

Forecasting the US Presidential Election: A Poll-of-Polls Approach Using Linear Models By Using Data of Pollsters From Five Thirty Eight In 2024*

Ziqi Zhu

Yuanchen Miao

Claire Chang

October 23, 2024

This paper presents a multiple linear regression model for forecasting the outcome of the 2024 U.S. Presidential Election. Using nationwide polling data, we predict the percentage of support for U.S. presidential candidates based on key factors such as sample size, pollster quality, and timing. Our model aggregates these predictions to simulate potential election outcomes. The results reveal how specific poll attributes affect the accuracy of forecasts, offering valuable insights for predicting the likelihood of various electoral scenarios.

Table of contents

1	Introduction	2
2	Data	3
2.1	Overview	3
2.2	Measurement	3
2.3	Outcome variables	3
2.4	Predictor variables	3
2.4.1	Sample Size	3
2.4.2	Poll Score	4
2.4.3	Numeric Grade	4
2.4.4	Transparency Score	4
2.4.5	Days Since Start	4

*Code and data are available at: https://github.com/zzq20010617/2024_USelection_prediction

3	Models	4
3.1	MLR Model set-up	4
3.1.1	Result	5
3.2	Bayesian Model set-up	6
3.2.1	Result	7
4	Prediction	7
4.1	Result	9
5	Discussion	9
5.1	Crossing point in the prediction	9
5.2	Weaknesses and next steps	10
	Appendix	11
.1	Idealized Survey and Methodology	11
.2	Trafalgar group’s methodology overview and evaluation	11
A	Model details	12
A.1	Assumption check for MLR models	12
A.2	Significant check for MLR models	13
A.3	Posterior predictive check	16
A.4	Diagnostics	16
	References	17

1 Introduction

This paper examines the development of a multiple linear regression (MLR) model to predict the percentage of support (pct) for U.S. presidential candidates based on polling data. The data includes a range of predictors, such as sample size, pollster quality factors, and time-related factors. The focus of the analysis is on nationwide polling data, aggregated from different pollsters, to create a robust model for forecasting election outcomes. The goal is to provide a clearer understanding of how various poll attributes influence polling results and to derive insights that can predict election outcomes.

The primary estimand in this study is the expected percentage of support (pct) for a U.S. presidential candidate, given polling data and relevant predictor variables like rating of the pollster and sample size of the poll.

Our Results shows that

This analysis is useful in understanding how various polling factors contribute to predicting election results. By modeling the relationship between polling result and factors, the study enhances the ability to forecast election outcomes based on public opinion data. This model

provides a practical tool for researchers, political strategists, and analysts to assess the reliability of polling data and its implications for elections.

2 Data

2.1 Overview

Our data is download from [Five Thirty Eight](#)(FiveThirtyEight, n.d.), the website gathered survey data from different pollsters of 2024 US president election, We use the statistical programming **language R** (R Core Team 2023), and packages **lubridate**(Grolemund and Wickham 2011), **dplyr**(Wickham et al. 2023), **tidyverse**(Wickham et al. 2019), **model summary**(Arel-Bundock 2022) to process the data. Also, we use **tibble** (Müller and Wickham 2023) to create the table for simulate table and save it as csv files. Graphs are created by using ggplot2 in from package **tidyverse**(Wickham et al. 2019)Following **Telling stories with data**(Data 2024), we consider using multiple linear model to predict and forecasting the result of US president election in 2024. We are creating different graphs and models for the candidates from different parties in general state and combine all models and graphs together to create a new scatter plot to predict and forecasting the US president election in 2024.

2.2 Measurement

The dataset is about public opinions about U.S. presidential candidates, which are captured through polling.

2.3 Outcome variables

Our primary outcome variable is **support percentage (pct)**, which represents the percentage of respondents who support each candidate. This variable is modeled as the response variable in the MLR analyse.

2.4 Predictor variables

2.4.1 Sample Size

Sample size is an important predictor because it influences the reliability of a poll. Larger sample sizes tend to reduce sampling error, providing more accurate reflections of voter sentiment. In our MLR model, the sample size is log-transformed to account for diminishing returns—larger polls do not necessarily offer proportionally better accuracy.

2.4.2 Poll Score

The poll score is a measure of the error and bias we can attribute to a pollster, negative number is better.

2.4.3 Numeric Grade

This variable is an aggregate score of the poll based on a numeric scale. It serves as an indicator of the pollster's historical performance and the methodology used. 3 is the maximum and some pollsters have no rating.

2.4.4 Transparency Score

Transparency score measures how openly a pollster reports their methodology and results. A higher transparency score suggests that the pollster has disclosed key details about how the poll was conducted, improving trust in the poll results.

2.4.5 Days Since Start

This variable represents the number of days since the beginning of the election polling period. It helps capture the dynamic nature of voter preferences over time, accounting for shifts in public opinion as the election date nears.

3 Models

The goal of our modelling strategy is try to capture the trend of support percentage for different parties in 2024 Presidential election as much as possible. In the following section we briefly describe the models we used to investigate. Background details and diagnostics are included in Appendix A.

3.1 MLR Model set-up

We first filtered our data by different parties, as the single linear model but include other 3 parties(GRE, IND, and LIB, Conservative only has two entries in cleaned nationwide data and has 0 pct so its been removed), and fit linear model for each of them. The formula for Multiple linear regression is as follow

$$\text{pct}_i = \beta_0 + \beta_1 \log(\text{sample_size}_i) + \beta_2 \text{pollscore}_i + \beta_3 \text{numeric_grade}_i \quad (1)$$

$$+ \beta_4 \text{transparency_score}_i + \beta_5 \text{days_since_start}_i + \epsilon_i \quad (2)$$

Table 1: Explanatory models of pct based on pollscore, sample size, and time

	DEM	REP	GRE	IND	LIB
(Intercept)	37.33 (1.67)	50.12 (1.62)	3.34 (0.87)	-5.76 (3.68)	4.13 (1.70)
log(sample_size)	0.38 (0.18)	-0.23 (0.18)	-0.19 (0.11)	2.42 (0.45)	-0.12 (0.21)
pollscore	-0.72 (0.38)	-3.33 (0.37)	0.02 (0.13)	-1.16 (0.58)	-1.00 (0.30)
numeric_grade	0.20 (0.49)	-3.55 (0.47)	-0.01 (0.17)	-1.80 (0.81)	-1.07 (0.39)
transparency_score	-0.34 (0.07)	0.07 (0.06)	0.00 (0.03)	0.13 (0.14)	0.01 (0.06)
days_since_start	0.03 (0.00)	0.01 (0.00)	0.00 (0.00)	-0.03 (0.00)	0.00 (0.00)
Num.Obs.	1235	1235	270	603	105
R2	0.357	0.121	0.155	0.179	0.133
R2 Adj.	0.355	0.117	0.139	0.172	0.090
AIC	6780.2	6711.8	552.1	3623.3	269.7
BIC	6816.0	6747.6	577.3	3654.1	288.3
Log.Lik.	-3383.085	-3348.901	-269.070	-1804.655	-127.852
RMSE	3.74	3.64	0.66	4.83	0.82

Response Variable is pct. Predictors are log(sample_size), pollscore, numeric_grade, transparency_score, and days_since_start, detail can be find in Section 2.3, and Section 2.4

3.1.1 Result

Our results for the MLR model are summarized in Table 1.

We first check the residual vs. fitted plot as shown in Appendix A, and see no obvious violation of assumptions. According to the summary, DEM and REP models are based on the largest number of observations (both 1235). Smaller parties have fewer observations, which may lead to less robust models. The r^2 shows that the model explain the highest amount of variance on Democratic Party than others with value 0.357, and other parties less than 0.2. Based on the summary of models for DEM and REP Section A.2, we are not sure if those predictor that

measure quality of polls are significant or not, so we keep them for now in Bayesian model. Prediction of MLR models is shown in Figure 1, which we can see blue line which is Democratic party has a steeper slope and start to take the lead between day 100 and day 200.

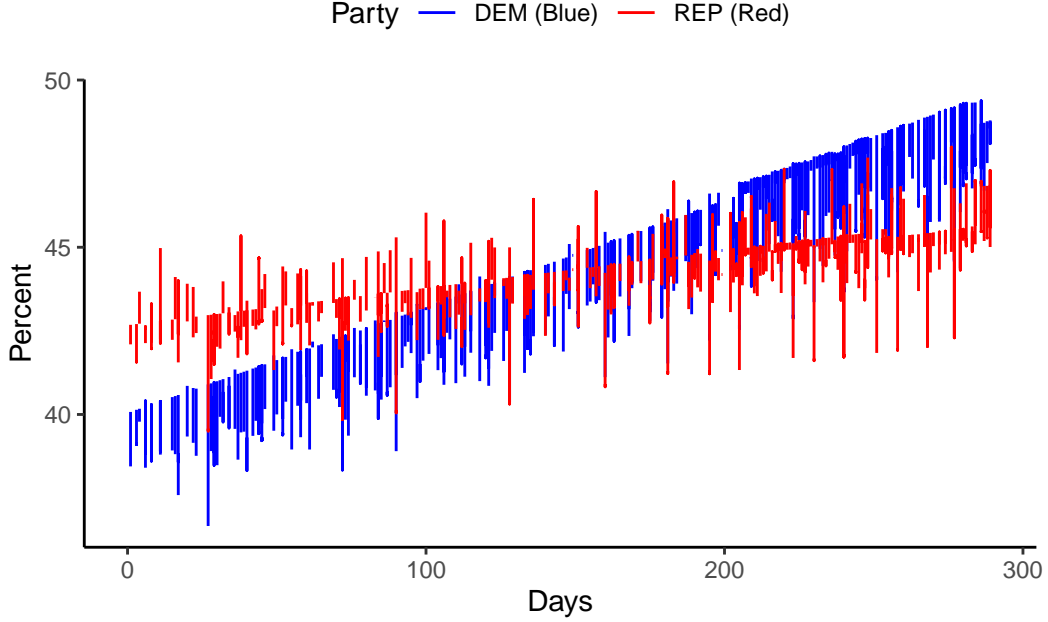


Figure 1: Prediction with MLR models for DEM and REP party

3.2 Bayesian Model set-up

The Bayesian model is build from the MLR models from above section. With the inspiration from the example R code, we introduce the random effects for Pollsters, and we use a logistic link to predict the probability of support, modeling count data which is in binomial distribution. We choose to use a default prior and enable autoscale, the formula is shown below

$$\log\left(\frac{p_i}{1-p_i}\right) = \beta_0 + \beta_1 \log(\text{sample_size}_i) + \beta_2 \text{pollscore}_i + \beta_3 \text{numeric_grade}_i \quad (3)$$

$$+ \beta_4 \text{transparency_score}_i + \beta_5 \text{days_since_start}_i + \alpha_j \quad (4)$$

$$\alpha_j \sim \text{Normal}(0, \sigma_{\text{pollster}}) \quad (5)$$

$$\beta_0 \sim \text{Normal}(0, 2.5) \quad (6)$$

$$\beta_1 \sim \text{Normal}(0, 2.5) \quad (7)$$

$$\beta_2 \sim \text{Normal}(0, 2.5) \quad (8)$$

$$\beta_3 \sim \text{Normal}(0, 2.5) \quad (9)$$

$$\beta_4 \sim \text{Normal}(0, 2.5) \quad (10)$$

$$\beta_5 \sim \text{Normal}(0, 2.5) \quad (11)$$

$$\sigma_{\text{pollster}} \sim \text{Exponential}(1) \quad (12)$$

y_i represent the number of individuals in the sample that support the Democratic party (this corresponds to `num_party`). The response is modeled as binomial: $y_i \sim \text{Binomial}(\text{sample_size}_i, p_i)$, where p_i is the probability that an individual in poll i supports the Democratic party, and sample_size_i is the total number of individuals surveyed in poll i .

3.2.1 Result

Summary of Bayesian model is shown in Table 2, and the ppcheck is in appendix Figure 5.

4 Prediction

To get a prediction by the Bayesian model we have, we first spline fit the model, and then create a data frame and generate posterior predictions by the new data frame. ## Spline Fit for Bayesian model set-up

We run the model in R (R Core Team 2023) using the `rstanarm` package of (Goodrich et al. 2022). We use priors $\text{Normal}(0, 5)$ to allow for more flexibility of predictors effects with “`autoscale = TRUE`”. Then we create a data frame like Table 3, we set the quality scale of predict data to the best to perform a poll that has high quality. the result is shown as Section 4.1

Table 2

	(1)
(Intercept)	−0.69
log(sample_size)	0.05
pollscore	−0.11
numeric_grade	−0.08
transparency_score	0.00
days_since_start	0.00
Sigma[pollster × (Intercept),(Intercept)]	0.01
Num.Obs.	1235
ICC	0.9
Log.Lik.	−10 595.196
ELPD	−10 883.4
ELPD s.e.	322.8
LOOIC	21 766.8
LOOIC s.e.	645.6
WAIC	21 795.1
RMSE	0.03

Table 3

	pollscore	sample_size	numeric_grade	transparency_score	days_since_start
1	−1.5	1200	3	10	0.000000
2	−1.5	1200	3	10	2.929293
3	−1.5	1200	3	10	5.858586
4	−1.5	1200	3	10	8.787879
5	−1.5	1200	3	10	11.717172
6	−1.5	1200	3	10	14.646465
	pollster				
1	TIPP				
2	TIPP				
3	TIPP				
4	TIPP				
5	TIPP				
6	TIPP				

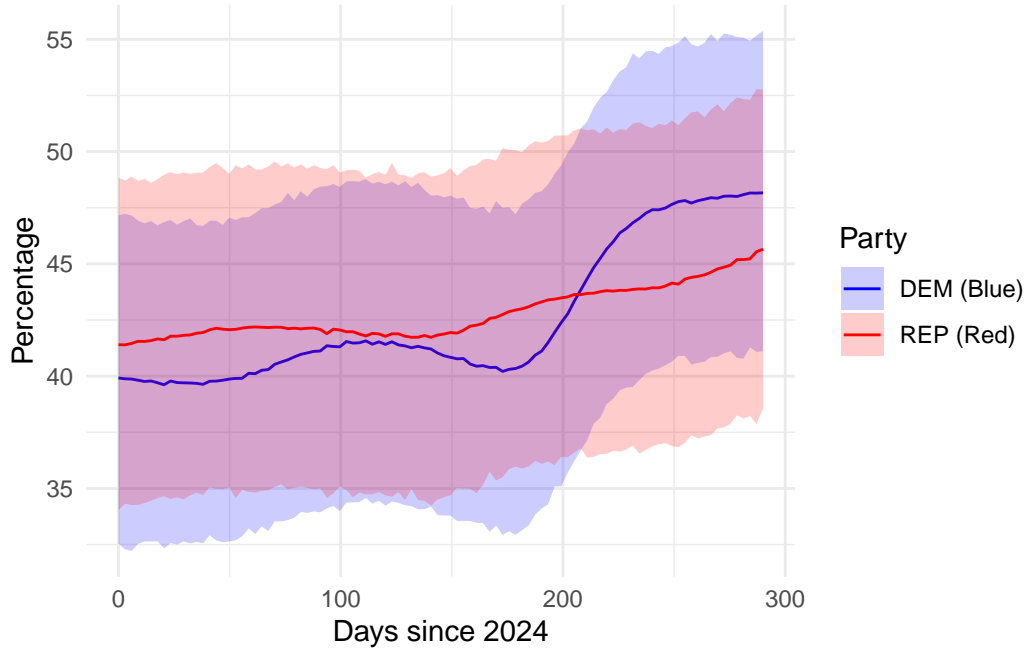


Figure 2: Poll Percentage over Time with Spline Fit for DEM and REP party

4.1 Result

The prediction graph Figure 2 contains prediction of spline fit model for two parties, blue line represents the DEM party's predicted percentage of support, and the red line represents the REP party's predicted percentage of support over time (measured in days since the start of 2024). We can tell that the predict support of DEM party is taking the lead by around 2 percent until the last day of the dataset.

5 Discussion

5.1 Crossing point in the prediction

From the prediction of spline fit model Figure 2, There is a clear increase in the support for the DEM party starting around day 200, while the REP party's support remains more stable with only a slight increase over time. The increase in support for the DEM party in the predictions could be attributed to Harris's rise as the Democratic candidate after Biden's exit from the race On July 21, 2024, which is about 230 days from start of 2024.

5.2 Weaknesses and next steps

Weaknesses and next steps should also be included.

Appendix

.1 Idealized Survey and Methodology

.2 Trafalgar group’s methodology overview and evaluation

The Trafalgar Group conducts polls ranging from major political campaigns to marketing surveys. The organization’s 0.7 pollster rating indicates moderate accuracy and reliability compared to other pollsters. The population consists of all eligible voters in the U.S., while the sampling frame includes registered voters across states, segmented by demographics like age, gender, party affiliation, and ethnicity. Trafalgar Group typically samples likely voters, focusing on those most likely to participate in upcoming elections. Trafalgar Group recruits its sample using a mix of interactive voice response (IVR), live phone calls, text messages, online panels, and email surveys. Trafalgar Group employs stratified random sampling to ensure proportional representation of subgroups like political party, gender, and region. This method improves accuracy by reflecting the electorate’s demographic and political makeup but may risk over-stratification, giving smaller voter groups disproportionate influence and potentially skewing results. Trafalgar Group handles non-response by using weighting to adjust the sample, ensuring respondent demographics match population parameters. This method reduces bias by accounting for underrepresented groups, though heavy reliance on post-survey adjustments may introduce new biases, particularly when there are large discrepancies in response rates across demographics.

A Model details

A.1 Assumption check for MLR models

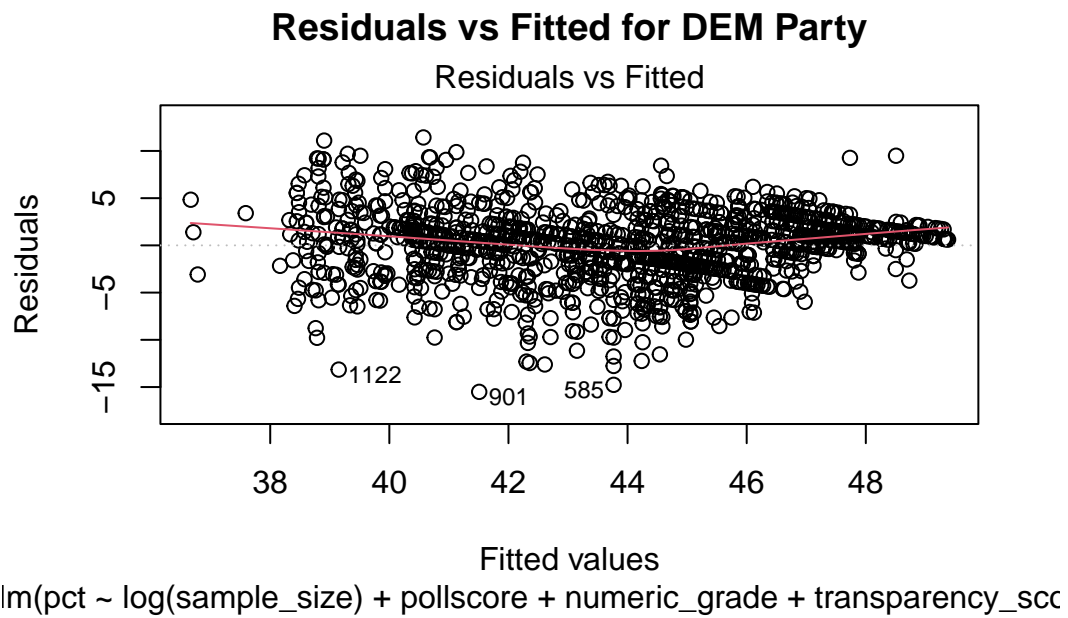


Figure 3: Residuals vs Fitted for DEM Party

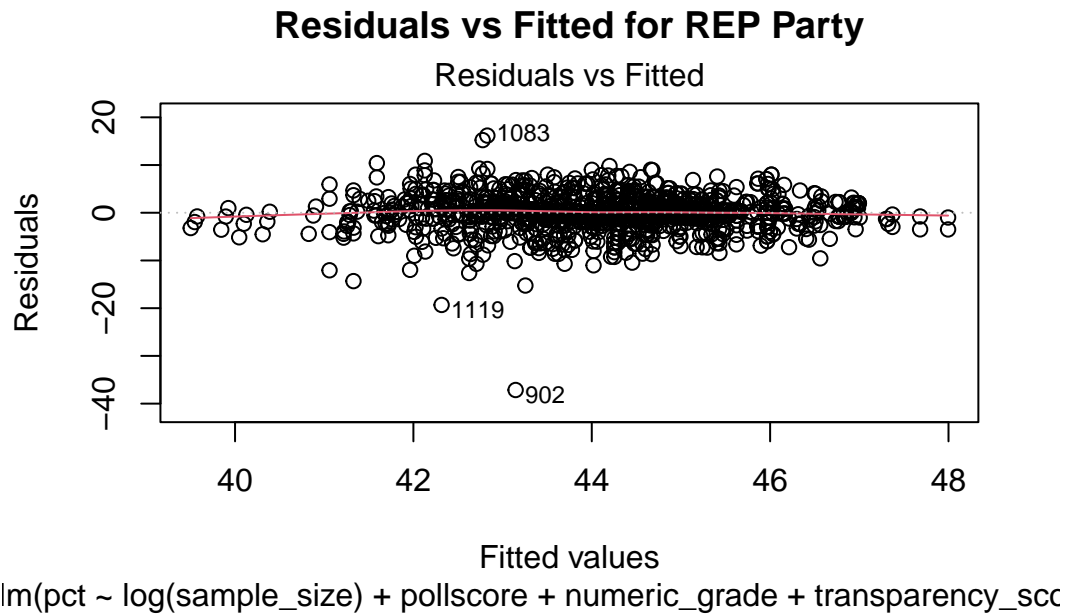


Figure 4: Residuals vs Fitted for REP Party

A.2 Significant check for MLR models

see Table 4, and Table 5

Table 4: Summary of DEM data fit by MLR model

Call:

```
lm(formula = pct ~ log(sample_size) + pollscore + numeric_grade +
    transparency_score + days_since_start, data = party_data)
```

Residuals:

	Min	1Q	Median	3Q	Max
	-15.5091	-2.0181	0.4498	2.1124	11.4283

Coefficients:

	Estimate	Std. Error	t value	Pr(> t)
(Intercept)	37.33222	1.66922	22.365	< 2e-16 ***
log(sample_size)	0.37656	0.18292	2.059	0.0397 *
pollscore	-0.72219	0.37612	-1.920	0.0551 .
numeric_grade	0.19976	0.48656	0.411	0.6815
transparency_score	-0.33589	0.06612	-5.080	4.35e-07 ***
days_since_start	0.03163	0.00132	23.971	< 2e-16 ***

Signif. codes: 0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

Residual standard error: 3.754 on 1229 degrees of freedom

(10 observations deleted due to missingness)

Multiple R-squared: 0.3572, Adjusted R-squared: 0.3546

F-statistic: 136.6 on 5 and 1229 DF, p-value: < 2.2e-16

Table 5: Summary of REP data fit by MLR model

Call:

```
lm(formula = pct ~ log(sample_size) + pollscore + numeric_grade +
    transparency_score + days_since_start, data = party_data)
```

Residuals:

Min	1Q	Median	3Q	Max
-37.146	-1.297	-0.054	1.861	16.171

Coefficients:

	Estimate	Std. Error	t value	Pr(> t)
(Intercept)	50.124690	1.623651	30.872	< 2e-16 ***
log(sample_size)	-0.230845	0.177923	-1.297	0.195
pollscore	-3.333315	0.365853	-9.111	< 2e-16 ***
numeric_grade	-3.547745	0.473277	-7.496	1.25e-13 ***
transparency_score	0.066665	0.064311	1.037	0.300
days_since_start	0.011144	0.001284	8.682	< 2e-16 ***

Signif. codes: 0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

Residual standard error: 3.652 on 1229 degrees of freedom

(10 observations deleted due to missingness)

Multiple R-squared: 0.1206, Adjusted R-squared: 0.1171

F-statistic: 33.72 on 5 and 1229 DF, p-value: < 2.2e-16

A.3 Posterior predictive check

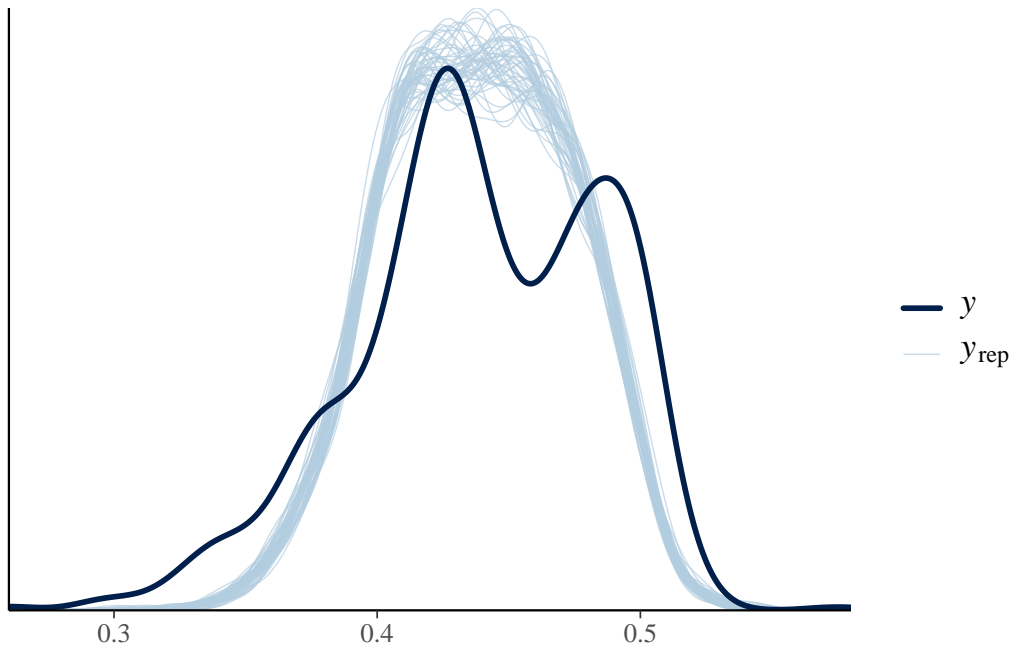


Figure 5: Posterior prediction check for Bayesian model

A.4 Diagnostics

Checking the convergence of the MCMC algorithm

References

- 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>.
- Data, Telling Stories with. 2024. “Telling Stories with Data.” <https://tellingstorieswithdata.com/>.
- FiveThirtyEight. n.d. “FiveThirtyEight Presidential General Election Polls - National (2024).” <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/>.
- Grolemund, Garrett, and Hadley Wickham. 2011. “Dates and Times Made Easy with lubridate.” *Journal of Statistical Software* 40 (3): 1–25. <https://www.jstatsoft.org/v40/i03/>.
- Müller, Kirill, and Hadley Wickham. 2023. *Tibble: Simple Data Frames*. <https://CRAN.R-project.org/package=tibble>.
- R Core Team. 2023. *R: A Language and Environment for Statistical Computing*. Vienna, Austria: R Foundation for Statistical Computing. <https://www.R-project.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>.
- 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>.