

Predicting the 2024 US Presidential Election: A Polling-Based Forecast*

Trump will get support about 47.55% by pollster and may lose the election

Tianrui Fu Yiyue Deng Jianing Li

November 3, 2024

This paper uses polling data to predict the outcome of the 2024 U.S. presidential election. By applying generalized linear regression and Bayesian models, we analyze the factors related to the polls, including polling organizations and the ratings associated with those organizations, on Donald Trump’s support rate. The research results indicate a significant relationship between the variables, with Donald Trump’s support rate showing an upward trend. Additionally, we further discuss the advantages and disadvantages of the two models, as well as potential methods for improvement.

Table of contents

1	Introduction	1
2	Data	2
2.1	Overview	2
2.2	Cleaning Data	2
2.3	Measurement	2
2.4	Outcome variables	3
2.5	Predictor variables	5
2.5.1	Pollster	5
2.5.2	Poll Score	6
2.5.3	Transparency Score	6
2.5.4	Numeric Grade	7
2.5.5	End date	8
2.5.6	Relationship between pollster, pollscore and end date	8

*Code and data are available at: <https://github.com/FrankFU323/U.S.-presidential-election.git>

2.5.7	Relationship between transparency score and numeric grade	9
3	Model	10
3.1	Model Selection	10
3.2	Model set-up	10
3.2.1	Model 1 – GLM	11
3.2.2	Model 2 – Bayesian model for Trump	11
3.3	Model justification	11
4	Results	12
4.1	Result of model for analysis data	12
4.2	Result for model prediction	14
5	Discussion	16
5.1	First discussion point	16
5.2	Second discussion point	16
5.3	Third discussion point	17
5.4	Weaknesses and next steps	17
A	Appendix	18
A.1	Methodology Overview and Evaluation of YouGov Polling	18
A.2	Idealized Survey Methodology	19
A.3	Clean data	22
A.4	Variable details	23
A.5	Results of model	23
A.5.1	Table for summary of model results	23
A.5.2	Bayesian model for Harris data set	25
A.5.3	Predictions for both Trump and Harris	26
B	Model details	28
	References	29

1 Introduction

Overview paragraph

Estimand paragraph

Results paragraph

Why it matters paragraph

Telegraphing paragraph: The remainder of this paper is structured as follows. Section [A](#)....

2 Data

2.1 Overview

We use the statistical programming language R Core Team (2023), the template from Alexander (2023) completed by following packages tidyverse Wickham et al. (2019), dplyr Wickham et al. (2023), rstanarm Goodrich et al. (2022), arrow Richardson et al. (2024), modelr Wickham (2023), modelsummary Arel-Bundock (2022), ggplot2 Wickham (2016), here Müller (2020), kableExtra Zhu (2024) and knitr Xie (2023) to complete this analysis. Our data about the latest polling outcomes is from the website FiveThirtyEight (2024). It has 52 variables and 15891 observations in this data set, including pollster, poll score, etc.

Table 1: The example of chosen variable dataset

pollster	pollscore	numeric_grade	transparency_score	end_date	pct
Marist	-0.9	2.9	7	2024-08-04	48.0
Emerson	-1.1	2.9	7	2024-07-23	47.5
Beacon/Shaw	-1.1	2.8	9	2024-09-24	46.0
AtlasIntel	-0.8	2.7	6	2024-10-17	49.9
Quinnipiac	-0.5	2.8	9	2024-09-22	48.0
Marquette Law School	-1.1	3.0	10	2024-08-01	47.0
YouGov	-1.1	3.0	9	2024-08-31	49.0
Emerson	-1.1	2.9	7	2024-09-28	49.9
Marist	-0.9	2.9	7	2024-09-05	48.0
MassINC Polling Group	-0.8	2.8	7	2024-09-18	40.0

2.2 Cleaning Data

We have cleaned the data set and the detailed procedure see from the appendix Section [A.3](#).

2.3 Measurement

In the realm of political polling, measurement is critical as it transforms abstract voter sentiments into quantifiable estimates that can influence electoral outcomes. The fundamental challenge lies in converting individual opinions—such as a voter’s intention to support a particular candidate—into structured data that can effectively forecast the number of electoral college votes that candidate might secure. This transformation begins with the design of the polling methodology, where specific phenomena in the real world, such as voter preferences, are captured through surveys.

This dataset originates from FiveThirtyEight (2024), a trusted source known for its rigorous standards in polling data collection, ensuring the widest possible coverage of voters and comprehensive collection and disclosure of all relevant information. This includes details like the pollster’s name, identification number, survey dates, and associated trust levels. In our analysis dataset, which consists of 492 observations related to Donald Trump’s polling outcomes, we have identified 38 relevant variables that play a significant role in understanding electoral support. Key among these are the numeric grade, transparency score, poll score, pollster, end date, and percentage of support for Trump. Each variable is meticulously measured to ensure that they accurately reflect the sentiments of likely voters.

The dataset presents results as the proportion of votes each candidate receives according to each pollster, with this approach adjusting over time to account for changes in voter sentiment, which can shift rapidly due to the approach of Election Day, candidate speeches, and other factors. By applying these measurement principles, we systematically convert nuanced public opinion into a structured dataset that not only captures voter intentions but also provides a reliable foundation for predictive modeling.

Ultimately, this rigorous measurement approach allows us to build a Bayesian model to assess how different pollsters might influence the final outcome of the U.S. presidential election, thus bridging the gap between individual voter opinions and electoral predictions.

The explain of each used variable please see in the Section [A.4](#).

2.4 Outcome variables

The Figure [1](#) shows the distribution of polling data over time, with each color representing a different polling organization. We can see that as time progresses from August to October, there is a steady increase in the amount of polling data, peaking around early October. This suggests that multiple organizations have contributed to polling data on a consistent basis over this period.

In the Figure [2](#), it reveals a central peak around the 45-50% range, indicating that most of the values for this variable are concentrated in this area. The shape of the plot suggests a right-skewed distribution, with relatively few data points falling below 40% or above 50%. This gives an overview of the central tendency and variability of the polling percentages.

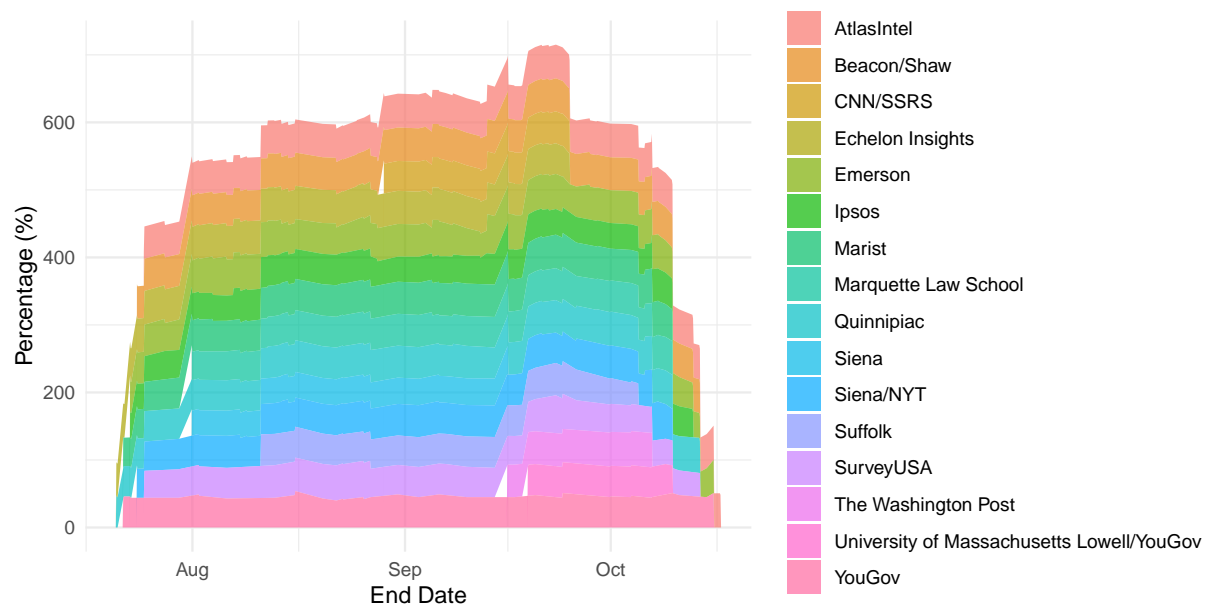


Figure 1: Percentage of Each Pollster Over Time

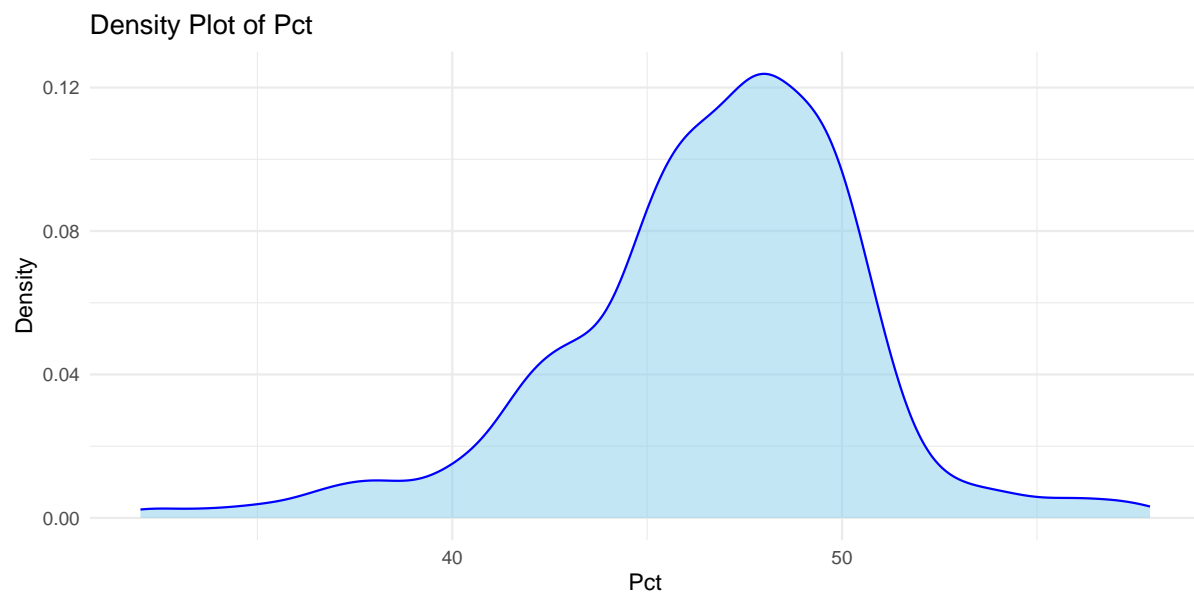


Figure 2: Density plot of pct

2.5 Predictor variables

2.5.1 Pollster

In Figure 3, the predictor variable is the count of polls conducted by each pollster. The bar chart displays the number of polls released by different polling organizations, ordered from highest to lowest frequency. The bars range in color from light blue to dark blue, indicating the frequency of polls conducted by each organization. Siena/NYT and YouGov are the most active pollsters, with 94 and 59 polls conducted, respectively. Emerson and Beacon/Shaw also have high polling frequencies, with 58 and 46 polls. In contrast, several pollsters, like Christopher Newport University and Data Orbital, conducted only one poll, indicating their minimal activity in comparison.

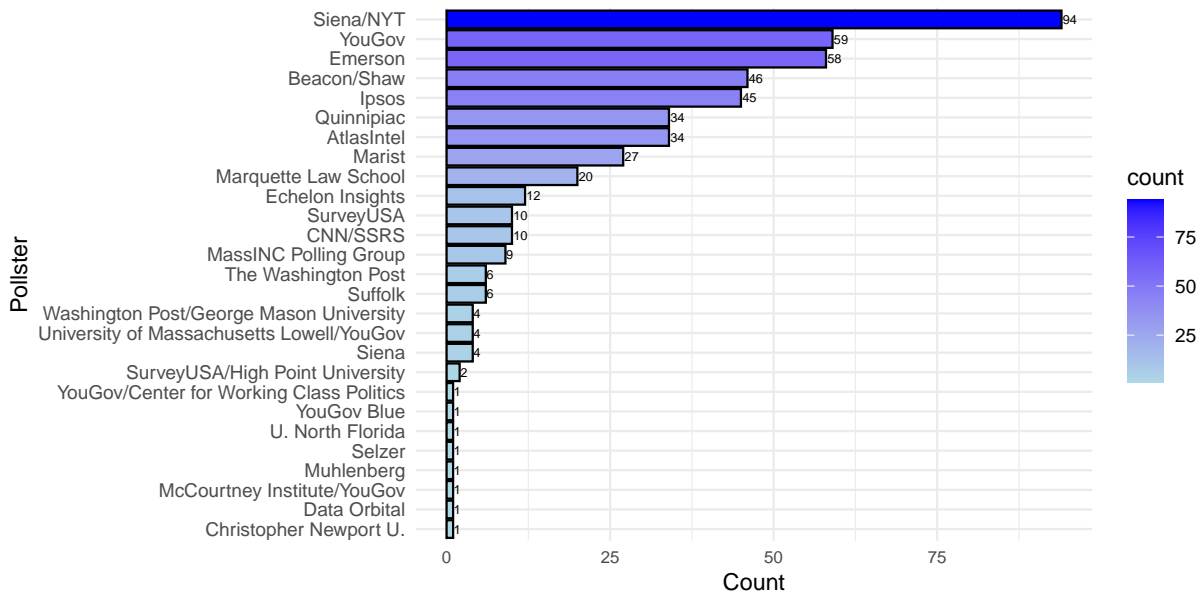


Figure 3: Frequency of each pollster

2.5.2 Poll Score

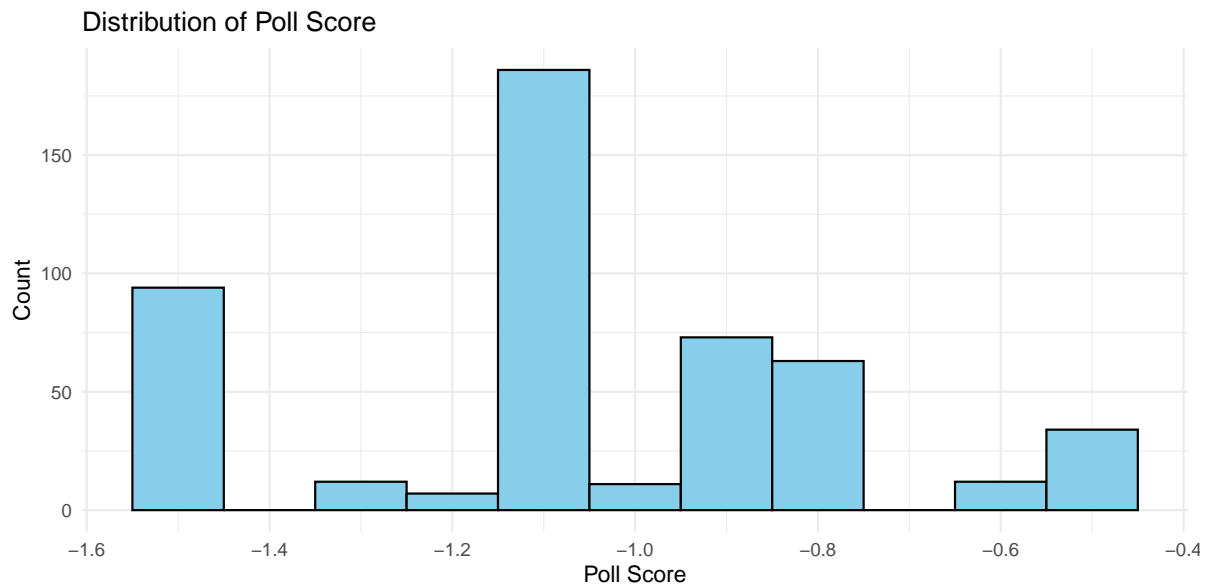


Figure 4: Distribution of Poll Score

2.5.3 Transparency Score

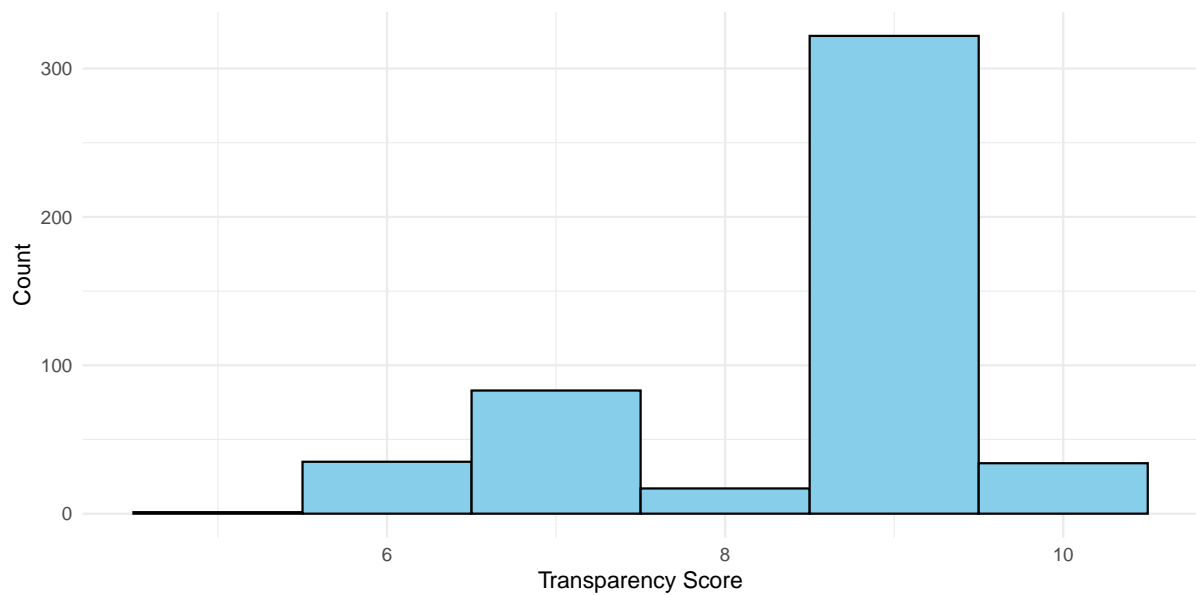


Figure 5: Distribution of Transparency Score

2.5.4 Numeric Grade

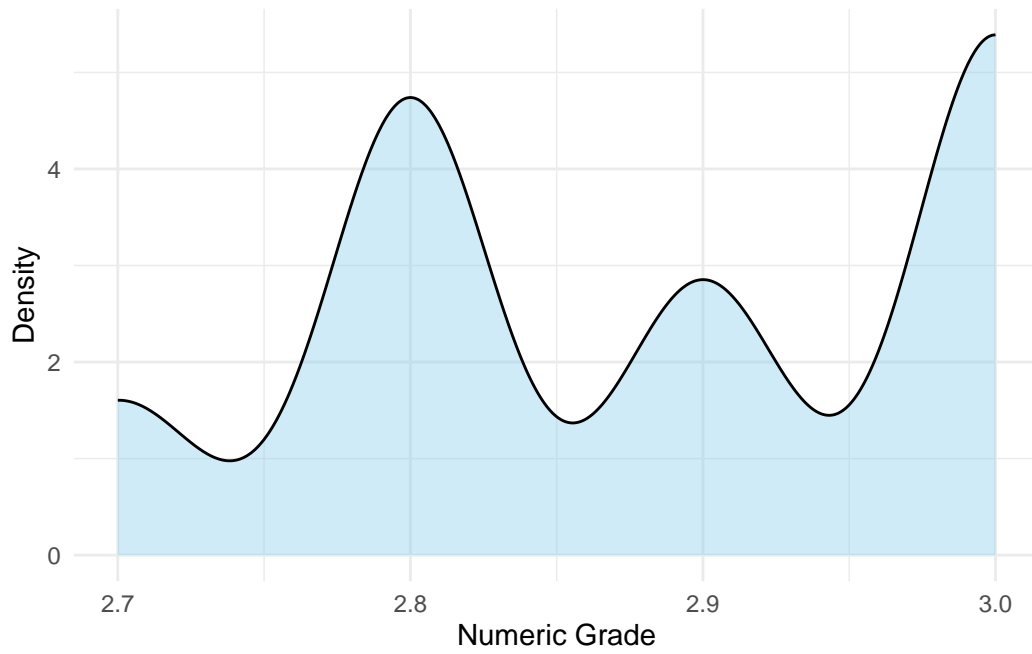


Figure 6: Density Plot of Numeric Grade

2.5.5 End date

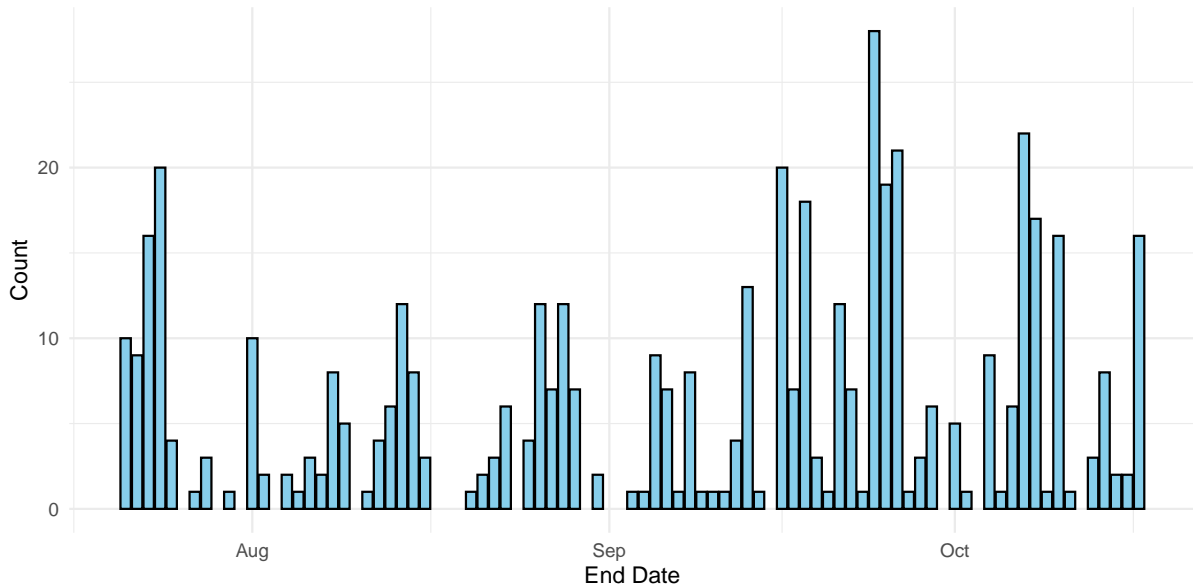


Figure 7: Frequency of Records by End Date

2.5.6 Relationship between pollster, pollscore and end date

In Figure 8, the predictor variable is the average poll score for each pollster over time, from August to October. This heat map visualizes the poll scores of each organization on different dates, with color intensity indicating the score levels—darker colors represent lower scores. The pollsters with the most consistently low average poll scores (darker colors) include Selzer, Siena/NYT, and Marquette Law School. These pollsters frequently show results lower than others over time, suggesting they may consistently lean toward one direction in their results. In contrast, pollsters like YouGov Center for Working Class Politics and YouGov Blue display lighter colors, indicating relatively higher scores across their polls.

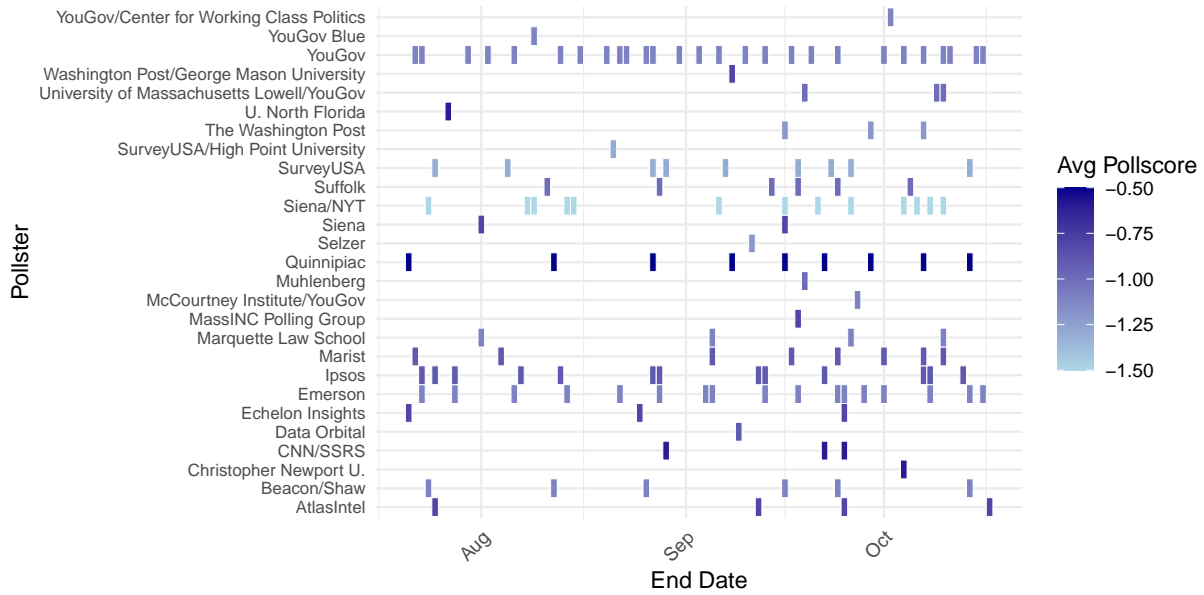


Figure 8: Average Pollscore by Pollster and End Date

2.5.7 Relationship between transparency score and numeric grade

In Figure 9, there is a clear positive relationship between transparency score and numeric grade. Pollsters with higher transparency scores, such as those scoring around 10, tend to achieve higher numeric grades, close to 3.0. Conversely, pollsters with lower transparency scores tend to have lower numeric grades, closer to 2.7. This pattern suggests that greater transparency is associated with better overall pollster ratings.

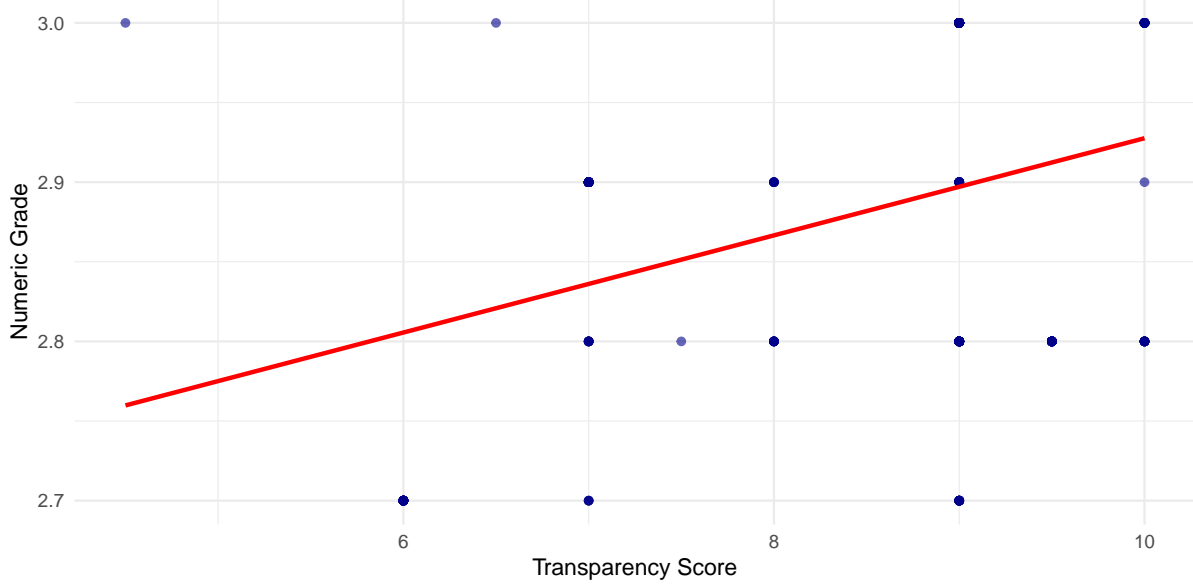


Figure 9: Scatter Plot of transparency score and numeric grade

3 Model

3.1 Model Selection

Our model aims to predict the changes in Donald Trump’s probability of winning the 2024 U.S. election over time while examining the influence of key factors, such as pollster reliability scores, transparency, and specific pollster effects, on polling results. For model selection, we initially considered a generalized linear regression to provide a straightforward analysis of the relationships between these factors. However, given the limitations of generalized linear regression in handling uncertainty and dynamic data, we ultimately chose a Bayesian model. The Bayesian approach enables the integration of prior information and dynamically updates predictions, offering greater stability and accuracy under high uncertainty. Below is a brief overview of our model.

Background details and diagnostics are included in [Appendix B](#).

3.2 Model set-up

In this Bayesian framework, we assume a normal distribution of poll results around a mean affected by key predictors: numeric grade, transparency score, and pollscore. Define y_i as the percentage of Donald Trump. Then β_i represents the numeric grade, γ_i represents the transparency score, and δ_i represents the pollscore.

3.2.1 Model 1 – GLM

$$y_i = \alpha + \beta_1 \cdot \text{numeric_grade}_i + \beta_2 \cdot \text{transparency_score}_i + \beta_3 \cdot \text{pollscore}_i \quad (1)$$

$$+ \sum_{j=1}^N \gamma_j \cdot \text{pollster}_j + \delta \cdot \text{end_date}_i + \epsilon_i \quad (2)$$

where: - α is the intercept, representing the average poll level. - β_1 , β_2 , and β_3 are the regression coefficients for numeric grade, transparency score, and pollscore, respectively. - γ_j indicates the fixed effect of each pollster. - δ is the regression coefficient for the end date. - ϵ_i is the error term.

3.2.2 Model 2 – Bayesian model for Trump

$$y_i | \mu_i, \sigma \sim \text{Normal}(\mu_i, \sigma) \quad (3)$$

$$\mu_i = \alpha + \beta_1 \cdot \text{numeric_grade}_i + \beta_2 \cdot \text{transparency_score}_i + \beta_3 \cdot \text{pollscore}_i \quad (4)$$

$$+ \sum_{j=1}^N \gamma_j \cdot \text{pollster}_j + \delta_i \cdot \text{end_date}_i \quad (5)$$

$$\alpha \sim \text{Normal}(50, 10) \quad (6)$$

$$\beta_1, \beta_2, \beta_3 \sim \text{Normal}(0, 5) \quad (7)$$

$$\gamma_j \sim \text{Normal}(0, 5) \quad (8)$$

$$\sigma \sim \text{Exponential}(1) \quad (9)$$

where: - α represents an average poll result. - β_1 , β_2 , and β_3 capture the unique effect of their respective predictors on the poll percentage. - γ_j denotes a specific pollster effect. - δ_i accounts for time trends.

We run the model in R Core Team (2023) using the package of rstanarm Goodrich et al. (2022).

3.3 Model justification

The choice of a Bayesian model to analyze Donald Trump’s voting support rate stems from the need to effectively integrate uncertainty and leverage prior knowledge in estimating outcomes. Given that the dataset includes key predictors related to polling organizations—such as numeric_grade, transparency_score, pollster, poll_score, end_date, and pct—this approach allows for a comprehensive understanding of how these variables influence voting results.

The `numeric_grade` and `transparency_score` are crucial in assessing the reliability of polling organizations. By integrating these factors, the Bayesian model can provide a nuanced analysis of how the quality and transparency of polling organizations impact Trump’s support rate. This is particularly important in the context of public opinion polling, as perceptions of reliability can significantly affect the interpretation of results.

Moreover, the Bayesian framework enables the incorporation of prior distributions for model parameters, reflecting beliefs about their possible values before observing the data. This is especially beneficial in a domain where historical data and expert opinion can provide valuable insights.

Additionally, by including the `pollster` variable, the model accounts for the fixed effects of different polling organizations, allowing for a more targeted analysis that acknowledges the unique characteristics of each organization. The `end_date` variable helps to consider time trends, ensuring that the analysis remains relevant to the evolving political landscape, especially as elections approach.

However, there are areas for improvement within the model. First, refining the selection of variables related to voting support could enhance the model’s predictive power. For instance, incorporating socio-economic factors, demographic data, or shifts in voter behavior may provide a more comprehensive perspective. Additionally, the model’s predictive performance could be assessed through cross-validation or other model evaluation techniques to ensure the robustness and reliability of the results.

By employing this Bayesian model, our aim is to provide robust estimates of Trump’s support rate, considering both the statistical characteristics of the data and the inherent uncertainty related to polling, while exploring avenues for further model improvement.

4 Results

The results part is combined with the result for models of analysis data and the result of predictions by Bayesian model on November 5, 2024. In order to get to know whether the votes of Trump will win the election, we combined the data of Harris who have the most competitive candidate other than Trump to do predictions. It would help us be faster and more accurate to do the predict judgment.

4.1 Result of model for analysis data

In Figure 10, we applied a generalized linear regression model to predict the percentage of polls Trump according each pollster. Specifically, the support rate data points range from 40% to 55%, with most clustered between 45% and 50%. The orange trend line shows a subtle upward trend, indicating a gradual increase in Trump’s average support rate over time. For example, in early August, the support rate is around 45%, while by mid-October, the average

support rate shown by the trend line is close to 50%. This suggests that Trump’s support rate has risen by about 5% over these three months. For the model summary of Generalized Linear Regression and Bayesian Model, the table can see at Section A.5.

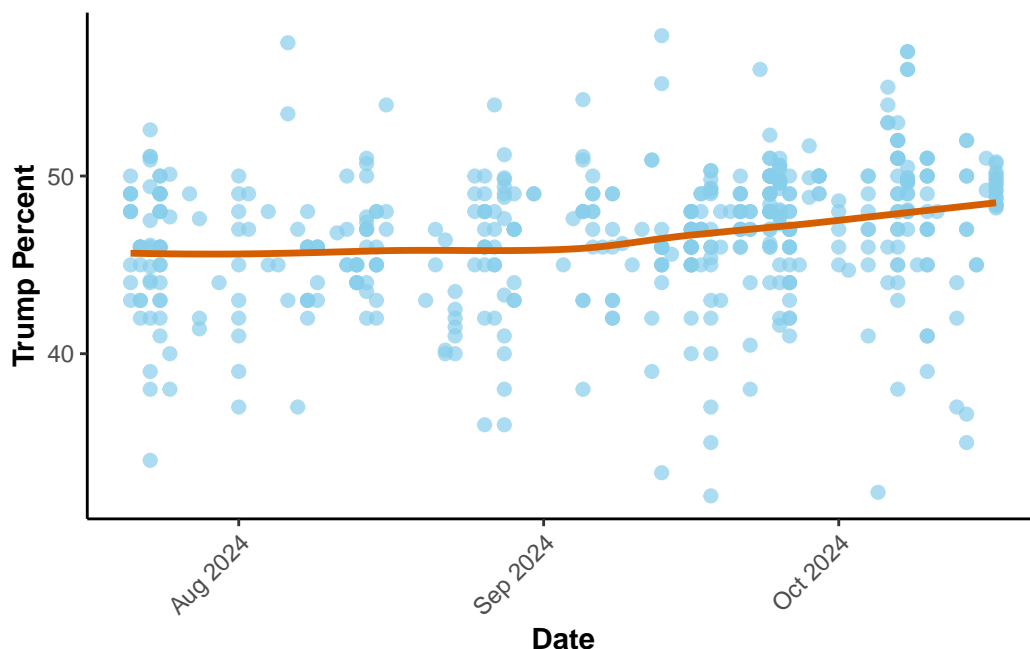


Figure 10: Donald Trump Support Over Time by Generalized Linear Regression

In Figure 11, we applied a bayesian model to predict the percentage of polls Trump according each pollster. Each data point is color-coded by pollster, with support rates ranging from 40% to 55%. The shaded gray area around the trend line represents the confidence interval, indicating the uncertainty in support rate variations. The overall trend line (blue) shows a slight upward trend, increasing from about 45% in August to nearly 50% in October. In particular, data points in early August are more dispersed, with some pollsters reporting low support rates (around 40%) and others reporting higher rates (over 50%). By October, most poll results have converged, with support rates centered between 45% and 50%. Notably, certain pollsters, like YouGov and Marist, provide either consistently higher or lower support rate estimates compared to others, indicating some bias or variation among polling organizations.

We also applied a Bayesian model to predict the percentage of polls Harris according each pollster. The figure can see it in Section A.5.

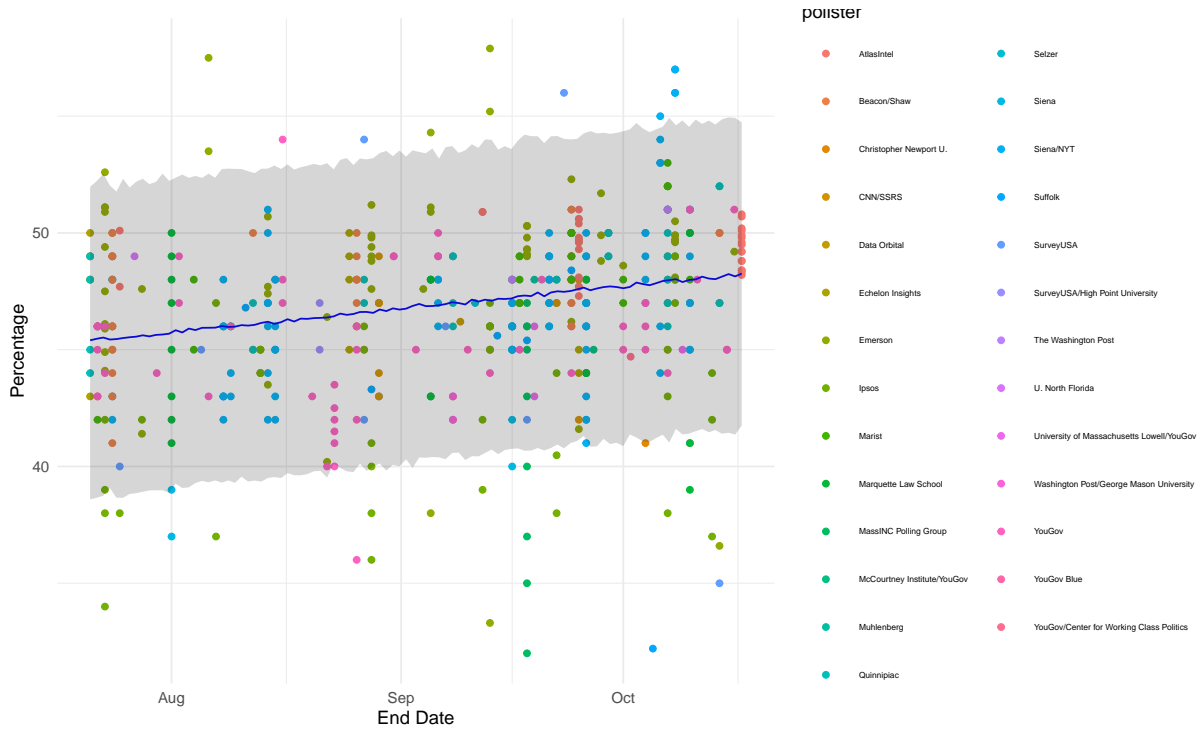


Figure 11: Poll Percentage over Time with Bayesian Fit for Trump

4.2 Result for model prediction

Figure 12 compares the predicted mean support rates for Trump and Harris across different polling organizations in 5th November of 2024. The vertical axis shows the predicted support rate (in percentage), and the horizontal axis lists the various pollsters. The red line represents Harris's predicted support rate, while the blue line represents Trump's.

For example, in predictions from Christopher Newport U. and MassINC Polling Group, Harris has a noticeably higher average support rate than Trump, exceeding 50%, while Trump's support rate is below 50%. In Siena's predictions, Trump's support rate is slightly lower than Harris's, reaching around 41%, while Harris's support is close to 53%. In Christopher Newport U.'s, MassINC Polling Group's, Siena's and Washington Post/George Mason University's predictions, Trump's support rate is all lower than Harris's and have a big gap.

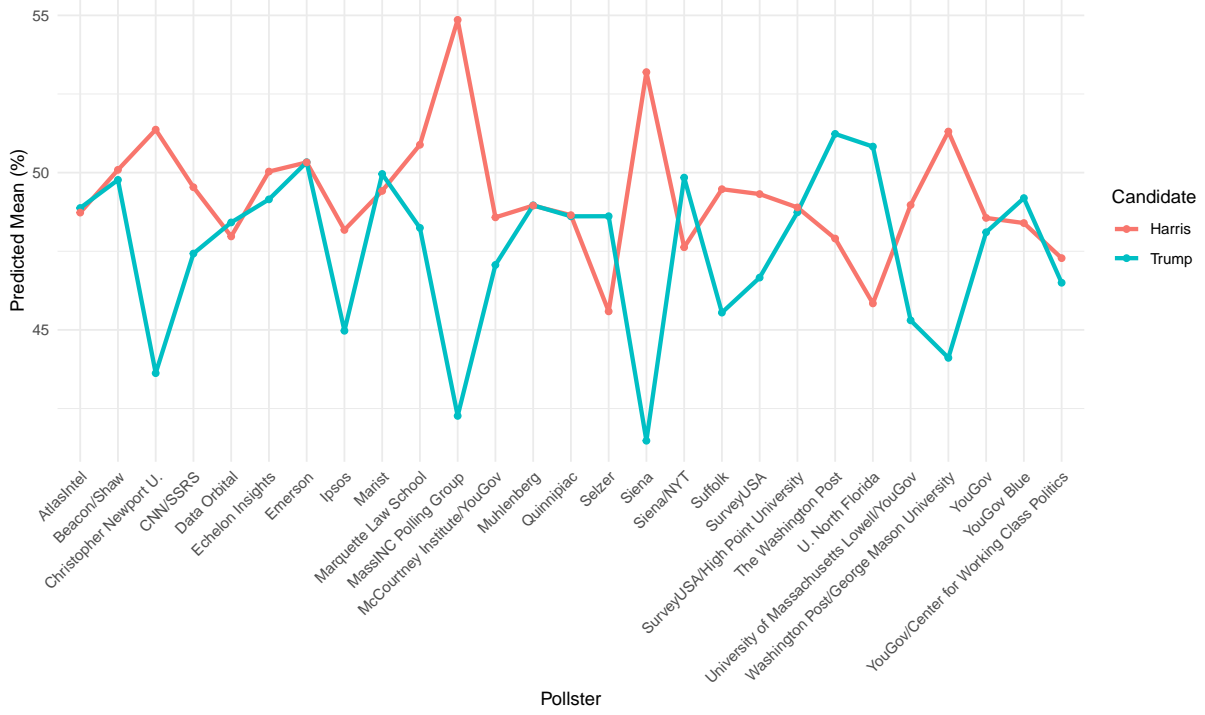


Figure 12: Predicted Mean Percentage by Pollster for Trump and Harris

In Table 2, we merged and compared the Bayesian model predictions of the vote shares for Trump and Harris on November 5, 2024, by pollster. The green background color represents the same pollster, indicating that the side with the larger predicted vote share is expected to win based on the data collected and summarized by that pollster. Conversely, red indicates the losing side. In Table 2, a total of 27 pollsters participated in the polling, with 11 pollsters predicting a Trump victory and 16 pollsters predicting a Harris victory. In Table 3, we calculated the mean support for Trump and Harris, estimating Trump’s probability of winning to be 41%. For the further more predictions about the Bayesian model, we can see at Section A.5.

Table 2: Predictions for both Trump and Harris by pollster

Predictions Summary			
	pollster	pred_mean.Trump	pred_mean.Harris
2.5%	AtlasIntel	48.87532	48.72877
2.5%1	Emerson	50.33160	50.32562
2.5%2	YouGov	48.10369	48.56045
2.5%3	Beacon/Shaw	49.76953	50.08884
2.5%4	Quinnipiac	48.61109	48.64800
2.5%5	SurveyUSA	46.66310	49.31814
2.5%6	Ipsos	44.97272	48.17852
2.5%7	Marist	49.95707	49.41884
2.5%8	Siena/NYT	49.84233	47.62689

2.5%9	University of Massachusetts Lowell/YouGov	45.30381	48.97090
2.5%10	Marquette Law School	48.24373	50.88540
2.5%11	The Washington Post	51.23118	47.90387
2.5%12	Suffolk	45.54956	49.47362
2.5%13	Christopher Newport U.	43.62446	51.36701
2.5%14	YouGov/Center for Working Class Politics	46.49916	47.28147
2.5%15	McCourtney Institute/YouGov	47.06573	48.57826
2.5%16	Echelon Insights	49.14967	50.03402
2.5%17	CNN/SSRS	47.42745	49.53568
2.5%18	Muhlenberg	48.96127	48.95174
2.5%19	MassINC Polling Group	42.26607	54.85720
2.5%20	Siena	41.47655	53.19529
2.5%21	Selzer	48.61415	45.58879
2.5%22	Data Orbital	48.41652	47.97229
2.5%23	Washington Post/George Mason University	44.11135	51.30620
2.5%24	SurveyUSA/High Point University	48.73813	48.89372
2.5%25	YouGov Blue	49.19037	48.39704
2.5%26	U. North Florida	50.82811	45.84207

Table 3: Summary of Predictions by mean and lead probability

Metric	Value
Trump Mean Support	47.55
Harris Mean Support	49.26
Trump Lead Probability	0.37

5 Discussion

5.1 First discussion point

If my paper were 10 pages, then should be at least 2.5 pages. The discussion is a chance to show off what you know and what you learnt from all this.

5.2 Second discussion point

Please don't use these as sub-heading labels - change them to be what your point actually is.

5.3 Third discussion point

5.4 Weaknesses and next steps

Weaknesses and next steps should also be included.

A Appendix

A.1 Methodology Overview and Evaluation of YouGov Polling

Introduction

YouGov is an online polling company that uses a combination of non-probability sampling and a model known as multilevel regression with post-stratification (MRP) to estimate voter intentions. This appendix reviews YouGov’s approach, including sample recruitment, data collection, response handling, questionnaire design, and MRP modeling. It also assesses the strengths and limitations of YouGov’s methods.

1. Population, Frame, and Sample Composition

YouGov’s target population for its election polling includes American adults, with a focus on registered voters. It uses an online panel as its frame, drawing participants from a pool of volunteers who sign up and provide demographic information. This setup allows YouGov to adjust samples to align with the population it seeks to represent. Additionally, YouGov uses the TargetSmart voter file to verify that the sample aligns with national voter demographics. During election periods, YouGov increases its sample size, beginning with nearly 100,000 responses and adding another 20,000 in September and October to update its model.

2. Sample Recruitment and Representativeness

YouGov recruits panel members online through ads and partnerships with other websites. This approach allows any American adult with internet access to join. Members provide basic demographic information, which YouGov uses to select participants and apply statistical weighting to reflect the population accurately. Panelists earn points for participating, redeemable for small rewards, which encourages engagement and enhances data quality. However, since recruitment is internet-based, some groups without reliable internet access, such as rural or low-income populations, may be underrepresented.

3. Sampling Approach and Methodological Trade-offs

YouGov’s approach involves non-probability sampling, meaning that not everyone in the population has an equal chance of being selected. To address this, YouGov applies the MRP model, dividing respondents into subgroups based on characteristics like age, gender, race, education, region, and political affiliation. Each subgroup is weighted to match its share in the population. This method is efficient but has trade-offs, as MRP may not fully capture smaller states or hard-to-reach groups.

4. Non-response Handling and Quality Control Measures

YouGov uses strict quality control measures to manage non-response. Surveys include checks for speed and consistency, removing low-quality responses. Panelists who repeatedly fail these checks are excluded, helping to reduce bias and improve data quality. These quality controls are essential for the MRP model, which relies on accurate data to produce reliable estimates.

5. Questionnaire Design

YouGov’s questionnaires are neutral and straightforward, with randomized question order to reduce bias. They include various question types, such as text, images, and audio, helping respondents understand the context. Options like “prefer not to say” are included for sensitive questions, supporting respondent privacy. However, since the surveys are conducted online, some groups without internet access may be excluded, and shorter questionnaires may limit the depth of data collected.

6. MRP Model for Vote Estimation

YouGov’s MRP model estimates voter intentions through three stages: estimating likelihood to vote, predicting support for a candidate among likely voters, and aggregating results to calculate overall support. By matching responses to the TargetSmart voter file, the model generates estimates at both national and regional levels. The MRP model also helps YouGov track changes in voter intentions over time.

7. Strengths and Limitations of the YouGov Approach

YouGov’s approach has several strengths. The MRP model allows for accurate predictions at regional and national levels, and the use of voter files enhances sample representativeness. Repeated interviews with panelists allow YouGov to observe shifts in voter intentions over time. However, limitations include potential representativeness issues due to non-probability sampling, particularly in small states and among underrepresented groups. The internet-based survey method also may not fully capture populations without reliable internet access.

Conclusion

Through a combination of MRP modeling, quality control, and a large online panel, YouGov provides a structured approach for election polling. Although challenges remain in achieving full representativeness and internet coverage, YouGov’s approach offers useful information on U.S. voter intentions. The MRP model has proven effective in past elections and serves as a practical tool for future polling and tracking trends.

A.2 Idealized Survey Methodology

The budget for this survey is \$100,000, aimed at predicting the outcome of the 2024 United States Presidential Election. This methodology includes stratified sampling, multi-platform recruitment, data validation, and multi-wave data aggregation to ensure representative and reliable data.

Sampling Approach: Stratified Random Sampling

To obtain a representative sample, we use stratified random sampling with a sample size of 5,000 respondents. Sampling is stratified by age, gender, education, and geographic region to ensure broad coverage across voter demographics. The age groups include 18-29, 30-44, 45-64,

and 65 and above. Gender is categorized as male, female, and other, while education is classified as high school or below, bachelor’s degree, and master’s degree or above. The geographic region includes all U.S. states. Stratified sampling ensures that the sample accurately reflects voter characteristics, providing a solid foundation for subsequent analysis.

Recruitment: Multi-Platform and Interactive Engagement

Recruitment is conducted through a multi-platform strategy to ensure wide coverage among voters. Targeted ads are deployed on Google, Facebook, and Twitter to attract respondents with specific demographic characteristics. The ad content is concise and highlights anonymity and the research purpose to encourage participation. In addition, we use Random Digit Dialing (RDD) to contact older adults and residents in remote areas who may be less accessible online, ensuring their participation in the survey. To increase the completion rate, a small reward system is implemented, with one out of every 100 participants receiving a \$5 to \$10 incentive.

Survey Platform: Google Forms and Phone Outreach

Google Forms is used as the primary data collection platform, enabling respondents to complete the survey on a computer or mobile device with ease. We have set Google Forms to restrict submissions to one per Google account to prevent duplicate responses. For those who may not have easy online access, telephone outreach will be conducted, especially targeting older adults and rural voters, to ensure the diversity and inclusivity of the sample.

Data Validation: Post-Stratification Weighting

To ensure data quality, post-stratification weighting will be applied during the analysis phase based on demographic characteristics. This adjustment aligns the sample structure with the national voter distribution, reducing bias and increasing the accuracy of our findings. Multi-layered data validation measures enhance the reliability and representativeness of the results.

Poll Aggregation: Multi-Wave Polling and Adjustments

To track changes in voter sentiment over time, we will conduct multiple waves of data collection and aggregation. The survey will be administered in multiple rounds throughout the election cycle, with each round spaced 3-4 weeks apart. Each wave of data will be collected and analyzed independently, then aggregated using weighted averages to smooth out single-instance fluctuations and help identify long-term trends, providing a reliable basis for predictions.

Budget Allocation: Phased Spending and Testing

To maximize budget efficiency, we will use a phased spending strategy. An initial 30% of the advertising budget will be allocated to testing across platforms to determine effectiveness, with the remaining 70% directed to the most successful channels. The budget breakdown is as follows: social media advertising (\$40,000) for broad and targeted outreach, phone outreach (\$20,000) for RDD calls to less accessible populations, small rewards (\$15,000) to encourage completion, and data cleaning and analysis (\$25,000) for validation, weighting, and aggregation across multiple waves.

Survey Content and Link

Survey Title: 2024 United States Presidential Election Survey

Survey Link: <https://forms.gle/EzyHp3zuX8Cu6Ep8A>

Survey Questions:

1. What is your age?

- ☐ 18-29
- ☐ 30-44
- ☐ 45-64
- ☐ 65 and above

2. What is your gender?

- ☐ Male
- ☐ Female
- ☐ Other
- ☐ Prefer not to say

3. What is your highest level of education?

- ☐ High school or below
- ☐ Bachelor's degree
- ☐ Master's degree or above

4. Which state do you currently reside in?

(Dropdown menu with state names)

5. Do you plan to vote in the 2024 election?

- ☐ Yes
- ☐ No
- ☐ Not sure

6. If the election were held today, which candidate would you be more likely to support?

- o Kamala Harris (Democratic Party)
- o Donald Trump (Republican Party)
- o Other
- o Undecided

7. Which of the following issues would have the greatest impact on your decision in the 2024 election? (Select one)

- o Economic stability and growth
- o Access to quality healthcare
- o Education reform and funding
- o Environmental sustainability and climate action
- o Addressing social and racial inequalities
- o Other (please specify): _____

8. If a future candidate proposed a policy that aligns perfectly with your concerns, would you consider changing your voting intention?

- o Yes
- o No
- o Not sure

9. If you would like to participate in the reward draw, please provide your email address below. Your contact information will only be used to notify winners. We will reach out to winners via email.

Email: _____

A.3 Clean data

It starts by reading the dataset and standardizing column names. Next, it removes irrelevant columns like `sponsor_ids`, `sponsors`, and etc and filters for high-quality polls (with a rating of 2.7 or higher) that specifically track Trump’s support. National polls missing state information are labeled as “National,” and dates are formatted and filtered to include only those on or after July 21, 2024 (when Trump declared his candidacy). Finally, the percentage supporting Trump is converted to an actual count for further modeling.

A.4 Variable details

We have used the variable `numeric_grade`, `transparency_score`, `pollster`, `poll_score`, `end_date` and `pct` in the building of Bayesian model. Among these variables in the data set after analysis,

- The `numeric_grade` is the numeric rating given to the pollster to indicate their quality or reliability from 2.7 to 3.0.
- The `transparency_score` is the grade for how transparent a pollster is, calculated based on how much information it discloses about its polls and weighted by recency from 4.5 to 10.0.
- The `pollster` is the name of the polling organization that conducted the poll included Marquette Law School, CNN/SSRS and etc.
- The `poll_score` is the numeric value representing the score or reliability of the pollster in question from -1.5 to -0.5 and negative numbers of poll score are better.
- The `end_date` is the date of polling ends.
- The `pct` is the percentage of the vote or support that the candidate received in the poll and keep integer.

A.5 Results of model

A.5.1 Table for summary of model results

In Table 4, we conducted a summary of three models, including a generalized linear regression and a Bayesian model for Trump's data, as well as a Bayesian model for Harris's data. This summary includes the γ_j values for each pollster and monitoring data for the models, such as AIC, BIC, etc.

Table 4: Explanatory models of flight time based on wing width and wing length

[!h]

	Model Summary		
	Model by glm	Model by bayes with Trump	Model by bayes with Harris
(Intercept)	-524.927 (126.816)	-590.212	-450.720
numeric_grade	-14.386 (14.848)	-3.356	0.225
transparency_score	0.903 (0.638)	-0.164	-0.213
pollscore	-3.028 (15.053)	1.121	0.842
Beacon/Shaw	-2.833 (5.143)	0.976	1.387
Christopher Newport U.	-9.186 (4.621)	-5.294	2.641
CNN/SSRS	-4.208 (3.383)	-1.343	0.792
Data Orbital	-3.602 (4.436)	-0.426	-0.776
Echelon Insights	-3.313 (2.258)	0.389	1.337
Emerson	1.010 (3.603)	1.561	1.543
Ipsos	-7.483 (2.857)	-3.927	-0.555
Marist	0.162 (2.423)	1.111	0.734
Marquette Law School	-3.213 (4.838)	-0.561	2.052
MassINC Polling Group	-7.280 (1.695)	-6.566	6.224
McCourtney Institute/YouGov	-1.354 (4.725)	-1.745	-0.218
Muhlenberg	-2.040 (4.528)	0.196	0.205
Quinnipiac	-1.513 (4.667)	-0.189	-0.058
Selzer	-2.952 (6.632)	-0.220	-3.090
Siena	-9.684 (1.843)	-7.344	4.433
Siena/NYT	-2.174 (9.330)	0.976	-1.115
Suffolk	-4.753 (3.179)	-3.290	0.656
SurveyUSA	-5.810 (7.583)	-2.163	0.519
SurveyUSA/High Point University	-5.876 (8.815)	-0.224	0.219
The Washington Post	0.576 (5.526)	2.315	-0.867
University of Massachusetts Lowell/YouGov	-6.065 (3.707)	-3.502	0.238
Washington Post/George Mason University	-8.943 (2.580)	-4.783	2.659
YouGov	-2.338 (4.355)	-0.764	-0.322
YouGov Blue	-1.128 (5.416)	0.335	-0.338
End Date	0.030 (0.006)	0.032	0.025
U. North Florida		1.990	-2.956
YouGov/Center for Working Class Politics		-2.360	-1.418
Num.Obs.	492	492	471
R2	0.332	0.318	0.248
R2 Adj.		0.258	0.163
AIC	2575.9		
BIC	2701.8		
Log.Lik.	-1257.932	-1264.371	-1203.153
ELPD		-1289.0	-1226.3
ELPD s.e.		25.3	23.9
LOOIC		2577.9	2452.6
LOOIC s.e.		50.6	47.8
WAIC		2575.6	2450.5
RMSE	3.12	3.38	3.22

^a This table shows the regression models with custom variable names.

A.5.2 Bayesian model for Harris data set

In Figure 13, we applied a bayesian model to predict the percentage of polls Harris according each pollster. Each data point is color-coded by pollster, with support rates ranging from 37.5% to 65%. The shaded gray area around the trend line represents the confidence interval, indicating the uncertainty in support rate variations. The overall trend line (blue) shows a slight upward trend, increasing from about 46% in August to nearly 48% in October. Compared to Trump's model, the increasing trend of Harris is lower than Trump's. By October, most poll results have converged, with support rates centered between 45% and 50%.

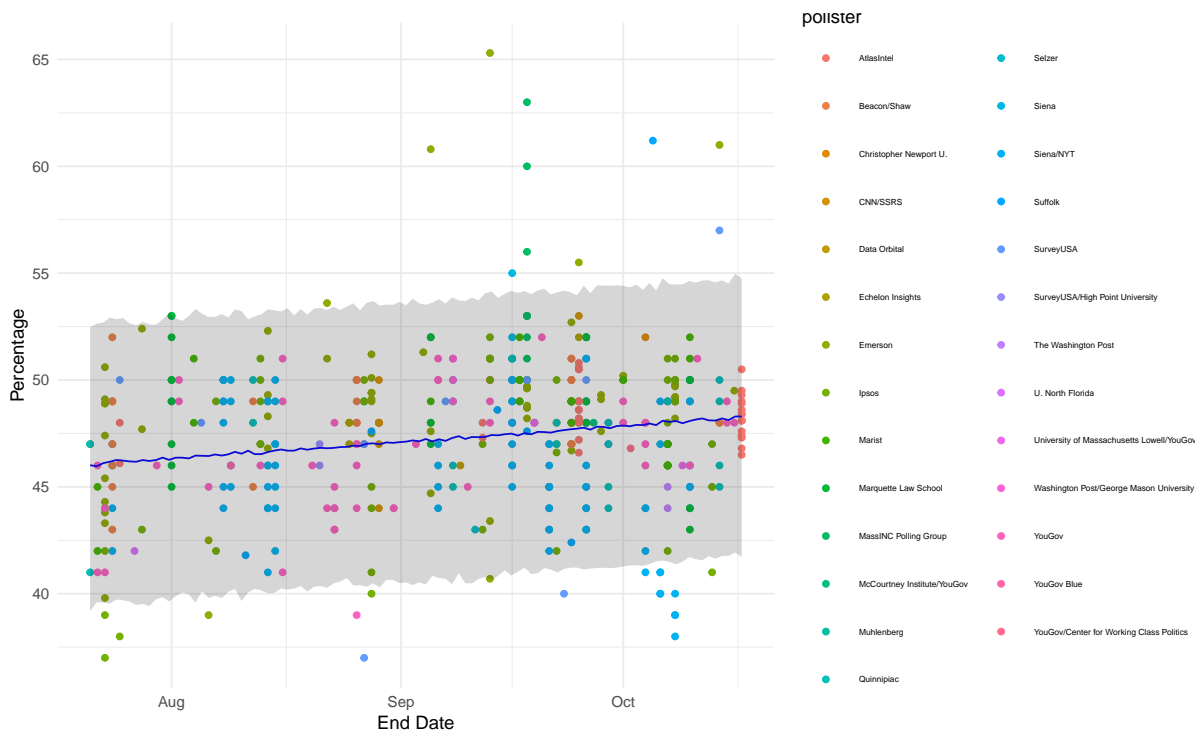


Figure 13: Poll Percentage over Time with Bayesian Fit for Harris

A.5.3 Predictions for both Trump and Harris

In Table 5, we summarized the Bayesian model predictions for Trump and Harris on November 5, 2024, based on 27 pollsters, including the mean, lower bound, and upper bound.

Table 5: predictions for trump and harris by pollster

Predictions Summary					
	Pollster	Candidate	Predicted Mean	Lower Bound	Upper Bound
2.5%	AtlasIntel	Trump	48.94631	42.43347	55.72567
2.5%1	Emerson	Trump	50.34653	43.55557	56.75106
2.5%2	YouGov	Trump	48.05967	41.58112	54.52504
2.5%3	Beacon/Shaw	Trump	49.77089	43.25543	56.09705
2.5%4	Quinnipiac	Trump	48.62878	41.39774	55.91628
2.5%5	SurveyUSA	Trump	46.63862	39.84705	53.37863
2.5%6	Ipsos	Trump	44.92498	38.11308	51.31979
2.5%7	Marist	Trump	49.95651	43.28239	56.33919
2.5%8	Siena/NYT	Trump	49.85318	42.98870	56.77792
2.5%9	University of Massachusetts Lowell/YouGov	Trump	45.39482	38.35665	52.46243
2.5%10	Marquette Law School	Trump	48.18798	41.37041	54.72301
2.5%11	The Washington Post	Trump	51.16407	44.55568	58.18185
2.5%12	Suffolk	Trump	45.51106	38.68998	52.46700
2.5%13	Christopher Newport U.	Trump	43.50302	35.10869	51.95017
2.5%14	YouGov/Center for Working Class Politics	Trump	46.55850	38.01149	55.29104
2.5%15	McCourtney Institute/YouGov	Trump	47.05294	38.49958	55.38905
2.5%16	Echelon Insights	Trump	49.29397	42.32005	56.09514
2.5%17	CNN/SSRS	Trump	47.53324	40.13830	54.76963
2.5%18	Muhlenberg	Trump	49.05747	40.66915	57.71098
2.5%19	MassINC Polling Group	Trump	42.41278	35.40659	49.36277
2.5%20	Siena	Trump	41.57369	34.37409	48.91954

2.5%21	Selzer	Trump	48.56498	39.93915	56.77115
2.5%22	Data Orbital	Trump	48.36348	40.21363	56.85730
2.5%23	Washington Post/George Mason University	Trump	44.10653	36.66518	51.46213
2.5%24	SurveyUSA/High Point University	Trump	48.70149	40.94166	56.33244
2.5%25	YouGov Blue	Trump	49.12226	41.09789	57.42313
2.5%26	U. North Florida	Trump	50.92116	42.51759	59.69250
2.5%27	AtlasIntel	Harris	48.70733	42.22387	55.22802
2.5%28	Emerson	Harris	50.21708	43.74753	56.80086
2.5%29	YouGov	Harris	48.35765	41.75536	54.68130
2.5%30	Beacon/Shaw	Harris	50.15770	43.81949	56.35262
2.5%31	Quinnipiac	Harris	48.55050	40.99461	55.68766
2.5%32	SurveyUSA	Harris	49.35902	42.66895	56.22528
2.5%33	Ipsos	Harris	48.17484	41.60572	54.52572
2.5%34	Marist	Harris	49.41276	43.06539	55.92794
2.5%35	Siena/NYT	Harris	47.59682	40.91047	54.38905
2.5%36	University of Massachusetts Lowell/YouGov	Harris	48.96015	41.72479	55.89682
2.5%37	Marquette Law School	Harris	50.80365	44.10780	57.30961
2.5%38	The Washington Post	Harris	47.88813	40.95054	54.63634
2.5%39	Suffolk	Harris	49.45653	42.78044	56.06545
2.5%40	Christopher Newport U.	Harris	51.34027	42.84452	59.87448
2.5%41	YouGov/Center for Working Class Politics	Harris	47.34676	38.65365	55.84531
2.5%42	McCourtney Institute/YouGov	Harris	48.48063	39.97428	56.81290
2.5%43	Echelon Insights	Harris	50.00029	42.88039	56.88593
2.5%44	CNN/SSRS	Harris	49.52549	42.23237	56.86531
2.5%45	Muhlenberg	Harris	48.99389	40.95485	57.16249
2.5%46	MassINC Polling Group	Harris	54.88675	47.98038	61.71728
2.5%47	Siena	Harris	53.17665	46.23865	60.26020
2.5%48	Selzer	Harris	45.50769	37.16547	53.77141

2.5%49	Data Orbital	Harris	47.89524	39.77833	56.10508
2.5%50	Washington Post/George Mason University	Harris	51.43515	44.06587	58.59393
2.5%51	SurveyUSA/High Point University	Harris	48.89925	41.28931	56.53991
2.5%52	YouGov Blue	Harris	48.38442	40.21236	56.84120
2.5%53	U. North Florida	Harris	45.76879	37.05674	54.45347

B Model details

To maintain readability while demonstrating model robustness, we include the following in the appendix:

- **Prior Justification and Sensitivity Analyses:** Alternative priors and their justifications are provided, alongside a sensitivity analysis to examine the impact of these priors on the posterior distributions.
- **Model Validation and Out-of-Sample Testing:** Validation metrics, such as RMSE calculations, out-of-sample testing, and test-train splits, offer evidence of the model’s predictive accuracy, complementing in-sample performance.
- **Alternative Models and Comparison:** An analysis of simpler and more complex models is included, explaining the rationale for selecting this Bayesian model based on performance metrics and interpretability.

References

- Alexander, Rohan. 2023. *Telling Stories with Data: With Applications in r*. Chapman; Hall/CRC.
- Arel-Bundock, Vincent. 2022. “modelssummary: Data and Model Summaries in R.” *Journal of Statistical Software* 103 (1): 1–23. <https://doi.org/10.18637/jss.v103.i01>.
- FiveThirtyEight. 2024. “2024 National Presidential Poll Results.” <https://projects.fivethirtyeight.com/polls/president-general/2024/national/>.
- Goodrich, Ben, Jonah Gabry, Imad Ali, and Sam Brilleman. 2022. “rstanarm: Bayesian applied regression modeling via Stan.” <https://mc-stan.org/rstanarm/>.
- Müller, Kirill. 2020. *Here: A Simpler Way to Find Your Files*. <https://CRAN.R-project.org/package=here>.
- R Core Team. 2023. *R: A Language and Environment for Statistical Computing*. Vienna, Austria: R Foundation for Statistical Computing. <https://www.R-project.org/>.
- Richardson, Neal, Ian Cook, Nic Crane, Dewey Dunnington, Romain François, Jonathan Keane, Dragoş Moldovan-Grünfeld, Jeroen Ooms, Jacob Wujciak-Jens, and Apache Arrow. 2024. *Arrow: Integration to 'Apache' 'Arrow'*. <https://CRAN.R-project.org/package=arrow>.
- Wickham, Hadley. 2016. *Ggplot2: Elegant Graphics for Data Analysis*. Springer-Verlag New York. <https://ggplot2.tidyverse.org>.
- . 2023. *Modelr: Modelling Functions That Work with the Pipe*. <https://CRAN.R-project.org/package=modelr>.
- Wickham, Hadley, Mara Averick, Jennifer Bryan, Winston Chang, Lucy D’Agostino McGowan, Romain François, Garrett Grolemund, et al. 2019. “Welcome to the tidyverse.” *Journal of Open Source Software* 4 (43): 1686. <https://doi.org/10.21105/joss.01686>.
- Wickham, Hadley, Romain François, Lionel Henry, Kirill Müller, and Davis Vaughan. 2023. *Dplyr: A Grammar of Data Manipulation*. <https://CRAN.R-project.org/package=dplyr>.
- Xie, Yihui. 2023. *Knitr: A General-Purpose Package for Dynamic Report Generation in r*. <https://yihui.org/knitr/>.
- Zhu, Hao. 2024. *kableExtra: Construct Complex Table with 'Kable' and Pipe Syntax*. <https://CRAN.R-project.org/package=kableExtra>.