

A Poll-of-Polls Forecast for the 2024 U.S. Presidential Election: Kamala Harris Emerges as a Leading Candidate with Constant Support Rates Hovering Around 50%*

Multiple Linear Regression and Bayesian Models Accounting for Pollster and Geographical Disparities in Support Trends

Yunkyung Ko

October 29, 2024

In the anticipation of the 2024 US presidential election, this paper intends to analyze and predict the level of support that Kamala Harris will gain. The support polls for Harris remained relatively constant around 47%. However, there were some variations detected when accounting for pollster and geographical differences, with Siena/NYT and Texas standing out with its large fluctuation in support rates. Our objective is to propose a reliable electoral prediction considering those variations help global stakeholders take reasonable preemptive measures. [!!! NEED MODIFICATIONS FOR THE MAIN RESULTS PART !!!]

Table of Contents

1	Introduction	2
2	Data	3
2.1	Overview	3
2.2	Measurement	4
2.3	Outcome variable	4
2.4	Predictor variables	6
2.5	Correlation Between Predictor Variables	10

*Code and data are available at: https://github.com/koyunkyung/us_election_2024.

3	Model	13
3.1	Model set-up	14
4	Results	16
4.1	Results from the relation between the prediction variable and outcome variable of the analysis data	16
4.2	Results from the prediction model	18
5	Discussion	18
5.1	Why Harris could beat her polls	21
5.2	Pollsters herding around false consensus	21
5.3	State-Level Differences	21
5.4	Weaknesses and next steps	21
A	Appendix	23
B	Additional data details	23
B.1	Pollster Methodology Overview and Evaluation	23
B.2	Idealized Methodolgy	23
B.3	Idealized Survey	24
B.4	Survey Demo	24
C	Model details	28
C.1	Posterior predictive check	28
C.2	Diagnostics	29
	References	30

1 Introduction

The US presidential election is an event that receives a lot of international attention due to its far-reaching implications beyond the country’s borders. The effects of US elections are not only related to economy and international relations around the world, but also link to social and environmental issues such as climate change (Bijune and Ha 2024). In the anticipation of the 2024 US electoral competition, this paper is aimed at predicting possible outcomes of the election by analyzing the level of support that Kamala Harris will gain.

We forecast the support of Kamala Harris based on the polling results at the national and state levels, and apply a linear regression and a Bayesian approach. The main parameter of interest is the proportion of vote or support that Harris received in surveys, which is traced over time. By considering the effect of changes in poll-making organizations and geographical distinctions, our objective is to correct for variation across different voter bases with the pooling the polls approach (Jackman 2024).

Our initial linear model examining the support for Harris over time suggests that the rate of support remained relatively stable around 47%, with no significant increase or decrease. However large variability was detected, as seen in the spread of points around the fitted line. Consequently, pollster-specific effects and state-level random effects were added to the model but resulted in even higher variations. Some pollsters or states consistently reported higher or lower support for Harris compared to others, proposing the significance of considering distinct pollster and state environment when interpreting electoral polling results.

These prognostics provide much more than merely forecasting the election in question. Hinting the trajectory the US might take in matters of foreign policy, economics, and global politics, predictions enable stakeholders worldwide to take preemptive measures for the resulted changes of a newly elected government (CSIS 2024). As such, this study not only contributes to the domestic political discourse but also provides a valuable tool for global actors seeking to navigate the uncertainty surrounding the 2024 US presidential election (CSIS 2024).

The paper is structured as follows. Initially, Section 2 and Section 3 explores the data and methodology used, including filtering and modeling techniques applied to the polling data. Following that, Section 4 presents the results from the linear and Bayesian models, while the next section Section 5 discusses the broader implications of these findings. Finally, the paper concludes with remarks on future directions for research and applications of these models Section 5.

2 Data

2.1 Overview

We use the statistical programming language R (R Core Team 2023) to analyze US presidential polling data from FiveThirtyEight (FiveThirtyEight 2024), focusing on support for Kamala Harris. The dataset includes a wide range of poll results from various national and state-level polls, with key variables such as pollster, sample size, percentage of support for Harris, and end date of the poll. Following the guidance of Alexander (2023), we compiled the results of each opinion poll over a period of time and compared them taking into account the methodological peculiarities of polling by pollsters and geographical scope of the conducted polls.

To ensure data quality, we filtered the dataset to include only polls that measured Kamala Harris' support, with a numeric grade of pollster 2.7 or higher for reliability. We also limited the analysis to polls conducted after July 21, 2024, when Harris officially declared her candidacy, and excluded pollsters with fewer than 30 polls to focus on those with sufficient data for robust results.

In performing the analysis, we utilized several R packages. `tidyverse` (Wickham et al. 2019) was used for data manipulation and visualization and `rstanarm` (Goodrich et al. 2024),

`modelsummary` (Arel-Bundock 2022) was respectively used for Bayesian modeling and generating model summaries. For visualizing results, `ggplot2` (Wickham 2016) was used and `kableExtra` (Zhu 2024) helped format tables for presentation. These packages provided a framework for efficient data processing, modeling, and reporting.

(Data last updated on 24 Oct 2024.) - [!!! MODIFY AT FINAL !!!]

2.2 Measurement

[!!! MODIFY OBSV. NUMBERS AFTER DATA UPDATES !!!] - ERASE AT FINAL

The original dataset sourced from FiveThirtyEight (FiveThirtyEight 2024) aggregates a wide range of poll results (16506 observations, based on dataset available at Oct 29). The polls conducted by various polling organizations capture voter preferences by taking a representative sample of the electorate and asking for the voters' candidate of choice. Surveys were conducted at the state and national levels, providing the wide perspective on public feelings across the country.

Each poll represents a predictor of an actual event, namely voter opinion at a particular moment. Nevertheless, like all survey results, the raw data is susceptible to many potential limitations including the following: sampling error, variation in polling methods, distortion because of inappropriate survey responses such as missing data or response from respondents who misunderstood the questions (Alexander 2023).

While applying several filters to the original dataset such as restricting to those with a numeric grade of 2.7 or higher or pollsters with more than 30 polls improves data reliability, certain limitations still exist. Selection bias and sampling error remains as a concern, since polls always represent only part of the population. Differences in the way different organizations conducted their polling might introduce more inconsistencies. Finally, by focusing our attention on post-declaration polls only, we exclude earlier trends that could add more insight into how Harris' support has evolved over time.

2.3 Outcome variable

2.3.1 The Proportion(%) of Support that a Candidate Received in the Poll

The main variable of interest that we aim to forecast is the 'pct' variable, which represents the proportion of vote or support that a candidate received in the poll. Table 1 and Figure 1 respectively shows the summary statistics and distribution of the 'pct' variable in the original dataset (FiveThirtyEight 2024). Table 2 and Figure 2 shows the summary statistics and distribution of the same variable, but in the filtered dataset that only comprises of the supporting votes for Harris from relatively high-quality polling organizations. Comparing the summary statistics for the raw data (Table 1) and filtered data (Table 2), higher numbers were derived

from data filtered only by Harris supporters. Also, Figure 2 illustrates that a significant number of polls indicate support levels ranging from 40% to 50%, which suggests a stable yet not substantial endorsement. This proposes that Harris possesses a reliable foundational support, although her capacity to obtain a majority remains ambiguous.

Table 1: Summary statistics for the proportion(%) of support that candidates received in the poll

mean	median	min	max	sd	n
33.7	42	0	70	18.15	16506

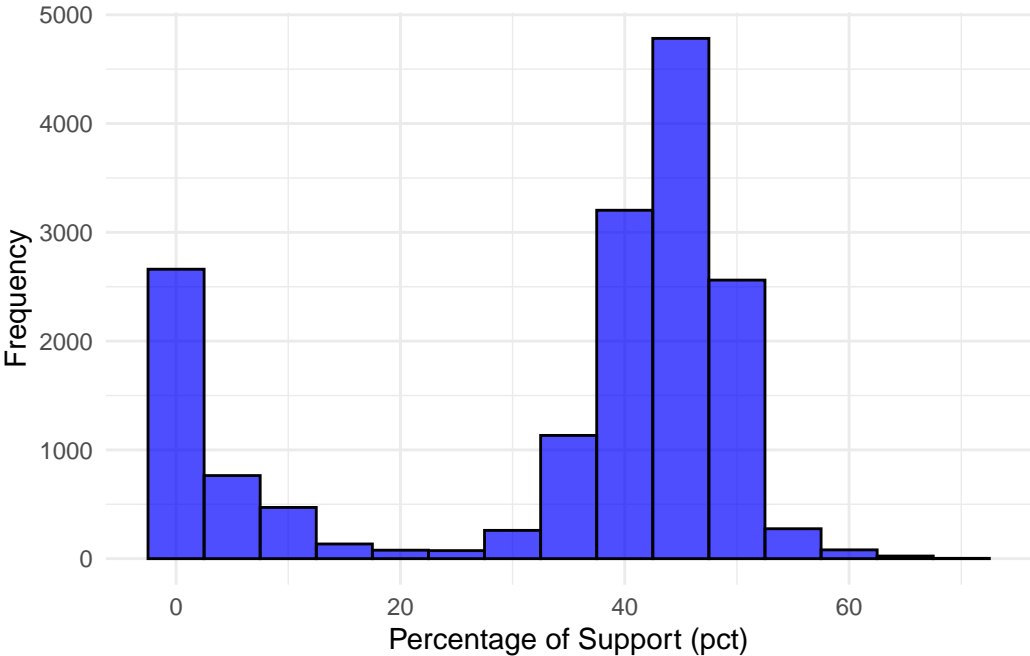


Figure 1: Distribution of the proportion(%) of support that candidates received in the poll

Table 2: Summary statistics for the proportion(%) of support that Harris received in high-quality polls

mean	median	min	max	sd	n
47.46	48	36.5	65.3	3.75	266

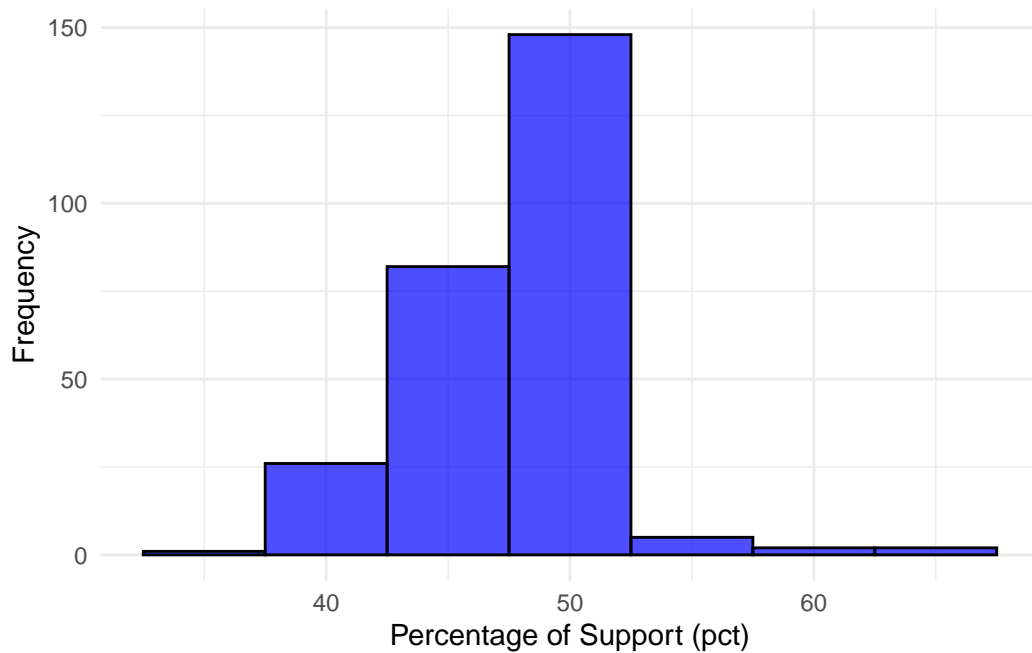


Figure 2: Distribution of the proportion(%) of support that Harris received in high-quality polls

2.4 Predictor variables

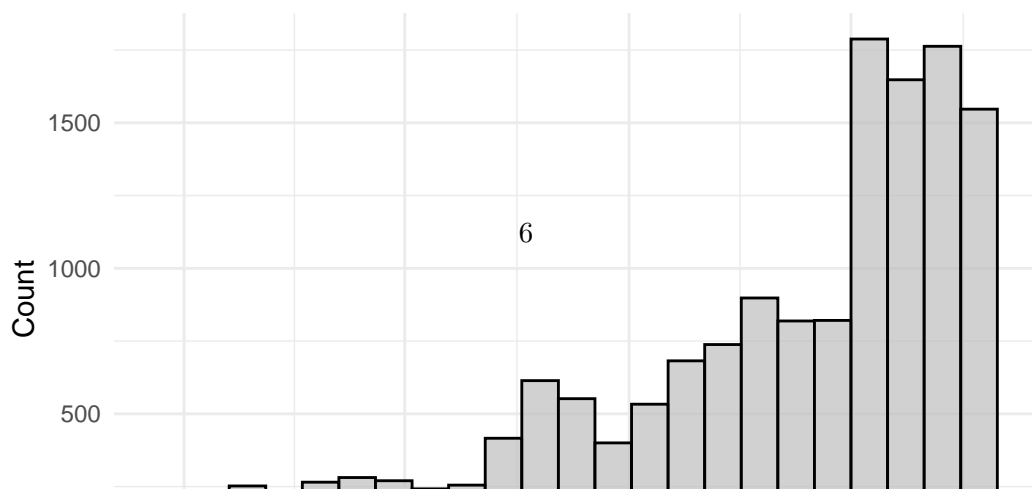
2.4.1 The Date the Poll was Concluded

[!!! MODIFY DATE RANGE AFTER DATA UPDATES !!!] - ERASE AT FINAL

The ‘end_date’ variable representing the time the poll was concluded was put into account to keep track of how support for a candidate changes in time. The reported end dates in the original dataset (FiveThirtyEight 2024) ranges from 1 January, 2023 to 9 September, 2024 (Table 3). Figure 3 shows that the polling data is more concentrated on survey results conducted in the recent period.

Table 3: Summary statistics for the date the poll was concluded

Min	Max
1/1/23	9/9/24



Harris. So, the date variable for filtered data ranges from 23 July, 2024 to 14 October, 2024 (Table 4). Figure 4 shows that overall, polling is conducted regularly but intensifies around specific dates.

Table 4: Summary statistics for the date the high-quality polls for Harris was concluded

Min	Max
2024-07-23	2024-10-26

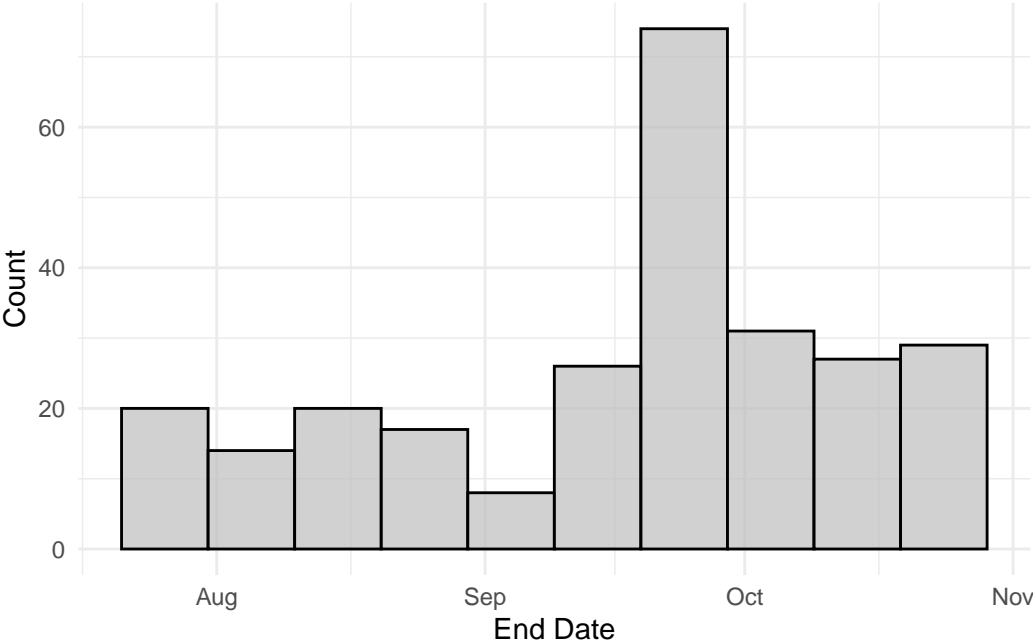


Figure 4: Distribution of the date the high-quality polls for Harris was concluded

2.4.2 Pollster and State

The ‘pollster’ and ‘state’ variable were selected to consider the effect of changes in poll-making organizations and geographical distinctions. The two variables respectively represent the polling organization that conducted the poll and the US state where the poll was conducted or focused.

Table 5 shows that the original dataset (FiveThirtyEight 2024) contains 222 distinct pollsters and 54 distinct states. After filtering for high-quality polls and assigning ‘other’ for states with fewer than 60 polls, the analysis data contains 3 distinct poll-making organizations and 19 geographical distinctions as shown in Table 6.

Table 5: Number of distinct polling organizations and US states where the poll was conducted

Pollster	State
229	54

Table 6: Number of distinct high-quality polling organizations and US states where more than 60 polls for Harris were conducted

Pollster	State
6	19

The distribution of polling counts for different pollsters in Figure 5 suggests that the analysis data is dominated by a few pollsters, particularly Siena/NYT. Depending on their polling methodology, the general results may have potential biases. A detailed analysis of the polling methodology and possible errors of the organization will be covered in Section B.

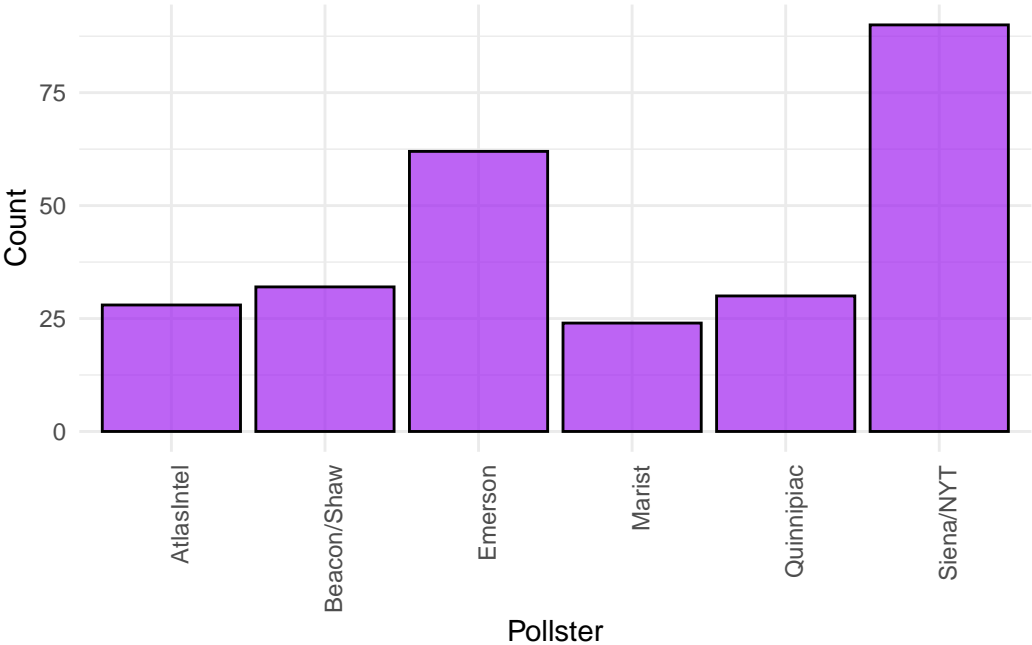


Figure 5: Distribution of polling organizations where high-quality polls for Harris were conducted

Figure 6 displays the distribution of polls across different states in the analysis data. Pennsylvania, Arizona, and Georgia are the top 3 states with high number of polls while states like

Montana, New Mexico, and Maryland have much fewer polls. Note that a significant number of national or unspecified state-level polls are aggregated in this analysis data regarding the high count in ‘Other’ category. The concentration of polls in certain states further suggests a strategic focus on areas likely to impact the election outcome (11Alive 2024).

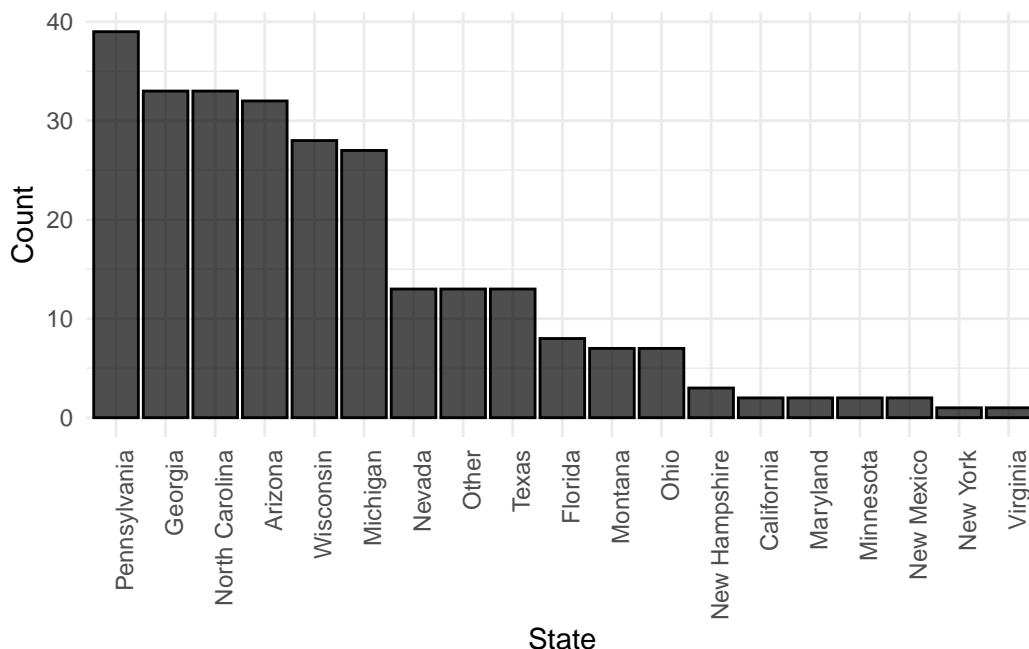


Figure 6: Distribution of US states where the more than 60 polls for Harris were conducted

2.4.3 Pollscore Measuring the Validity of Polling Questions

This variable was another factor we had put into consideration to check whether the validity of the polling questions affects the polling results. The ‘pollscore’ variable represents the score or reliability of the pollster in question. The numeric values are the error and bias that can be attributed to the pollster, which means negative numbers are better. Table 7 and Figure 7 suggests that while the majority of the polls are moderately to highly qualitative in the original dataset, a fraction of the polls with low-quality or no scores could add noise or uncertainty to the analysis.

Table 7: Summary statistics for the reliability scores of pollsters

mean	median	min	max	sd	n
-0.38	-0.3	-1.5	1.7	0.7	16506

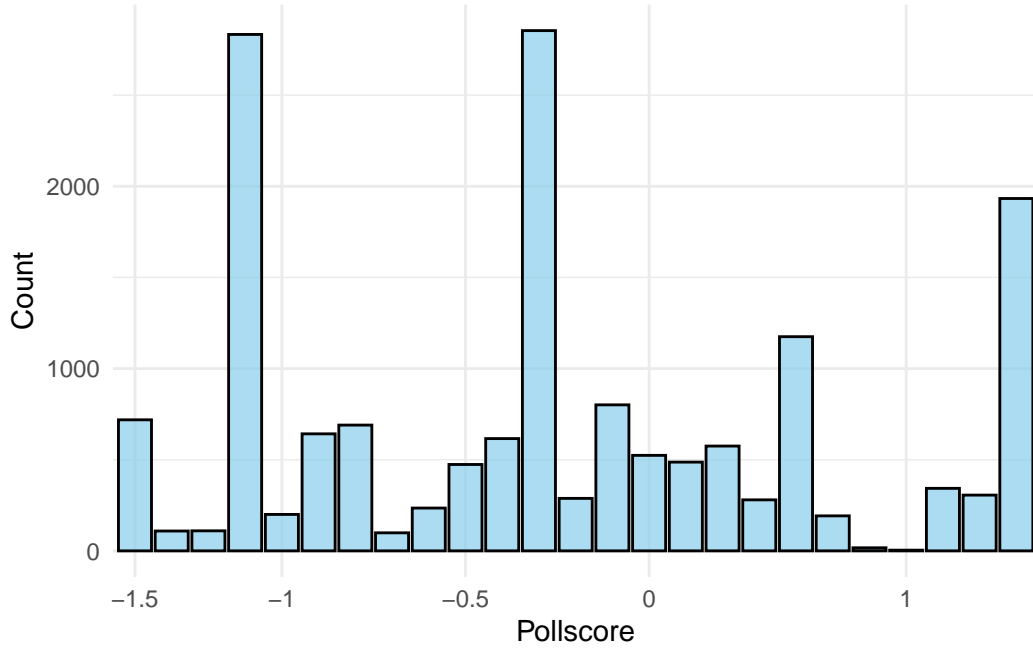


Figure 7: Distribution of the reliability scores of pollsters

After the filtering to polling data of high-quality polling organizations, we can find that the overall value and standard deviation of pollscores went down in Table 8. This implies that the polling data narrowed down to the responses from more reliable survey questions. Figure 8 also indicates that the data cleaning process effectively excluded less reliable sources, which can enhance the robustness of subsequent analyses.

Table 8: Summary statistics for the reliability scores of high-quality pollsters used for analysis

mean	median	min	max	sd	n
-1.12	-1.1	-1.5	-0.5	0.33	266

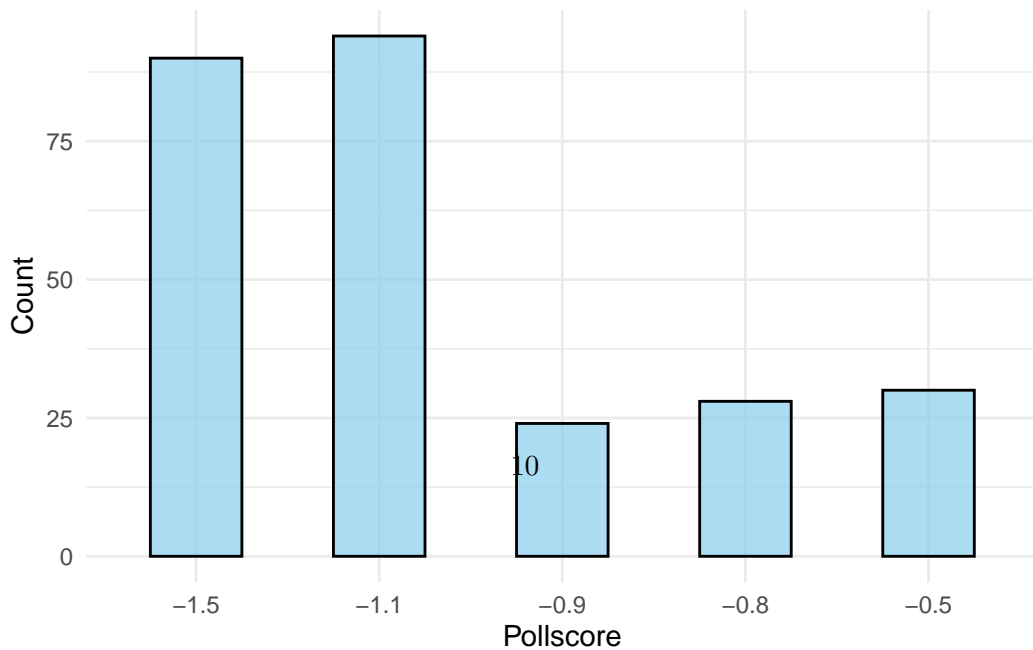


Figure 8: Distribution of the reliability scores of high-quality pollsters used for analysis

2.5.1 End Date and State

Figure 9 shows notable polling concentrations on certain states such as Pennsylvania, Ohio, and Nevada in week 2024-37 (September 9, 2024) or 2024-38 (September 15, 2024). This suggests that these states are key battlegrounds or areas of strategic focus during the election period. Moreover, polling activities are not consistent across all weeks and there are clear peaks in polling activity in week 2024-38 (September 15, 2024), which may correspond to significant political events, debates, or media focuses.

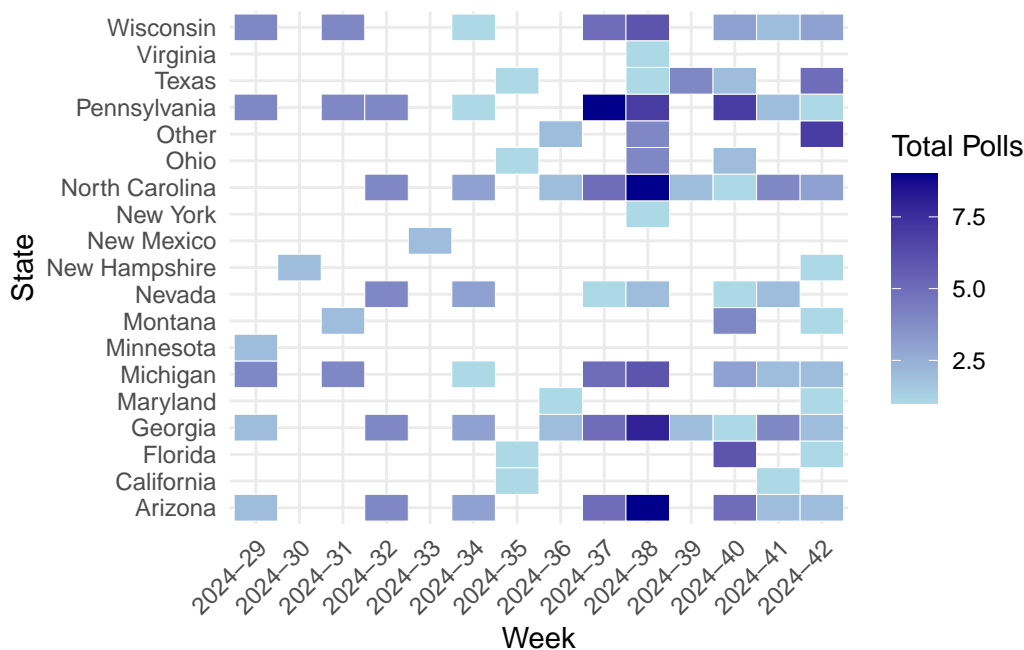


Figure 9: Polling concentration over time by state *Note:* Strategic focus of polling is directed in battleground states such as Pennsylvania and Ohio in September 15, 2024

2.5.2 End Date and Pollster

Figure 10 shows that Siena/NYT, which made up the largest proportion of polls in our analysis dataset, has concentrated polling efforts in week 2024-37 (September 9, 2024) and week 2024-38 (September 15, 2024). Notably active periods, possibly around key election events, might introduce temporal bias and overemphasize methodologies of Siena/NYT in that specific time period. Emerson shows overall consistency of polling activities across several weeks indicating steady involvement, while other polling organizations show more sporadic or minimal involvement.

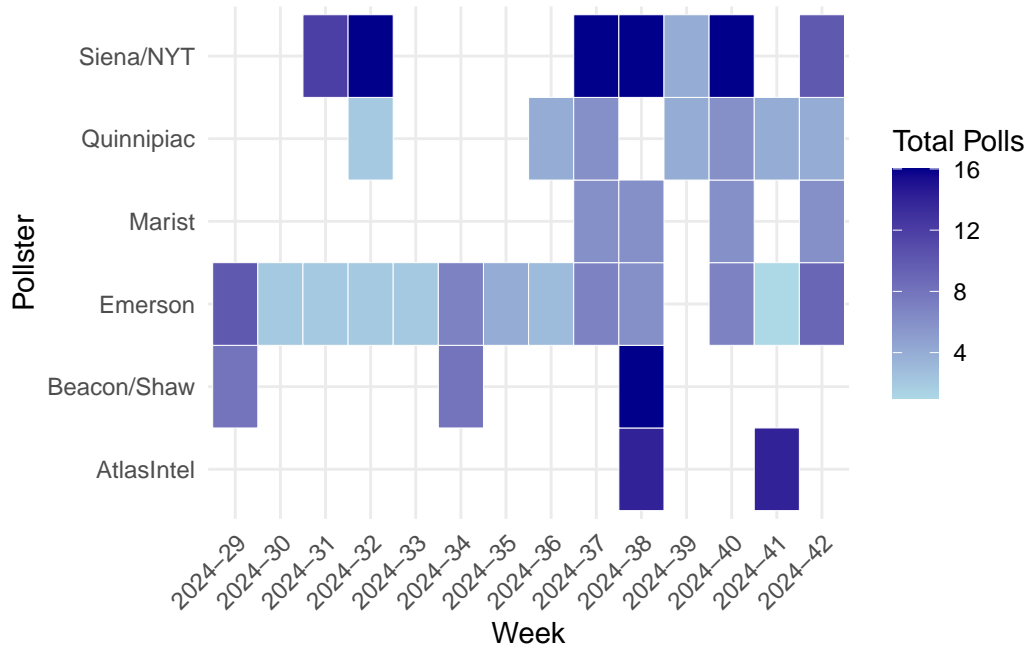


Figure 10: Polling concentration over time by pollster *Note:* Concentrated polling efforts of Siena/NYT in the first half of September, 2024 show temporal bias and potential overemphasizing of their polling methodologies around key election events.

2.5.3 Pollster and State

Figure 11 shows that Siena/NYT has the strongest state-level presence in key battleground states like Michigan and Arizona. Others like AtlasIntel and Beacon/Shaw has minimal polling coverage in targeted states such as Nevada and Georgia. This variation in coverage could influence the overall analysis in that some states might receive disproportionately more attention from certain pollsters.

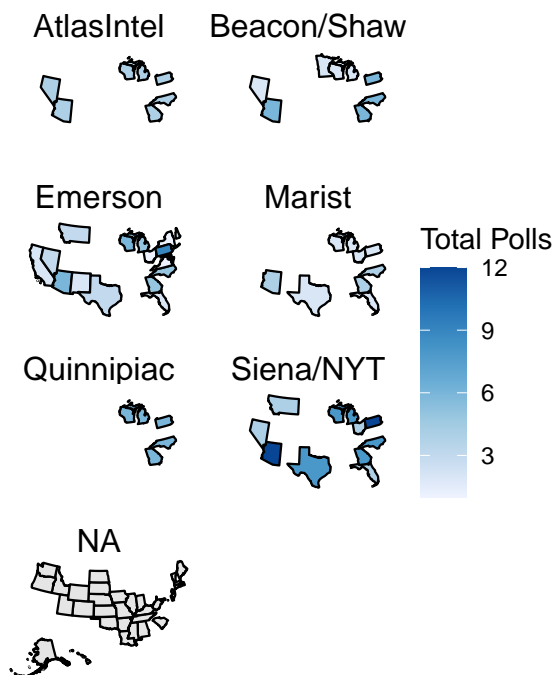


Figure 11: Polling activity by pollster and state *Note:* Siena/NYT has strongest state-level polling coverage in key battleground states

3 Model

The goal of our modelling strategy is to estimate Kamala Harris's support percentage in the 2024 US election polls, accounting for potential variations over time, as well as across different pollsters and states. The model balances complexity with interpretability, incorporating both linear and Bayesian frameworks to capture patterns in the data. Background details and diagnostics are included in Appendix C.

3.1 Model set-up

Define y_i as the percentage of support that Kamala Harris receives in poll i . We begin with two simple linear models and progress to more complex Bayesian models that account for hierarchical structures.

The following models outline our approach:

3.1.1 Linear Model by Date

$$y_i = \beta_0 + \beta_1 \cdot \text{end_date}_i + \epsilon_i \quad (1)$$

$$\epsilon_i \sim \text{Normal}(0, \sigma^2) \quad (2)$$

where: - y_i is the percentage of support for Harris in poll i , - β_0 is the intercept, - β_1 represents the effect of the poll's end date, - ϵ_i is the error term.

3.1.2 Linear Model by Date and Pollster

$$y_i = \beta_0 + \beta_1 \cdot \text{end_date}_i + \gamma_{p[i]} + \epsilon_i \quad (3)$$

$$\epsilon_i \sim \text{Normal}(0, \sigma^2) \quad (4)$$

where: - $\gamma_{p[i]}$ is a fixed effect for pollster p conducting poll i (e.g., Siena/NYT).

3.1.3 Bayesian Model with Random Intercept for Pollster

$$y_i | \mu_i, \sigma \sim \text{Normal}(\mu_i, \sigma) \quad (5)$$

$$\mu_i = \beta_0 + \beta_1 \cdot \text{end_date}_i + \gamma_{p[i]} \quad (6)$$

$$\gamma_p \sim \text{Normal}(0, \sigma_\gamma) \quad (7)$$

where: - γ_p is a random effect for pollster p .

3.1.4 Bayesian Model with Random Intercept for Pollster and State

$$y_i | \mu_i, \sigma \sim \text{Normal}(\mu_i, \sigma) \quad (8)$$

$$\mu_i = \beta_0 + \beta_1 \cdot \text{end_date}_i + \gamma_{p[i]} + \delta_{s[i]} \quad (9)$$

$$\gamma_p \sim \text{Normal}(0, \sigma_\gamma) \quad (10)$$

$$\delta_s \sim \text{Normal}(0, \sigma_\delta) \quad (11)$$

where: - $\delta_{s[i]}$ is a random effect for state s .

The Bayesian models are fit using **rstanarm** in R. The priors used are weakly informative: - $\beta_0 \sim \text{Normal}(0, 10)$ - $\sigma \sim \text{Exponential}(1)$

3.1.5 Model justification

Different pollsters and states induce variations in polling results, as pollsters may have distinct methodologies and states represent diverse voter bases. Incorporating random effects for both pollsters and states allows us to improve the robustness of the model.

These models are run through the **rstanarm** package (Goodrich et al. 2024) in R (R Core Team 2023), which makes Bayesian modeling available through the use of Stan's strong inference engine. To validate the models, RMSE and WAIC have been made use of to check the goodness of fit; Bayesian models with reduced RMSE and WAIC outperform linear models. We use weakly informative priors; for example, $\beta_0 \sim \text{Normal}(0, 10)$ and $\sigma \sim \text{Exponential}(1)$. This reflects our initial uncertainty but prevents overfitting. The priors were chosen conservatively to ensure that the model remains consistent.

Model diagnostics, including posterior predictive checks and convergence diagnostics, were carried out to ensure the reliability of the results. The Bayesian models converged successfully, as indicated by $\hat{R} = 1$ for all parameters.

The main assumption in these models is that the pollster and state effects can be treated as random. This assumes that the effects are normally distributed across pollsters and states, which may not always be accurate. Additionally, the model assumes that polling data is representative of the actual electorate, an assumption that can be violated if polls are biased or have non-random sampling issues. Despite these limitations, the hierarchical structure allows us to capture important variability, making the model suitable for predicting Harris's support. Future improvements could involve incorporating time-varying effects or exploring interactions between pollsters and states.

4 Results

4.1 Results from the relation between the prediction variable and outcome variable of the analysis data

Figure 12 shows an initial increase in support following Harris' announcement, peaking around mid-September. However, this initial momentum stabilized with minor decrease in support as the election gets closer. Moreover, the scattered data points around the summary line show high variance between individual polls.

```
`geom_smooth()` using method = 'loess' and formula = 'y ~ x'
```

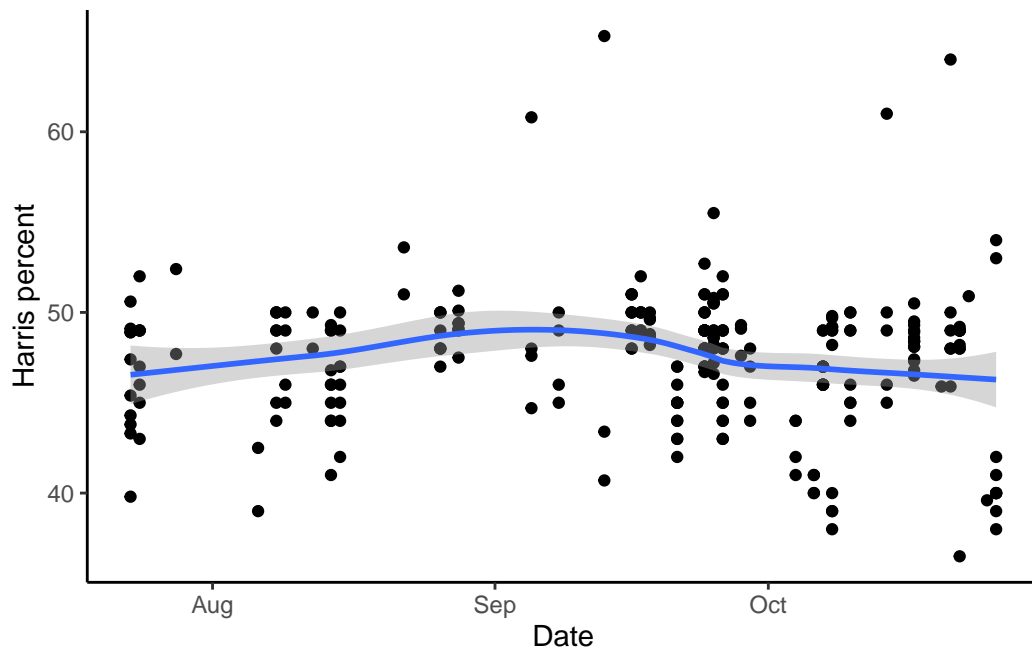


Figure 12: Polling votes for Harris over time *Note:* The date starts from after Harris officially announced her campaign for 2024 U.S. presidential election on July 21, 2024.

When Figure 13

```
`geom_smooth()` using method = 'loess' and formula = 'y ~ x'
```

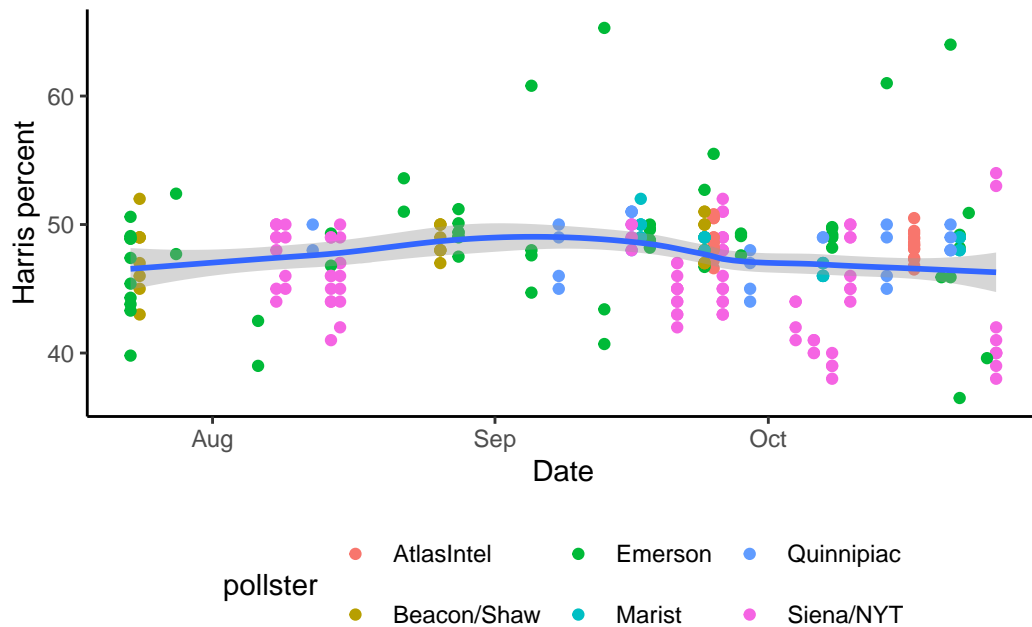



Figure 13: Polling votes for Harris over time by pollster *Note:* The date starts from after Harris officially announced her campaign for 2024 U.S. presidential election on July 21, 2024. Only filtered pollsters with high-quality polling questions are shown in the following.

```
Warning: Failed to fit group -1.
Failed to fit group -1.
Failed to fit group -1.
Failed to fit group -1.
Caused by error in `smooth.construct.tp.smooth.spec()`:
! A term has fewer unique covariate combinations than specified maximum degrees of freedom
```

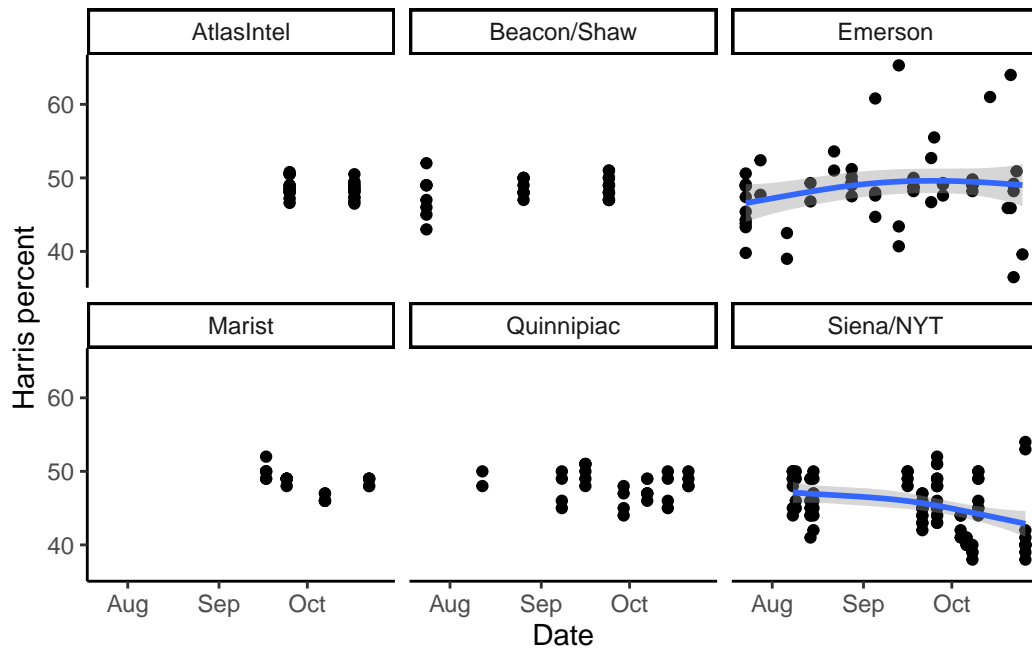


Figure 14: Polling votes for Harris over time by pollster (facets) *Note:* The date starts from after Harris officially announced her campaign for 2024 U.S. presidential election on July 21, 2024. Only filtered pollsters with high-quality polling questions are shown in the following.

4.2 Results from the prediction model

Prediction results derived from our model frameworks are summarized in Table 9 and Table 10. [!!! ADD MORE WORDS !!!]

5 Discussion

[!!! NEED MODIFICATION !!!]

Table 9: Linear models of support percentages for Harris based on date and pollster

	Linear by Date	Linear by Date, Pollster
(Intercept)	72.30 (231.63)	−53.70 (223.39)
end_date	0.00 (0.01)	0.01 (0.01)
pollsterEmerson		0.35 (0.85)
pollsterSiena/NYT		−2.68 (0.81)
Num.Obs.	172	172
R2	0.003	0.132
R2 Adj.	−0.018	0.092
Log.Lik.	−486.041	−474.968
ELPD	−489.5	−479.7
ELPD s.e.	16.7	16.3
LOOIC	979.0	959.4
LOOIC s.e.	33.4	32.5
WAIC	979.0	959.3
RMSE	4.07	3.81

Table 10: Bayesian models of support percentages for Harris based on pollster and state

	Bayesian with Pollster	Bayesian with Pollster, State
(Intercept)	−0.09 (0.06)	−0.03 (0.07)
Sigma[pollster × (Intercept),(Intercept)]	0.01 (0.01)	0.00 (0.00)
Sigma[state × (Intercept),(Intercept)]		0.06 (0.02)
Num.Obs.	172	172
ICC	0.9	0.9
Log.Lik.	−1030.101	−764.406
ELPD	−1038.3	−797.7
ELPD s.e.	82.9	26.3
LOOIC	2076.6	1595.4
LOOIC s.e.	165.9	52.6
WAIC	2076.4	1590.4
RMSE	0.04	0.02

5.1 Why Harris could beat her polls

Kamala Harris’ polling numbers show stability, a common trend for well-known candidates whose electorate solidifies in its opinion over time. Articles from The New York Times (Goldmacher and Weisman 2024) note that “once voters form opinions about candidates, those views rarely fluctuate significantly unless a major event occurs, such as a scandal or a high-profile policy shift. In the case of Harris, her polling numbers have been consistent, only changing slightly since she announced in 2024 (Figure 4). This suggests that her supporters are more or less consistent and unlikely to change much as the election draws closer.

5.2 Pollsters herding around false consensus

The polling averages for Harris and Trump don’t seem to show a big difference.

Various polling entities, including Siena/NYT and Quinnipiac, use distinct methodologies, resulting in potential systematic discrepancies within their outcomes. According to Michael Wines, the inherent biases of each pollster result in different interpretations of the elections prediction (Wines 2024). For instance, Siena/NYT’s live-interview method often skews slightly older and more conservative in its responses compared to online polling methods used by other organizations (Institute 2024). This variability makes clear that when interpreting for polling results, the source of the poll must be included at minimum.

5.3 State-Level Differences

Certain states, particularly battleground states like Florida, Pennsylvania, and Georgia, show significantly different polling results compared to national trends due to regional issues and voter priorities (Piper 2024). These reflect state-specific issues, such as local economies, immigration, and access to healthcare. Surveys carried out within these states often show wider variation in candidate support, which makes state-level analysis necessary for more accurate predictions. Incorporating state-specific factors into electoral models allows a better representation of the diverse political landscapes seen across the United States, hence improving the precision of prediction in these key areas.

5.4 Weaknesses and next steps

Future studies could focus on the incorporation of time-varying effects to capture dynamic shifts in public opinion, especially at times of large campaign events. Interaction terms between pollster and state can be included in the model, since some polling organizations may be more effective or influential in specific regions. This can give further explanation to how pollsters and regional dynamics affect overall support for a candidate. Moreover, extending the model

to incorporate voter demographics, such as age, gender, and education, could show which segments of the population are driving changes in support.

While our model provides a foundation to understand Harris' polling support, it will take further refinements are necessary to enhance its accuracy and applicability. The combination of hierarchical modeling with Bayesian methods is effective at accounting for the heterogeneity across pollsters and states, yet there remain considerable directions of inquiry to explore at the intersection of polling methodology, voter behavior, and regional electoral dynamics.

A Appendix

B Additional data details

B.1 Pollster Methodology Overview and Evaluation

The New York Times/Siena College polling partnership, the polling organization that accounted for the majority of polls in our analysis (Figure 5), conducts polls tailored for specific elections, such as state or national races (“New York Times/Siena Poll Methodology - June 2020” 2020). Their sample size typically includes 600 to 1000 likely voters per poll, with oversampling in battleground states to capture regional nuances (“New York Times/Siena Poll Methodology - June 2020” 2020). The methodology uses random-digit dialing (RDD) for landlines and mobile phones to ensure representative sample coverage across demographics. In addition, online surveys are administered to complement phone-based responses, ensuring broader accessibility (“New York Times/Siena Poll Methodology - June 2020” 2020). The stratified random sampling approach is employed, where the population is divided into strata (based on demographic variables like race, education, and geography), and a random sample is drawn from each stratum (Alexander 2023). This allows for precision in reflecting the political leanings and key demographic shifts in specific regions.

The organization intends to enhance transparency in how public opinion is assessed, ensuring that questions are carefully designed to represent contemporary political discussions, and that the terminology is polished through a process of iterative testing to achieve clarity. They devote extensive resources to cognitive testing to ensure question wording reflects what the public thinks (Institute 2024). Their polling methodology stands out in that its strategic focus on using representative samples reflect political leanings and demographics of a region for more contextual and precise polling. Siena/NYT has its reputation for accurately predicting key battleground state outcomes during previous elections, such as Florida in 2016(News 2024).

The limitations of Siena/NYT’s methodology are the challenge of polling itself. Since polling is a “snapshot in time”, the results can fluctuate based on recent political events or campaign dynamics. Additionally, while the effort to represent a broad demographic is laudable, there are still issues with nonresponse bias in polling-particularly among the hard-to-reach voter or voters suspicious of polling organizations themselves (Center 2023).

B.2 Idealized Methodology

The proposed methodology for forecasting the 2024 U.S. presidential election with a budget of \$100,000 would be designed as follows. First, a stratified random sampling method will be employed that allows for the capture of the demographic elements such as age, gender, race, and education level. This could alleviate bias and make the sample more representative of the population (Center 2023). The data collection process will encompass both telephone

surveys and internet polling, thereby effectively engaging a diverse range of voters (Center 2023). Survey respondents would include older populations via conventional methods and younger, technologically intellectual individuals through digital platforms. This multifaceted approach enhances the reliability and inclusiveness of the gathered data.

In addition to robust sampling, the methodology incorporates weighting methods that account for groups that are underrepresented, ensuring that the outcomes are not skewed by sampling errors (Center 2023). The final model would use Bayesian hierarchical modeling, which allows for more flexible modeling of uncertainty and variation across states, pollsters, and other external factors. These models, along with out-of-sample testing and cross-validation, enable accurate prediction sensitive to the dynamics of real-world changes, including political events (Center 2023). The inclusion of external factors, such as major political events, debates, or sudden economic shifts, would help the model remain responsive to rapid changes in voter sentiment.

B.3 Idealized Survey

The proposed survey questionnaire design is provided in the following link: <https://forms.gle/zQ8iJyPk3HhJYrK>

B.4 Survey Demo

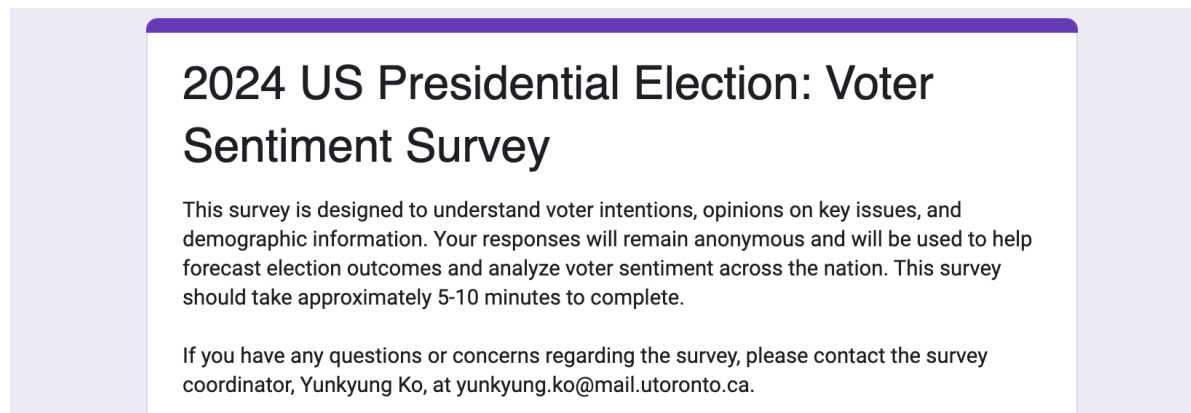
The image shows a screenshot of a survey introduction page. The page has a light purple background. At the top, there is a dark purple header bar. Below the header, the title "2024 US Presidential Election: Voter Sentiment Survey" is displayed in a large, bold, black font. Underneath the title, a paragraph of text explains the purpose of the survey: "This survey is designed to understand voter intentions, opinions on key issues, and demographic information. Your responses will remain anonymous and will be used to help forecast election outcomes and analyze voter sentiment across the nation. This survey should take approximately 5-10 minutes to complete." At the bottom of the page, there is a line of text providing contact information: "If you have any questions or concerns regarding the survey, please contact the survey coordinator, Yunkyoung Ko, at yunkyoung.ko@mail.utoronto.ca."

Figure 15: Survey Intro

Demographics

What is your age? *

☐ 18-24

☐ 25-34

☐ 35-44

☐ 45-54

☐ 55-64

☐ 65-74

☐ 75-84

☐ 85+

☐ 기타: _____

What is your gender *

☐ Male

☐ Female

☐ Prefer not to say

☐ 기타: _____

Which of the following best describes your race/ethnicity? *

☐ White

☐ African American

☐ Hispanic or Latino

☐ Asian

☐ Prefer not to say

☐ 기타: _____

What is your highest level of education? *

☐ High school or less

☐ College

☐ Bachelor's degree

☐ Graduate degree

☐ Prefer not to say

☐ 기타: _____

Figure 16: Survey Questions 1

Voting Intentions

If the 2024 U.S. presidential election were held today, for whom would you vote? *

- ☐ Kamala Harris
- ☐ Donald Trump
- ☐ Other
- ☐ Undecided

Candidate Favorability

Please rate your favorability to the following candidates on the scale below.

How do you view Kamala Harris? *

	1	2	3	4	5	
Very Unfavorable	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	Very Favorable

How do you view Donald Trump? *

	1	2	3	4	5	
Very Unfavorable	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	Very Favorable

Figure 17: Survey Questions 2

Issues of Interest

How important are the following issues to you when deciding who to vote for?
(Please rate on a scale from 1-5, with 1 being "Not important" and 5 being "Very important")

	1	2	3	4	5
Economy	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>
Healthcare	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>
Climate Change	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>
Education	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>
Immigration	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>

If none of those options provided above are of your main interest when deciding who to vote for, what is it?

내 답변

Likelihood to vote

On a scale of 1-10, how likely are you to vote in the 2024 U.S. presidential election? *

	1	2	3	4	5	6	7	8	9	10	
Likely to Vote	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	Not Likely to Vote

Figure 18: Survey Questions 3

Thank You for Your Participation

Thank you for completing the survey! Your input is valuable and will help provide a clearer understanding of voter sentiment leading up to the 2024 U.S. presidential election. If you have any further questions or concerns, feel free to contact Yunkyung Ko at yunkyung.ko@mail.utoronto.ca.

Figure 19: Survey Final Thanks

C Model details

C.1 Posterior predictive check

In the first posterior predictive check (Figure 20a), we compare the observed data with replicated data generated from the posterior distribution. This shows that the model is able to replicate the overall distribution of the observed data, with the replicated curves (light blue) closely following the true data (dark blue line). This indicates that the model fits the data well in terms of capturing the main pattern or trend (Stan Development Team 2023).

In the second plot (Figure 20b), the replicated data which had both the pollster and state variable as random intercepts, shows relatively closer approximation to the true data distribution. The narrowing of uncertainty in the posterior relative to the prior indicates the impact of the data on refining the model's predictions. This reassures that the model fits the data reasonably well and that the prior information has been appropriately updated by the observed data (Stan Development Team 2023).

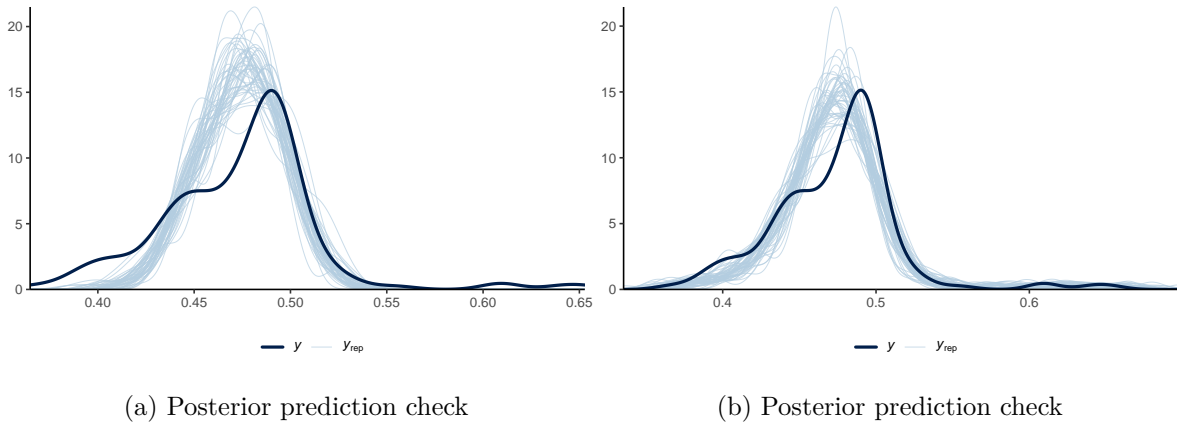


Figure 20: Examining how the Bayesian model fits, and is affected by, the data

C.2 Diagnostics

Figure 21a is a trace plot. The sampled values for posterior distribution of intercept parameter across iterations of the MCMC algorithm shows good convergence (Gabry, Češnovar, and contributors 2021b). The lines for the parameter appear to be stable and fluctuating around a central value without any clear trends or patterns. This suggests that the MCMC algorithm has likely converged, and the posterior samples are representative of the target distribution.

Figure 21b is a Rhat plot. The Rhat value is approximately 1.0 for the intercept, which shows that the variance within and between multiple chains have converged. An Rhat value close to 1 indicates that the chains have mixed well and are drawing from the same distribution while values significantly greater than 1 would indicate that further iterations are needed (Gabry, Češnovar, and contributors 2021a). This suggests that the Bayesian models for both “pollster” and “state” have likely converged, and the results derived from these models are reliable.

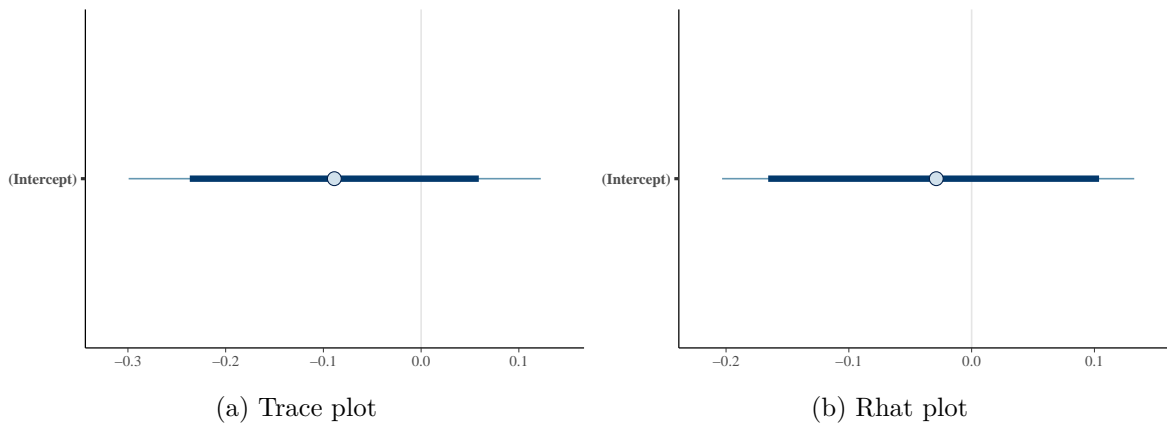


Figure 21: Checking the convergence of the MCMC algorithm

References

- 11Alive. 2024. “What Does a Battleground State Mean? What Are the Swing States?” 2024. <https://www.11alive.com/article/news/politics/elections/what-does-a-battleground-state-mean-what-are-the-swing-states/85-89793c86-c6b3-4394-a1f9-e4cf6136f35e>.
- Alexander, Rohan. 2023. *Telling Stories with Data*. Chapman; Hall/CRC. <https://tellingstorieswithdata.com/>.
- Arel-Bundock, Vincent. 2022. “modelsummary: Data and Model Summaries in R.” *Journal of Statistical Software* 103 (1): 1–23. <https://doi.org/10.18637/jss.v103.i01>.
- Bijune, Aiste, and Lan Ha. 2024. “US 2024 Election: Implications for the Global Economy.” *Euromonitor International*. <https://www.euromonitor.com/article/us-2024-election-implications-for-the-global-economy>.
- Center, Pew Research. 2023. “How Public Polling Has Changed in the 21st Century.” *Pew Research Center*. <https://www.pewresearch.org/methods/2023/04/19/how-public-polling-has-changed-in-the-21st-century/>.
- CSIS. 2024. “The Global Impact of the 2024 u.s. Presidential Election.” <https://features.csis.org/2024-us-election-global-impact/>.
- FiveThirtyEight. 2024. “Dataset: US Presidential General Election Polls.” https://projects.fivethirtyeight.com/polls/data/president_polls.csv.
- Gabry, Jonah, Šimon Češnovar, and contributors. 2021a. *MCMC Diagnostics in Bayesplot*. <https://mc-stan.org/bayesplot/reference/MCMC-diagnostics.html>.
- . 2021b. *MCMC Trace Plots in Bayesplot*. <https://mc-stan.org/bayesplot/reference/MCMC-traces.html>.
- Goldmacher, Shane, and Jonathan Weisman. 2024. “Harris Gains in Polling Against Trump, According to National Poll.” *The New York Times*. https://www.nytimes.com/2024/10/08/us/politics/harris-trump-poll-national.html?campaign_id=60&emc=edit_na_20241008&instance_id=0&nl=breaking-news®i_id=193003212&segment_id=179873&user_id=2d04b9e4643940584b7476aee2111e62.
- Goodrich, Ben, Jonah Gabry, Imad Ali, and Sam Brilleman. 2024. “rstanarm: Bayesian applied regression modeling via Stan.” <https://mc-stan.org/rstanarm/>.
- Institute, Siena College Research. 2024. “New York Times/Siena College National Poll - October 2024.” 2024. <https://scri.siena.edu/2024/10/08/new-york-times-siena-college-national-poll-3/>.
- Jackman, Simon. 2024. “Pooling the Polls over an Election Campaign.” chrome-extension://efaidnbmnnnibpcajpcglclefindmkaj/https://cdn-uploads.piazza.com/paste/ket0p3a9re9628/8a55272a479e1aaeabbd01b3030aa1d38c4a7fc323a9c3fee3bd5fec2fdca44/Pooling_the_polls_over_an_election_campaign.pdf.
- “New York Times/Siena Poll Methodology - June 2020.” 2020. <chrome-extension://efaidnbmnnnibpcajpcglclefindmkaj/https://int.nyt.com/data/documenttools/nyt-siena-poll-methodology-june-2020/f6f533b4d07f4cbe/full.pdf>.
- News, Siena College. 2024. “A Perfect Partnership in Polling.” 2024. <https://www.siena.edu/news/story/a-perfect-partnership-in-polling/>.

- Piper, Jessica. 2024. “Trump Maintains Narrow Lead in Key Battleground States, NYT Poll Finds.” *Politico*. <https://www.politico.com/news/2024/05/13/trump-narrow-lead-battleground-00157541>.
- R Core Team. 2023. *R: A Language and Environment for Statistical Computing*. Vienna, Austria: R Foundation for Statistical Computing. <https://www.R-project.org/>.
- Stan Development Team. 2023. *Posterior Predictive Checks*. <https://mc-stan.org/docs/stan-users-guide/posterior-predictive-checks.html>.
- Wickham, Hadley. 2016. *ggplot2: Elegant Graphics for Data Analysis*. Springer-Verlag New York. <https://ggplot2.tidyverse.org>.
- 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>.
- Wines, Michael. 2024. “Polling Problems: How Margins of Error Could Mislead in the 2024 Election.” *The New York Times*. <https://www.nytimes.com/2024/10/14/us/elections/poll-problems-margin-of-error.html>.
- Zhu, Hao. 2024. *kableExtra: Construct Complex Table with ‘kable’ and Pipe Syntax*. <https://CRAN.R-project.org/package=kableExtra>.