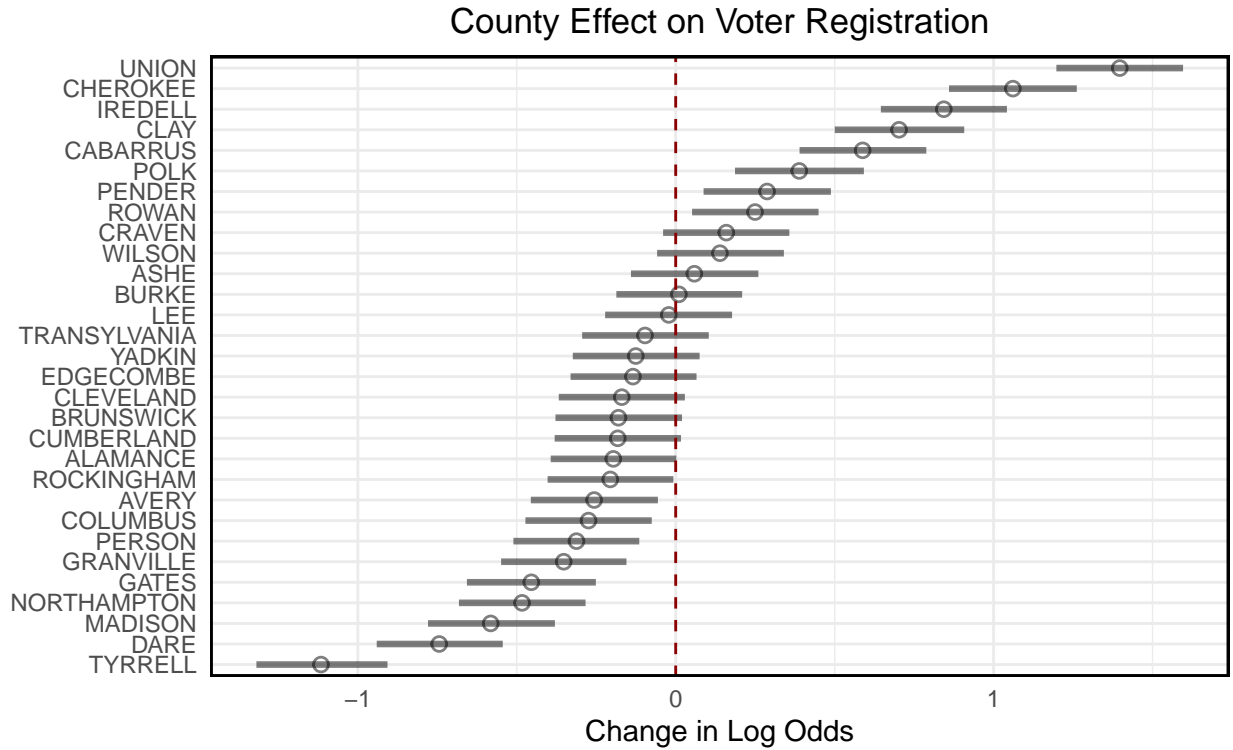


Figure 4: County Level Intercept Estimates from Model II

5.3 *How do registration rates differ by gender for each party?*

Using Model II we can look at the gender effect within each political party from random slope estimates and assess the relative change in log odds from the baseline gender (Female). Figure 5 shows 95% CRs and point estimates. Males across all political parties are less likely to register than corresponding females and republican males register at the highest rates.

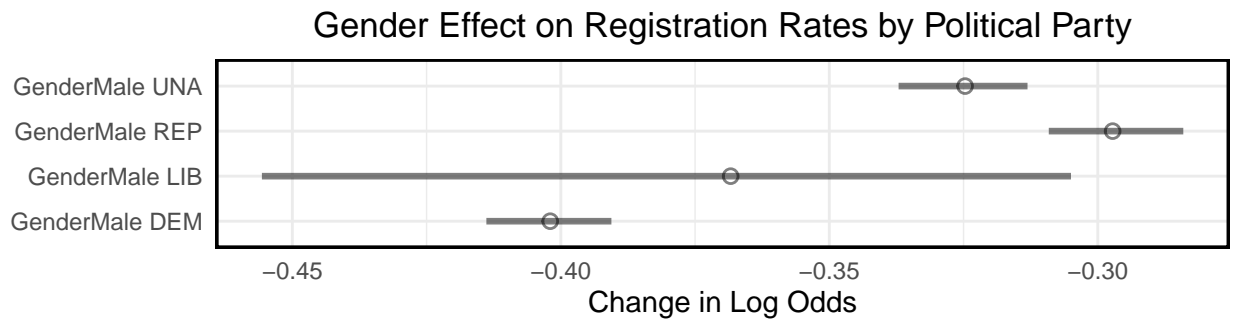


Figure 5: Random Slope Estimates (Model II) for Gender across Political Party

5.4 Registration differences by age group for each party?

Random slope estimates from Model III are the registration effects for different age groups relative to a baseline (18-25). Figure 6 shows 95% CRs and point estimates for different political parties and age groups. Across age categories republicans are more likely to register than other political parties and registration rate tends to increase with age.

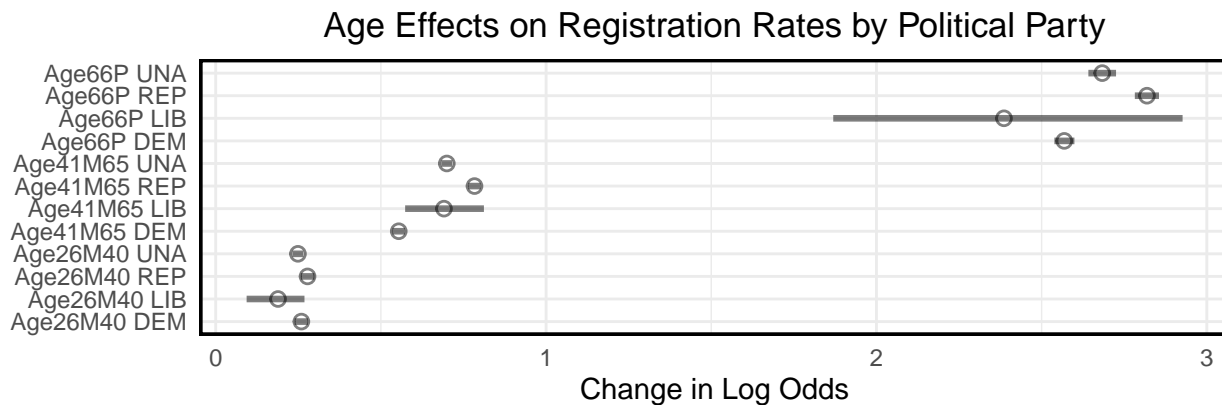


Figure 6: Random Slope Estimates for Age across Political Party

5.5 Predictive Accuracy

Comparing predicted voters to actual registrations using the out-of-sample data can assess model validity. Figure 7 displays predictions for each demographic group within one county and dashed lines denote average error. Surprisingly, Model I has the best predictive performance, with average error of roughly 10 individuals, but median error of less than one.

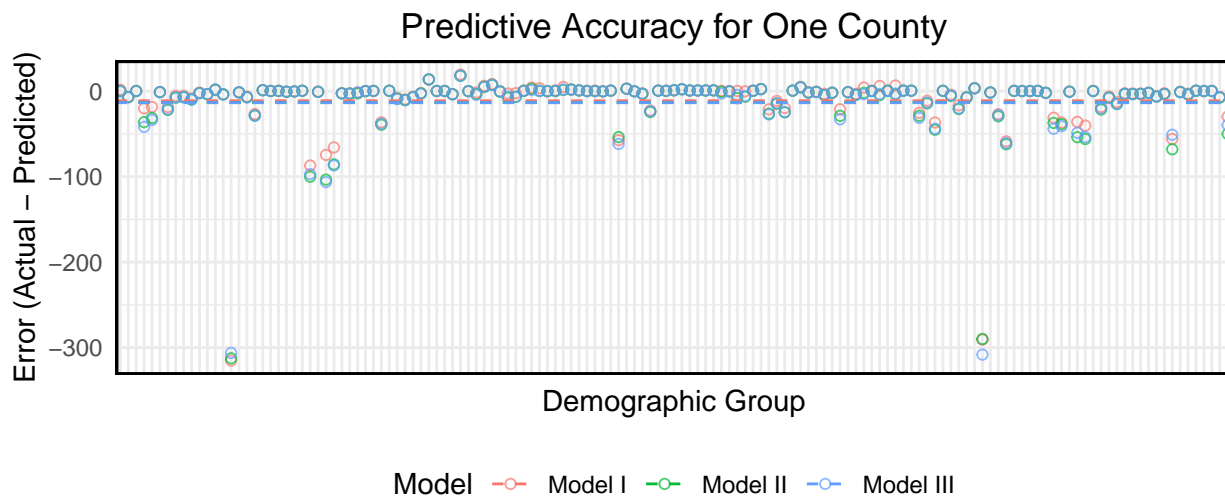


Figure 7: Predictive Accuracy across Demographic Groups

Table 3 displays the overall predictive performance for the out-of-sample test set. The average absolute error denotes the difference, in number of individuals, between the model predictions and the observed number of registered voters. The median error is substantial less than the average and indicates some predictions are significantly wrong.

Table 3: Prediction Summaries

	Avg. Absolute Error	Median Absolute Error
<i>Model I</i>	38.58	2.56
<i>Model II</i>	35.92	2.69
<i>Model III</i>	35.91	2.69

Table 4 shows the 95% credible region summaries for each model prediction. All models display substantially worse performance on the out-of-sample test set than the coverage level. Model I predictions also have dramatically wider intervals than other models.

Table 4: Confidence Interval Summaries

	Confidence Interval Coverage	Avg. CI Width
<i>Model I</i>	0.8	150.88
<i>Model II</i>	0.6	16.4
<i>Model III</i>	0.61	16.46

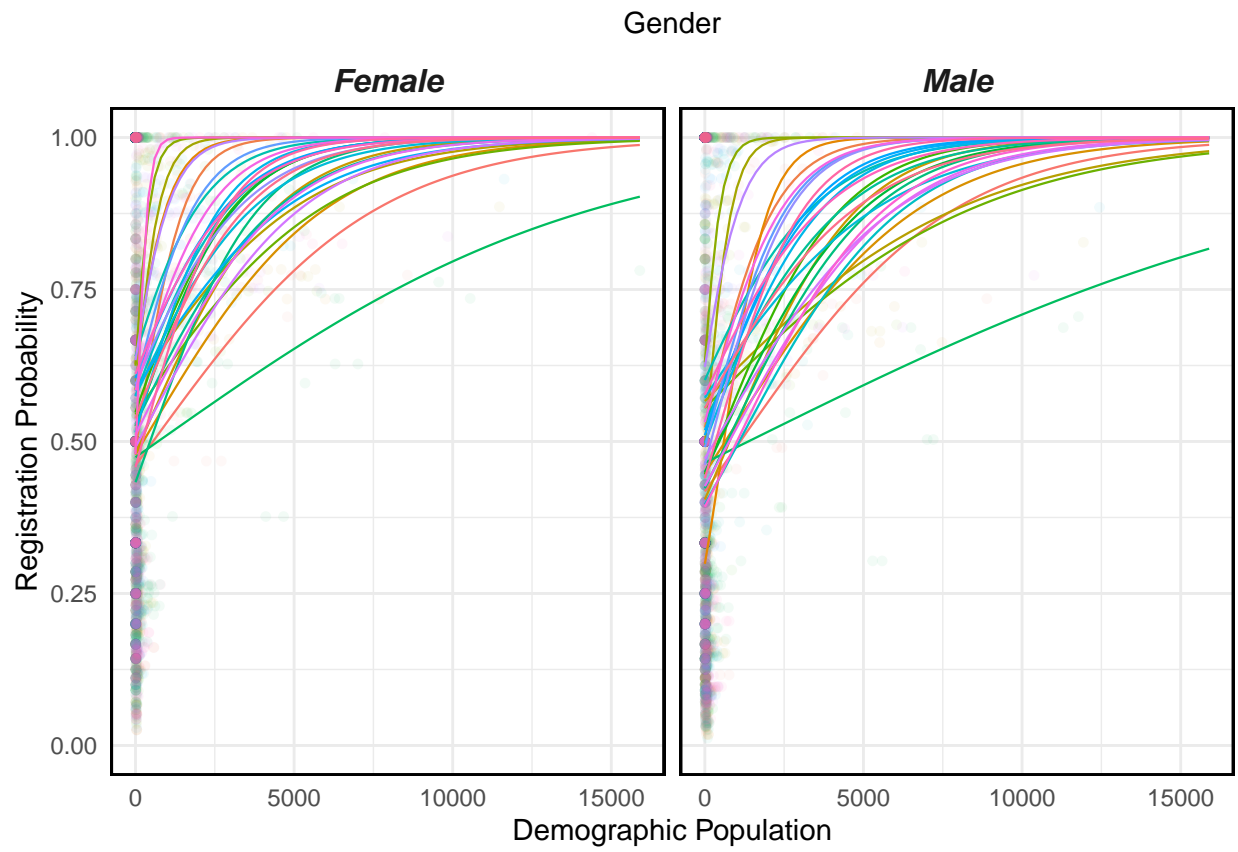
6 Limitations & Conclusion

Population estimates in 2016 are approximated from the 2010 Census which may introduce unaccounted for variability if significant population changes occurred over the six year period. There are multiple instances of “invalid” observations with more registered voters than demographic population estimates, which should be impossible. To minimize variability a weighted sub-sample of 30 counties was drawn for estimation, with weights inversely proportional to the fraction of invalid observations. This addresses detectable differences, however, does not account for the possibility that all population estimates differ over a six year period. In the future, one could look to include 2020 Census estimates, when they become available in 2023, and impute more realistic population measures using estimates from the 2010 and 2020 Census. The sub-sample also reduces the size of the dataset, adding more variability, but required to reduce the MCMC sampling time which can still take over 12 hours on the reduced set. To alleviate variance concerns in fixed effect estimates, one could leverage the entire dataset with a bootstrap resampling procedure and run multiple models in parallel. An additional limitation is the unknown population estimates for the political affiliation of unregistered voters as the Census does not contain this information. These totals were imputed based on the party affiliations of registered voters within each demographic category as a rough approximation, but introduces another source of error.

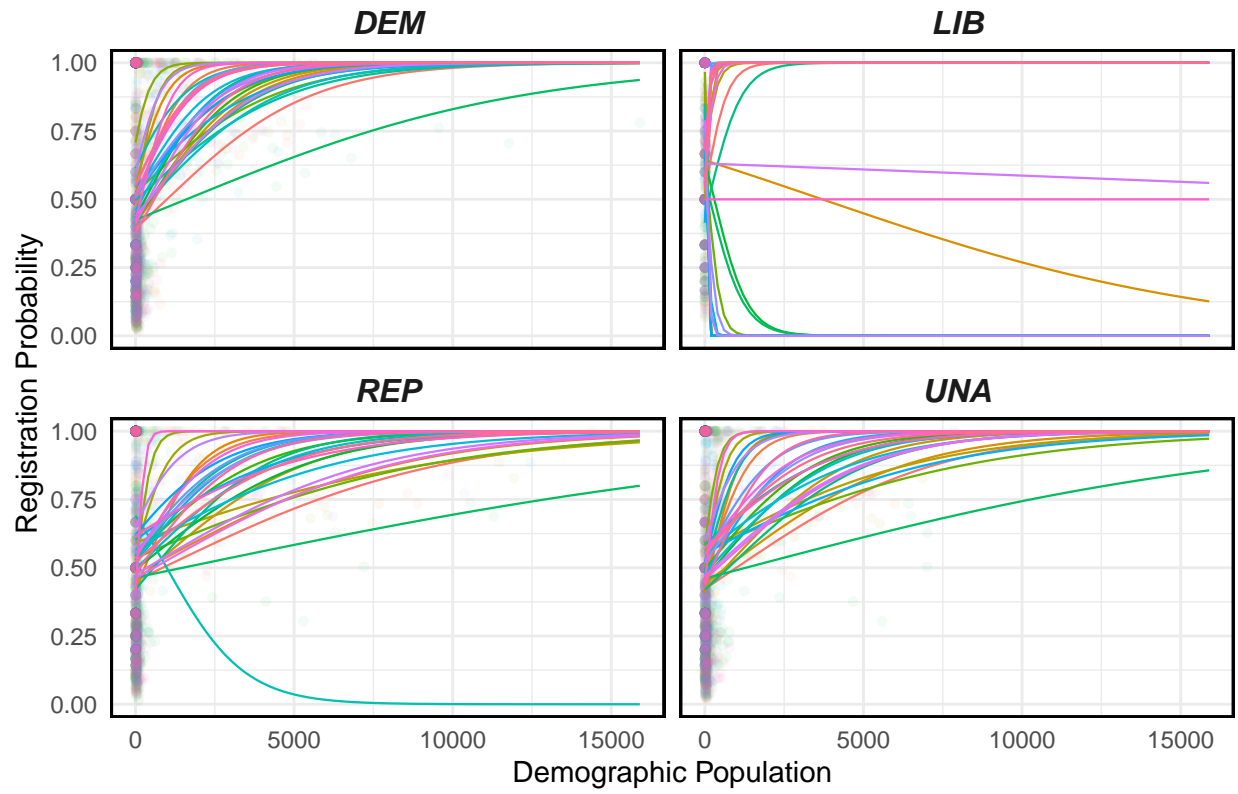
Historical voting records and Census population estimates were used to evaluate demographic differences in registration tendencies. Three Bayesian multilevel models were created to address specific questions of interest using different effect structures. The full dataset was sub-sample to improve computation time, but may lead to higher variance estimates. Preliminary results suggest significant differences in registration tendencies across demographic categories in North Carolina.

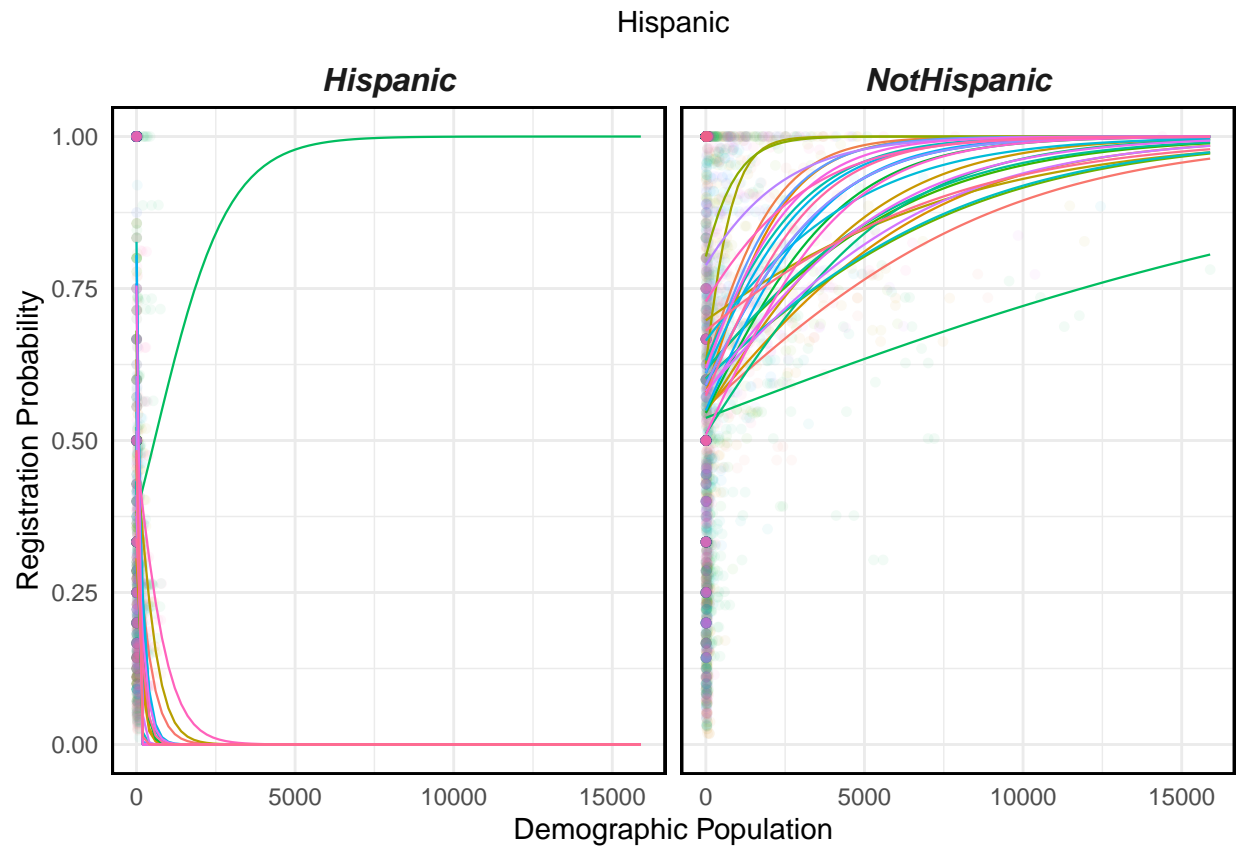
7 Appendix

Additional Covariate Plots

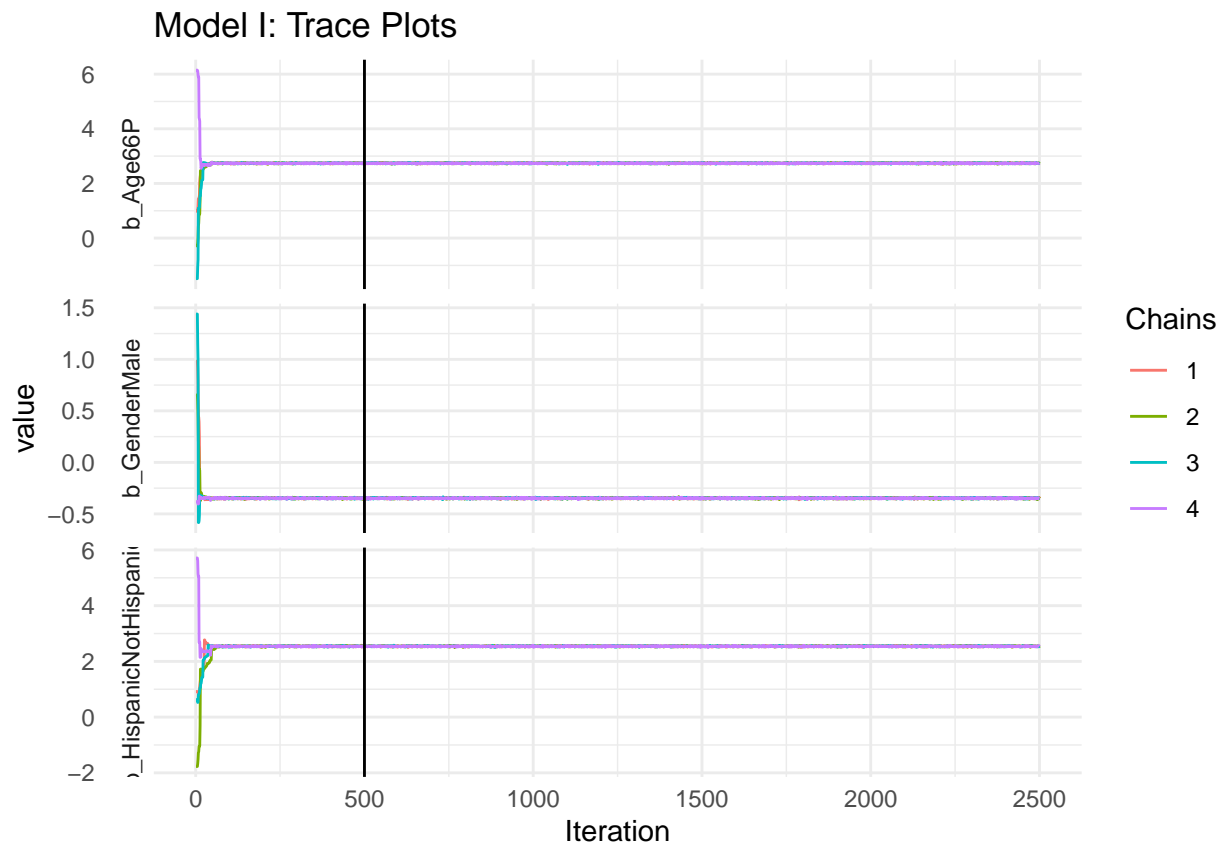


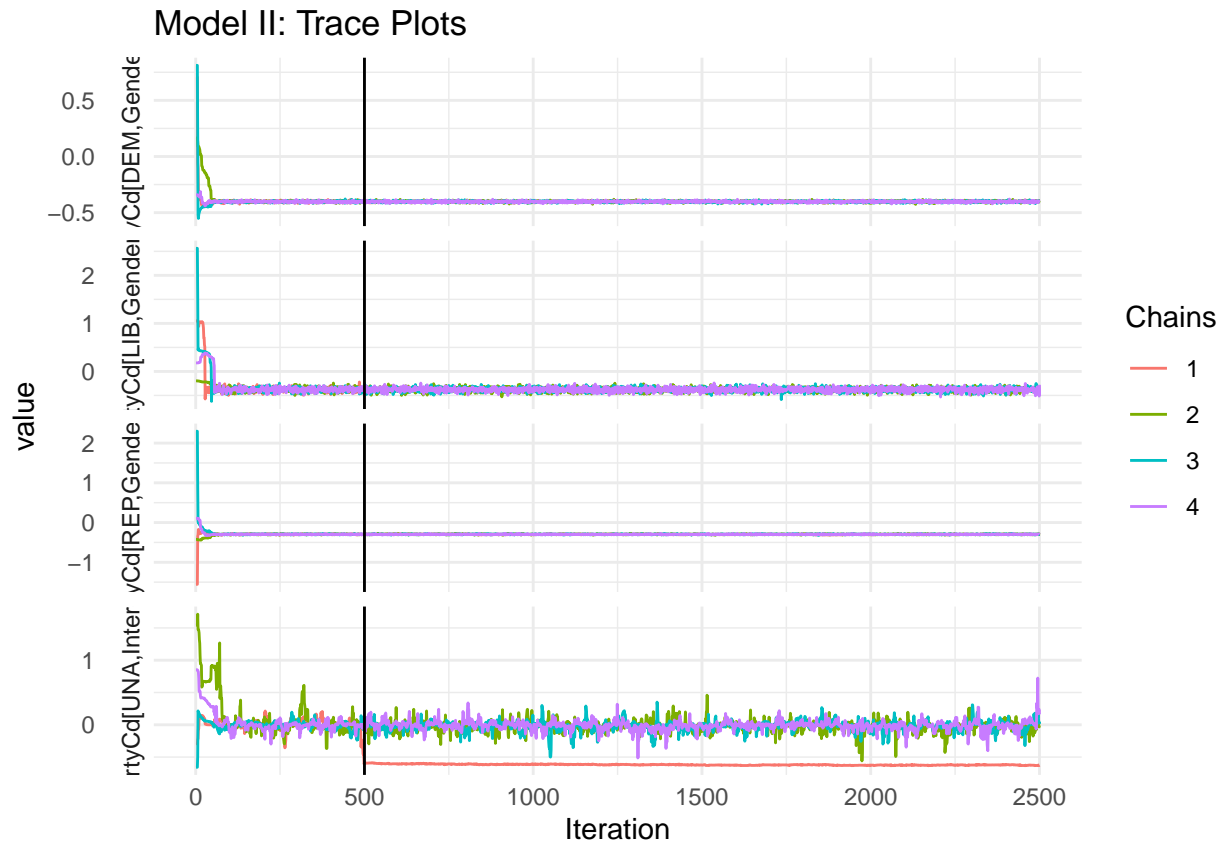
Party Affiliation

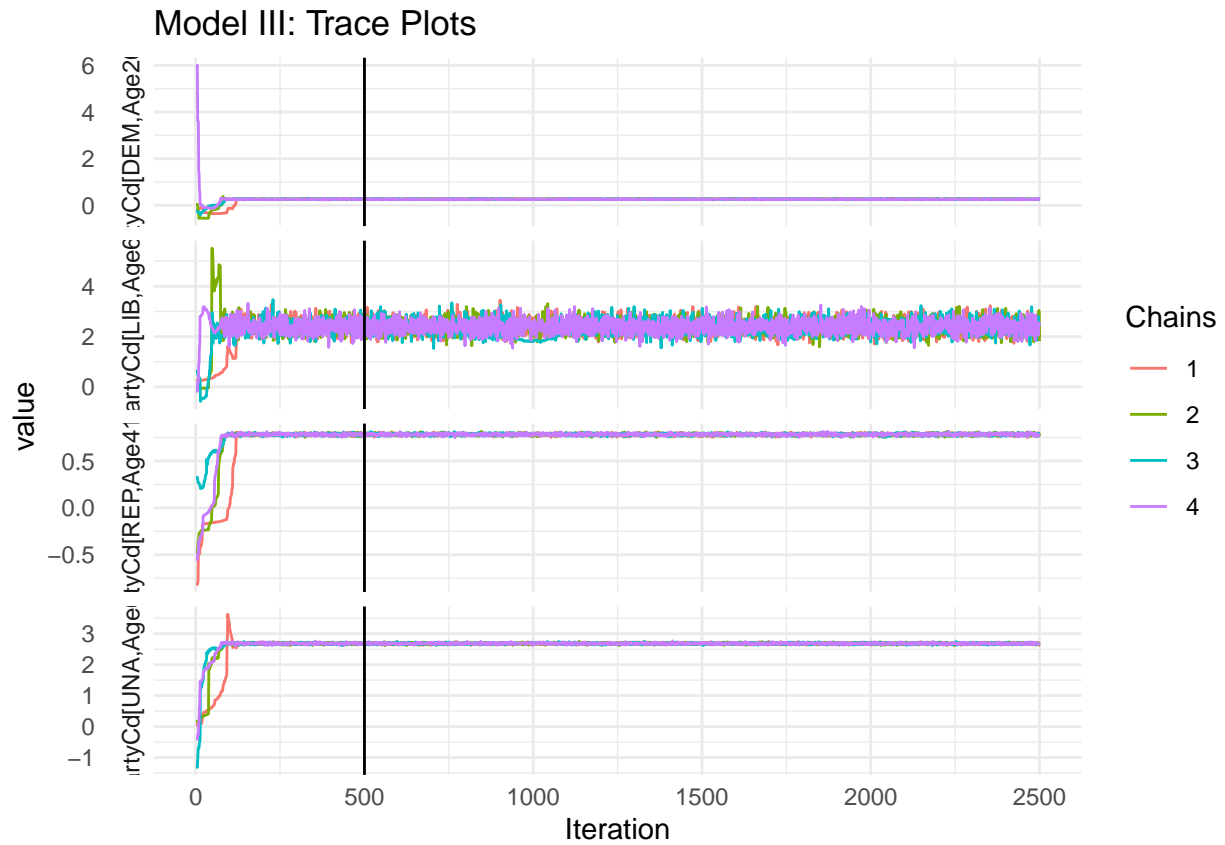




Trace Plots







County Sample List

Geography
ALAMANCE
ASHE
AVERY
BRUNSWICK
BURKE
CABARRUS
CHEROKEE
CLAY
CLEVELAND
COLUMBUS
CRAVEN
CUMBERLAND
DARE
EDGECOMBE
GATES
GRANVILLE
IREDELL

Geography

LEE
MADISON
NORTHAMPTON
PENDER
PERSON
POLK
ROCKINGHAM
ROWAN
TRANSYLVANIA
TYRRELL
UNION
WILSON
YADKIN
