

MAT5314 Project 1: Data Visualization

Teng Li(7373086)
Shiya Gao(300381032)
Chuhan Yue(300376046)
Yang Lyu(8701121)

Introduction

A data set of the 2016 US election polls was given. In this project we aim to understand the data structure by creating various visualizations.

First, we would like to give a brief introduction to the U.S. election system, because it is crucial to understand the background of the data. Voters in each state vote to choose the President of the United States. The candidate who wins the majority of the votes will receive all the electoral votes in that state. Then the sum of the electoral votes in each state is calculated. The total number of electoral votes is 538. The candidate who wins half of the votes plus 1 will win and become the new President of the United States.

Secondly, we want to analyze the key factors for the candidate's victory. Since each state has a different number of electoral votes, it is crucial for electors to win in several key states. The reason is that if a candidate wins a certain state, he will win all the electoral votes in that state. So there will be tight competition in states with more votes.

The data set was published by FiveThirtyEight to illustrate the reliability and quality of each pollster to which a letter grade ranging from A+ to D- was given.

Method

We use various R packages to present the data set and to plot the graphs.

Result

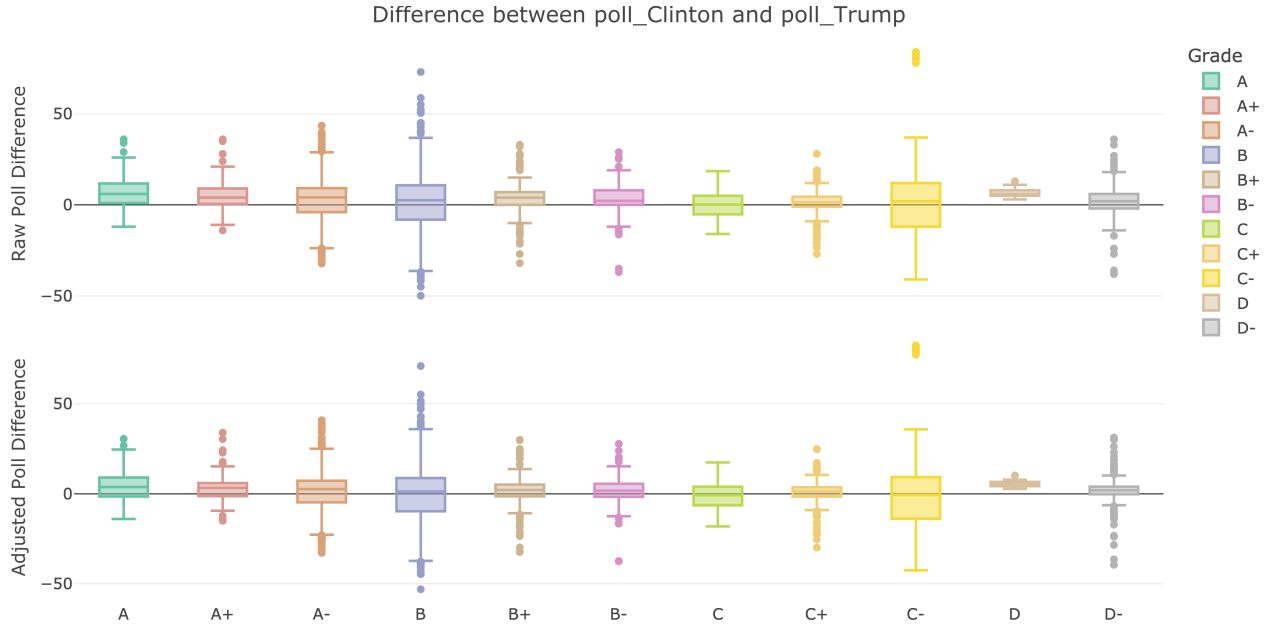
We first created a data variable definition table to give an initial understanding of the data. As one can see, there were a few variables with missing values:

Note that the poll results for Johnson and McMullin had lots of missing values. In particular, Johnson had 33.48% raw poll result and 33.48% adjusted poll result missing, and McMullin had 99.29% and 99.29% missing. Due to the fact that these two candidate didn't make to the final election, we chose to ignore their data in some of the analysis.

Since the final two candidates in Election 2016 are Clinton and Trump, we plotted box plots of the difference of their poll results, one for the raw poll and one for the adjusted poll. We saw that there's little difference between the distribution of the raw and the adjusted data. However, the mean of each grade of the adjusted poll result was a little closer to zero than that of the raw poll result. This indicated that the adjustment that FiveThirtyEight made was an improvement because the raw poll difference of each grade was mostly above zero, which clearly showed that the poll result was more in favour of Clinton yet Trump was the final winner of the election.

Table 1: Data Variable Definition

Variables	Size	Type	Example	Number.Unique	Number.Missing	Comment
state	4208	character	U.S., New Mexico, Virginia	57	0	The name of the state (or national) where the election is held
startdate	4208	character	2016-11-03, 2016-11-01, 2016-11-02	352	0	Start date of poll
enddate	4208	character	2016-11-06, 2016-11-07, 2016-11-05	345	0	End date of poll
pollster	4208	character	ABC News/Washington Post, Google Consumer Surveys, Ipsos	196	0	Organization name that conducts or analyzes opinion polls
grade	4208	character	A+, B, A-	11	429	Grade assigned by Fivethirtyeight to pollster
samplesize	4208	integer	2220, 26574, 2195	1767	1	Sample size of polls for each pollster
population	4208	character	lv, rv, a	4	0	Type of population being polled
rawpoll_clinton	4208	numeric	47, 38.03, 42	1312	0	Poll Percentage for Hillary Clinton
rawpoll_trump	4208	numeric	43, 35.69, 39	1385	0	Poll Percentage for Donald Trump
rawpoll_johnson	4208	numeric	4, 5.46, 6	585	1409	Poll Percentage for Gary Johnson
rawpoll_mcmullin	4208	numeric	NA, 24, 27.6	17	4178	Poll Percentage for Evan McMullin
adjpoll_clinton	4208	numeric	45.20163, 43.34557, 42.02638	4200	0	Adjusted percentage for Hillary Clinton
adjpoll_trump	4208	numeric	41.7243, 41.21439, 38.8162	4204	0	Adjusted percentage for Donald Trump
adjpoll_johnson	4208	numeric	4.626221, 5.175792, 6.844734	2211	1409	Adjusted percentage for Gary Johnson
adjpoll_mcmullin	4208	numeric	NA, 24, 27.70142	31	4178	Adjusted percentage for Evan McMullin

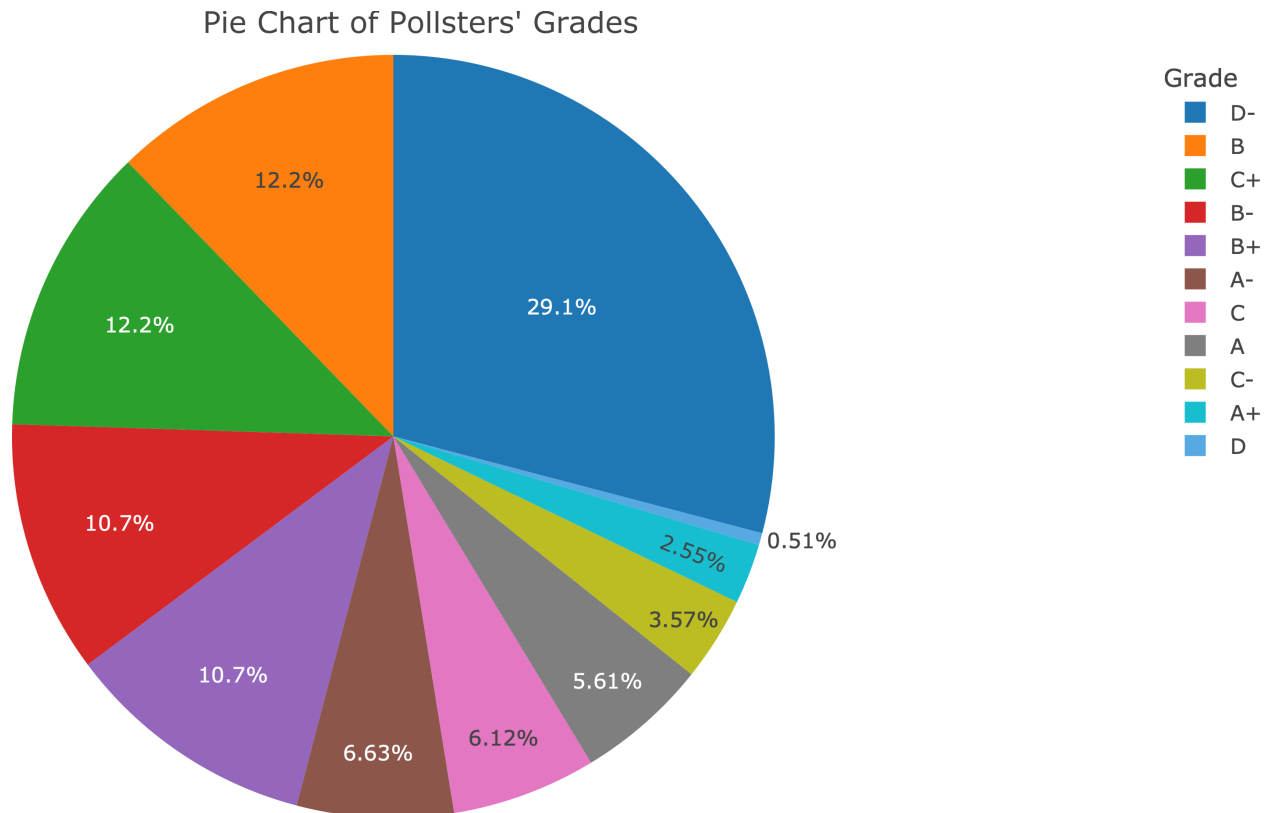


In the second step, we want to look at the percentage of different grades of pollsters. Therefore, we extract the “pollster” and “grade” columns from the original data frame, merge the duplicate rows and re-sort the data according to the grade.

```
## # A tibble: 196 x 2
##   grade pollster
##   <chr> <chr>
## 1 A     Behavior Research Center (Rocky Mountain)
## 2 A     Fairleigh Dickinson University (PublicMind)
## 3 A     Fox News/Anderson Robbins Research/Shaw & Company Research
## 4 A     Marist College
## 5 A     Marquette University
```

```
## 6 A      Muhlenberg College
## 7 A      National Journal
## 8 A      Public Policy Institute of California
## 9 A      Research & Polling, Inc.
## 10 A     Siena College
## # i 186 more rows
```

From the “Pie Chart of Grades” below, we could see that the pollsters with B grades account for nearly 50% of the total, C and A grades account for about 30% and 25% respectively. Pollsters rated D only account for less than 1%. Furthermore, grade B and C+ have the largest proportions, at about 17.3%, followed by B+ and B-, both at about 15.1%.



Polls trend by time

Using the end date as the standard, we drew a scatter plot of the changes in the support rates of the four candidates over time.

Figure 1 indicates that Clinton and Trump’s approval ratings are significantly higher than those of Johnson and McMullin at the overall level. Another interesting demonstration is the divergence of support as the vote draws to a close. This shows that as the voting deadline approaches, people’s intentions are more inclined to one of Clinton or Trump. Voting is also more polarized. A low poll rate for one candidate also implies that other candidates may have received high poll rates. This creates differences in pollsters or states.

Comparison of candidates’ polls in each state

We want to process the metadata by counting the polls received by each of the four candidates in each state. This result is easier to obtain by multiplying the given size and the proportion. Regardless of the various pollsters, we combine the number of polls for each candidate received in each state although the polls may come from different pollsters. In this case, we treat NA as 0.

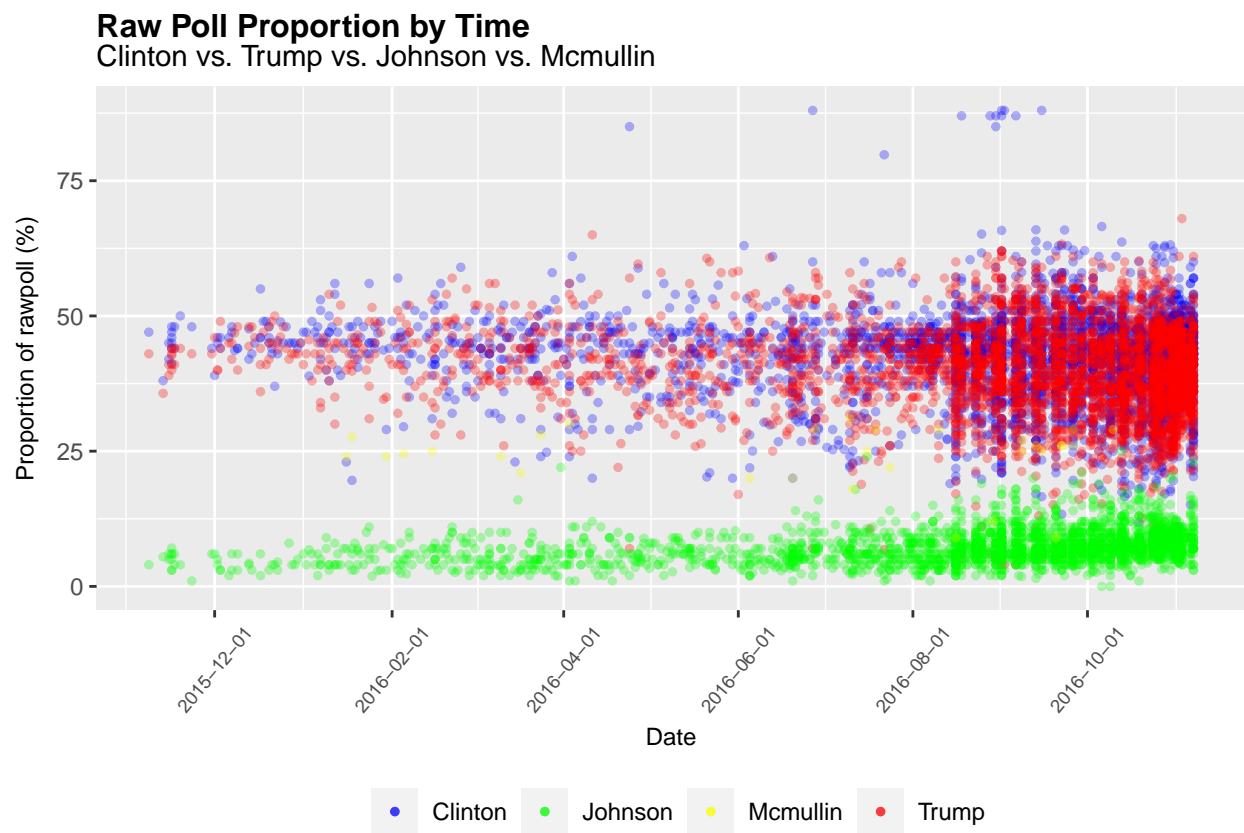


Figure 1: Raw Poll Proportion by Time.

The visualization of the poll proportion of the four candidates in each state as follows.

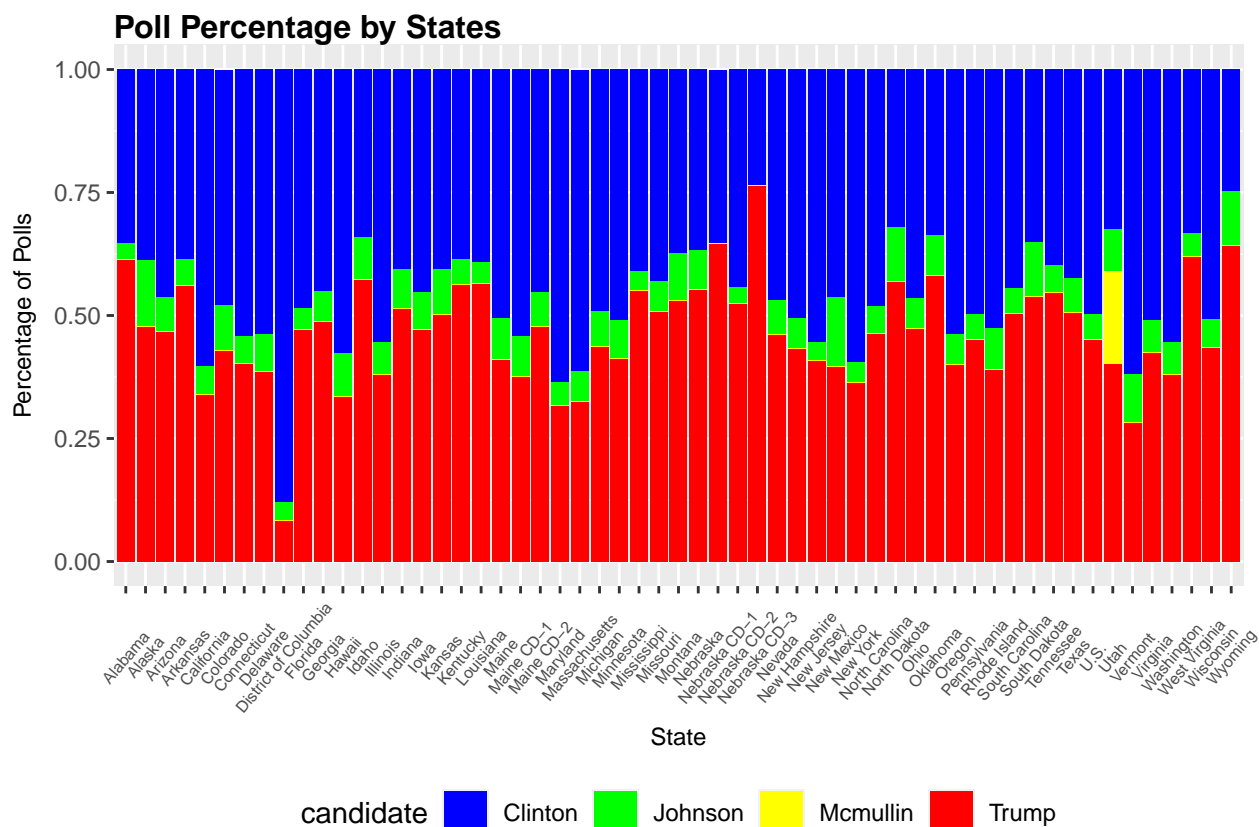


Figure 2: Poll Percentage by States.

Figure 2 shows the poll proportions of the four candidates in each state distinguished by colors. We can clearly observe which candidate is likely to win all the electoral votes in each state, which is helpful in estimating the outcome of the presidential election. From the chart above, we conclude that Clinton and Trump are clearly ahead of Johnson and McMullin. Furthermore, Clinton is clearly ahead of Trump in California, DC, Hawaii, Illinois, Maryland, Massachusetts, New Jersey, New York, Oregon, Rhode Island, Vermont, and Washington. Another side, in Alabama, Alaska, Idaho, Indiana, Iowa, Kansas, Kentucky, Louisiana, Mississippi, Missouri, Montana, Nebraska, North Dakota, Ohio, Oklahoma, South Carolina, South Dakota, Tennessee, Texas, Utah, West Virginia, Wyoming, Trump clearly leads Clinton. This helps us tally who would win all of the state's electoral votes.

Discussion

Conclusion