# Data 102 Final Project: Predicting and Evaluating Election Outcomes

Julianna Lee, Minoo Kim, Joshua Yoo, Alyssa Mar

December 2024

# Contents

# 1 Introduction

## 1.1 Data Overview

We decided to combine two data sets together. The first dataset we used is the FiveThirtyEight 2018 Primary Candidate Endorsements and we combined it with the Federal Election Commission 2018 Campaign Financing Data. We decided to use an additional financial dataset because the Primary Candidate Endorsement dataset was not sufficient enough to be used for our predictions. We would essentially only be predicting on binary variables but with the addition of Total Contribution for each candidate, we would better be able to predict whether they won the primary election or not.

Our data was collected in an observational manner. The Primary Candidate Endorsements data consists of census data that shows candidates who have appeared on the ballot in 2018 in the Republican and Democratic primaries. The Federal Election Commission data also consists of census data that shows financial campaign data for each candidate.

The Federal Election Commission data doesn't exclude any groups from the data but the Primary Candidate Endorsements data only includes data from Democratic and Republican candidates. All other parties are excluded from the Primary Candidate Endorsements dataset.

Each row corresponds to an individual candidate who appeared on the 2018 ballot for Senate, House, or governor, excluding races where a Democratic or Republican incumbent was running. This candidate-specific granularity allows us to interpret our findings at the level of individual candidates. There is no selection bias, measurement error, or convenience sampling concerns regarding the context of this data. Since one dataset only includes Democratic and Republican candidates, we removed the other parties from our analysis and conclusions.

For the Primary Candidate Endorsements dataset, the Republicans don't have any Partisan Lean, only the Democratic candidates did. We assume that the missing entry meant that the partisan lean of the district just wasn't collected for the Republican candidates. For our first research question about trying to predict whether a candidate won the primary or not, we filled in the missing partisan lean values with 0. For our second research question, we removed the Republican rows and used performed the causal inference for only Democratic candidates.

In order to combine the two datasets, we had to clean the candidate name columns to be equal to each other. We used regex and string replacement to get the name columns to be in the format (First Name Last Name). This allowed us to see the financial campaign data for each candidate in addition to their endorsement data. Additionally, we dropped columns that we didn't need, and kept only the ones we used for prediction and causal inference. The last step of preprocessing that we did was converting all the categorical variables into binary variables. For instance, for the column for the candidate's political party, we attributed 0 to be Democrat and 1 to Republican. By converting all the columns to numerical data, we are able to better use the data for our decision tree and causal inferences.

## 1.2 Research Question 1

**Question:** How can we predict winning primary elections from campaign financing and primary election data?

The real world decisions that could be made using the results would be predicting a potential outcome of someone's election campaign with their current features. This would be a helpful gauge as to how the candidate is currently doing and if their current funding/race is projected to turn into a favorable outcome.

The methods we chose are a good fit for this research question because they help us predict an outcome variable based on several predictors. Specifically, we'll use logistic regression as the GLM since it's designed for binary outcomes, and a decision tree as the nonparametric method. Decision trees are especially useful for capturing more complex, nonlinear relationships between variables, which logistic regression might miss because of its linearity assumptions. For our features, we'll use the primary percentage and party support columns (both expressed as percentages), along with financial data on Total Contribution for each candidate. By combining these datasets, we'll focus our predictions specifically on election outcomes for the year 2018.

One limitation of the decision tree is that it is prone to overfitting the data. If we add too many features, it may not do well because it will reduce the ability to generalize our analysis. One limitation of the logistic regression model is that it assumes a linear relationship between the predictors and the log-odds of the

outcome. If the actual relationships are highly nonlinear, this method may under perform and not fit the data well.

## 1.3 Research Question 2

**Question:** Does endorsement status (whether or not the candidate is endorsed) cause an increase in the primary percentage of the candidate?

This could apply to real life situations if campaign managers are seeing if it is worth being endorsed. By better understanding the causal link between endorsements and primary vote outcomes, it provides actionable insights for candidates, organizations, and voters in optimizing electoral strategies.

Causal inference is a good fit for this research question because we are trying to see whether or not endorsement status directly affects primary percentage. The technique we plan to use is Inverse Propensity Weighting which accounts for confounders. By using the endorsement status as our target variable and the confounders as our predictors, we can weight the propensity scores of the two groups to estimate the ATE.

Inverse Propensity Weighting requires the unconfoundedness assumption to be met. If the model does not include all relevant confounders, it could lead to biased estimates.

## 1.4 Prior Work

1. Norpoth, Helmut. *Post-Mortem of the 2024 Forecast.* Primary Model, PrimaryVote.com

In *Post-Mortem of the 2024 Forecast,* Helmut Norporth tackles why and how Donald Trump won the 2024 presidential elections despite the model favoring Kamala Harris by substantial margins. The paper explores several factors why the model failed to predict correctly. This question is relevant to our research because even critically acclaimed models can fall short in real-world scenarios. Factors such as assumptions of stability, voter sentiment, campaigning strategies, and potential overfitting were deemed as causes for the incorrect prediction. These factors evidence the world behind what data can capture. Historical election patterns are reliable predictors, but unforeseen political powers may have reversed expectations. There may have been shifts in last-minute voter sentiments due to current events. The model that Primary Vote used did not account for campaigning strategies. These causes may have led to potential overfitting. Similarly, our model may face challenges in capturing the nuanced and unpredictable elements of the political climate because of our limited feature selections. Although we are predicting primary candidates, there are similarities between the features and methods used. Primary Model heavily relies on polling averages, whereas our model relies on various features such as endorsement, Total Contribution, and partisan lean. Different features limit both models when applied to the unforgiving and unpredictable nature of the world and people.

2. Stening, Tanner. *Do Political Endorsements Still Matter?* Northeastern Global News, 26 May 2022, northeastern.edu.

In "Do Political Endorsements Still Matter," Tanner Stening examines how political endorsements positively and/or negatively impact the candidate. This is relevant to our research question because it asks a very similar question of endorsement impact on candidate races. In his article, Tanner uses Donald Trump's endorsement as an example. According to Costas Panagopoulos, endorsements can make a difference in political races, however, there are more fundamental variables that determine a winner. Endorsements can have both positive and negative effects. Trump's endorsements, for example, have been well-received by his loyal followers, but not as influential among those outside his core base. Therefore, according to Tanner, the effects of endorsements vary between people and, more importantly, timing. Taking all this into account, it will be interesting to see how endorsements play a role in our model within our limited dataset.
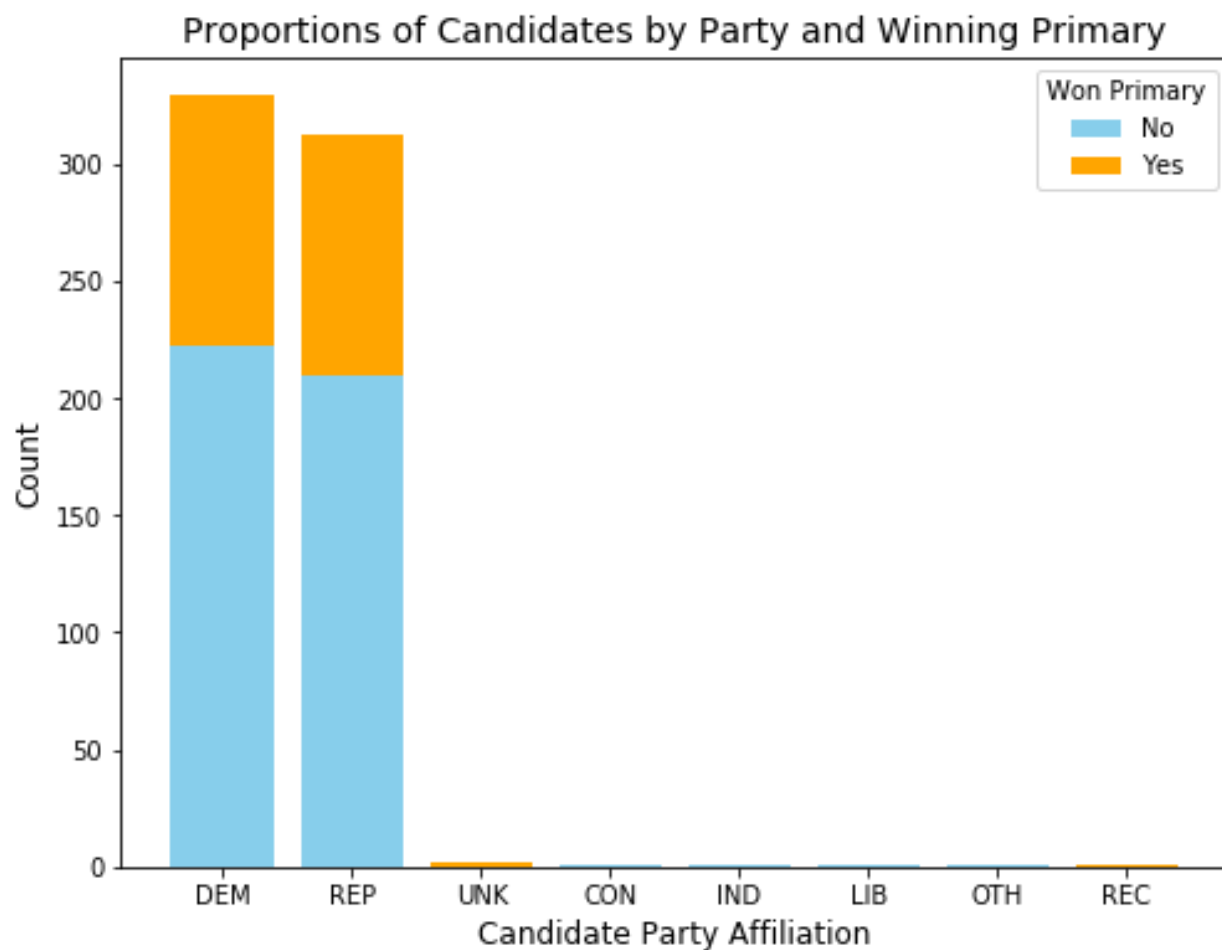
## 1.5 EDA

### 1.5.1 Figure 1



Figure 1: Proportion of Candidates by Party

There seems to be a similar number of Democratic and Republican candidates in the dataset to equally represent each party. The proportion of candidates winning a primary or not also seems to be approximately even between these two parties. This could suggest that party affiliation does not play a huge role in winning a primary, but rather, has more to do with other factors of each candidate. Therefore, this relates to our first research question regarding party affiliation and predicting whether or not a candidate won a primary.
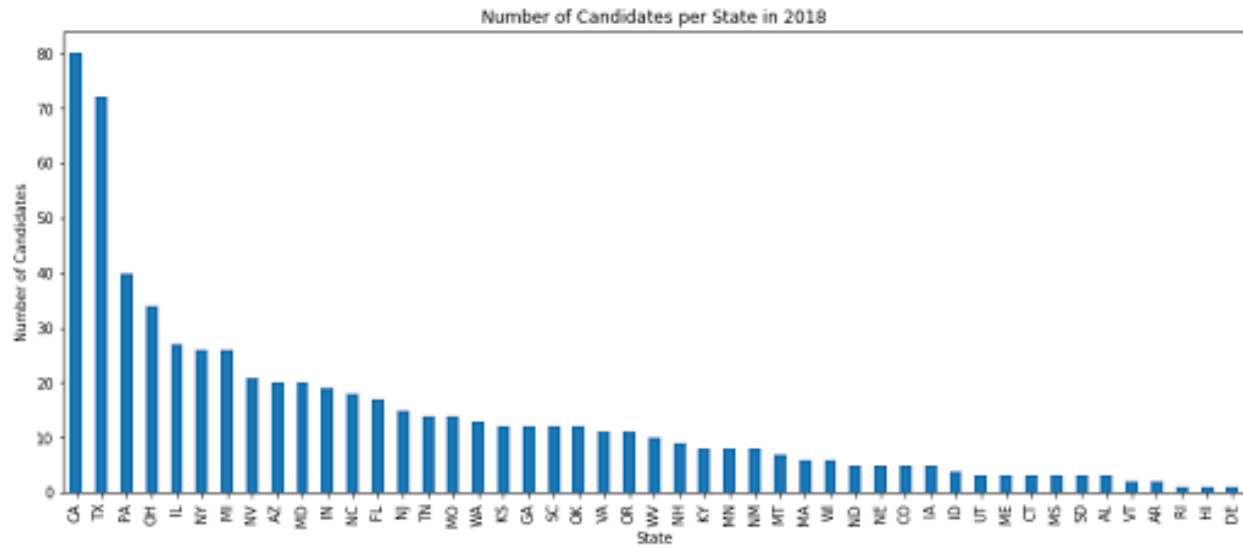
### 1.5.2 Figure 2



Figure 2: Number of Candidates per State

One thing that we noticed is that the state with the highest number of candidates is California, followed by Texas. The states with higher populations tended to have a higher number of candidates. Also, not all states are represented on this bar graph, 3 states do not have candidates. This is relevant to our second research question: 'Does endorsement status (whether or not the candidate is endorsed) cause an increase in the district's lean? ' For states with more candidates, like California and Texas, this question may be harder to answer because there may be more confounding variables across all the candidates.
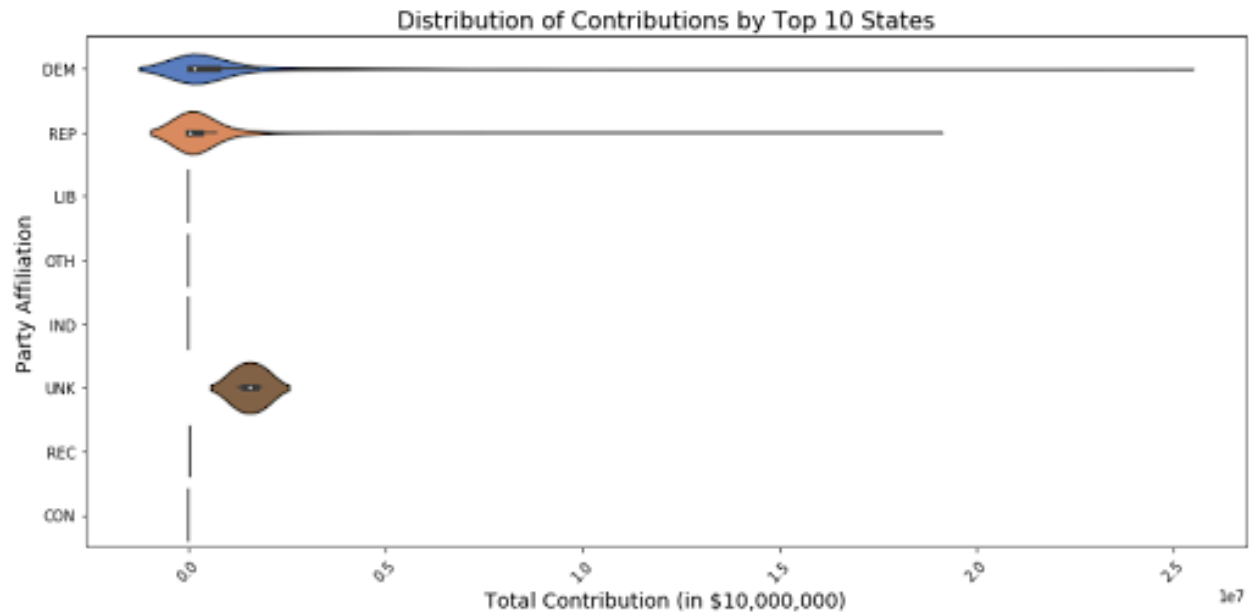
### 1.5.3 Figure 3



Figure 3: Distribution of Contributions

Seeing the distribution of the contributions for each party suggests that the Democratic and Republican parties have more candidates that are receiving contributions compared to the other parties. This is also shown by the wide range of values in the violin plots for the two parties, whereas the other parties only have a distribution of zero. This relates to both research questions where we can see how the contributions are distributed across both party affiliation and state. This helps show the nuance in how contributions are divided between different candidates; we can take this into account when trying to predict whether they won the primary because we can see from the plot that the Democratic and Republican parties are the main parties that get contributions in general.
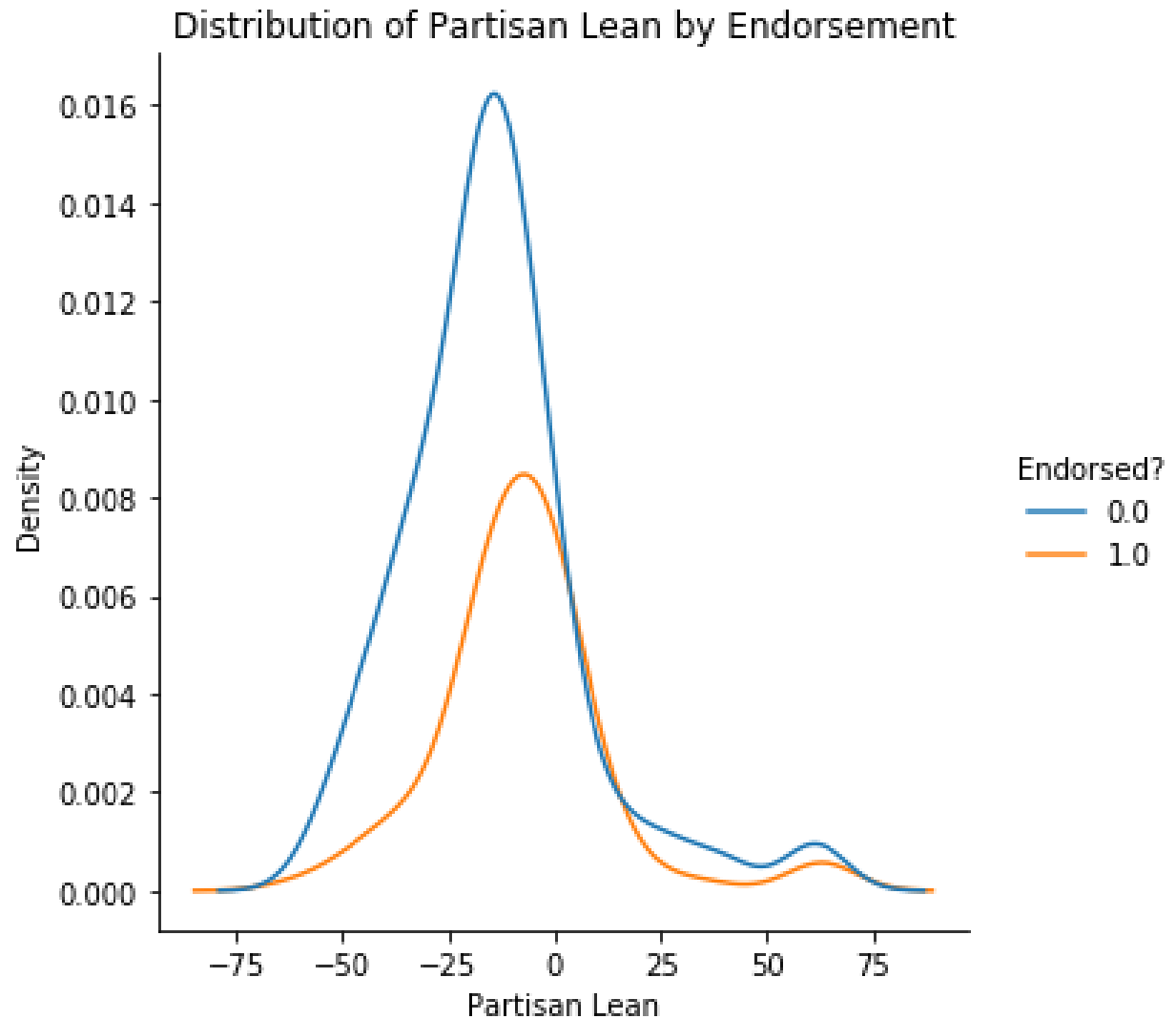
Figure 4: Distribution of Partisan Lean

The plot shows that being endorsed and not endorsed both have a similar mean below zero. However, we can still see how not being endorsed has a higher concentration around that mean. We also notice a slight increase in the density of partisan lean between values of 50 and 75. This relates to our second research question because this graph examines the two variables we are attempting to find a causal relationship between. From this graph, we can initially hypothesize that there may not be a causal relationship between endorsement and partisan lean.
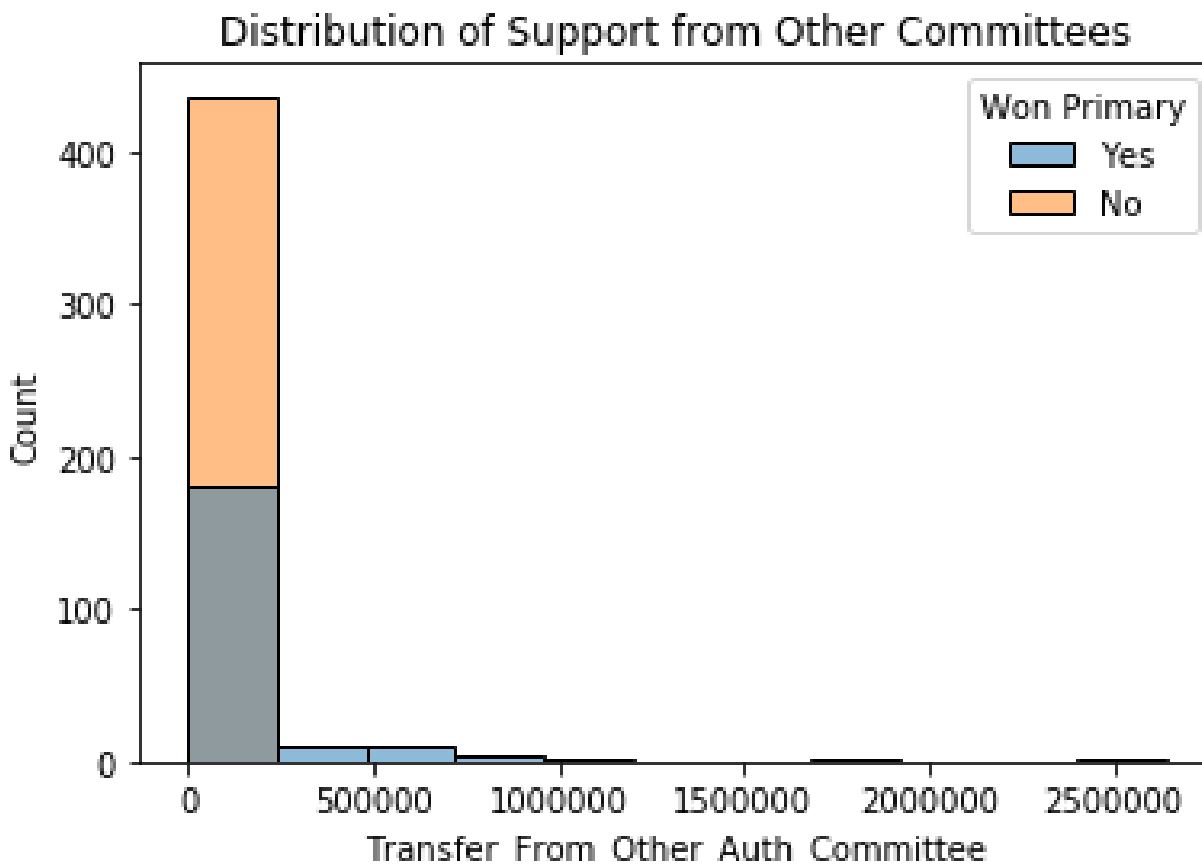
Figure 5: Distribution of Support

We can see that most people who did not win the primary election had no transfers from other committees, whereas the few people who did get transfers seem to have won the primary. This relates to the first research question in trying to predict whether they win the primary or not. It is still difficult to determine whether or not the candidate won because a reasonable number of people who won the primary did not get the transfer.

# 2  QUESTION 1: Prediction with GLMs and Nonparametric Methods

**Question:** How can we predict winning primary elections from campaign financing and primary election data?

## 2.1  Methods

We are trying to predict whether a candidate won the 2018 primary election based on state, race type, primary status, partisan lean, race, total contribution, candidate party affiliation, and endorsement status.

We combined the candidate information dataset with a financial campaign dataset to provide further insight into what might factor into a candidate's election outcome. We chose to use financial campaigning

data because candidates who spend more on their campaigning may be more well-known in the community because they have more resources to do so.

We are using a frequentist GLM with a logistic link function because we are trying to predict a binary outcome. We picked a frequentist GLM over a Bayesian GLM because we did not have any indicative prior information on how our features correlate to winning the primary. Additionally, forming a prior over all our features would be too complex, and would not share any substantial information.

We assumed that each candidate's information in the dataset is independent from other candidates and the predictors we used are linearly independent from each other. Additionally, we assumed that all the endorsement indicators in the dataset were represented by the *Total Contribution* column, of which caused us to remove the organizational endorsement columns. Endorsements were represented in Total Contribution, since any contribution greater than 0 would indicate an endorsement.

We are using a decision tree because it helps us better interpret the way the model is making decisions, showing us what splits the tree is making based on the features.

Since a lot of the columns in the data frame had categorical variables, in order to use a decision tree, we had to convert these categorical variables to numbers. For instance, for the 'Cand_Party_Affiliation' column, we assigned 0 to Democrats and 1 to Republicans. Additionally, for columns with more categorical variables like 'State', we chose to approach it a bit differently and didn't assign 1-50 to each state because it could possibly overfit the data. Instead, we assigned 1 to small states (states with less than 10 candidates), 2 to medium states (having greater than 10 but less than 20 candidates), and 3 to large states (having greater than or equal to 20 candidates). Additionally, there were several columns with 'NaN's so we filled those with 0, where we deemed the value 0 to not contribute information that might indicate someone winning a primary election. One example of this can be seen through the 'Partisan Lean' column. 0 in that column represented a neutral lean or no lean at all so we assumed that for the candidates that had 'NaN' for that column.

Because we are classifying a binary outcome, decision trees would be useful for splitting points to create homogenous groups of winning or losing a primary election at the leaves. The election dataset is relatively small and we wanted to be able to interpret the results and visualize where the splits were occurring. We assumed that the data was not corrupted by noise, meaning that training a random forest was not necessary.

We evaluated the performance of the GLM with the Chi-Squared Statistic given by statsmodels (goodness of fit measure). This value is the average ratio of your model's discrepancy and the outcomes. We would expect this to be around or better than 435-7 (n = 435 points, d = 7 features).

We evaluated the performance of the decision tree with the root mean squared error rate. We will also look at the accuracy, precision and recall of the decision tree to understand how it classifies true and false values.

## 2.2 Results

$$\text{Accuracy} = \frac{TP + TN}{TP + TN + FP + FN} \tag{1}$$

$$\text{Precision} = \frac{TP}{TP + FP} \tag{2}$$

$$\text{Recall} = \frac{TP}{TP + FN} \tag{3}$$

Looking at our decision tree, the accuracy (1) on the test set was .78. For those that did not win the primary (predicted 0), the model has a precision (2) of .85 and a recall (3) of .82. For those that did win primary (predicted 1), the model has a precision of .67 and recall of .71. Our model didn't seem to suffer from overfitting since our training RMSE error was .15 while our test error was higher at .46. The precision (2) of the model shows us the ratio of correctly predicted positives, while the recall (3) shows us the ratio of true positives that were correctly predicted. These statistics are apt for our model, since it is more necessary to rightly predict winning an election over not predicting a positive outcome.

When fitting the model into the GLM with a binomial family, we get a log-likelihood of -28.375 an r-squared of 0.6892, and Chi Squared value of 179. A key predictor of the model is Primary % and Total

Contribution. Primary % has a p-value less than 0.001 and a positive coefficient of 0.1767. Total Contribution also has a low p-value of 0.003 and a positive but small coefficient.

The confidence intervals of State and Total Contribution are narrow and provide insights into their estimated effects. The confidence intervals of State is $[-0.164, 0.450]$; the confidence interval for Total Contribution is $[7.87 \times 10^{-7}, 1.55 \times 10^{-6}]$. They have relatively narrow confidence intervals meaning that the model is confident about the effects of these variables when predicting whether or not someone wins primary. It is difficult to estimate the effect of Total Contribution since we did not scale the numbers in accordance with the other numbers (values are in 100s of thousands), so although the coefficient is small, the effect may be large multiplied by a larger number (total contribution).

For every \$1 contributed, the log odds of winning increase by a very small amount of $1.168 \times 10^{-6}$.

## 2.3   Discussion

To compare our GLM to our Decision Tree, we looked at the RMSE values between the two models, which used the same features. Comparing these values, the GLM did better, with a test RMSE of .345, whereas the Decision Tree had a test RMSE of .439. We believe the GLM performed better even though the GLM had fewer features than our Decision Tree since we were able to provide a link function and by using fewer features, provided less space for overfitting. The Decision Tree may have overfit with the large number of features and large difference in the training (.16) and test (.43) RMSE. In terms of confidence in our model, we believe that the features chosen by the regression and AIC model are strong, but not necessarily sufficient to determine whether or not a candidate would win the primary election in a future election. We also are hesitant on the lack of data in our training set, especially considering that each contestant's race can depend on many variables, of which have different weights for each person.

For the GLM, the model was measured using Chi Squared Statistic, of which can be used for absolute measurement, where the standard of a model's Chi Squared value is determined by subtracting the number of data points by the number of features (n-d), of which a well fit model has a Chi Squared value less than n-d. The number itself is calculated based on how far deviated the true and predicted values are from each other. This value takes into account the problem of overfitting since it looks at the number of features relative to the number of data points by subtracting n-d. Looking at our selected GLM, we had a Chi Squared statistic on our training set of 613. This is greater than our value of $435 - 4 = 431$, which indicates that our model did not fit our training data well.

For the Decision Tree, we looked at accuracy, precision, and recall of our training set to see how well the model predicts the training set. Accuracy tells us how many labels were correctly identified (regardless of the Won or Did Not Win label), of which we obtained .97 accuracy. We obtained a precision of 1 and recall of .96, of which was high since we did not set a max-depth parameter. Precision tells us how many data points that the model labeled positive were actually positive while recall tells us the proportion of how many of the actual True labels we labeled True.

Initially, we had guessed that financial parameters, like Total Contribution, would play a large role in predicting whether a candidate wins the primary election. However, given its relatively small coefficient value of $1.168 \times 10^{-6}$, it seems to have minimal effect, and thus, we can conclude that this feature has a weak relationship with the outcome. The predictor that had a strong relationship to the outcome was the party of the candidate with a coefficient of .89. These variables have a large positive effect on predicting the outcome. Looking at all 4 variables, it seems like the Candidate's State and Party had the largest influence on the outcome, whereas our numerical variables had little effect.

Looking at our Decision Tree, it is nearly impossible to determine what relationships between variables existed. This is because we did not set a max depth for our Decision Tree, and we inputted various numerical variables, which makes interpretability less valuable.

Random noise is the major limitation for our nonparametric model as we assumed there was none. In the instance that our data was corrupted, this could decrease our model accuracy drastically, in which case we could implement a random forest model instead to account for training on separate datasets (bootstrapping) and random features. This way we would be able to make better generalizations.

Our GLM model does not use binary variables (aside from the Party the candidate is in), and our Partisan Lean and Total Contribution features have high densities of values at 0. This may be helpful in determining which candidates didn't do as well, but could also be conflated with missing values. There is the possibility

that the simplicity of our model only allows us to encompass certain parts of winning an election, failing to account for additional factors.

Having information on the candidates that the candidate of interest is running against would be beneficial since votes will either go to the candidate of interest or the competing candidates. Additionally, demographic and views of the candidate were only contained for Democrat candidates, whereas it would have been more beneficial to a wider group of candidates (of other parties) as well as the Democratic party had we had this data on all candidates.

The uncertainty in our model appears to be relatively normal with bounds that are not too tight, but don't span large ranges of numbers. To further look at the uncertainty in our model, we can look at the standard error measurements of various features. We can see that the standard error of all our features is less than .5, which can indicate that our error is well encompassed relative to the number of samples we have. The uncertainty may stem from the sparseness of data, of which we had many values imputed into our dataset. This made it difficult to have substantial data that we could predict our model on. Additionally, we can see increase in uncertainty because we chose not to consider data regarding the candidate's beliefs/stances, of which affects people's votes in reality.

# 3 Question 2: Causal Inference

**Question:** Does endorsement status (whether or not the candidate is endorsed) cause an increase in the percentage of primary votes a candidate receives?

## 3.1 Methods

Endorsement status was the variable we used as the treatment, with values of being endorsed or not endorsed. Thus, the control variable was not being endorsed.

Because we are observing a potential causal relationship between endorsement status and the percentage of votes a candidate receives in their primary election, the outcome variable was primary percentage.

The confounders consist of the following variables: state, race type, primary status, race, total contribution, and candidate party affiliation. These all affect how many primary election votes a candidate receives as well as how much a candidate gets endorsed during the election.

The unconfoundedness assumption holds because the selected variables account for all confounding pathways. Demographic aspects of the candidate are covered in the state, race, and candidate party affiliation variables as voter preferences vary across states (such as partisan leanings) and the state as a whole can encompass the many political and demographic differences that may influence endorsement likelihood and vote percentage. This is further reflected in the other two variables mentioned. For the election itself, race type, primary status, and Total Contribution sufficiently account for the context associated with how the candidate runs the race that may influence both the treatment and outcome variables. For example, for Total Contribution, we may notice that higher contributions can impact the percentage of votes in a certain direction.

To adjust for confounders when estimating causal effects, we used inverse propensity weighting. By reweighting observations (each candidate) based on their propensity scores, we can balance confounders to ensure that the distributions in both endorsed and not endorsed groups are more comparable, removing confounding bias.

There are no colliders in the dataset as no data is a result of the treatment or outcome.
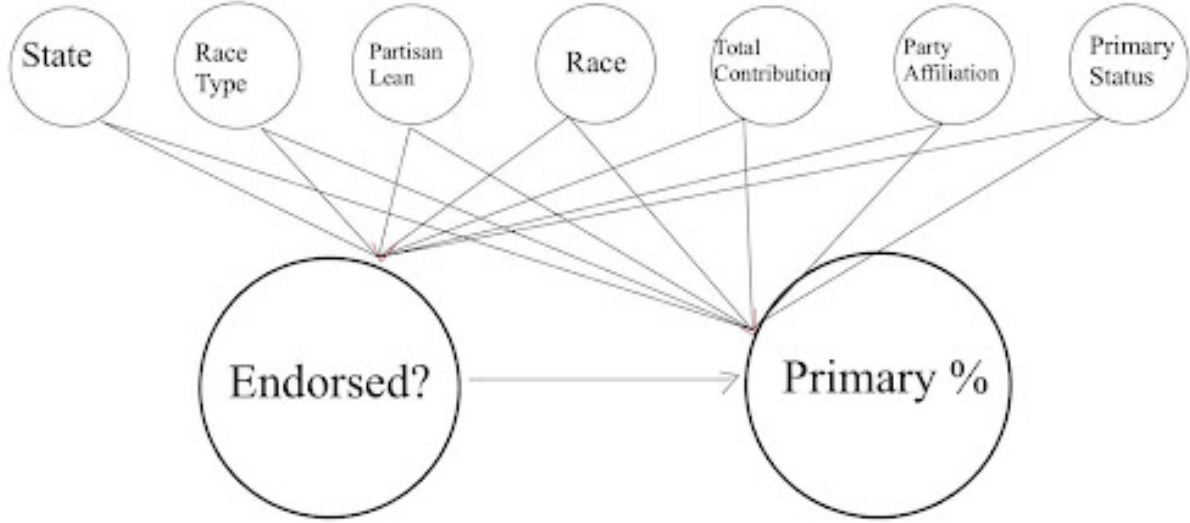
Figure 6: DAG

## 3.2 Results

After using inverse propensity weighting to achieve the average treatment effect estimate, we concluded that being endorsed causes a decrease in primary election vote percentage by around 23.77%. The propensity scores were calculated using a fitted logistic regression model and taken from the predicted probabilities of receiving the treatment given the confounders. For both candidate types that were and weren't endorsed, the propensity scores were all above 0.5. This indicates that all candidates have a relatively high likelihood of receiving the treatment and not receiving the treatment, given the covariates, which could suggest that a certain covariate dominates the logistic regression model. As there were some propensity scores in the 0.9 to 1 range, we also trimmed the scores to see if our estimator contained a lot of variance. However, this did not produce an estimate significantly different from our previous ATE estimate.

Although the negative causal effect is contrary to our initial hypothesis, there may be many reasons as to why endorsements cause lower vote percentage. The logistic model, for example, may not fully represent the interaction between the treatment, confounding, and outcome variables, which could lead to uncertainty in calculating propensity scores. There may also be selection bias in our data that is leading to the negative causal effect that we observed as selection of candidates in the data may not be random.

## 3.3 Discussion

Because we used a logistic regression model, if the true relationship between the treatment and confounders is non-linear, the estimated propensity scores may not be calculated accurately, leading to biased IPW results. Additionally, because we observed that the propensity scores were all above 0.5, this may suggest that certain confounders have a dominating influence over other confounders, and thus, misrepresent the probability of receiving the treatment or not. Although we did trim scores to account for some variability, there still may be high variance due to extreme weights, which could also cause the ATE estimate to be inaccurate. Another limitation might be that the data itself may be concentrated in specific subgroups. Because we saw a negative causal relationship between the treatment and outcome, there could be candidates in the data that have vastly differing data values.

Metrics like media coverage or endorsement times may be useful for answering this causal question as this may give additional insight into how elections are conducted. We could also use factors like endorsement rates by organizations to use as instrumental variables to introduce randomness, and thus, solve the issue of confoundedness. Because endorsements rates only affect whether a candidate gets endorsed and does not directly affect primary vote percentages, is independent of the confounders, and has a nonzero causal effect on candidates getting endorsements, this variable could be considered an instrumental variable.

We are fairly confident that there is a causal relationship between endorsement status and primary vote percentage, because given the confounders, the IPW was an unbiased estimate of the ATE, which was nonzero. This nonzero estimate may also be due to the way endorsements could target weaker/less-preferred candidates who need more financial support, and thus, create a more negative result in vote percentage.

# 4 Conclusion

## 4.1 Outcomes Summary

While understanding the contents of our data, we noticed that the Democratic Data had a lot more descriptive data about the candidate that the Republican Data set did not have. This led us to select features that both Democratic or Republican candidates shared, and for Partisan Lean that was not recorded for Republicans was imputed to 0. After cleaning our data, we ran a GLM with a logistic link function and a Decision Tree to predict whether or not candidates won the 2018 Primary Election. For our GLM, we used the AIC score to test models with different features, of which we chose to use all non-binary features [Total_Contribution, Partisan_Lean, State, Cand_Party_Affiliation] aside from the Candidate's Party Affiliation (Cand_Party_Affiliation). For our Decision Tree, we ran LASSO regression to see if there were any non-contributing features, and then ran our decision tree on 7 features ['State', 'Race Type', 'Partisan Lean', 'Race', 'Total_Contribution', 'Cand_Party_Affiliation', 'Endorsed?'] with no maximum_depth. Comparing the models using RMSE we saw that our GLM outperformed our decision tree with RMSE scores of .345 and .439 respectively. We attribute this difference in performance to an overfitting of the Decision Tree since the RMSE for the training and test set of the Decision Tree varies noticeably. We then ran causal inference to see how being endorsed affected Primary % (the percentage of votes a candidate got in the 2018 Primary Election). Even in light of accounting for confounders using Inverse Propensity Weighting, we saw an unexpected decrease in causal effect of Endorsement Status on Primary %. We had expected being endorsed to cause an increase in support shown through votes in the Primary %, but suspect this difference to be due to a difference in the support of those who have money and are able to support financially and the more general population.

## 4.2 Critical Evaluation

We did not account for multiple years of elections while training our data since our Democratic and Republican data was only based off of 2018. There was FEC data for multiple years, but we were limited by the party data. Additionally, we did not have much information on the stances of the candidates on certain issues, relative to what the general population was feeling about these policies and issues. If we had known this alignment either through some sort of score measuring candidate views relative to population views, or information on what platforms both Democratic and Republicans supported, this could have been more

effective. Lastly, with the variability of the data and issues that are dealt with each year, it would be difficult to accurately predict elections in the long run.

The missing domain knowledge in our data treatment was what we were to do with the missing values that were in our dataset. We also did not know how to condense categorical variables such as states and the sparseness in the data of endorsements. A question we would ask a domain expert is "What information could help us account for missing value, whether that is values to impute or additional data that could accurately represent values that were missing such as Partisan Lean?" The answer would have helped us better deal with many data points that were not complete in our model, giving a fuller picture on the distribution of these missing values and increasing robustness in our model. Additionally, we might gain further insight into how to select futures and what data to look for on other websites. Essentially, we did not know what effect each feature had on predicting the model which made it difficult to choose relevant features. Another helpful piece of knowledge would be how to account for all the confounders in our causal inference problem. We were not able to come up with any further confounders, of which was difficult to know the accuracy of our inference.

Our conclusions from our model are somewhat robust, but there is additional room for improvement. We realize that our goal of only selecting a binary outcome for our GLM and decision tree can make it easier for the model to not be as in depth into the realities of elections. To improve this model, we would possibly choose another outcome such as Primary % so that the model can predict a wider variety of outcomes. The modeling choice may cause bias in focusing on high scores, whereas winning an election only requires a % above 50.

Our results are somewhat generalizable to short term elections that may occur in the future. We would be able to apply this model to the entire country since we included data on all 50 states, but it may be accurate for elections within bigger states. Additionally, we are able to apply our findings to understand the relationship between endorsements and how they apply to election outcomes, though this is subject to variation depending on the endorsers.

## 4.3    Recommendations

A study that can build upon this work could look deeper into what endorsements caused a negative causal relationship to primary %. This could look into additional confounders needed to be clarified through domain knowledge. The study may be able to go further in depth into what endorsements are being accounted for this causal problem and compare causal effects on endorsements in 2018 with other year outcomes.

A call to action would be for candidates to consider support beyond endorsements. Not that there is no benefit from endorsements, rather there may be higher benefit in allocating resources to benefit those who may not be able to endorse. This could include taking into consideration the viewpoints of this group of people or allocating more time to these people to gain their support. Additionally, the people apart from endorsement groups likely do not have as strong of sentiment toward one candidate over the other, of which could be effective.

The impacts of this could be allocating time and resources into things that may not be guaranteed to be efficient, as is endorsements. It also does not build as much momentum nation or state wide, rather focuses on smaller groups, which candidates may not be able to do. This also would call into question how candidates gain information on these people and what information they are using to determine whether or not certain voters fall under this plan of action. This may involve using personal data to determine such things such as socioeconomic status and the location that voters reside in. This implementation, though ideally effective, would be difficult to implement and would still require lots of resources on the backend to be able to afford such action.