

Case Study 3: Election Prediction

Bob Ding, Becca Erenbaum, Grace O’Leary, Rena Zhong

11/1/2020

1. Introduction

The outcome of the 2016 election not only stunned the nation, but also sent shockwaves through the statistical and polling communities. Over the course of the election year in 2016 up until the week of the election, poll predictions of Hillary Clinton’s likelihood of beating Donald Trump ranged from 71%- 99% probability [1]. So when Trump beat these odds, the polling industry lost a lot of trust from the general public [2].

This small, specialized industry that is the political polling business, has a great deal of influence on how the election is portrayed in the news media, voter decisions, and candidate partisan policy initiatives. [1]. Million dollar decisions on advertising and campaign strategy are dictated by polls.

Polls conducted early in the election year are a weak predictor of election outcomes because the general public has paid less attention to the race and is less knowledgeable of the candidates’ platforms. Voters that tend to sway between parties often report that they are undecided, however these are arguably the most important opinions in predicting the election outcomes. [3] As the election nears, polls become more accurate, but that is where historical prediction models have been lacking – they fail to take into account the differential accuracy in poll results. Historical models, regression based models that rely on outcomes from past elections and structural factors, predict election outcomes at a single point in time which renders high levels of uncertainty [4].

Drew A. Linzer developed a dynamic bayesian forecasting model to predict the U.S. presidential election at the national and state level that combines these historical models with everchanging poll updates. Linzer’s model uses hierarchical specification to handle states being polled on different days and takes into account sampling errors of the polls and national campaign effects [4].

In this report we aim to:

- Predict the outcome of the presidential election and the electoral college vote using the Linzer model to predict swing state outcomes in combination
- Predict whether the US Senate remains in Republican control by using an adaptation of the Linzer model and FiveThirtyEight Senate poll data for each state
- Predict the outcomes of all 13 NC Congressional elections using Linzer model with input from FiveThirtyEight Senate poll data for North Carolina along with our model that predicts who will vote in North Carolina
- Predict the outcome of the NC Senate election and the associated uncertainty using the Linzer model.

Taking into account the anomaly that was the 2016 election, we chose to use the Linzer model to best account for slight and unexpected changes leading up to election day. In Section 2 of this report we will discuss datasets used for all tasks, including a brief exploratory data analysis. In Section 3, we formulate models and methodologies that answer all the research questions. Section 4 will present model diagnostic and sanity check validation of prediction. Then, in section 5 we present the major results of our analysis. Section 6 will be focused on conducting sensitivity analysis to test data imputation hypothesis and prior choices.

2. Data Source and EDA

2.1 Description of Data

In order to answer these questions, we used a total of four datasets: senate polls, house polls, 2020 US presidential election polls, and North Carolina voter registration history snapshot dataset. Both the senate and house polls dataset comes from the fivethirtyeight website [5], while the presidential election polls dataset comes from the Economist website [6]. Senate and house polls have 38 variables and 4061 and 2655 observations respectively. The presidential election polls have 1447 observations and 17 variables which are: poll state, pollster, sponsor, start and end date, entry time, number of observations, population, method, Biden, Trump, Biden margin, other, undecided, URL, include, and notes.

The North Carolina voter history dataset contains information on the 2016 elections and about how voters voted (method, election, county, party affiliation, precinct). This will help us model and identify potential voters in 2020 (Interim report). On top of the created model, we used the 2020 voter registration snapshot dataset to identify potential voters, and use this information as additional input to model house polls results.

2.2 Exploratory Data Analysis

To get a basic understanding of the four datasets we were working with, we went through the data to understand what we were working with. In the senate and the house polls, we explored the variables of: methodology of the poll, the state in which the poll was conducted, which party the candidate was, and the percentage of the poll. In the presidential election polls, we explored similar variables of: method of poll and state, in addition to the margin between Biden and Trump. Because all three of these datasets had similar variables, we are able to compare them to each other. In Figure A, we can see how the methods of polling were distributed between the house, senate, and presidential polls. In Figure B, we can see how each poll's distribution of which state was polled, and in Figure C, we can see how the house and senate polls' candidate party are different from each other. Finally, in Figure D, we can see the distribution of the margin between Biden and Trump.

To explore the North Carolina voter history dataset, we looked at variables such as method, county, and party affiliation of the voter. We can see in Figure E that the majority of North Carolina voters are from two counties, Wake and Mecklenberg. Party affiliation and the method that the voters used are included in our appendix.

3. Model Formulation

3.1 Model for Question 1

In this section, we build a model to answer our first research question: Predict the outcome of the presidential election and the electoral college vote using the Linzer model to predict swing state outcomes in combination. We follow the route of Linzer and build our model. Specifically, the motivation is that we want to model Joe Biden's percentage of EC vote in the state specific level of detail. We believe that throughout time, each state's preference for Joe Biden is randomly varying and follows a random walk process. Besides, each state's uncertainty of their attitude towards Joe Biden is also state-widely different. We make the assumption that such uncertainty is constant throughout time. Therefore, we create the following extension of Linzer Model.

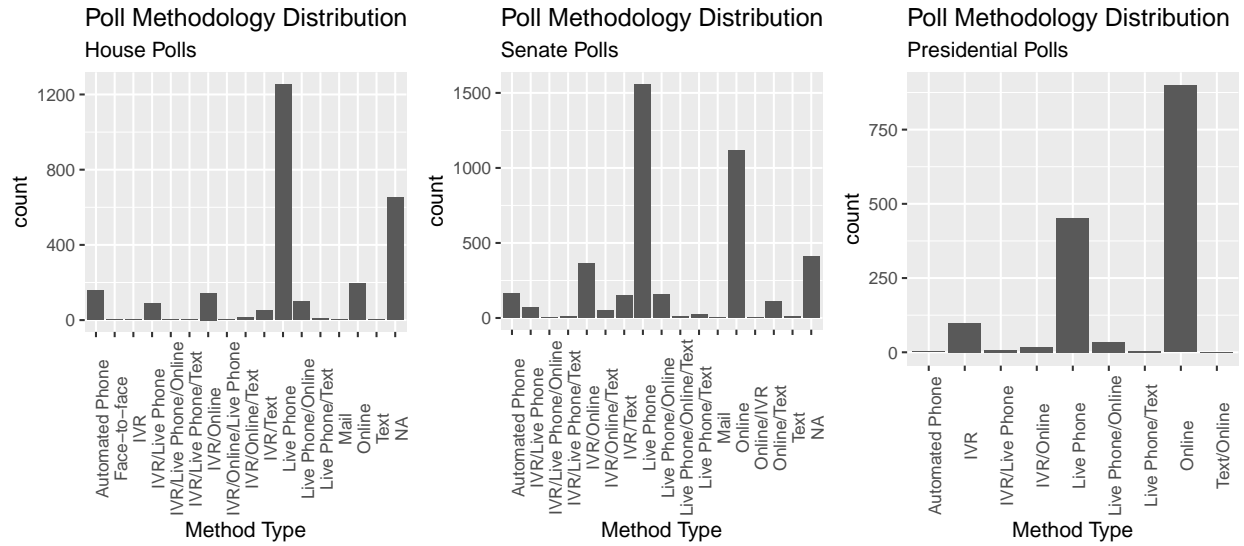


Figure 1: Figure A: Poll Methodology

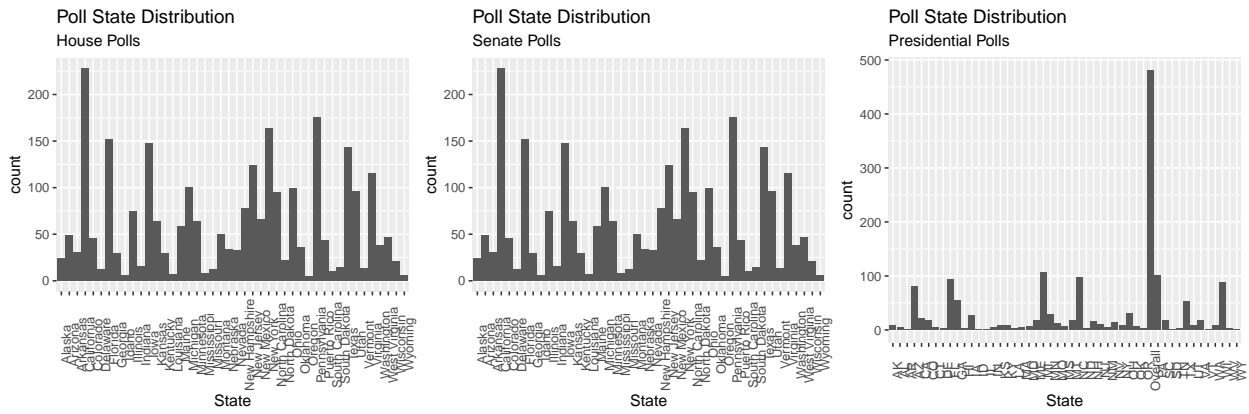


Figure 2: Figure B: Polls State Distribution

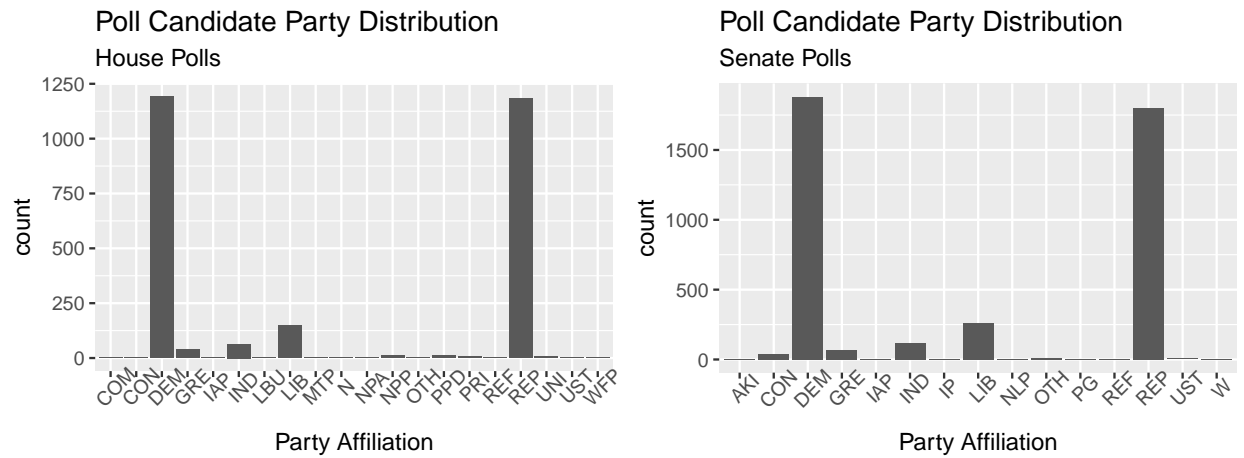


Figure 3: Figure C: Candidate Party Polls

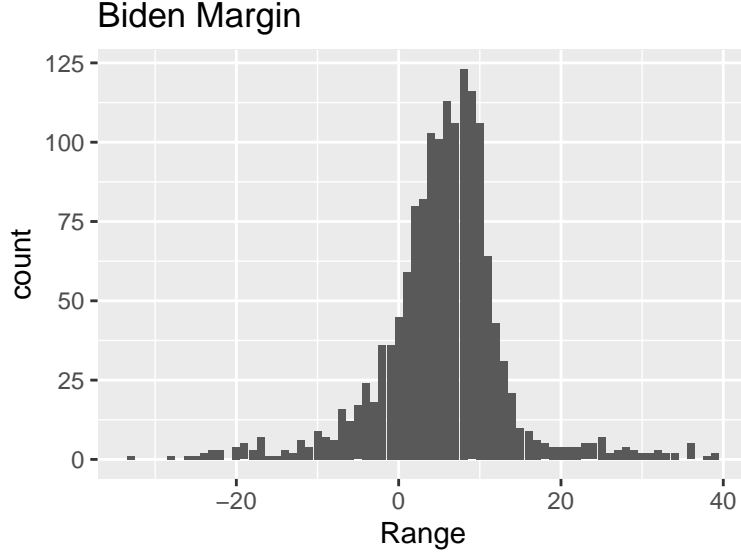


Figure 4: Figure D: Biden Margin

$$\begin{aligned}
y_k &\sim N(\beta_{i[k]t[k]}, (\sigma_y^2)_{i[k]}) \\
\text{for } t > 1 : \beta_{it} &\sim N(\beta_{i,t-1}, (\sigma_\beta^2)) \\
\text{for } t = 1 : \beta_{i1} &\sim N(\mu_0, \sigma_0^2) \\
(\sigma_y^2)_{i[k]} &\sim \text{InvGamma}(\nu_y, \tau_y) \\
(\sigma_\beta^2) &\sim \text{InvGamma}(\nu_\beta, \tau_\beta) \\
\mu_0 &\sim N(50, 17) \\
\sigma_0^2 &\sim \text{InvGamma}(\frac{1}{2}, \frac{1}{2}) \\
\nu_y &\sim \text{Uni}(0, 100) \\
\tau_y &\sim \text{Uni}(0, 100) \\
\nu_\beta &\sim \text{Uni}(0, 100) \\
\tau_\beta &\sim \text{Uni}(0, 100)
\end{aligned}$$

Where k index the polls, i index the states, and t index the date. $i[k]$ represents the k^{th} poll's corresponding state index, and $t[k]$ represents the k^{th} poll's corresponding date index. We model samples from the economist poll data of 2020 and therefore are able to sample the posterior distribution of 1) Joe Biden's chance of winning, and 2) total electoral vote. The sampled posterior distribution y (percentage of supportance within each state) will be calculated by comparing to rounding by 50% to simulate the "winner takes all" procedure. After such adjustment, we can simply multiply the number of electoral college votes of each state to obtain a posterior distribution of electoral college votes of the entire nation. Thus, 2) can be obtained. By comparing to 270 again provides the probability of Joe Biden being elected.

3.2 Model for Question 2, 4

In this section, we're answering question 2 and 4 jointly: Predict whether the US Senate remains in republican control and predict the outcome of the NC Senate election. Similarly, using the above model with exactly the same notation, and by switching Economist dataset to FiveThirtyEight senate polls data, the exact

inference procedure can be produced. After modeling and predicting the Republican's support percentage, we can predict whether the Republican or Democrat party wins the state majority. This procedure includes the state of North Carolina. Furthermore, by summing the posterior samples of Republican senators in each state on election day, we can obtain the posterior distribution of US Senate remains in Republican Party's control.

3.3 Model for Question 3

In this section, we answer research question 3: Predict the outcomes of all 13 NC Congressional elections. The dataset we've used is the FiveThirtyEight house poll data. However, the dataset's sample sizes on each NC congressional district are small. This is due to the fact house polls happen less frequently compared to either senate polls or presidential polls. Besides, not all congressional districts are competitive because there is no term limit for representatives. FiveThirtyEight has only recorded poll history data for competitive districts including districts 2,3,7,8,9,11,13. Thus, correlation between non-competitive districts cannot be captured without inputting extra datasets. Thus, we concatenated this section's modeling procedure with the output of the interim reports' "who vote" results.

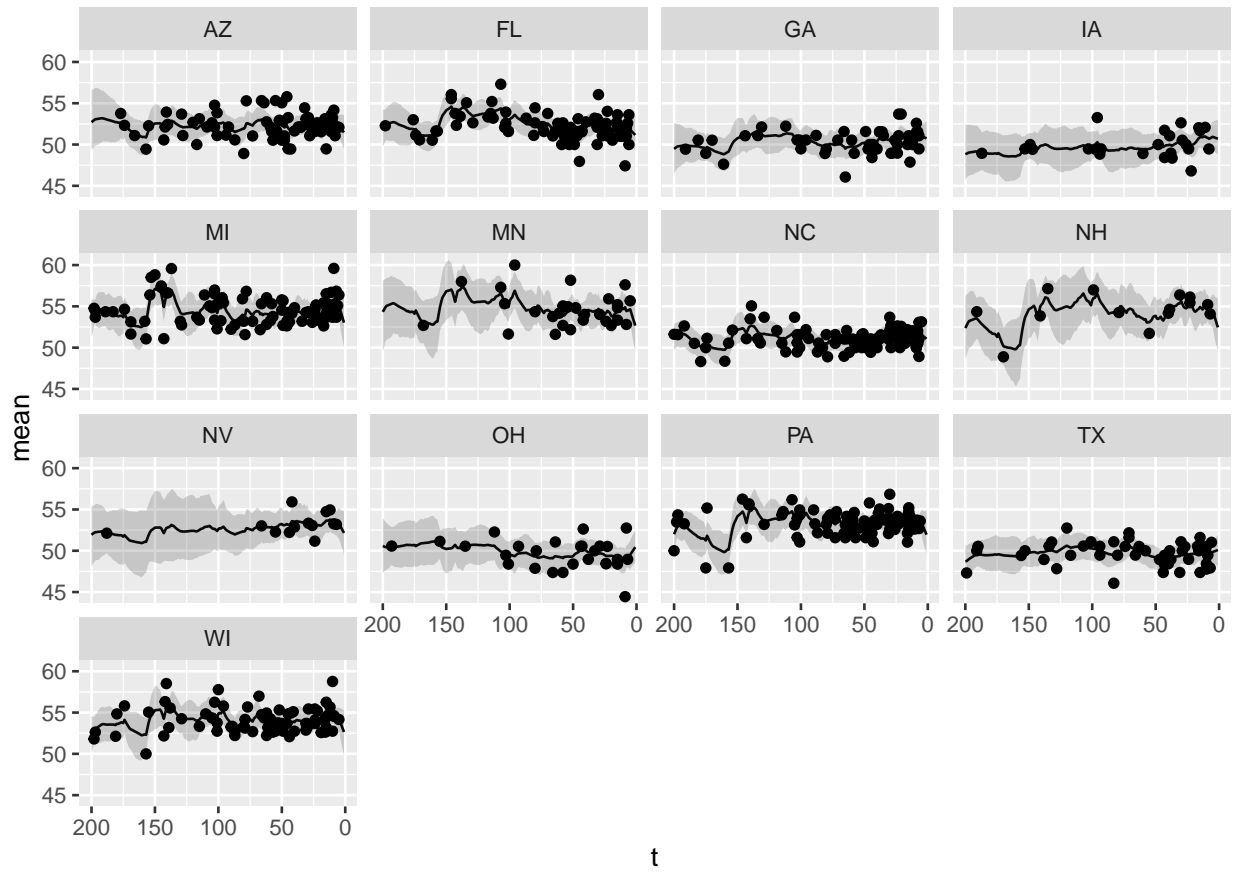
Using the model to predict who will vote in 2020, we ran this binary output model on all the registered voters in NC to predict whether they'll vote or not. In the NC voter registration profile 2020 snapshot dataset, there is a column indicating party affiliation of each potential voter. Thus, we can predict the percentage of Republican voters for each congressional district and use this percentage as an imputed value for polls percentage. In this way, we're essentially imputing polling results for non-competitive districts, which can be also fed into the Linzer Model defined below. Such imputation is reasonable but bold. Thus, sensitivity analysis of imputed percentage will be further explored later in section *****

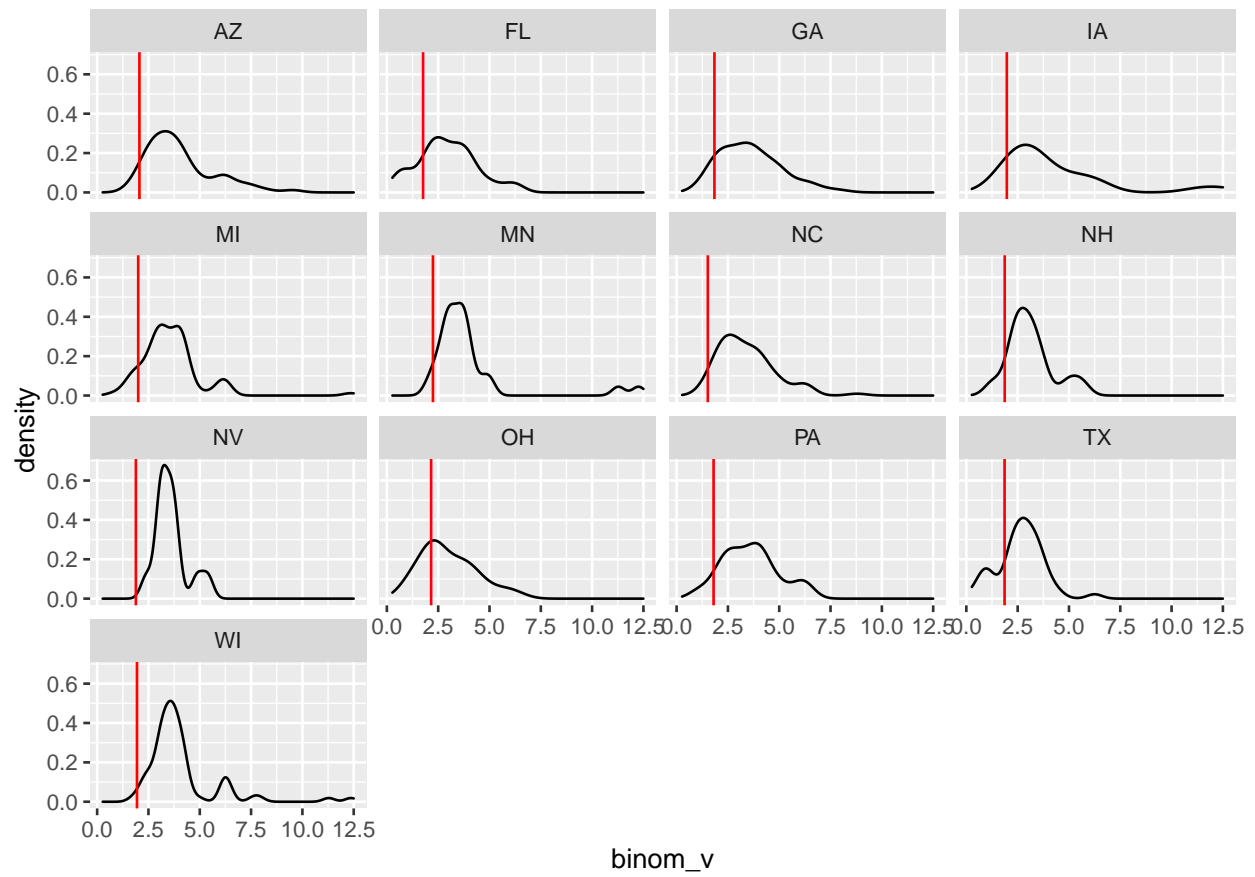
Below is the model we're using for this task:

$$\begin{aligned}
y_k &\sim N(\beta_{i[k]t[k]}, (\sigma_y^2)_{i[k]}) \\
\text{for } t > 1 : \beta_{it} &\sim N(\beta_{i,t-1}, (\sigma_\beta^2)) \\
\text{for } t = 1 : \beta_{i1} &\sim N(\mu_0, \sigma_0^2) \\
(\sigma_y^2)_{i[k]} &\sim \text{InvGamma}(\nu_y, \tau_y) \\
(\sigma_\beta^2) &\sim \text{InvGamma}(\nu_\beta, \tau_\beta) \\
\mu_0 &\sim N(50, 17) \\
\sigma_0^2 &\sim \text{InvGamma}(\frac{1}{2}, \frac{1}{2}) \\
\nu_y &\sim \text{Uni}(0, 100) \\
\tau_y &\sim \text{Uni}(0, 100) \\
\nu_\beta &\sim \text{Uni}(0, 100) \\
\tau_\beta &\sim \text{Uni}(0, 100)
\end{aligned}$$

Where k index the polls or the imputed polls, i index the congressional district, and t index the polls date. $i[k]$ represents the k^{th} poll's corresponding congressional district index, and $t[k]$ represents the k^{th} poll's corresponding date index.

```
## Warning: Removed 2044 rows containing missing values (geom_point).
```





```
## [1] 0.6173
```

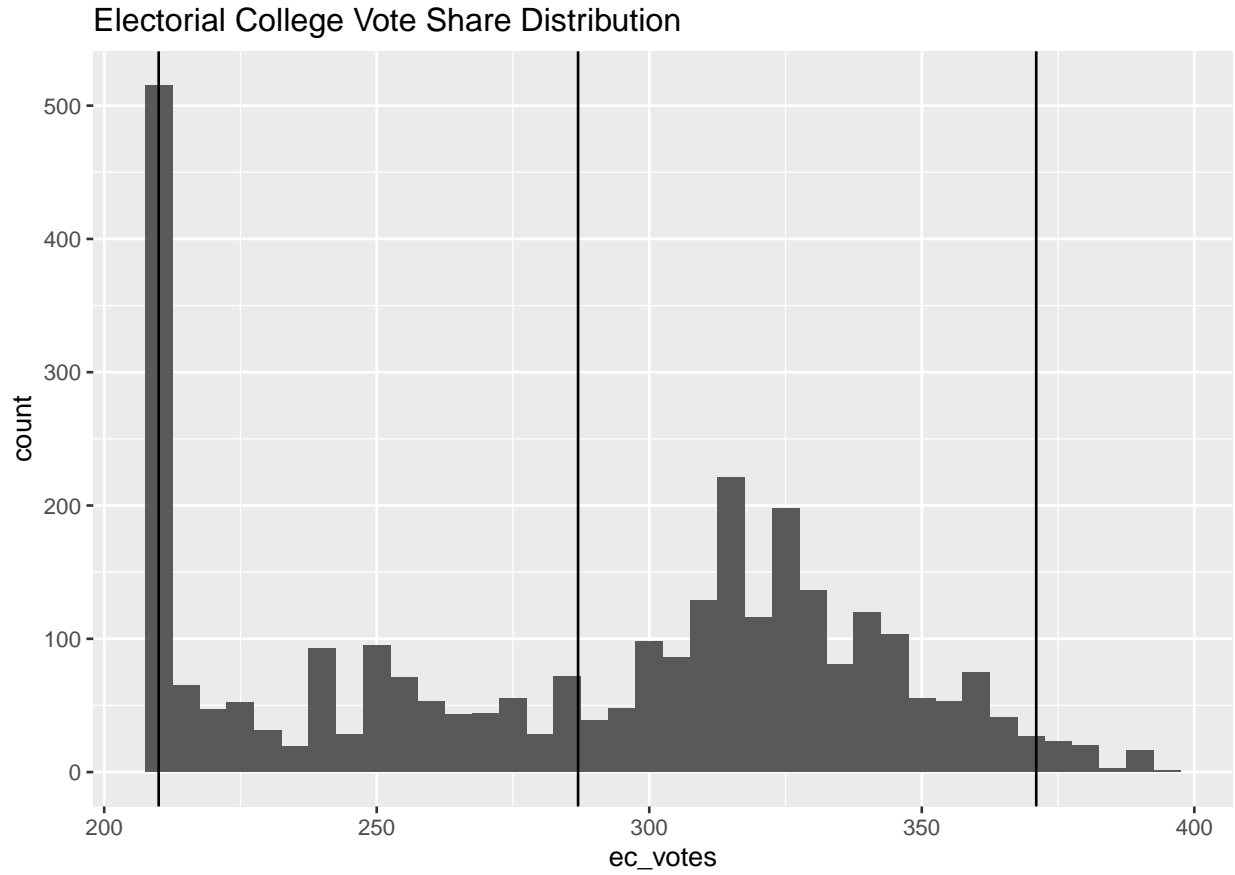


Table 1: EC Vote Total (Biden)

	x
2.5%	210
97.5%	371

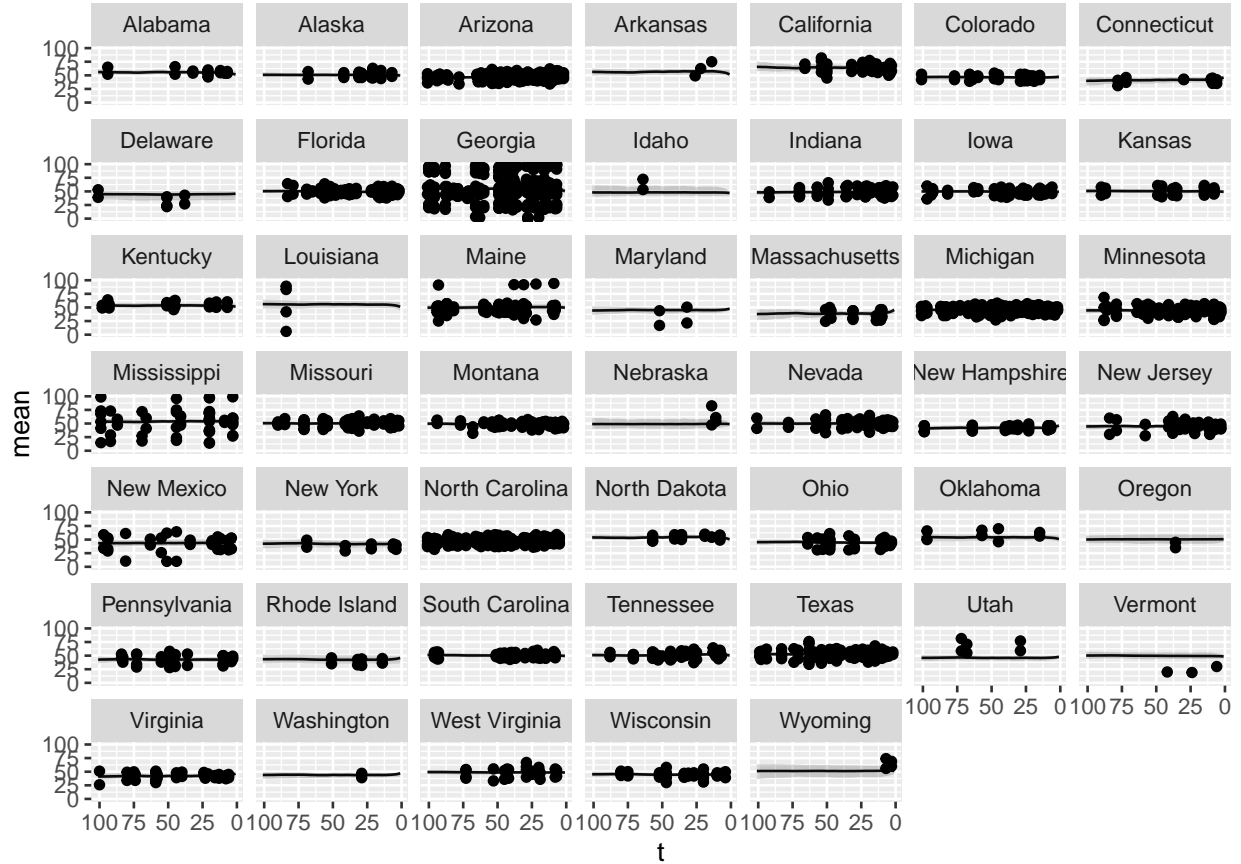
Table 2: Swing State Winning Probability

	Winning Probability
AZ	0.2977
NC	0.1497
MI	0.6717
WI	0.6453
MN	0.6170
TX	0.0187
FL	0.2180
PA	0.4550
GA	0.0990
OH	0.0583
NV	0.5280
NH	0.5627
IA	0.1180

Table 3: Swing State Share Percentage Interval Estimate (Biden)

	2.5%	50%	97.5%
AZ	51.50	53.44	55.56
NC	51.26	53.05	54.91
MI	52.12	54.97	58.27
WI	51.83	54.73	57.31
MN	51.58	54.66	57.79
TX	50.21	52.12	53.88
FL	51.16	53.07	55.49
PA	51.62	53.84	56.88
GA	50.90	52.70	54.81
OH	50.54	52.46	54.39
NV	51.53	54.12	56.85
NH	51.70	54.29	57.85
IA	50.81	52.62	55.02

Warning: Removed 4059 rows containing missing values (geom_point).



Warning: Removed 2 rows containing non-finite values (stat_density).



Table 4: Senator All State Winning Probability

Winning Probability	
Michigan	0.3887
Minnesota	0.3507
Arizona	0.8253
North Carolina	0.6490
Iowa	0.7870
Virginia	0.4643
Georgia	0.8623
Texas	0.9380
Alaska	0.8710
New Hampshire	0.3940
Alabama	0.9340
Kentucky	0.9607
Wyoming	0.8270
Montana	0.6840
Kansas	0.8210
Maine	0.8827
South Carolina	0.8317
Mississippi	0.9297
New Jersey	0.4300
Tennessee	0.8817
Nebraska	0.7777
Massachusetts	0.3837

	Winning Probability
Arkansas	0.9003
Oklahoma	0.8580
Colorado	0.4403
New Mexico	0.4257
West Virginia	0.7060
Oregon	0.7920
Delaware	0.4687
Idaho	0.4590
Louisiana	0.8573
Florida	0.9140
Nevada	0.8103
Indiana	0.7287
Missouri	0.9433
Pennsylvania	0.3850
California	0.9903
Ohio	0.4493
New York	0.3773
Wisconsin	0.4673
Connecticut	0.3823
Vermont	0.6643
North Dakota	0.8903
Rhode Island	0.4450
Washington	0.4180
Utah	0.5340
Maryland	0.5493

Table 5: Senator All State Vote Share Percentage Interval Estimate

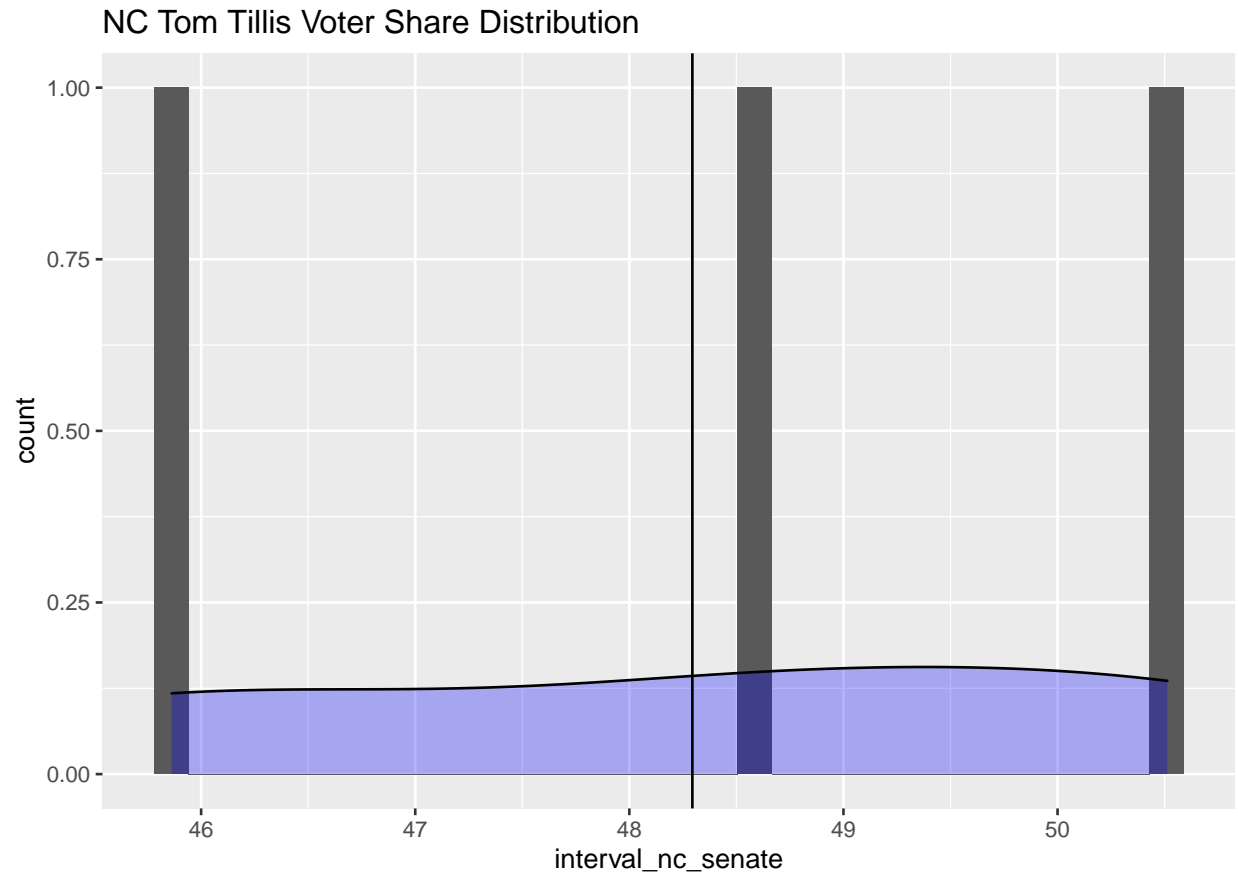
	2.5%	50%	97.5%
Michigan	44.32	47.37	49.92
Minnesota	40.97	46.50	49.90
Arizona	46.58	49.11	50.98
North Carolina	45.86	48.51	50.51
Iowa	45.57	49.19	52.89
Virginia	39.41	47.73	50.55
Georgia	46.01	49.70	55.67
Texas	47.45	50.27	54.77
Alaska	46.72	49.42	54.03
New Hampshire	40.60	47.05	50.29
Alabama	47.14	50.24	57.60
Kentucky	47.66	50.87	57.13
Wyoming	45.66	49.50	60.53
Montana	45.68	48.58	50.31
Kansas	46.12	49.22	52.67
Maine	46.66	49.44	54.15
South Carolina	46.11	49.33	52.64
Mississippi	47.07	49.82	56.90
New Jersey	43.39	47.60	50.42
Tennessee	46.19	49.70	55.02
Nebraska	45.40	49.15	52.26
Massachusetts	37.77	46.87	50.02

	2.5%	50%	97.5%
Arkansas	46.73	49.84	60.14
Oklahoma	45.92	49.72	57.52
Colorado	42.94	47.50	52.42
New Mexico	40.00	47.37	50.23
West Virginia	45.72	48.77	51.34
Oregon	45.69	49.40	58.71
Delaware	34.71	47.77	50.71
Idaho	42.09	47.67	50.51
Louisiana	44.74	49.68	59.66
Florida	47.25	49.42	51.94
Nevada	46.49	49.06	51.65
Indiana	45.15	48.80	51.29
Missouri	47.65	49.45	52.44
Pennsylvania	42.14	47.13	50.00
California	48.35	51.43	63.26
Ohio	43.42	47.66	51.20
New York	36.66	47.13	50.08
Wisconsin	44.08	47.76	50.60
Connecticut	38.90	46.95	50.41
Vermont	44.37	48.62	53.02
North Dakota	46.16	49.92	54.88
Rhode Island	35.65	47.30	50.79
Washington	41.02	47.50	50.20
Utah	43.74	48.17	50.89
Maryland	43.84	48.20	50.25

Table 6: North Carolina Vote Share Interval Estimate

	x
2.5%	45.86
50%	48.51
97.5%	50.51

```
## 'stat_bin()' using 'bins = 30'. Pick better value with 'binwidth'.
```



```
## [1] 0.649
```

```
## [1] 0.856
```

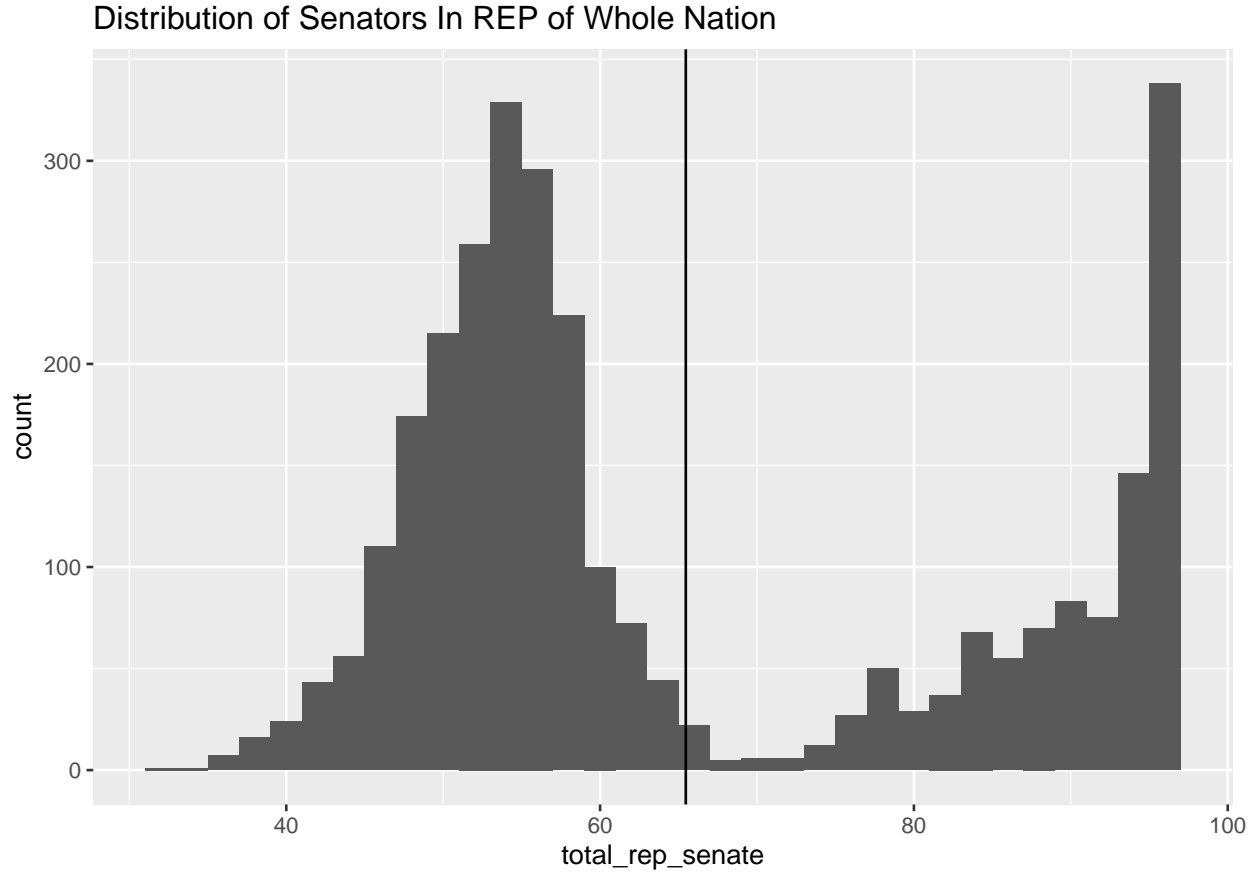
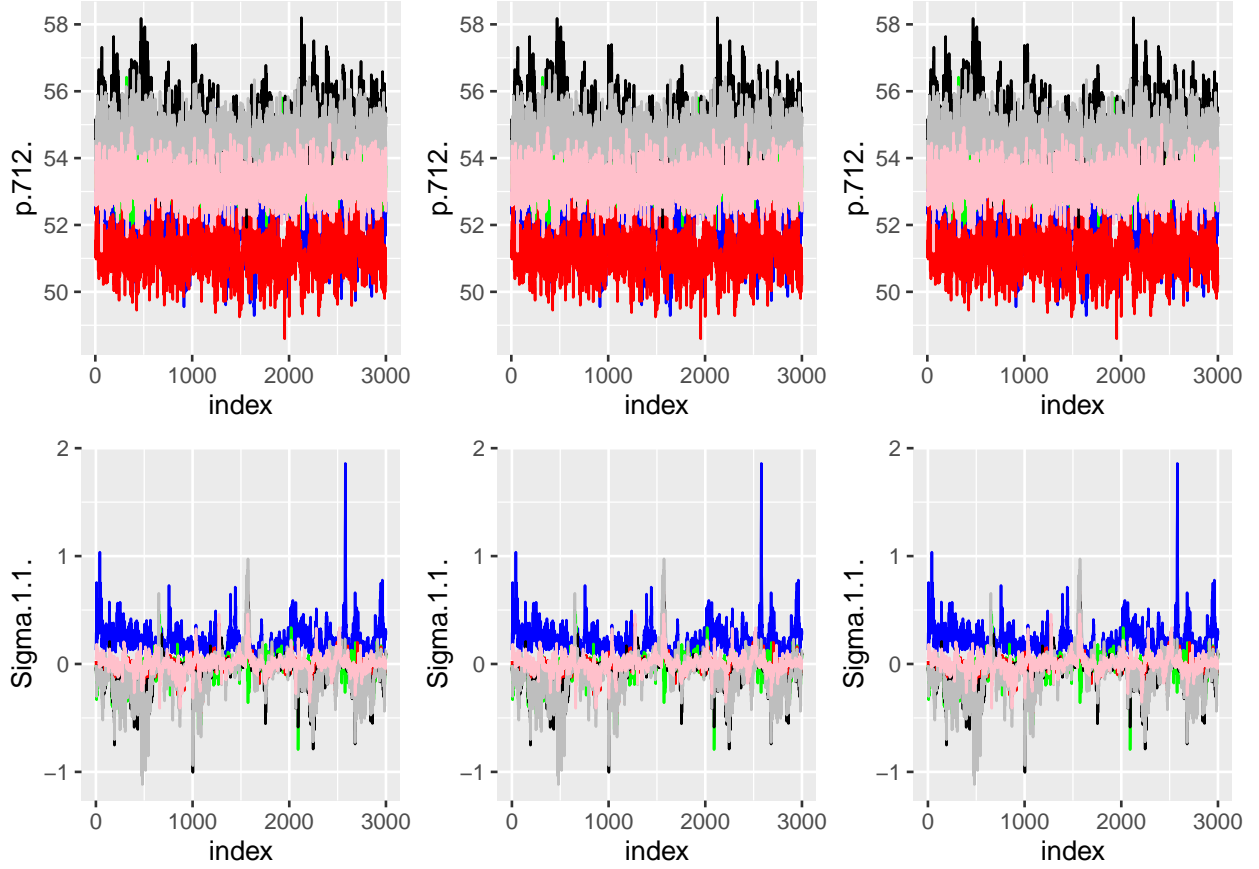


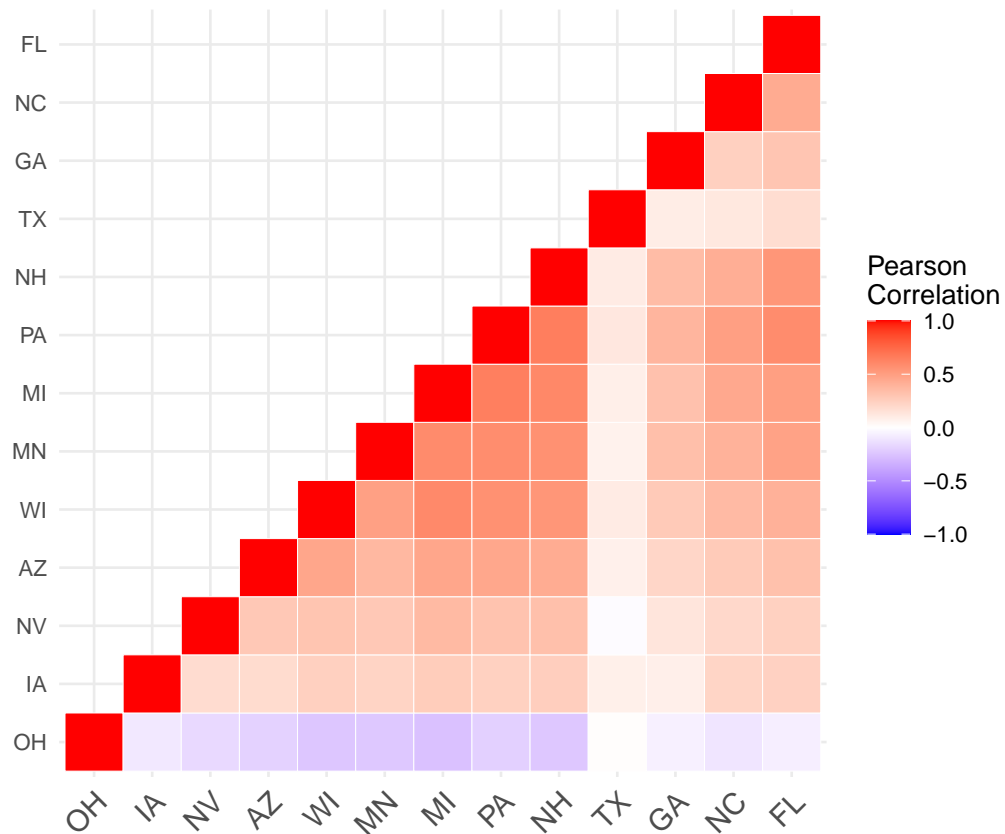
Table 7: House All Congressional District Winning Probability (REP)

Winning Probability	
District 9	0.4267
District 11	0.6737
District 8	0.4390
District 3	0.6677
District 13	0.6470
District 2	0.4227
District 7	0.6667
district 1	0.0000
district 4	0.0000
district 5	1.0000
district 6	0.0000
district 10	1.0000
district 12	0.0000

Table 8: house All State Vote Share Percentage Interval Estimate (REP)

	2.5%	50%	97.5%
District 9	45.99	49.63	54.85
District 11	41.05	55.87	65.48
District 8	31.81	48.59	56.99
District 3	44.35	53.36	63.10
District 13	32.75	52.93	63.68
District 2	43.77	49.34	55.57
District 7	29.28	61.71	68.49
district 1	26.19	26.85	27.68
district 4	21.92	22.79	23.66
district 5	68.32	68.56	68.79
district 6	45.59	45.72	45.86
district 10	69.83	70.02	70.20
district 12	36.13	36.27	36.41





```
mean(as.data.frame(jags_sims_mv$BUGSoutput$summary)$Rhat)
```

```
## [1] 1.042
```

```
mean(as.data.frame(jags_sims_senator$BUGSoutput$summary)$Rhat)
```

```
## [1] 1.409
```

```
mean(as.data.frame(jags_sims_house$BUGSoutput$summary)$Rhat)
```

```
## [1] 3.727
```