

# Final Written Report: What makes someone more likely to vote?

Approximately Normal: Jenny Huang, Bella Larsen, Jessie Bierschenk, Aidan Gildea

17 November 2020

## Introduction and Data

In our regression analysis, we will be assuming the role of a nonprofit political organization aiming to design the most effective GOTV (“get out the vote”) effort possible. Our primary focus is to increase voter turnout, however, as an organization funded solely by donations, we need to make sure our outreach efforts are as financially efficient as possible. Therefore, it is critical that we are sending our mailing literature primarily to those we believe may not vote – as this will help turnout voters who may not of voted otherwise. To do so, we will conduct a regression analysis that investigates the different factors and characteristics that appear to be involved or related to voters in the U.S. We are using data on voting behaviors in the U.S. over the past 14 years [1]. These data are sourced from IPUMS (an organization that provides census and survey data) and the American Statistical Association (ASA) [2, 1]. These data include 28 variables on more than 640,000 voters in the U.S. The data collected contains voter characteristics such as age, geographic location, sex, race, marital status, employment, citizenship, ethnicity, education, and voting history and tendencies [1]. This information is particularly relevant in exploring what factors may be related to U.S. - Americans voting or not.

Motivation for our research comes from previous efforts to predict voting behaviors. In an MIT Election data study, researchers describe how understanding voter turnout is important when observing the particular tendencies of certain groups of people as well as factors that motivate individual U.S. citizens to vote [3]. The study highlights general assumptions of voter turnout predictions, noting how higher turnout rates tend to be related to individuals with the following traits: married, white, female, higher education, higher income, older age [3]. The article also addresses how reform may be able to increase voter turnout [3]. In another study done by Harvard graduate student Anthony George Fowler, voter turnout and its implications and repercussions are further examined in the U.S. as well as Australia and Mexico [4]. The study explores the 2008 U.S. election and addresses partisan gaps, voter knowledge (how politically-informed a voter is), and race as main variables of interest in exploring voter tendencies in the U.S. [4]. Both of these studies provide motivation for further and continued investigation into voter data and statistics – especially for our efforts to understand what populations usually do not make it to the voting booth.

In anticipation of our GOTV effort, we are interested in predicting whether a person voted or not based on a list of predictor variables including sex, age, marital status, veteran status, citizenship status (native born or naturalized citizen), whether or not someone is Hispanic or Latinx, employment status and more (described in more detail in Section 2). Our proposed research question is: do voter turnout rates depend on these predictors? Which predictors are more impactful than others? We are also interested in looking at voter turnout over time. We will use the predictor (year) to see if there are changes in voter turnout by demographics over time, or perhaps compare models from different time periods to determine how voter trends have changed over time. Our organization’s initial hypothesis is that the significant predictors of voter turnout will include age, level of education, whether they voted in previous elections, and race. Based on historical patterns, people in older age categories tended to vote more than people in younger age categories and people with higher levels of education tended to vote more frequently than people with lower levels of education [5]. Finally, if a person has voted previously, we predict they will be more likely to vote again compared to someone who has not voted previously. Taking these hypotheses into consideration, we will explore which factors are most significant in relation to voting attendance. Our findings will hopefully allow

us to identify populations that are statistically less likely to vote, informing who we target with our GOTV literature in the future.

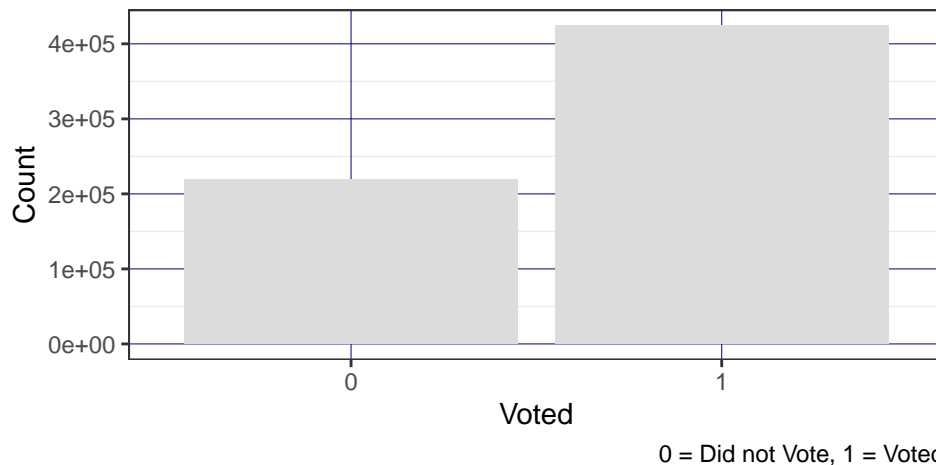
While our GOTV initiative hopes to turnout people who we determine are unlikely to vote, our statistical analysis will use a broad margin to determine who should receive our informational materials. Every vote matters, so we want to make sure no one gets left behind!

We will begin our EDA by visualizing the relationship between the response variable describing whether or not someone voted and several of the other variables of particular interest.

First, we will simply look at the distribution of those who voted throughout the last 8 years of elections.

## Visualizing the Distribution of Voting Status

*More people reportedly voted than did not vote*



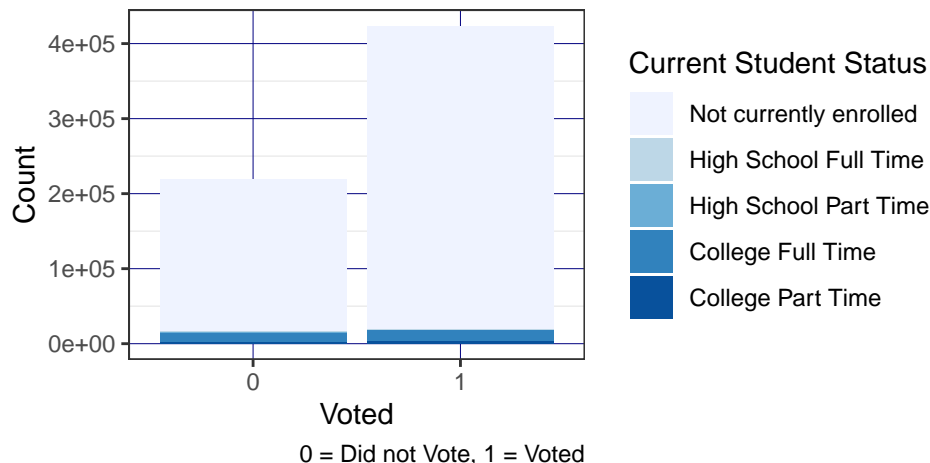
*see theme code inspiration at reference [6] see scale fill code inspiration at reference [7]*

From the barplot above, it is clear the more individuals in the data set voted (voted = 1) than did not (voted = 0).

Many political nonprofits engage with college campuses, so we want to analyze whether or not being a student influences the frequency of voting. We will explore this preliminarily by visualizing the distribution of if school aged individuals (18-24) voted or not – categorized by their current student level. This is seen in the bar plot below.

## Voting Distribution of 16–24 Year Olds

*Examining relationship between student status and voting*



see *scale fill brewer code inspiration from reference [6]*

From the bar plot, it is evident that a majority of these individuals were not currently enrolled. This may be a result of a general national trend, but we want to investigate if it is the result of a larger proportion of older individuals within in the range of ages between 16-24. We will investigate this by analyzing those who are not currently enrolled in school within this age range.

Table 1: Proportion of Respondents Not Currently Enrolled In School By Age

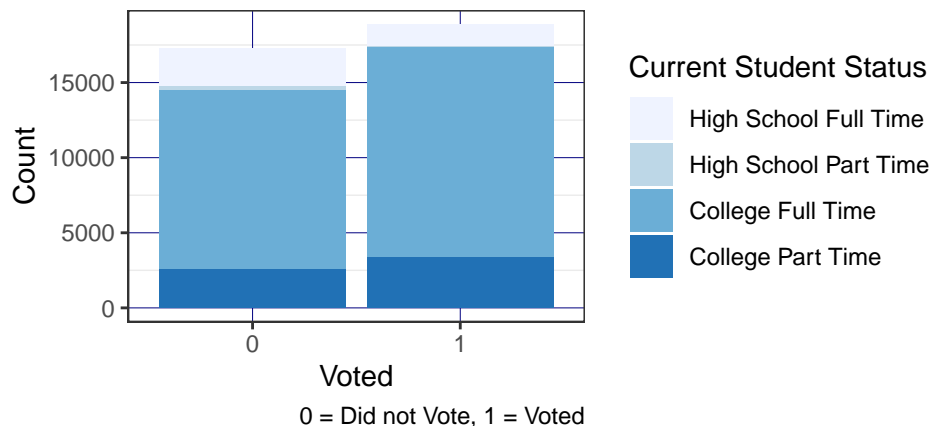
AGE	n	prop
18	2327	0.065
19	3671	0.102
20	4367	0.122
21	4760	0.133
22	5955	0.166
23	6991	0.195
24	7791	0.217

From the table above, it is apparent that more than 40% of those not currently enrolled in school are 23-24 years old. This could be a potential reason for why this age range includes so many who are not currently enrolled as a student.

To more meaningfully analyze the relationship between being a student and if they vote or not, we adjusted our visualization to only include those currently enrolled in some level of education. This is seen in the visualization below.

## Voting Distribution of 16–24 Year Olds Enrolled in School

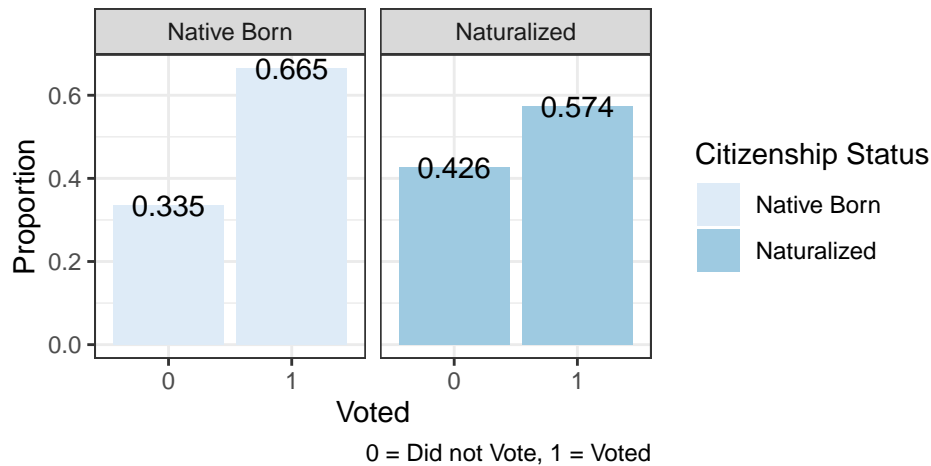
*Examining relationship between student status and voting*



This visualization shows that eligible voters between the age 16-24 enrolled in some education at the time were mainly full time college students. More full time and part time college students that were eligible to vote did vote compared to those that did not vote. The opposite is true for high school students: more full time and part time high school students that were eligible to vote did not vote compared to those that did vote.

## Voting distribution based on citizenship status

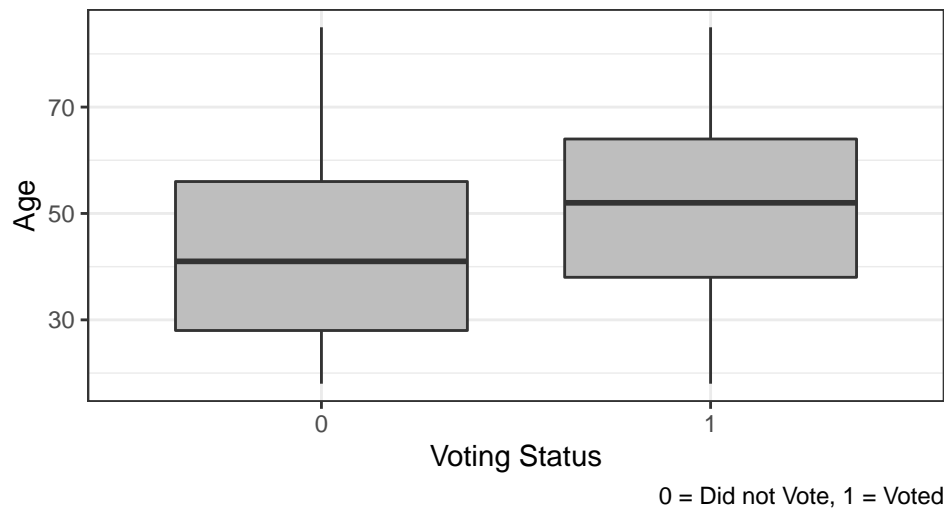
Examining relationship between citizenship status and voting



*see geom\_text() code inspiration in reference [8]*

This visualization shows that for both native born and naturalized individuals, more citizens that were eligible to vote did vote compared to those that did not vote; however, the proportion is much greater for native born citizens than for naturalized. For a similar plot comparing respondents by veteran status, please consult Appendix A.

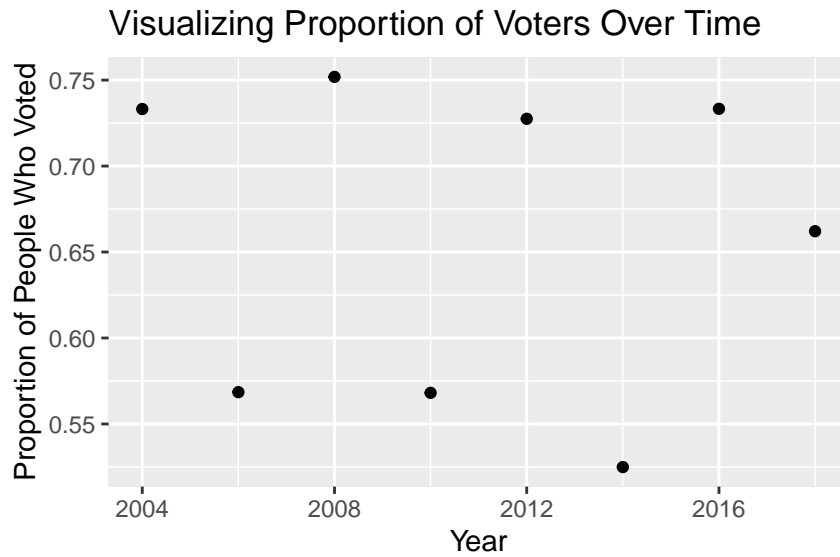
## Relationship between age and voting



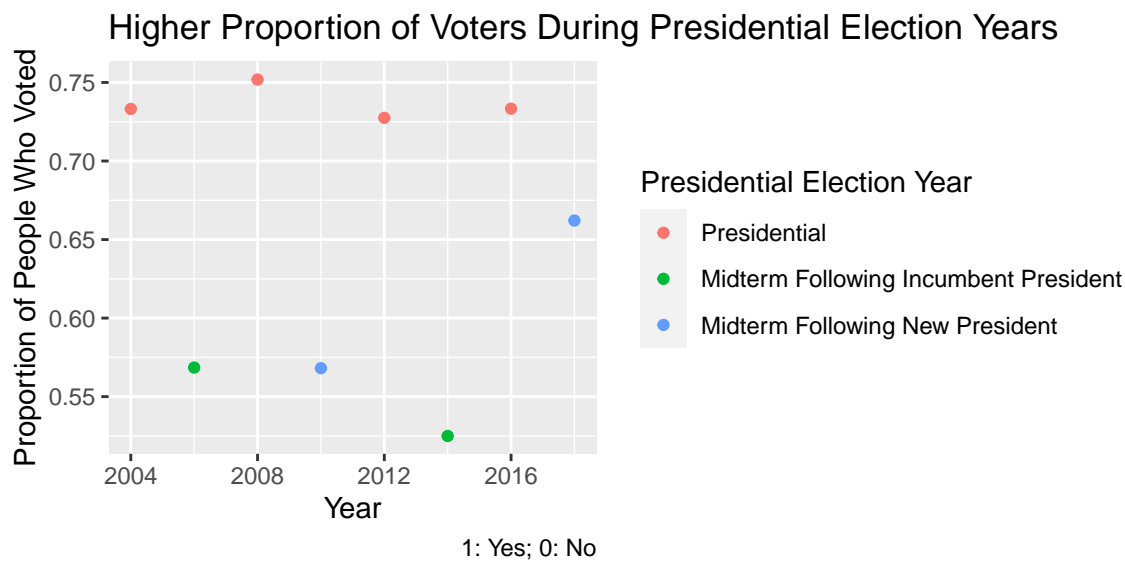
This box plot shows that eligible voters that did vote were generally older than eligible voters that did not vote.

We are also interested in looking at how voter turnout has changed over the years.

We notice that the proportion of people who voted fluctuates depending on whether the year falls on a presidential election. In the trend of the proportion of voting over time, we see a clear divide between the years when there is a presidential elections versus when there is not. In the future, we may decide to add the variable "Election Year" as an interaction term with year as a divide between years that fall on an election.



We decided to look at the scatterplot of voter turnout over time broken down by the status of the election (Presidential election, midterm after incumbent president, or midterm after new president) to see how it may differ depending on the election.



From the above scatterplot, we confirmed that there is higher voter turnout during presidential elections compared to midterm elections. In addition, there is equal or higher voter turnout for midterm elections following the election of a new president compared to midterm elections following the election of an incumbent president, and an especially high voter turnout in the midterm election after Trump’s election in 2016.

Finally, it is important to acknowledge that our any missingness in our data, as it may influence the outcome of our regression analysis. The initial data sourced from the This Is Statistics: Fall Data Challenge used non-uniform values for the different variables, therefore we cleaned the data to be more intelligible for our analysis. While cleaning, we did impute “NA” values by either combining them with other categories and/or removing them to improve our analysis. We believe this will not negatively affect our analysis, yet will move forward taking it into consideration.

After analyzing these preliminary visualizations and observations, we will now begin to build our model.

# Methodology

## Model Selection

### Select a random subset of the data to create model.

In order to make our model, we have decided to take a random sample of 10,000 to be sure that the model selection is accurate and not obscured by too many observations.

Using the random sample of 10,000 observations, we began the model selection process. See Appendix B for the full backward selection model output.

Table 2: Model Resulting From Backward Selection

term	estimate	std.error	statistic	p.value
(Intercept)	16.369	11.317	1.446	0.148
sexMale	-0.078	0.048	-1.628	0.104
marstDivorced/Separated	-0.702	0.071	-9.924	0.000
marstNot Married/Other	-0.474	0.057	-8.336	0.000
citizenNaturalized	-0.684	0.103	-6.637	0.000
employedYes	0.406	0.060	6.768	0.000
highest_educationHigh School Degree/GED	-1.357	0.074	-18.430	0.000
highest_educationSome College	-0.630	0.081	-7.744	0.000
highest_educationSome High School	-2.220	0.099	-22.398	0.000
highest_educationAssociate Degree	-0.524	0.095	-5.494	0.000
highest_educationMasters Degree	0.287	0.125	2.304	0.021
highest_educationProfessional Degree	0.487	0.282	1.727	0.084
highest_educationDoctorate Degree	0.403	0.275	1.464	0.143
highest_educationNone/Unknown	-3.161	0.653	-4.838	0.000
current_studentHigh School Full Time	1.177	0.305	3.865	0.000
current_studentHigh School Part Time	0.299	1.134	0.264	0.792
current_studentCollege Full Time	0.444	0.123	3.626	0.000
current_studentCollege Part Time	0.573	0.248	2.312	0.021
raceBlack	0.599	0.084	7.142	0.000
raceAsian or Pacific Islander	-0.372	0.133	-2.791	0.005
raceNative American	-0.060	0.202	-0.295	0.768
race2 or more races	0.258	0.205	1.256	0.209
AGE	0.040	0.002	22.706	0.000
Presidential_Election_StatusMidterm Following Incumbent President	-0.989	0.057	-17.347	0.000
Presidential_Election_StatusMidterm Following New President	-0.669	0.063	-10.696	0.000
YEAR	-0.008	0.006	-1.440	0.150

The final model included YEAR + sex + current\_student + citizen + employed + race + marst + Presidential\_Election\_Status + AGE + highest\_education.

The final model is:

The backward selection based on AIC took out the variables metro, veteran, and hispanic\_status (see Appendix B).

Interpretation of coefficients of interest:

Baseline: Female, married, native born, not employed, Bachelors Degree is highest education, not currently enrolled in school if between the age of 16-24, White, and the time of voting is midterm following the election

of an incumbent president.

We expect the odds of an eligible voter voting for a divorced/separated eligible voter to be 0.50 ( $\exp(-0.702)$ ) times the odds of individuals that are female, married, native born, not employed, have Bachelors Degree as their highest education, not currently enrolled in school if between the age of 16-24, White, and the election is a midterm election following the election of an incumbent president, after factoring in all other voter characteristics/information.

For each additional year in age, we expect the odds of an eligible voter voting to be 1.04 ( $\exp(0.040)$ ) times the odds of individuals that are female, married, native born, not employed, have Bachelors Degree as their highest education, not currently enrolled in school if between the age of 16-24, White, and the election is a midterm election following the election of an incumbent president, after factoring in all other voter characteristics/information.

We expect the odds of an eligible voter voting when it is a presidential election year to be 2.689 ( $\exp(0.989)$ ) times the odds of individuals that are female, married, native born, not employed, have Bachelors Degree as their highest education, not currently enrolled in school if between the age of 16-24, White, and the election is a midterm election following the election of an incumbent president, after factoring in all other voter characteristics/information.

## Checking Model Conditions

### Linearity

According to the empirical logit plot (see Appendix C), there is a linear relationship between the empirical logit and the predictor variable age. Hence linearity is satisfied for AGE.

### Randomness

It is possible that randomness is not satisfied because our data is from the census survey, which may not be random (ie might select for people who have time to fill it out). However, there is no reason to believe that this will not generalize to the US population as a whole in a significant way, particularly due to the large sample size.

### Independence

Independence may be violated because geographic location may influence voting due to factors such as (residuals by state ID). Hence, we will look at misclassification rate by region.

Here, we create a confusion matrix with a decision threshold of 0.5.

We will join the augmented data with region in order to look at misclassification rates by region.

Below is a table of the misclassification rates by region.

Table 3: Missclassification Rates by Region

region	voted	predicted_voted	n	prop
Midwest	0	Will Vote	81460	0.191
Midwest	1	Will Not Vote	32314	0.076
Northeast	0	Will Vote	116798	0.195
Northeast	1	Will Not Vote	44736	0.075
South	0	Will Vote	131479	0.191
South	1	Will Not Vote	51284	0.075
West	0	Will Vote	93400	0.192
West	1	Will Not Vote	37109	0.076
NA	0	Will Vote	8040	0.186
NA	1	Will Not Vote	3433	0.080

Consulted Census data for the region fips number corresponding to region name [9]

We then created a plot of misclassification rate by region to determine if independence is satisfied. See the Appendix C for the boxplot visualizing the misclassification rates by region.

Based on the misclassification rates, we have no reason to believe that the independence condition is not satisfied. The misclassification rates across regions are relatively similar, which suggests that there is not an issue of spatial correlation and that the independence condition is satisfied.

## Checking for influential points

### Cook's distance

We will also look for influential points using Cook's Distance. See the plot for Cook's distance in Appendix C.

According to the plot of Cook's Distance (see Appendix C), there are no influential points, so all points can be left in the model.

### Multicollinearity

All of the VIF values are under the threshold of 10 (see table in Appendix C), indicating that there is no evidence of multicollinearity in our data.

## Interaction Terms

We will add in several interaction terms of interest to us and use a drop-in-deviance test to see if they are meaningful predictors of the odds of someone voting.

The reduced model is the same as the above model titled "Model Resulting From Backward Selection," and the full model is the reduced model plus interactions terms for sex and employment, presidential election status and highest education level, age and presidential election status, sex and presidential election status, and race and presidential election status.

The following hypotheses will be used:

$H_0$ : the coefficients for the interaction between sex and employment, presidential election status and highest education level, age and presidential election status, sex and presidential election status, and race and presidential election status are all zero

$H_0$  :

$$\beta_{sex*employed} = \beta_{Presidential-Election-Status*highest\_education} = \beta_{AGE*Presidential-Election-Status} = \beta_{sex*Presidential-Election-Status} = \beta_{race*Presidential-Election-Status} = 0$$

$H_a$ : at least one of these coefficients for the interaction terms  $\neq$  zero

$\alpha = 0.05$

Table 4: Drop-In-Deviance Test Results For Interaction Terms

Resid..Df	Resid..Dev	df	Deviance	p.value
9974	10673.41	NA	NA	NA
9946	10592.36	28	81.054	0

The p-value (7.840e-07) is very small (less than the alpha level 0.05), so we can reject the null hypothesis. Thus, we conclude that the data provide sufficient evidence that the coefficients associated with the additional interaction terms are not equal to 0. Therefore, we should add them to the model.

Significant interaction terms: The effect of the election being a presidential election for an individual with the highest level of education as a High School Degree/GED is significant with a p-value of 0.011 assuming



an alpha level of 0.05. The coefficient is negative, indicating that this effect would mean that during a presidential election for an individual with the highest level of education as a High School Degree/GED, they are less likely to vote than individuals that are female, married, native born, not employed, have Bachelors Degree as their highest education, not currently enrolled in school if between the age of 16-24, White, and the election is a midterm election following the election of an incumbent president.

The effect of the election being a presidential election for an individual that is Black is significant with a p-value of 0.013 assuming an alpha level of 0.05. The coefficient is positive, indicating that this effect would mean that during a presidential election for an individual that is Black, they are more likely to vote than individuals that are female, married, native born, not employed, have Bachelors Degree as their highest education, not currently enrolled in school if between the age of 16-24, White, and the election is a midterm election following the election of an incumbent president

We will use a drop-in-deviance test to determine whether or not sex is a meaningful predictor of the odds of someone voting.

The following hypotheses will be used:

$H_0$ : the coefficients for the main effect for sex, the interaction between sex and employment, and sex and presidential election status are all zero. All of the coefficients associated with sex are equal to zero.

$H_0 : \beta_{sex} = \beta_{sex*employed} = \beta_{sex*Presidential-Election-Status} = 0$   $H_a$ : at least one of these coefficients for the coefficients associated with sex  $\neq$  zero

$\alpha = 0.05$

Table 5: Drop-In-Deviance Test Results For Sex

Resid..Df	Resid..Dev	df	Deviance	p.value
9949	10596.31	NA	NA	NA
9946	10592.36	3	3.945	0.267

From the drop-in-deviance test to include variable sex, the p-value (0.267) is greater than an alpha-level of 0.10, so we will exclude the main effect for sex and the interaction terms including sex from the model. This contradicts the model output from the backward selection, which is most likely due their different criteria for statistical significance. Backward selection makes decisions based on the AIC, while the drop-in-deviance test uses the p-value. A potential discrepancy between the two values may have resulted in the different determinations of the significance of sex. We have chosen to use the results of the drop-in-deviance test, and therefore, will be excluding sex from the final model.

## Results

Table 6: Final Model Output

term	estimate	std.error	statistic	p.value
(Intercept)	14.581	11.394	1.280	0.201
YEAR	-	0.006	-1.371	0.170
	0.008			
current_studentHigh School Full Time	1.120	0.307	3.643	0.000
current_studentHigh School Part Time	0.388	1.139	0.340	0.734
current_studentCollege Full Time	0.461	0.124	3.710	0.000
current_studentCollege Part Time	0.660	0.251	2.633	0.008
citizenNaturalized	-	0.103	-6.672	0.000
	0.690			
employedYes	0.397	0.060	6.621	0.000

term	estimate	std.error	statistic	p.value
raceBlack	0.876	0.132	6.613	0.000
raceAsian or Pacific Islander	-	0.190	-3.801	0.000
	0.722			
raceNative American	-	0.284	-1.134	0.257
	0.322			
race2 or more races	-	0.278	-0.687	0.492
	0.191			
marstDivorced/Separated	-	0.071	-9.996	0.000
	0.709			
marstNot Married/Other	-	0.057	-8.471	0.000
	0.483			
Is_Presidential_ElectionPresidential	1.522	0.213	7.162	0.000
AGE	0.033	0.002	14.700	0.000
highest_educationHigh School Degree/GED	-	0.117	-	0.000
	1.590		13.536	
highest_educationSome College	-	0.129	-6.511	0.000
	0.838			
highest_educationSome High School	-	0.144	-	0.000
	2.449		17.031	
highest_educationAssociate Degree	-	0.153	-3.938	0.000
	0.604			
highest_educationMasters Degree	0.648	0.237	2.738	0.006
highest_educationProfessional Degree	1.228	0.613	2.003	0.045
highest_educationDoctorate Degree	0.104	0.458	0.226	0.821
highest_educationNone/Unknown	-	0.806	-4.133	0.000
	3.330			
Presidential_Election_StatusMidterm Following Incumbent President	-	0.231	-1.412	0.158
	0.327			
Presidential_Election_StatusMidterm Following New President	NA	NA	NA	NA
highest_educationHigh School	0.455	0.178	2.555	0.011
Degree/GED:Presidential_Election_StatusMidterm Following Incumbent President				
highest_educationSome College:Presidential_Election_StatusMidterm Following Incumbent President	0.296	0.193	1.536	0.125
highest_educationSome High School:Presidential_Election_StatusMidterm Following Incumbent President	0.468	0.235	1.994	0.046
highest_educationAssociate Degree:Presidential_Election_StatusMidterm Following Incumbent President	0.091	0.229	0.399	0.690
highest_educationMasters Degree:Presidential_Election_StatusMidterm Following Incumbent President	-	0.315	-1.960	0.050
highest_educationProfessional Degree:Presidential_Election_StatusMidterm Following Incumbent President	0.618			
highest_educationDoctorate Degree:Presidential_Election_StatusMidterm Following Incumbent President	-	0.752	-1.441	0.149
highest_educationNone/Unknown:Presidential_Election_StatusMidterm Following Incumbent President	1.084			
	1.012	0.722	1.402	0.161
	-	109.766	-0.082	0.934
	9.035			

term	estimate	std.error	statistic	p.value
highest_educationHigh School Degree/GED:Presidential_Election_StatusMidterm Following New President	0.315	0.182	1.730	0.084
highest_educationSome College:Presidential_Election_StatusMidterm Following New President	0.382	0.198	1.932	0.053
highest_educationSome High School:Presidential_Election_StatusMidterm Following New President	0.334	0.244	1.370	0.171
highest_educationAssociate Degree:Presidential_Election_StatusMidterm Following New President	0.173	0.238	0.728	0.467
highest_educationMasters Degree:Presidential_Election_StatusMidterm Following New President	- 0.382	0.327	-1.168	0.243
highest_educationProfessional Degree:Presidential_Election_StatusMidterm Following New President	- 0.982	0.800	-1.226	0.220
highest_educationDoctorate Degree:Presidential_Election_StatusMidterm Following New President	0.019	0.633	0.030	0.976
highest_educationNone/Unknown:Presidential_Election_StatusMidterm Following New President	1.185	1.435	0.826	0.409
AGE:Presidential_Election_StatusMidterm Following Incumbent President	0.013	0.003	3.859	0.000
AGE:Presidential_Election_StatusMidterm Following New President	0.013	0.004	3.632	0.000
raceBlack:Presidential_Election_StatusMidterm Following Incumbent President	- 0.504	0.204	-2.473	0.013
raceAsian or Pacific Islander:Presidential_Election_StatusMidterm Following Incumbent President	0.337	0.302	1.116	0.264
raceNative American:Presidential_Election_StatusMidterm Following Incumbent President	0.274	0.496	0.553	0.580
race2 or more races:Presidential_Election_StatusMidterm Following Incumbent President	0.767	0.506	1.518	0.129
raceBlack:Presidential_Election_StatusMidterm Following New President	- 0.426	0.203	-2.101	0.036
raceAsian or Pacific Islander:Presidential_Election_StatusMidterm Following New President	0.890	0.301	2.954	0.003
raceNative American:Presidential_Election_StatusMidterm Following New President	0.705	0.492	1.434	0.152
race2 or more races:Presidential_Election_StatusMidterm Following New President	1.119	0.513	2.180	0.029

## Checking Model Conditions for Final Model

### Checking for influential points

#### Cook's distance

We will also look for influential points in our final model using Cook's Distance.

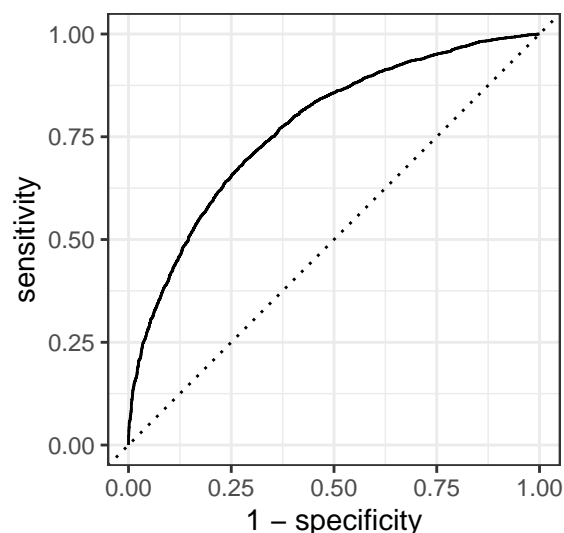
According to the Cook's Distance plot for the final model (see Appendix D), there are no influential points, so all points can be left in the final model.

## Multicollinearity

All of the VIF values are under the threshold of 10 (see Appendix D), indicating that there is no evidence of multicollinearity in our final model.

## Creating Classifier

We'll fit an ROC curve to help us determine a decision-making threshold.



Below we look at values within a range of thresholds in order to choose the threshold that yields high sensitivity and low values of (1 - specificity).

Table 7: ROC Curve Threshold Table

.threshold	specificity	sensitivity	pred_prob
0.419	0.338	0.925	0.603
0.419	0.339	0.925	0.603
0.420	0.339	0.925	0.603
0.420	0.339	0.925	0.603
0.420	0.339	0.925	0.603

According to our ROC curve and modeling objectives, we will choose a threshold of .42 because we want to lean towards having a higher false negative rate (type II error) than a false positive rate as it does not hurt to mail a few extra ballots to people who may already be planning to vote.

Any data point with a probability over 0.60 will be predicted to be in the “voted” category.

## Discussion

While completing this project provided excellent insight into the voting patterns in the United States, there are several limitations of this project that must be acknowledged. First, the original dataset hard-coded several variables with “unknown” and “missing” combined into one category. This could have obscured the effects of missingness in our data and contributed to decreasing the predictive accuracy of our model and results. For this reason, if given the opportunity to expand this project, we would like to explore missingness in our dataset and its implications with more depth to improve our model. Because several categorical variables already had levels that combined missing and unknown, we decided to code the other variables in the same way, effectively imputing a new category for missing that was used in the model. While this was the

best way to deal with missingness within the scope of this project, we would be interested in expanding this in the future and teasing out missing from unknown to create a more accurate model and improve our results.

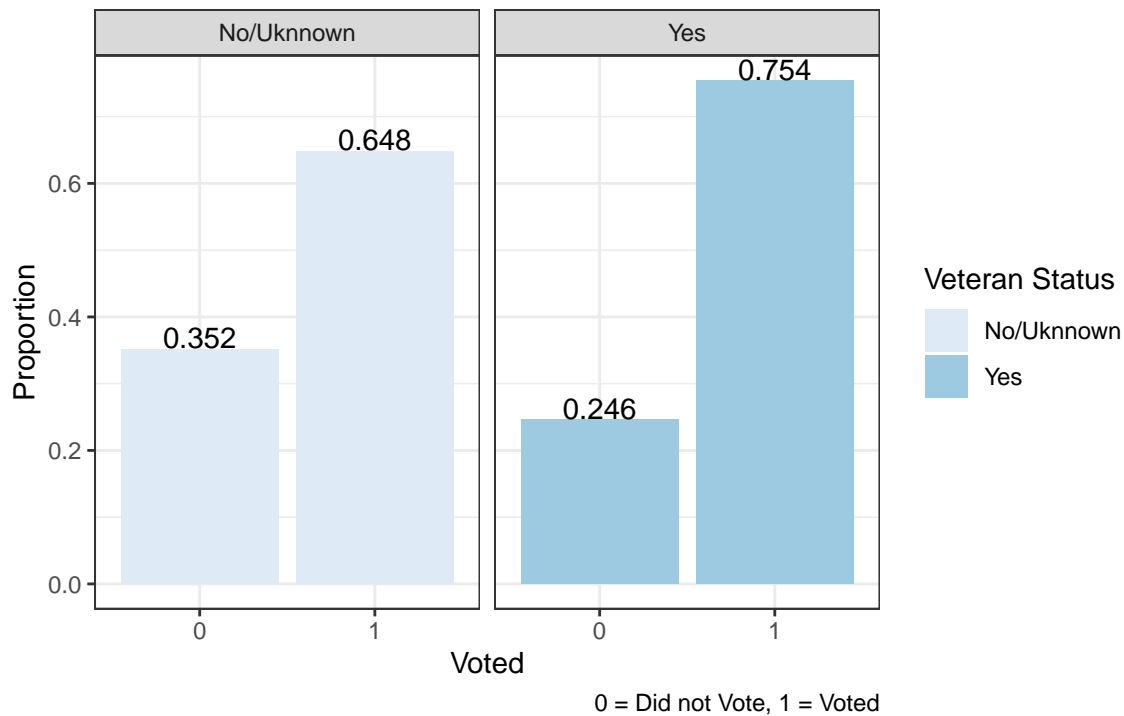
## References

- [1] “Fall Data Challenge | This Is Statistics.” n.d. Accessed November 16, 2020. <https://thisisstatistics.org/falldatachallenge/>.
- [2] University of Minnesota. “What Is IPUMS?” Text. IPUMS. February 7, 2019. <https://ipums.org/what-is-ipums>.
- [3] MIT Election Data + Science Lab. n.d. “Voter Turnout.” Accessed November 16, 2020. <https://plotly.com/~cwimpy/69/>.
- [4] Fowler, Anthony George. 2013. “Five Studies on the Causes and Consequences of Voter Turnout,” October. <https://dash.harvard.edu/handle/1/11156810>.
- [5] McDonald, Michael. n.d. “Voter Turnout Demographics - United States Elections Project.” Accessed November 16, 2020. <http://www.electproject.org/home/voter-turnout/demographics>.
- [6] “GGPlot Cheat Sheet for Great Customization - Articles - STHDA.” <http://www.sthda.com/english/articles/32-r-graphics-essentials/125-ggplot-cheat-sheet-for-great-customization/> (November 15, 2020).
- [7] “The A - Z Of Rcolorbrewer Palette You Must Know.” 2018. Datanovia. <https://www.datanovia.com/en/blog/the-a-z-of-rcolorbrewer-palette/> (November 15, 2020).
- [8] “How to Put Labels over Geom\_bar for Each Bar in R with Ggplot2.” Stack Overflow. <https://stackoverflow.com/questions/12018499/how-to-put-labels-over-geom-bar-for-each-bar-in-r-with-ggplot2> (November 15, 2020).
- [9] Prepared by the Geography Division of the U.S. Census Bureau. “Census Regions and Divisions of the United States.” [https://www2.census.gov/geo/pdfs/maps-data/maps/reference/us\\_regdiv.pdf](https://www2.census.gov/geo/pdfs/maps-data/maps/reference/us_regdiv.pdf).

## Appendix A

### Voting distribution based on veteran status

Examining relationship between veteran status and voting



*see geom\_text() code inspiration in reference [4]*

## Appendix B

### Backward Selection Output

```
## Start:  AIC=10729.96
## voted ~ metro + sex + marst + veteran + citizen + hispanic_status +
##         employed + highest_education + current_student + race + AGE +
##         Presidential_Election_Status + YEAR
##
##              Df Deviance   AIC
## - metro          1    10672 10728
## - veteran         1    10672 10728
## - hispanic_status 1    10673 10729
## - YEAR            1    10674 10730
## <none>            1    10672 10730
## - sex            1    10675 10731
## - current_student 4    10700 10750
## - citizen         1    10709 10765
## - employed        1    10718 10774
## - race            4    10735 10785
## - marst           2    10801 10855
## - Presidential_Election_Status 2    11008 11062
## - AGE            1    11197 11253
## - highest_education 8    11547 11589
```

```

##
## Step: AIC=10727.96
## voted ~ sex + marst + veteran + citizen + hispanic_status + employed +
##     highest_education + current_student + race + AGE + Presidential_Election_Status +
##     YEAR
##
##           Df Deviance   AIC
## - veteran           1    10672 10726
## - hispanic_status    1    10673 10727
## - YEAR               1    10674 10728
## <none>              10672 10728
## - sex               1    10675 10729
## - current_student    4    10700 10748
## - citizen            1    10709 10763
## - employed           1    10718 10772
## - race               4    10736 10784
## - marst              2    10803 10855
## - Presidential_Election_Status 2    11008 11060
## - AGE                1    11197 11251
## - highest_education   8    11554 11594
##
## Step: AIC=10726.03
## voted ~ sex + marst + citizen + hispanic_status + employed +
##     highest_education + current_student + race + AGE + Presidential_Election_Status +
##     YEAR
##
##           Df Deviance   AIC
## - hispanic_status    1    10673 10725
## - YEAR               1    10674 10726
## <none>              10672 10726
## - sex               1    10675 10727
## - current_student    4    10700 10746
## - citizen            1    10709 10761
## - employed           1    10718 10770
## - race               4    10737 10783
## - marst              2    10803 10853
## - Presidential_Election_Status 2    11008 11058
## - AGE                1    11219 11271
## - highest_education   8    11556 11594
##
## Step: AIC=10725.41
## voted ~ sex + marst + citizen + employed + highest_education +
##     current_student + race + AGE + Presidential_Election_Status +
##     YEAR
##
##           Df Deviance   AIC
## <none>              10673 10725
## - YEAR               1    10676 10726
## - sex               1    10676 10726
## - current_student    4    10702 10746
## - citizen            1    10717 10767
## - employed           1    10719 10769
## - race               4    10739 10783
## - marst              2    10805 10853

```



```
## - Presidential_Election_Status 2    11009 11057
## - AGE                          1    11233 11283
## - highest_education            8    11580 11616
```

## Appendix C

### Checking Model Conditions for the resulting model from the selection process

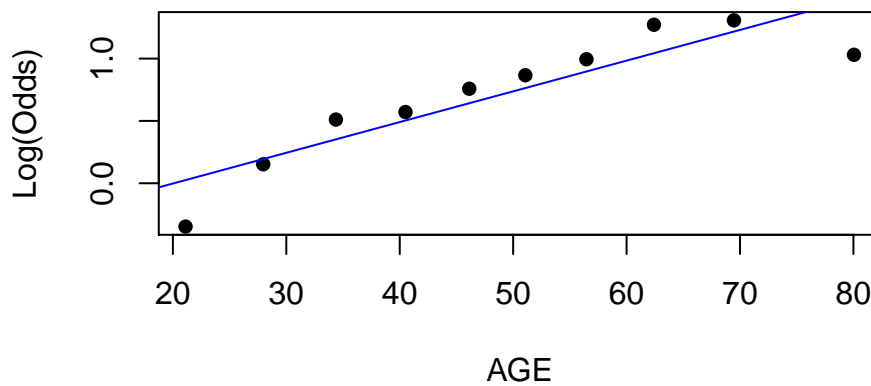
#### Linearity

Table 8: Empirical Logit Results For Sex

sex	voted	n	prop	emp_logit
Female	1	3582	0.67	0.710
Male	1	3029	0.65	0.621

Table 9: Empirical Logit Results For Marital Status

marst	voted	n	prop	emp_logit
Married	1	4145	0.735	1.018
Divorced/Separated	1	783	0.588	0.355
Not Married/Other	1	1683	0.556	0.226



#### Independence

##### Misclassification rate by region

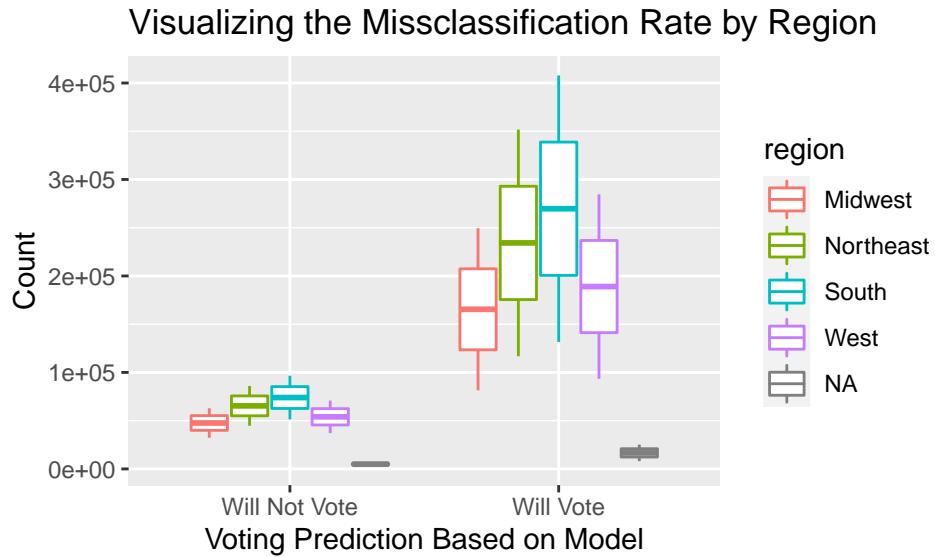
Table 10: Missclassification Rates by Region

region	voted	predicted_voted	n	prop
Midwest	0	Will Vote	81460	0.191
Midwest	1	Will Not Vote	32314	0.076
Northeast	0	Will Vote	116798	0.195
Northeast	1	Will Not Vote	44736	0.075
South	0	Will Vote	131479	0.191
South	1	Will Not Vote	51284	0.075
West	0	Will Vote	93400	0.192
West	1	Will Not Vote	37109	0.076

region	voted	predicted_voted	n	prop
NA	0	Will Vote	8040	0.186
NA	1	Will Not Vote	3433	0.080

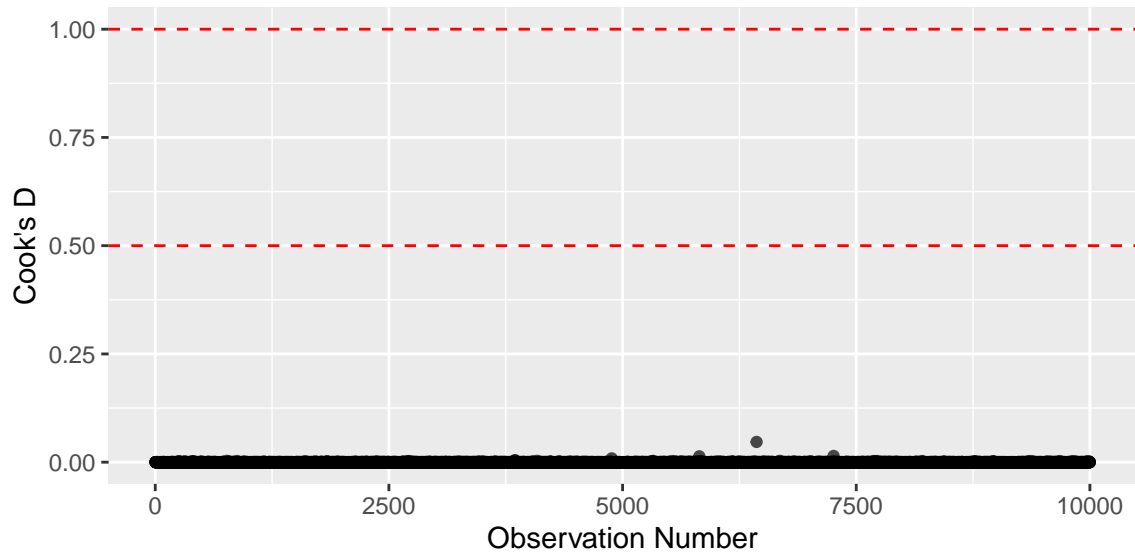
*Consulted Census data for the region fips number corresponding to region name [3]*

#### Boxplot of misclassification rate by region



#### Checking for influential points

##### Cook's distance plot



#### Multicollinearity: VIF Table

names	x
current_studentHigh School Part Time	1.006
raceNative American	1.009

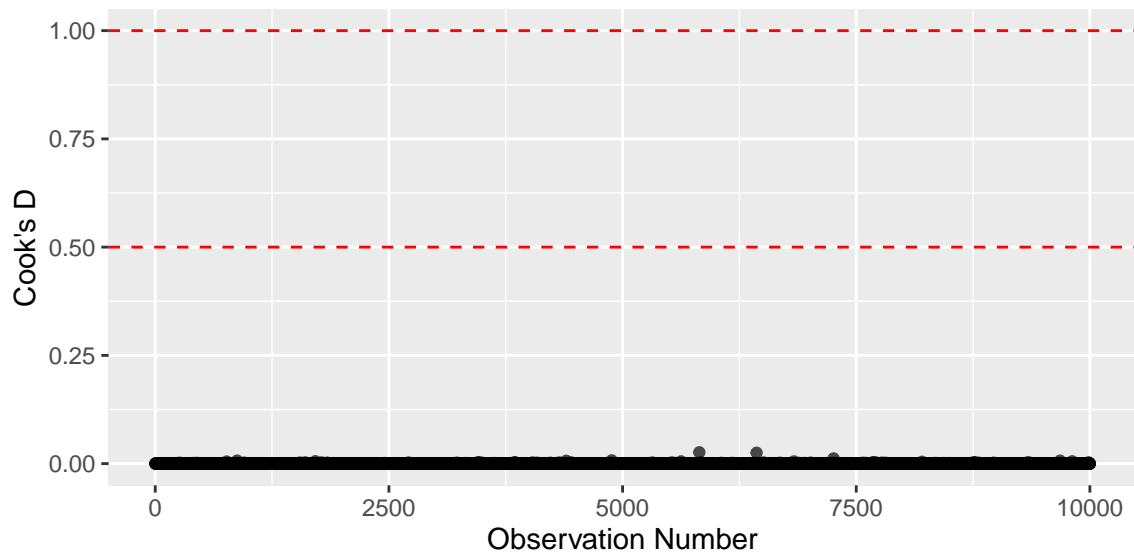
names	x
race2 or more races	1.009
highest_educationNone/Unknown	1.016
sexMale	1.023
current_studentCollege Part Time	1.038
highest_educationProfessional Degree	1.043
raceBlack	1.049
highest_educationDoctorate Degree	1.050
current_studentHigh School Full Time	1.086
marstDivorced/Separated	1.116
YEAR	1.168
Presidential_Election_StatusMidterm Following Incumbent President	1.176
raceAsian or Pacific Islander	1.206
citizenNaturalized	1.214
highest_educationMasters Degree	1.246
marstNot Married/Other	1.286
current_studentCollege Full Time	1.303
Presidential_Election_StatusMidterm Following New President	1.311
employedYes	1.433
highest_educationAssociate Degree	1.508
AGE	1.709
highest_educationSome High School	1.744
highest_educationSome College	1.950
highest_educationHigh School Degree/GED	2.216

## Appendix D

### Checking Model Conditions for the Final Model

#### Checking for influential points

##### Cook's distance plot for final model



##### Multicollinearity: VIF Table for final model

names	x
YEAR	NA
current_studentHigh School Full Time	NA
current_studentHigh School Part Time	NA
current_studentCollege Full Time	NA
current_studentCollege Part Time	NA
citizenNaturalized	NA
employedYes	NA
raceBlack	NA
raceAsian or Pacific Islander	NA
raceNative American	NA
race2 or more races	NA
marstDivorced/Separated	NA
marstNot Married/Other	NA
Is_Presidential_ElectionPresidential	NA
AGE	NA
highest_educationHigh School Degree/GED	NA
highest_educationSome College	NA
highest_educationSome High School	NA
highest_educationAssociate Degree	NA
highest_educationMasters Degree	NA
highest_educationProfessional Degree	NA
highest_educationDoctorate Degree	NA
highest_educationNone/Unknown	NA
Presidential_Election_StatusMidterm Following Incumbent President	NA
Presidential_Election_StatusMidterm Following New President	NA
highest_educationHigh School Degree/GED:Presidential_Election_StatusMidterm Following Incumbent President	NA
highest_educationSome College:Presidential_Election_StatusMidterm Following Incumbent President	NA
highest_educationSome High School:Presidential_Election_StatusMidterm Following Incumbent President	NA
highest_educationAssociate Degree:Presidential_Election_StatusMidterm Following Incumbent President	NA
highest_educationMasters Degree:Presidential_Election_StatusMidterm Following Incumbent President	NA
highest_educationProfessional Degree:Presidential_Election_StatusMidterm Following Incumbent President	NA
highest_educationDoctorate Degree:Presidential_Election_StatusMidterm Following Incumbent President	NA
highest_educationNone/Unknown:Presidential_Election_StatusMidterm Following Incumbent President	NA
highest_educationHigh School Degree/GED:Presidential_Election_StatusMidterm Following New President	NA
highest_educationSome College:Presidential_Election_StatusMidterm Following New President	NA
highest_educationSome High School:Presidential_Election_StatusMidterm Following New President	NA
highest_educationAssociate Degree:Presidential_Election_StatusMidterm Following New President	NA
highest_educationMasters Degree:Presidential_Election_StatusMidterm Following New President	NA
highest_educationProfessional Degree:Presidential_Election_StatusMidterm Following New President	NA
highest_educationDoctorate Degree:Presidential_Election_StatusMidterm Following New President	NA
highest_educationNone/Unknown:Presidential_Election_StatusMidterm Following New President	NA

names	x
AGE:Presidential_Election_StatusMidterm Following Incumbent President	NA
AGE:Presidential_Election_StatusMidterm Following New President	NA
raceBlack:Presidential_Election_StatusMidterm Following Incumbent President	NA
raceAsian or Pacific Islander:Presidential_Election_StatusMidterm Following Incumbent President	NA
raceNative American:Presidential_Election_StatusMidterm Following Incumbent President	NA
race2 or more races:Presidential_Election_StatusMidterm Following Incumbent President	NA
raceBlack:Presidential_Election_StatusMidterm Following New President	NA
raceAsian or Pacific Islander:Presidential_Election_StatusMidterm Following New President	NA
raceNative American:Presidential_Election_StatusMidterm Following New President	NA
race2 or more races:Presidential_Election_StatusMidterm Following New President	NA