**1. Data Overview**
**Primary Dataset Characteristics**

We primarily worked with the data given by FiveThirtyEight on the primary elections of 2018 along with the Federal Election Commission data on campaign donations in 2022. The first dataset–from FiveThirtyEight–were comprised of two different censuses, separated by political party, along with relevant questions specific to each party; for instance, when looking at Republican candidates, one of the critical fields that indicated funding and support was the column "NRA Endorsed." Each row in this set represents an individual candidate who has run for office in the primaries. The secondary dataset, which was a census of campaign financing information, grappled solely with donors, receivers, amounts, and times of contributions in the election process. Instead of displaying a granularity of individuals, the set's rows each represented a singular donation from an individual or campaign to a political candidate. We were specifically interested in swing states during the 2017-2018 period–we expected that running analysis on a non-election year but a critical period in US campaigning would provide interesting and unique insight into the inner workings of the election process.

**Additional Data Information**

Because we hoped to address the question of partisanship and party uncertainty in our research, we found an additional data set from Pew Research Center to ground our work. We chose to concentrate our analysis on the swing states from prior primary elections along with states historically known to lean towards the right or left. As such, our analysis, at least for the initial Generalized Linear Model section, was confined to the swing states Nevada, Arizona, Michigan, Pennsylvania, Georgia, and North Carolina, along with the more partisan states such as California or Texas. Because we wanted our conception of red, blue, and purple states to be representative of real information and data-driven, we used a dataset measuring party affiliation by state from the Pew Research Center to drive our analysis. For this dataset, each row simply represented the partisan lean of that specific state.

**Systematic Exclusion and Bias**

The FiveThirtyEight dataset appears to be quite comprehensive in its examination of Republican and Democratic datasets, however one element our group did note was that the vast majority of information was collected from campaign websites. As such, if a candidate chose to engage in a primarily offline campaign strategy, their information could be incorrect or excluded entirely in the dataset. This may be especially true for grassroots campaigns or campaigns that are primarily locally driven, and therefore may result in a slight bias towards more experienced, established candidates with prior funding.

Regarding the Federal Election Commission dataset, our group inspected the "Filings and Reports" section of the data source and determined that, while there should not be any systematic exclusions, the FEC has frequently struggled with false filings that may be misrepresentative of

actual donation amounts. The slight discrepancy may lead to some degree of measurement error in important metrics for our analysis, such as donation averages by party.

The Pew Research Center dataset makes no mention of any systematic exclusion, however, when we attempted to see the survey question they had asked, we were sent to a "404" error page. As such, we are unclear on if the wording of the question may have caused systematic exclusion or resulted in convenience sampling, drawing only responses from the most impassioned of survey participants.

**Data Collection Consent and Privacy**

Both the Federal Election Commission dataset and the Pew Research dataset required specific input from participants to aggregate the data, meaning that participants were aware of the collection. For both, as well, we can assume that respondents were aware of its intended use, as the Federal Election Commission serves a critical purpose in campaign financing data collection and the Pew Research Center asked a specific question through their resources as a think tank and research organization.

The FiveThirtyEight dataset, however, required no voluntary participation from the collection subjects and instead was primarily put together through web scraping and research conducted on public platforms. Because each of the candidates willingly put forth information about their campaigns via their website, one could argue that they had consented to any use of that information following. However, there is some uncertainty on established consent for this specific usage of their information.

Out of all the datasets, the Pew Research Center table is the most successful in anonymizing participants. None of the datasets, though, harness differential privacy as a technique. For both the Federal Election Commision set and the FiveThirtyEight data, their inherent value is grounded in being linked to specific individuals who have run for office or donated to a campaign. Accordingly, privacy measures do not make much sense in this context.

**Missing Data**

The dataset with the most missing information, perhaps, was the FiveThirtyEight dataset. When looking into the construction of the data, we discovered that this was because any information that was not publicly available about the specific endorsements of a candidate were simply logged as "NaN." Because the dataset hones in on quite a few separate endorsements for each candidate, there was a high likelihood that multiple columns for any given candidate would be marked as empty. In our analysis, because a lack of public mention of an endorsement likely meant that the candidate had not joined hands with that specific group, we simply filled "NaN" values with "No."

The Democratic dataset from FiveThirtyEight also had a particularly helpful column that measured partisan lean for each candidate. Since so much of our analysis hinged on partisanship for our research questions, it would have been helpful to have the same column present in the Republican dataset so that we could draw a direct comparison. In general, the lack of overlap
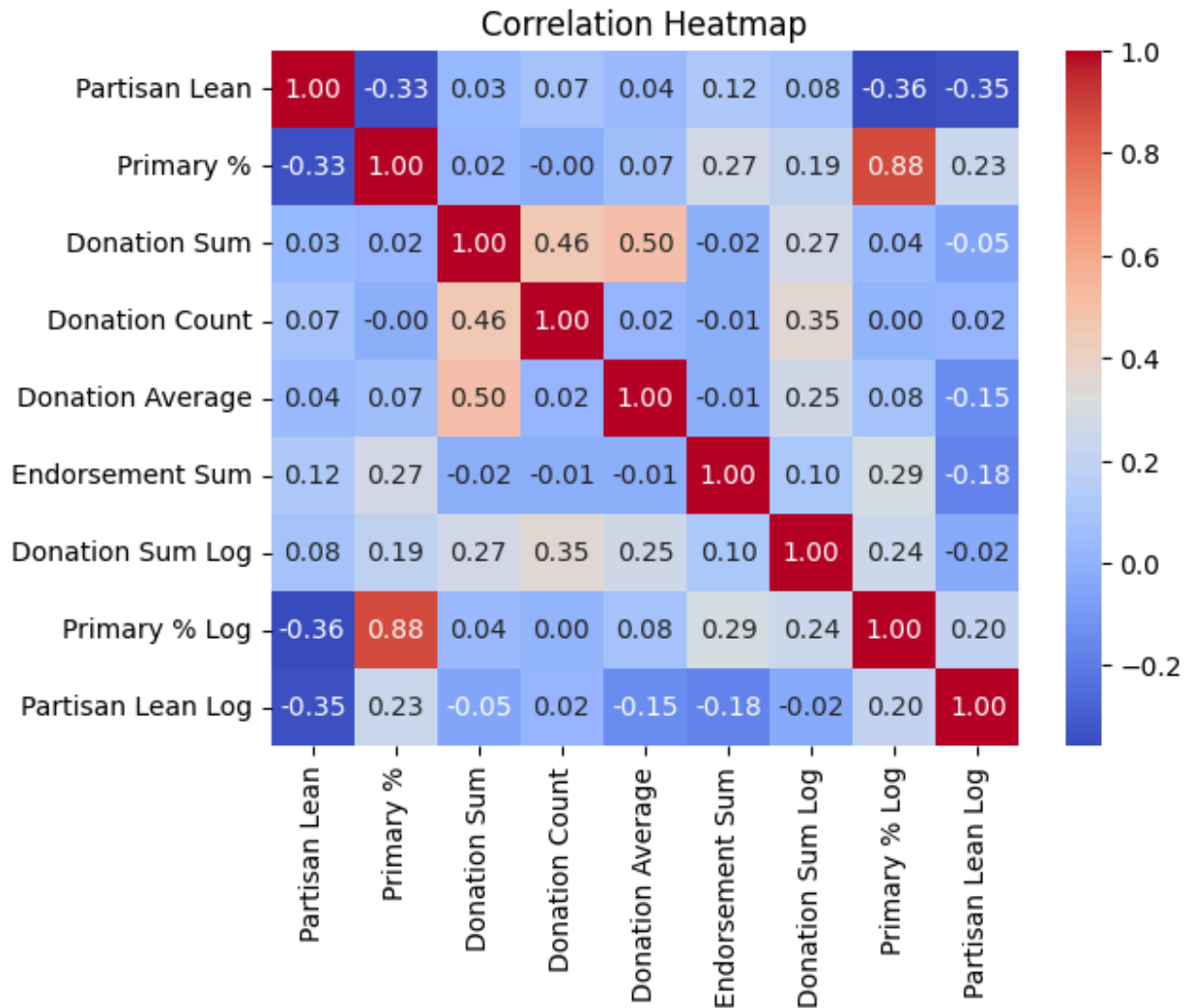
between the columns in the Republican and Democratic dataframes for the FiveThirtyEight data, specifically, was a consistent challenge and sometimes frustrating.
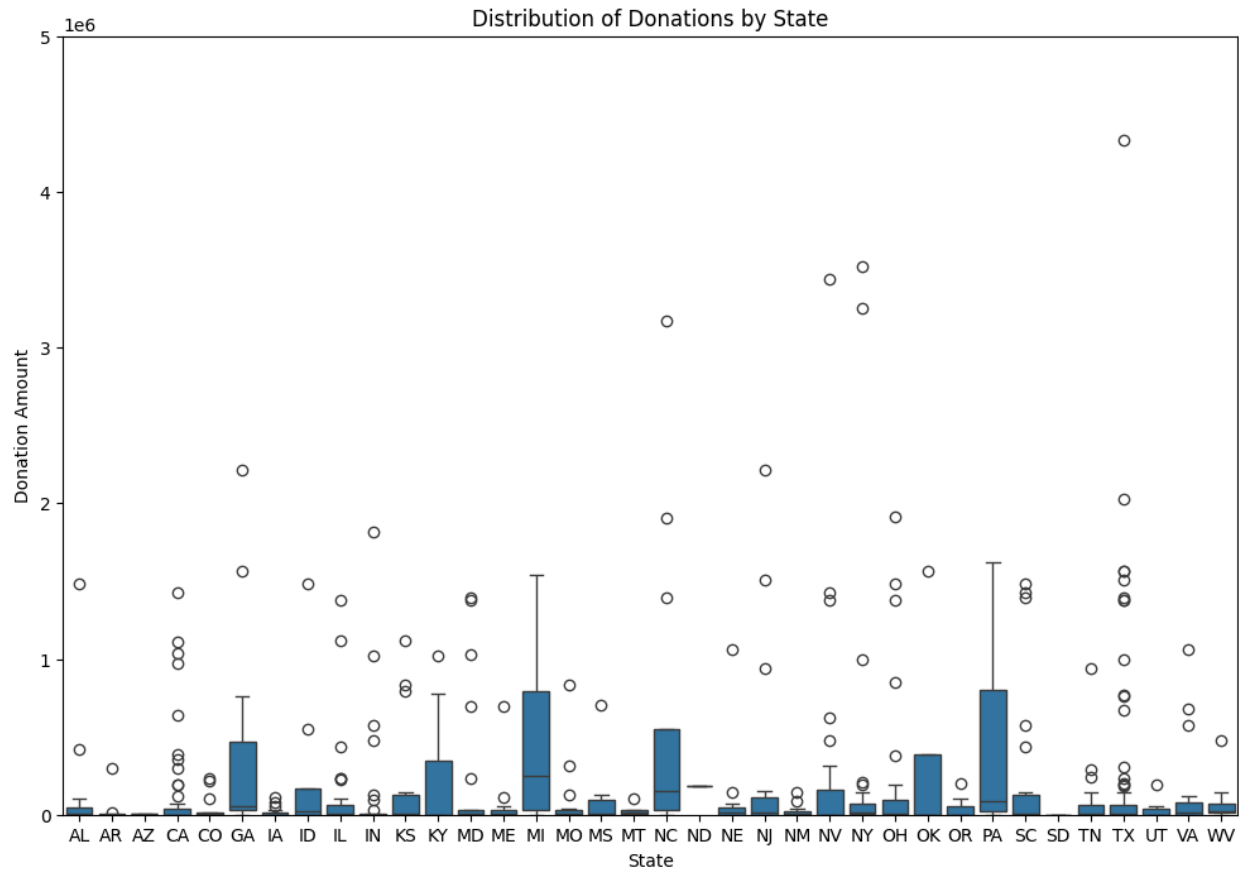
**Cleaning and Preprocessing**

For ease of comparison and reference in our dataset, we first ensured that each dataset referenced the candidates in question by their last names. We further added critical columns that were relevant to our analysis to the Republican and Democratic dataframes, including "Donation Sum," "Donation Count," "Donation Average," "Endorsement Sum," "Endorsement Count," and "Endorsement Average." In later sections–which were specific to each research question–we also took the log of columns that had especially large values to normalize them in relation to the datasets.
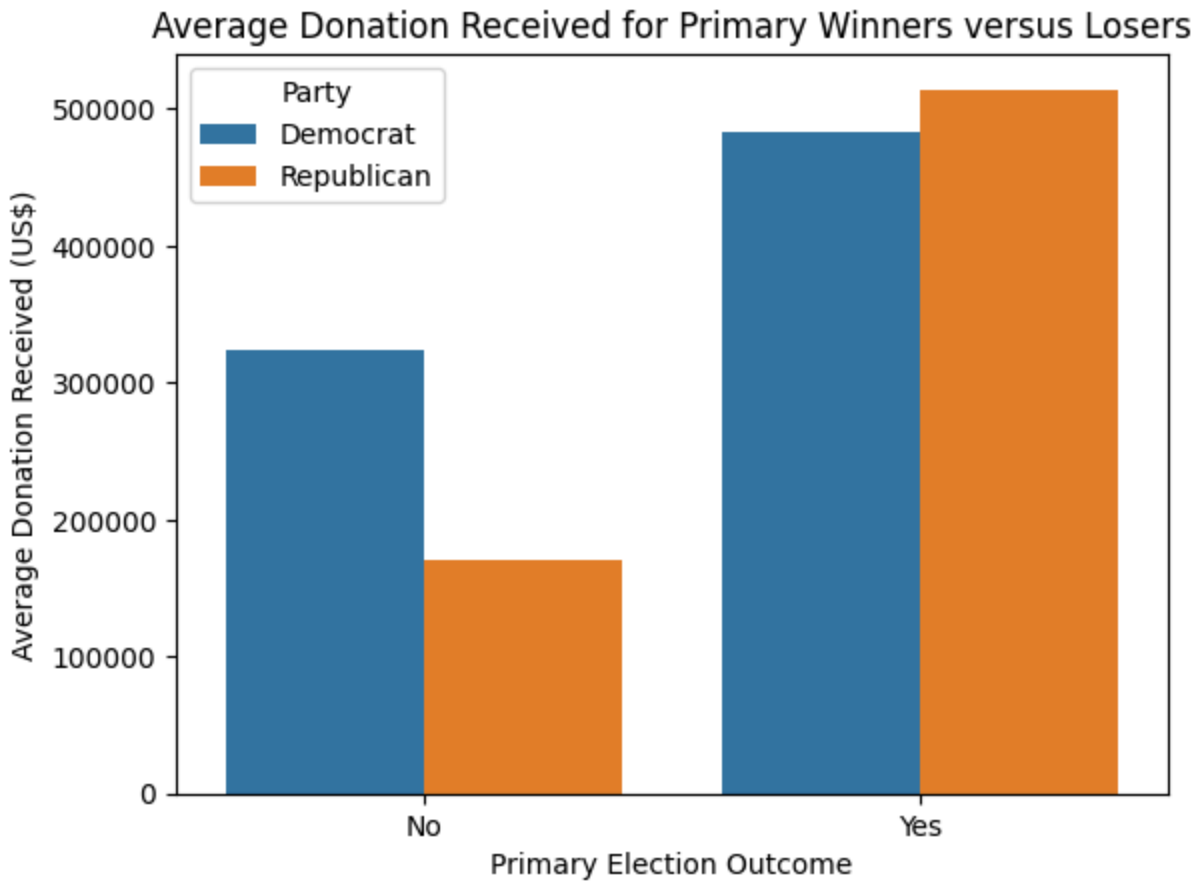
## 2. EDA

In the first stage of our Exploratory Data Analysis, we wanted to see the correlation each numeric feature had with each other. We did this to see potential relationships, as well as to avoid multicollinearity. From the heatmap below, there is a pretty strong relationship between Primary Percent and Partisan Lean, of -0.33. Most of the other correlations are not very strong.

## Correlation Heatmap

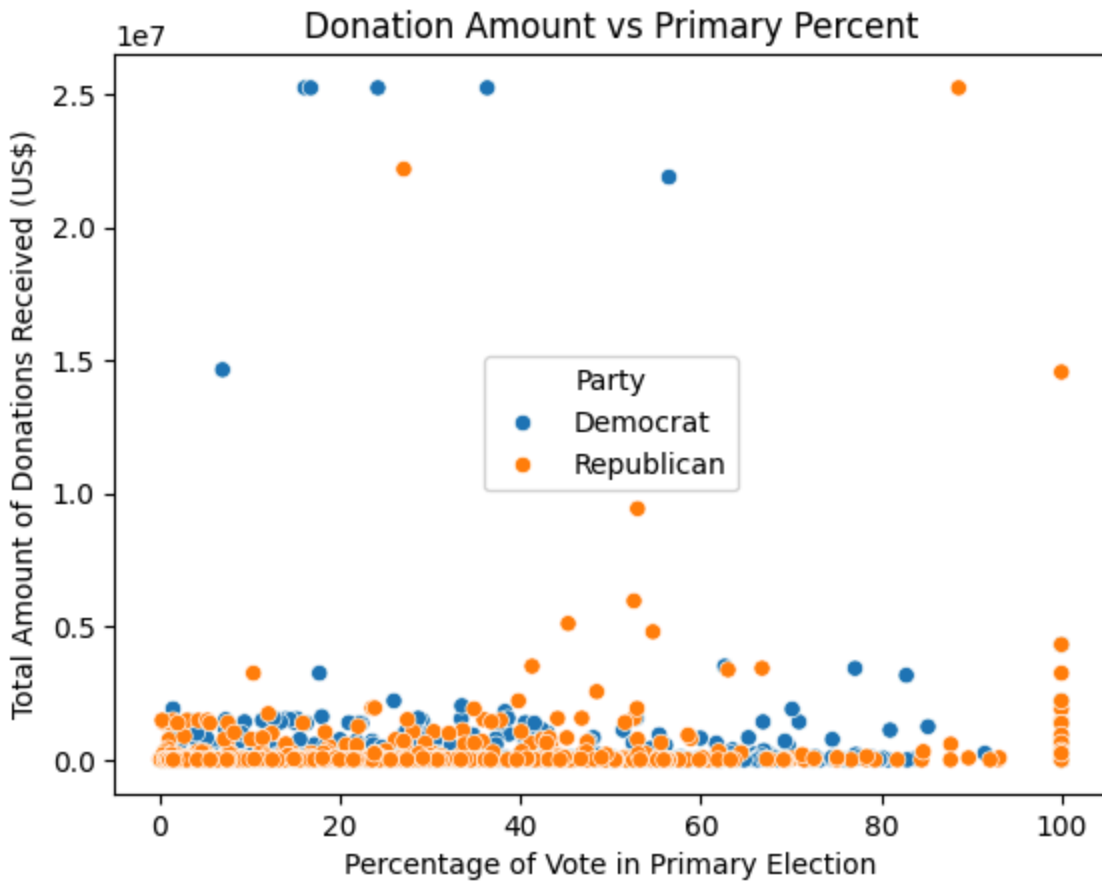| | Partisan Lean | Primary % | Donation Sum | Donation Count | Donation Average | Endorsement Sum | Donation Sum Log | Primary % Log | Partisan Lean Log |
|---|---|---|---|---|---|---|---|---|---|
| Partisan Lean | 1.00 | -0.33 | 0.03 | 0.07 | 0.04 | 0.12 | 0.08 | -0.36 | -0.35 |
| Primary % | -0.33 | 1.00 | 0.02 | -0.00 | 0.07 | 0.27 | 0.19 | 0.88 | 0.23 |
| Donation Sum | 0.03 | 0.02 | 1.00 | 0.46 | 0.50 | -0.02 | 0.27 | 0.04 | -0.05 |
| Donation Count | 0.07 | -0.00 | 0.46 | 1.00 | 0.02 | -0.01 | 0.35 | 0.00 | 0.02 |
| Donation Average | 0.04 | 0.07 | 0.50 | 0.02 | 1.00 | -0.01 | 0.25 | 0.08 | -0.15 |
| Endorsement Sum | 0.12 | 0.27 | -0.02 | -0.01 | -0.01 | 1.00 | 0.10 | 0.29 | -0.18 |
| Donation Sum Log | 0.08 | 0.19 | 0.27 | 0.35 | 0.25 | 0.10 | 1.00 | 0.24 | -0.02 |
| Primary % Log | -0.36 | 0.88 | 0.04 | 0.00 | 0.08 | 0.29 | 0.24 | 1.00 | 0.20 |
| Partisan Lean Log | -0.35 | 0.23 | -0.05 | 0.02 | -0.15 | -0.18 | -0.02 | 0.20 | 1.00 |

Next, we looked at how donations varied by state. Clearly, the most donations were in Michigan, North Carolina, Georgia, Oklahoma, and Pennsylvania. These are all battleground states, which means that there are more donations in states that are more contentious. It could be interesting to explore how being a battleground state impacts the relationship between donation amount and primary percent, as well as the potential impact of partisan lean.

Distribution of Donations by State

We wanted to explore the relationship between donation size and election outcome. In the visualization below, you can see that winners of their primary on average received much more in donation money. Clearly, this suggests a significant relationship between how much money a candidate receives, and their likelihood of winning. This is a relationship we will explore in our research question about the relationship between Donation Sum, Primary Percent, and Partisan Lean.

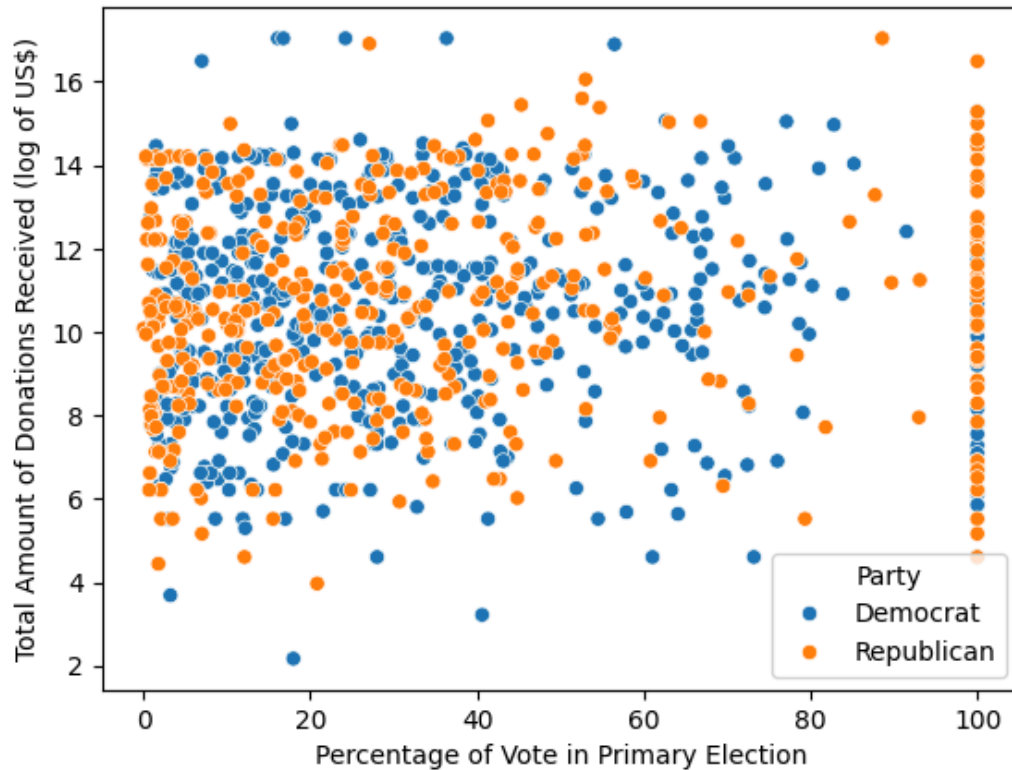Average Donation Received for Primary Winners versus Losers

During the process of Exploratory Data Analysis, we wanted to explore the relationship between Donation Amount and Primary Percent. Initially, we just did a scatterplot of these two features without any data cleaning, and got the following visualization:

Donation Amount vs Primary Percent

In order to make the relationship clearer, we make a visualization of only candidates that at least received some donations, since there were many that didn't receive any. We got the following visualization:

Donation Amount vs Total Donations (for Candidates Who Received Donations)

This visualization shows a potential relationship between Donation Amount and Total Donations, consistent across both Democrats and Republicans. This is relevant to our research question because one of our research questions is about the relationship between Donation Sum, Partisan Lean, and Primary Percent.

**3. Research Questions**

Our research questions are as follows:

1. How well does partisan lean and primary percent perform as predictors of the donation sum of a given candidate? And, how does conditioning by state impact their success as predictors?
    a. For this question, we decided that a GLM-based approach would be the most appropriate. By comparing coefficients of linear models trained across distinct states, we could shed critical insight on the relationship between partisan lean and election outcomes on a state-by-state basis. Our hope was that this would reveal the relationship between sociopolitical factors driven from a local level and election outcomes, which shape decision making for the nation.
    b. One potential limitation of this approach is the lack of clarity surrounding the distribution. By selecting the wrong inverse-link function, we could fundamentally overlook critical elements or correlations within the data.
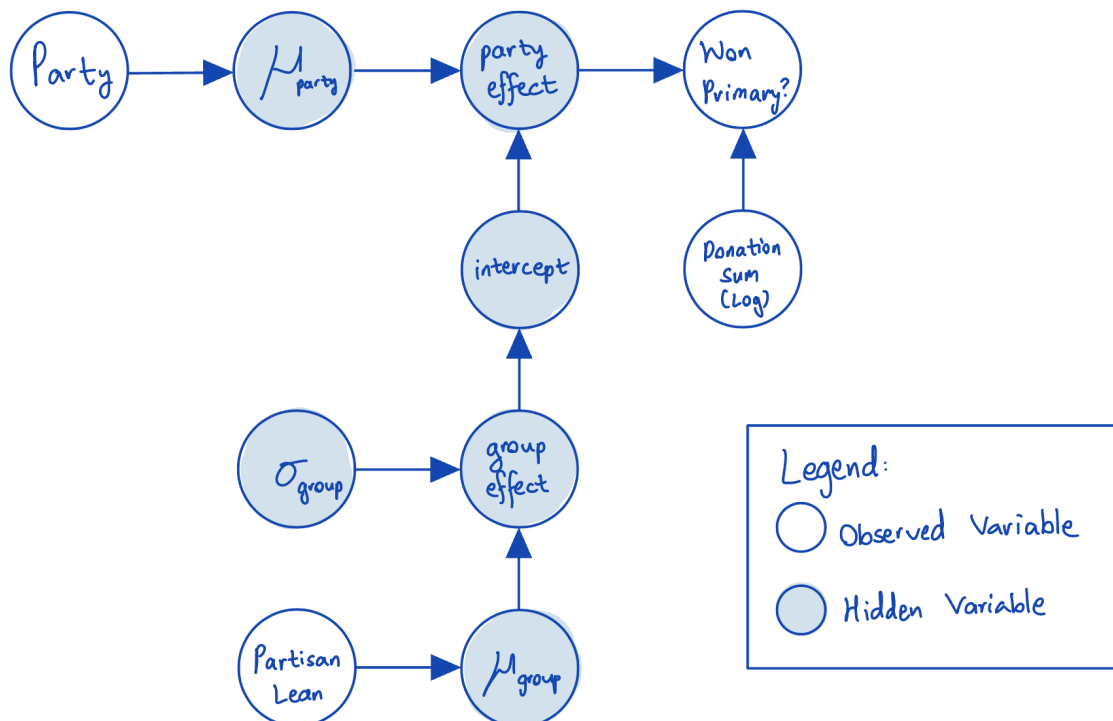
2. By comparing parameters between the two groups, can we isolate a non-negative relationship between election outcome and the amount that is donated to the campaign?
    a. This question aims to answer if grassroots campaigns fundamentally have an advantage or disadvantage over those funded by larger, more established PACs. We attempted to use Bayesian Hierarchical Modeling here because of the causal and varied nature of these relationships, which we felt a graphical model could encapsulate.
    b. One potential limitation here could be the lack or presence of instrumental and confounding variables–the relationship dynamics between variables in this question are quite complex, and it may be challenging to assert independence or causality.

## 4. Inference and Decisions
*Bayesian Model:*

- **Methods**
    - Draw a graphical model, clearly indicating which variables are observed. Provide descriptions of any hidden variables you're trying to estimate.



In our model, we observe the party affiliation of each candidate (`Party`), the partisan lean of the state (`Partisan Lean` in the graph, `Partisanship` in the model), and the sum of donations received on a logarithmic scale (`Donation Sum (Log)`).
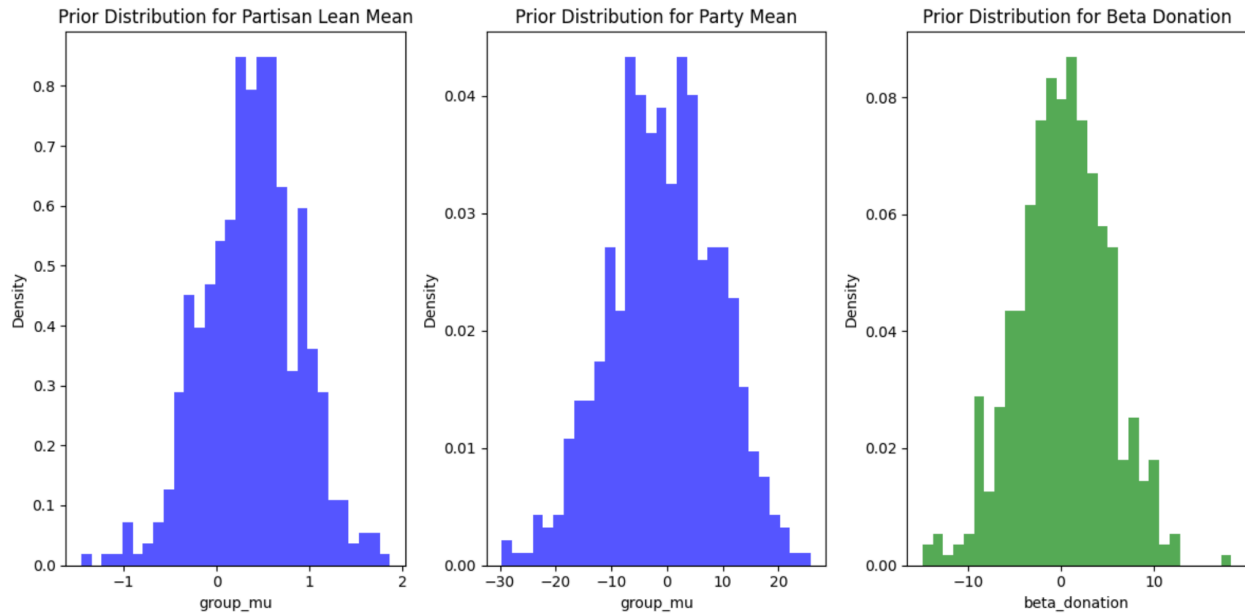
Hidden variables:

- **Group Mu (`group_mu`):** Average primary outcome for each partisan group.
- **Group Sigma (`group_sigma`):** Variability of primary outcomes within each partisan group.
- **Group Effect (`group_effect`):** Individual-level effect based on the partisan group.
- **Party Mu (`party_mu`):** Baseline mean primary outcome for each political party.
- **Party Effect (`party_effect`):** Party affiliation's effect on the election result.
- **Intercept:** Baseline constant for the primary election outcome.

  - Clearly describe the groups in your dataset, and explain why a hierarchical model is a good choice for modeling variability across the groups.

The `Partisan Lean` is a categorical variable, which was extracted from the `States.csv` dataset, with states with less than 6% difference in partisan support considered to be 'battleground' state (note that in the Bayesian hierarchical model, we use a different definition for 'battleground' states) and the rest considered Democrat- or Republican-leaning.

Assumption: Because depending on the two party's support in each state influences how easily a candidate can win, we chose to divide the states based on the partisan support (ie a Republican candidate doesn't have to spend as much on the primary to win their primary in a Democrat stronghold state: because there is less chance of winning the actual election, the competition is less fierce). As the number of contenders differ widely across states due to factors such as competitiveness and population, we have varying amounts of data for each group (or state), which makes a Bayesian model a better choice than a frequentist approach, allowing us to produce a reasonable model that doesn't heavily depend on the limitations of the dataset.

Prior Distributions:

1. **Group Mu: Normal Distribution**
   ○ The Normal Distribution accommodates continuous data with potentially symmetric variations. Its parameters are based on the overall mean and standard deviation of the binary primary outcomes, making it a natural choice for a prior that aggregates diverse partisan influences.
2. **Group Sigma: Half-Normal Distribution**
   ○ The Half-Normal Distribution is chosen to ensure a strictly positive standard deviation while maintaining a simple structure. This prior expresses that the variability within each partisan group is unknown but bounded, allowing us to model differences realistically.
3. **Group Effect: Normal Distribution**
   ○ It's modeled as a Normal Distribution centered around `group_mu` with variability controlled by `group_sigma` to reflect individual deviations within each partisan group. This allows capturing unique variations at an individual level while linking each candidate's outcome back to their partisan group's overall influence.
4. **Party Mu: Normal Distribution**
   ○ The Normal Distribution is chosen to represent the expected primary election outcome based on political affiliation, assuming a continuous and symmetric effect. It reflects inherent differences between political parties while allowing flexibility in deviations.
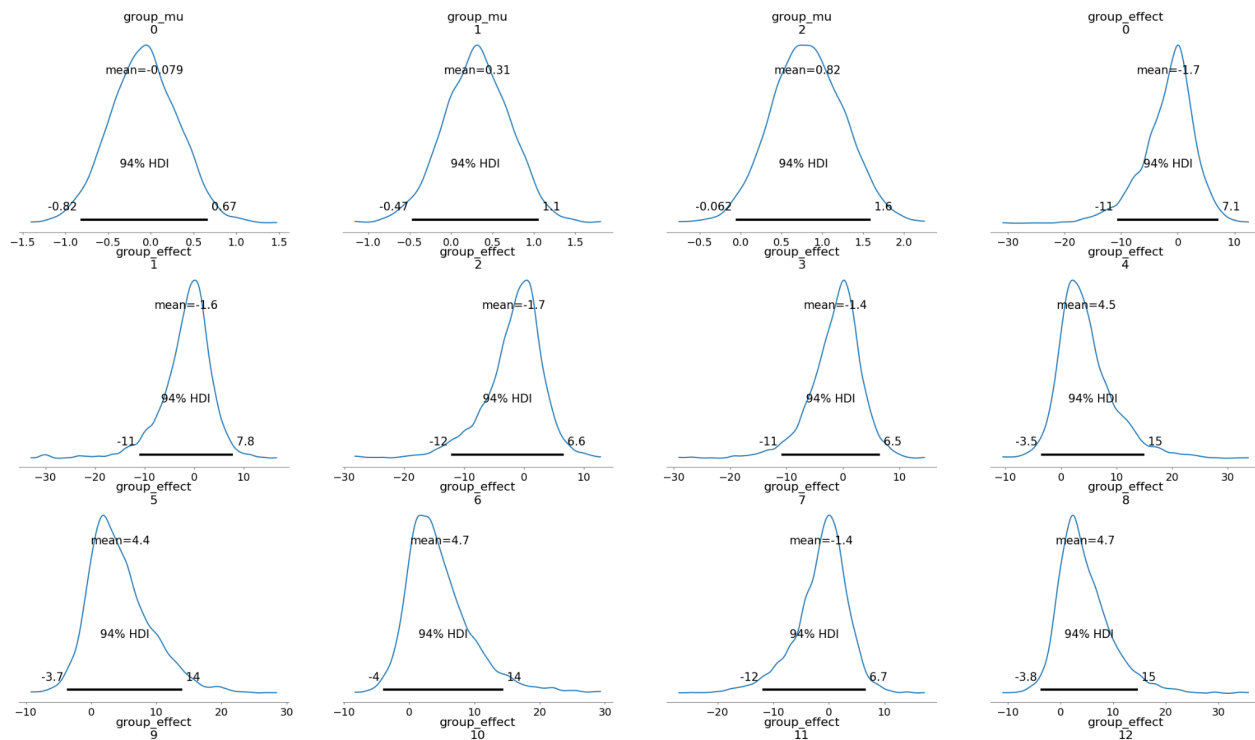5. **Party Effect: Normal Distribution**

○ The Normal Distribution is chosen to reflect each candidate's deviation from the baseline party average (`party_mu`). This normal prior ensures that individual candidates' outcomes are influenced by their party but still have variability.

6. **Beta Donation: Normal Distribution**
   ○ The Normal Distribution is chosen to represent uncertainty in the effect of donations on primary outcomes while assuming potential positive or negative impacts. Its broad prior variance reflects a lack of strong prior knowledge about this relationship.

7. **Intercept: Normal Distribution**
   ○ The Normal Distribution (centered at zero) implies a baseline expectation of no influence without additional predictors. It ensures that the model is adaptable to data-driven adjustments and allows capturing the baseline election results across all observations.

8. **Likelihood: Bernoulli Distribution**
   ○ A Bernoulli Distribution is suitable for representing a binary win/loss event. The predicted probability (from the logistic function) is used to model whether each candidate won or lost based on the various prior influences.

● **Results**

The hyperparameters for the prior distributions were chosen using an empirical Bayes approach, using the mean and standard deviation of the observed `Won Primary` column in the dataset to inform the prior mean (`group_mu`) and its variability (`group_sigma`). This choice is practical given the relatively large dataset and allows data-driven estimation of group-level priors. By incorporating these empirical summaries, the model efficiently captures observed patterns while still allowing for individual deviations within the full hierarchical structure.

○ Summarise and interpret your results. Are there any counterintuitive findings or surprises?

- The three group-level priors reflect different partisan groups, with posterior means of -0.079, 0.312, and 0.819. These indicate varying levels of influence on primary outcomes, with the highest group mean (0.819) indicating stronger positive group influence.
- The group effects show high variability, as indicated by large standard deviations. This suggests significant individual differences within each partisan group, implying that factors other than just partisanship influence primary outcomes.
- Party effects also exhibit high variability with some candidates showing extreme positive or negative deviations from their party's average. This reflects substantial heterogeneity among candidates and indicates that political affiliation alone doesn't determine election outcomes.
- The positive mean of `beta_donation` (0.417) suggests a significant positive relationship between donations and primary success. This aligns with the common understanding that higher campaign funding increases the likelihood of winning the primary.
- The intercept has a broad credible interval (HDI of -13.084 to 5.031), suggesting high uncertainty about the baseline probability of winning when other predictors are not considered.
- The posterior of `group_sigma` has a mean of 5.086, which signifies high variability across partisan groups. This further emphasizes substantial differences among individuals within each group.

The large standard deviations in both group and party effects were unexpected, indicating that neither partisan leaning nor party affiliation alone can predict primary outcomes well. Furthermore, the wide credible interval for the intercept points to a baseline election outcome that is highly uncertain without considering candidate-specific features.

In sum, the model suggests that while partisan and party influences are present, election results are significantly shaped by other (potentially unobserved/confounding) factors, such as individual candidate characteristics and donations.

The credible interval for `group_mu[0]`, representing the mean primary election outcome for the candidates in the Republican-leaning states:

This 94% credible interval means that there is a 94% probability that the true mean primary election outcome for this group lies between -0.824 and 0.668. The credible interval indicates considerable uncertainty around the expected election outcome for this partisan group. The negative lower bound suggests a non-zero possibility that the group could have an overall negative influence on primary election outcomes, while the positive upper bound reflects potential positive influence. In other words, the partisan lean of the state has effects ranging from -0.824 to 0.668 (with 94% probability), with the partisan lean giving positive effect on some individuals and negative effect on the others.

- **Discussion**

In our method, we focused on the party affiliation of the candidate, partisan lean of the state and the total amount of donation received by each candidate, but there are more factors playing in, such as the candidate's stance on issues ranging from abortion rights to gun control versus the population's support, the length of political experience etc that we didn't capture in our model.

Because of the limited features that contrasts with the complexity of reality, we ended up getting 970 different group effects and party effects with widely varying characteristics, reflecting the limitation with our approach.

When we initially used `Primary %` as our observed variable (the Y), the inference had trouble converging, which we assumed to have been caused by the differing primary election rules, competition and heterogeneity in the dataset. For instance, some candidates win with a 30% vote, while others lose with a 30% vote, all with different donation amounts and party dynamics. As such, we switched to using `Won Primary` which is a much simpler boolean outcome, which we could more easily capture with our model (say a Bernoulli distribution).

We had previously used the `Donation Sum` on a linear scale, but due to the uneven distribution and outliers, the model either failed to converge or produced unrealistic values that don't generalize to the rest of the candidates.

In addition, we also tried using a much simpler model based on the `Partisan Lean` of the state as the group. Although the model converged and produced some results, we believed that

there is a more interesting underlying relationship between the `Partisan Lean` and the `Party` affiliation of the candidate, which led us to our current model.

One of the primary limitations of our model was the shortage of features present for the initial model to be trained upon. As a result, instead of us being able to predict a nuanced, numerical value as our output, we were forced to constrict our model to a binary metric of winning or losing the primary election. By adding additional feature columns such as the size of the organization sending a donation, the frequency of donations, or other critical elements beyond simply the transactions themselves on a granular level in the Federal Election Commission dataset, we would likely be able to add more complex relationships into our Bayesian Hierarchical Model and represent them accordingly in our probabilistic conceptions.

### *GLM and Random Forest:*
- **Methods**

We used GLMs and Random Forest modeling to answer our first research question: Is there a relationship between state, partisan lean, and primary percent? And if so, how does the effect of partisan lean vary by state? The features we are using are State, Parisian Lean, and Primary Percent. We only used columns from our dataset with information on Democratic candidates, because there was no column for "Partisan Lean" in the Republican dataset. We want to explore how the effect of Partisan Lean on Primary Percent lessens or strengthens depending on the type of state we're looking at. We created multiple GLM models and Random Forest models to test the relationship when using Blue States, Very Blue States, Red States, Very Red States, Neutral States, No-Lean States, and Swing States. To separate states into these categories, we used supplemental data from Pew Research called "Party Affiliation by State" that showed what percent of people surveyed in each state identified as "Democrat", "Republican", or "No Lean".

For our GLM approach, we went the frequentist approach and made a linear regression model to help figure out our research question. In order to use the linear regression model, we assumed that there were some features in the dataset that were correlated to the donation sum in order to maximize our accuracy. We also assumed that the gamma distribution with an identity function would fit the GLM best because we assumed that there would be many small donations compared to a few large donations.

The nonparametric method we used was a Random Forest Regressor because it offers some of the highest accuracy when it comes to predictions, performs well when it comes to outliers, and deals with missing values well. In addition, it understands complicated, non-linear relationships well. Our data might have some outliers when it comes to donation sum so it is crucial our model can deal with outliers properly. Also the dataset had a lot of missing values. We can assume that the donation sum would have a wide range and potential outliers, so we chose the Random Forest Regressor as it deals with outliers the best. Depending on the model we

will use a RMSE for the random forest regression model and use the null deviance and deviance for the linear regression GLM.

- **Results**

Below is a table summarizing both the results of our Random Forest Model and our GLM Model. Our group added two columns so that we could interpret the results of our models, baseline RMSE for our random forest model, and null deviance for our GLM model. The baseline RMSE is the RMSE if the model predicts the mean every time. The null deviance gives the deviance when there are no predictors.

For the random forest models, they performed the worst on Battleground, Neutral, and No Lean states. For these groups of states, the test RMSE was actually higher than the baseline RMSE. It performed the best on the Very Blue states, which had a difference between test RMSE and baseline RMSE of 0.63. This difference is not high but it is better than the others. The rest had test RMSE that were within 0.2 of the baseline RMSE. For feature importance, each group had Primary Percent as a more important feature than Partisan Lean, with approximately 0.55 for Primary Percent and 0.45 for Partisan Lean for each, give or take about 0.05. From these results, it seems that our random forest model did not perform well at all.

Now for our GLM model, it seems to have performed better than our random forest model. The deviance of each model was consistently lower than the null deviance. The difference between deviance and null deviance was the greatest for Blue states, Neutral states, Battleground states, and Most No Lean States. It performed the worst on Very Red states. This is interesting, because it is almost the reverse of the Random Forest Model. However, since the Random Forest Model performed so poorly, we believe these results can be trusted more. Looking at the p-values, the p-value for Partisan Lean was above 0.05 for the Very Red, Blue, Red, and Neutral states. In fact, the p-value was 0.97 for Neutral states. The z-scores result mirror the p-values. Those that had a high p-value for Partisan Lean had a low z-score for the same feature. Based on the p-values, it is clear that some of the groups of states' results cannot be trusted, especially the Blue, Red, and Neutral states. Excluding those, the model seemed to have performed pretty well, with a solid difference between the null deviance and deviance for each group of states.

| State Group | Base RF RMSE | Test RF RMSE | RF Feature Importance, Partisan Lean | RF Feature Importance, Primary Percent | GLM Coef. Partisan Lean | GLM Coef. Primary Percent | GLM Null Deviance | GLM Deviance | GLM P-values | GLM Z-Score |
|---|---|---|---|---|---|---|---|---|---|---|
| Very blue | 5.18 | 4.55 | 0.424` | 0.576 | 0.0197 | -0.0425 | 4864.5 | 4172.2 | All <0.05 | All >2 |

| | | | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|
| Very red | 5.318 | 5.115 | 0.372 | 0.627 | 0.0318 | -0.0609 | 3656.1 | 3613.9 | 0.1 for Partisan Lean | 1.6 Part. Lean |
| Most no lean | 5.26 | 5.72 | 0.44 | 0.55 | 0.0097 | -0.0424 | 10,867 | 9048.4 | All <0.05 | All >2 |
| Blue | 5.315 | 5.28 | 0.43 | 0.57 | 0.0035 | -0.0292 | 13,279 | 9676.3 | 0.36 for Partisan Lean | 0.9 Part. Lean |
| Red | 5.47 | 5.43 | 0.39 | 0.6 | 0.0098 | -0.032 | 8116.6 | 7213.8 | 0.33 for Partisan Lean | 0.99 Part. Lean |
| Neutral | 5.54 | 5.61 | 0.47 | 0.52 | 0.0003 | -0.0149 | 8192 | 6215.9 | 0.97 for Partisan Lean | 0.04 Part. Lean |
| Battle-ground | 4.84 | 5.37 | 0.4 | 0.59 | -0.0081 | -0.014 | 3999.5 | 2114.5 | All <0.05 | All >1.9 |

Diving into the uncertainty of the GLM predictions, we can see that there is a high null deviance when it comes to the no lean, red, blue, and neutral states with a range of $\approx$ 8,000 - 13,000. Compared to the very blue, very red, and battleground states, there is a lot less deviance with a range of $\approx$ 3,500 - 5,000. This can be due to an underlying factor on people's beliefs, income, number of states, or population size.

- **Discussion**

The GLM seemed to perform better, since the measure our team used to assess its accuracy, deviance, performed consistently better than the error metric for the Random Forest model. Given the current performance, we definitely would not recommend the use of the Random Forest model for future datasets. The GLM might still be useful, though the deviance is still higher than the ideal desired value for a high accuracy.

When comparing the performance of the random forest regressions between each other, there seemed to be a low variance and high bias for the RMSE. As for the GLM, we can see that there is low bias and high variance from the null deviance. As stated before, the no lean, red, blue, and neutral states had a higher deviance when compared to the the very blue, very red, and battleground states.

Both models do indicate a relationship between Donation Sum, and both Partisan Lean and Primary Percent. For both models, they consistently had a greater feature importance for Primary Percent, and had a greater absolute value of the coefficient for that feature. This suggests that Primary Percent has a greater relationship with Donation Sum than Partisan Lean.

Additionally, Primary Percent consistently had a p-value below 0.05 and a Z-score above 2, whereas Partisan Lean sometimes had a p-value above that and a low Z-score. Overall, Primary Percent had a very slight but statistically significant negative correlation with Donation Sum, hovering between -0.014 and -0.06. Partisan Lean had a very slight positive correlation with Donation Sum, but it is only statistically significant for some of the states. Partisan Lean's coefficient was between -0.0081 and 0.0035.

       To add some depth to our analysis, we also chose to examine the results across groups. The strongest relationship between Donation Sum and the two features was with the Very Red states. This one had a Primary Percent coefficient of -0.06 and a Partisan Lean coefficient of 0.03. In our Random Forest model, the feature importances were 0.372 and 0.627 for Partisan Lean and Primary Percent respectively. However, the p-value for the coefficient for Partisan Lean is 0.1 in this case, indicating that the result for Partisan Lean is not statistically significant. The relationship was also strong with the Very Blue States, with coefficients 0.0197 and -0.0425 for Partisan Lean and Primary Percent respectively. The feature importances were also similar to that for the Very Red states. The relationship between the features and Donation Sum was the weakest for Blue, Red, Neutral and, surprisingly, Battleground states. These results indicate that partisanship of the state you're in indeed changes the relationship between Donation Sum and Primary Percent. The relationship between Partisan Lean and Donation Sum however, does not vary greatly across groups, but is greatest in Very Red states. All of these results suggest that the relationship between Partisan Lean and Donation Sum is weak. The relationship between Primary Percent and Donation Sum is also weak, but is statistically significant across all groups. It is also strong in the most partisan states, suggesting that in high partisan states, there is a stronger relationship between how much the candidate receives and donation, and how much they're likely to win their primary by.

       Some of the limitations of a random forest regression model are that it may be hard to see a relationship between the predictor and response variable and it requires a big dataset to precisely make accurate predictions. This is due to the fact that it is a black box model and outputs an output without really telling the reasoning. Furthermore with a larger sample size, random forest regression models can better capture the underlying patterns and relationships in the data. More data points provide a more comprehensive representation of the population, reducing the risk of overfitting and improving the model's ability to generalize to new, unseen data. A limitation for a linear regression model is that features used to predict need to be correlated to what you are trying to predict or else the model will not be accurate.

       Our models could likely be improved with additional data. There is likely an important relationship between how many endorsements a candidate gets and how much that candidate is donated to. This would be a confounding variable, that might be negatively impacting the accuracy of our results. Unfortunately, the dataset we currently have does not give us the total number of endorsements, but rather a select few of the endorsements, making the number too small to affect the model accurately.

**5. Conclusion**

After doing our analysis on research question 1, we concluded that partisan lean and primary percent are both sub-optimal features to use as predictors for donation sum. Partisan lean ended up being the worse feature out of the two. The generalizability of the results should be interpreted with caution and considered in the context of the specific dataset, since our model did not perform the best and every year is different in terms of candidacy and how the economy is doing. In order to do our research, we needed to merge the financial dataset of the candidates with the primary candidate endorsements. This allowed us to tie in the donation sum with the candidate and their endorsements. Looking into the future, we can start looking at other years and see if we have the same result.

The results of our Bayesian model indicates there's substantial variability among candidates within each group. The varying levels of influence among partisan groups suggest that political affiliation alone doesn't determine election outcomes. Surprisingly, both group and party effects show high variability, indicating that factors beyond partisanship and party affiliation significantly shape primary outcomes. The positive relationship between donations and primary success aligns with what we would expect. However, the broad credible interval for the intercept shows the high uncertainty about baseline election outcomes when other predictors are not considered. Additional data on factors like donation source, frequency, and other critical elements beyond transactional data could enhance the model's predictive power. Putting these features into the Bayesian Hierarchical Model could better represent the underlying relationships in primary election outcomes.

Given the results of out research, the best policy we could implement is campaign financing reform. Given the positive relationship between donations and primary success, there should be some regulations on compaign contributions to mitigate the influence of money in politics. Donations should reflect popularity among constituents rather than large organizations. Additionally, financial transparency should be better ensured when people run for office, as this would increase knowledge and trust for constituents.