# Predicting the 2024 U.S. Presidential Election

## A Poll-of-Polls Forecasting Approach

Mingrui Li

2024-11-03

**Abstract**

This project aims to generate a linear model to predict the outcome of the 2024 U.S. Presidential Election, using data from several polls combined with a linear regression model. The model, developed in Python, makes use of three key predictors of pollster reliability represented through numeric grade, sample size, and poll quality represented through pollscore to estimate Donald Trump's projected vote share against Kamala Harris. Our model predicts, using the aggregated and cleaned dataset, that Trump will receive 49.7% of the popular vote. We discuss here the methodology of our model and its strengths and limitations: the simplicity of the approach and the fact that it is data-driven are stressed, while obvious limitations, such as the absence of any state-level data or possible biases in the polls, are highlighted. Further developments will add in integration with state-by-state polling to complete electoral college projections and adjustments for factors such as voter turnout and swing state dynamics. This work underlines the relevance of aggregated polling data in election predictions and gives a solid, albeit preliminary, forecast of the 2024 election.

# Table of contents

# 1 Introduction

Election forecasting has always constituted one of those difficult but significant domains of political analysis, especially for high-stake elections like the forthcoming 2024 U.S. Presidential Election. This paper tries to forecast the popular vote shares by the leading candidates between Donald Trump and Kamala Harris through a poll-of-polls approach based on aggregative data from numerous reliable pollsters. We will go over three most important predictors of poll accuracy using a linear regression model in Python: reliability of the pollster, sample size, and poll quality rating. We'll select these variables so as to optimize model precision and to make sure that the polling data is representative and methodologically sound.

This analysis uses the cleaned and processed polling data without fringe candidates, standardized vote percentages between the two main candidates, to predict that Donald Trump will win about 49.7% of the popular vote. This approach is very simplified, as the model has strong points in terms of the data-driven methodological construction and the reliance on accessible predictors. However, this paper also acknowledges limitations to state-by-state polling data and traditional methods of polling, which might hold inherent biases.

This report proceeds as follows: Section 2 describes the procedure of data cleaning and model selection that should guarantee the identification of robust predictors. Section 3 gives an account of EDA and graphical representation of main features of this dataset. Section 4 set up a linear model, introduce assumptions, and check predictive accuracy by cross-validation and simulation. Sections 5 and 6 present the results and discuss, respectively, the predictive capabilities of the model, along with its limitations, emphasizing potential improvements that could be added in the future, such as state-by-state polling and voter turnout adjustments. It is our aim in this paper to provide an early but orderly and sound forecast of the Presidential Election of 2024.

# 2 Data Cleaning and Model Selection

We begin by using a dataset of polling data from the 2024 U.S. presidential general election to make an election forecast. The predictors include `numeric_grade`, which represents reliability pollster rating, `sample_size`, and internal ratings of poll quality represented as `pollscore`. We do some cleaning on the dataset, processing it into a format focused around the two main candidates, Donald Trump and Kamala Harris, then scaling the results so their vote percentages add up to 100%. It includes the cleaned predictors and the response variable, scaled_trump_pct, and is stored in analysis_data.parquet. (FiveThirtyEight, 2024)

The desired dataset consists of polls with different methodologies. The predictors include:

- **numeric_grade**: A measure of the pollster's reliability.
- **sample_size**: The number of respondents in the poll.

- **pollscore**: A rating of the poll's overall quality.

These three predictors, numeric grade, sample size, and pollscore—come because they represent important aspects of poll quality, reliability, and representation that are necessary to make a reasonably accurate forecast of the election outcome. Numeric grade is indicative of the credibility of the pollster and its track record. Higher values indicate greater reliability, and hence that pollster is more likely to produce data representative of public opinion. Sample_size represents the number of respondents, and larger samples reduce sampling error and increase confidence in the poll's findings. The pollscores provide a summary quality rating that measures other factors that could affect a poll's robustness, such as methodological rigor and control for bias. By including these other variables, the model captures strength and quality variables of the polling data that would serve useful in making a more believable forecast of the election outcome.

Since the raw data set consists of a huge amount of data with missing value, sufficient data cleaning is needed. During the process of data cleaning, we removed all data (polls) which does not have a value for all three included response variable above. Also, for the response variable (percentage of vote gained) Harris and Trump gain the majority of votes (above 95%) for almost every polls. Therefore, we removed all percentage of vote gained by other candadates except for Trump and Harris, and scaled the percentage gained by Trump and Harris so that they sums to 100%.

After data cleaning, we reduced the number of polls in raw data of 3227 to 965 in cleande data. The response variable is the **scaled percentage of Trump's vote**, calculated as:

$$\text{scaled\_trump\_pct} = \frac{\text{Trump\_pct}}{\text{Trump\_pct} + \text{Harris\_pct}}$$

# 3 Exploratory Data Analysis

We conduct a Exploratory Data Analysis (EDA) here.

```
First few rows of the data:
   numeric_grade  sample_size  pollscore  scaled_trump_pct
0            1.8       1729.0        0.4          0.506579
1            1.5       2050.0       -0.1          0.540230
2            0.5       1000.0        1.6          0.550562
3            1.8       1319.5        0.4          0.519737
4            1.4       1056.0        0.0          0.491892

Summary statistics:
       numeric_grade  sample_size   pollscore  scaled_trump_pct
count     965.000000   965.000000  965.000000        965.000000
mean        2.071710  1809.613718   -0.270777          0.496952
```

```
std        0.659245     2456.758810     0.695746             0.036848
min        0.500000      247.666667    -1.500000             0.290155
25%        1.800000      707.000000    -0.800000             0.482759
50%        1.900000     1000.000000    -0.300000             0.498483
75%        2.800000     1405.000000     0.200000             0.512685
max        3.000000    19442.500000     1.700000             0.676829

Missing values in each column:
numeric_grade       0
sample_size         0
pollscore           0
scaled_trump_pct    0
dtype: int64
```
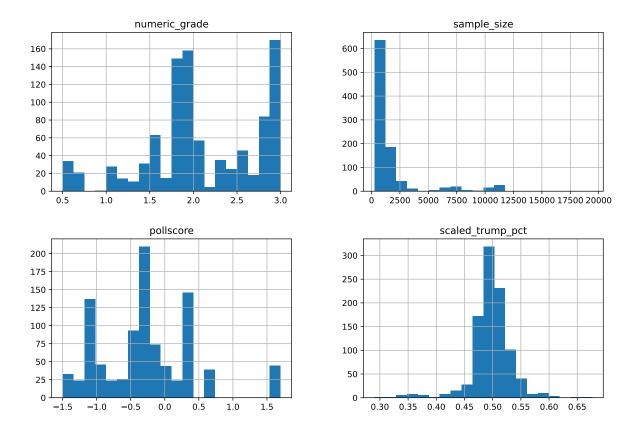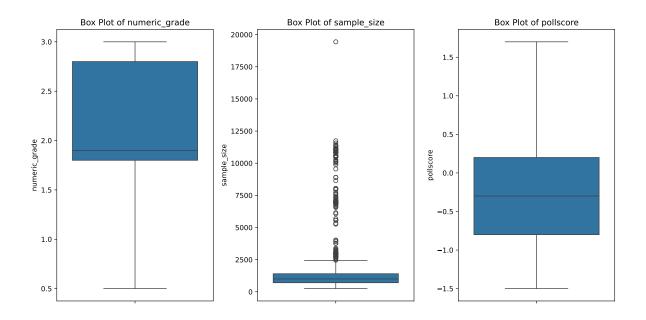
**Python Output 1**

Above we can see the first few rows of the dataset, summary statistics, and information on missing values. Summary statistics will give us an idea about the distribution of each variable, and we can immediately see that numeric_grade has a mean at approximately 2.07, while sample_size seems to have a high standard deviation, meaning a large range in sample sizes taken for the polls. The variables here are centered around zero, as can be seen from the summary statistics-pollscore-and the median for scaled_trump_pct is in the neighborhood of 0.5, indicating that Trump's scaled vote share usually hovers around 0.5. There are no NA values in this data set, meaning the data cleaning process is successful.

Distribution of Predictors and Response Variable

**Python Output 2**

The distribution of the variables is represented by using histograms for each variable. numeric_grade is distributed somewhat bimodally. There are a lot of polls rated near 2 and 3. The sample_size histogram is really skew, with a bunch of small sample sizes, but some sample sizes above 10,000. The histogram of pollscore is roughly normally distributed around -0.3, with some peaks that indicate variation in the quality scores of the polls. The histogram for scaled_trump_pct suggests it is generally between 0.45 and 0.55 with a little central tendency for Trump scaled percentage.

**Python Output 3**

Boxplots of numeric_grade, sample_size and pollscore The numeric_grade boxplot gives an impression of an approximately spread-out distribution with no serious outliers. The sample_size boxplot has outliers in a few large values to draw attention to the heterogeneity in poll sizes. The pollscore boxplot has a more central and tighter distribution mostly from about -1.5 to 1.5, with some consistency of good quality polls.

(Rossum & Foundation, 2024)

# 4 Model

The goal of our modeling strategy is to predict the scaled percentage of votes for Donald Trump. We use a simple **linear regression model** with three predictors: `numeric_grade`, `sample_size`, and `pollscore`. The linear model is represented as:

$$\text{Trump\_pct} = \alpha + \beta_1 \cdot \text{numeric\_grade} + \beta_2 \cdot \text{sample\_size} + \beta_3 \cdot \text{pollscore}$$

## 4.1 Model Assumption

Our linear regression model requires a number of assumptions to be met if our predictions are to be accurate and meaningful. First, we assume a linear relationship between our predictors

(numeric grade, sample size and pollscore) and our response variable, scaled Trump percentage. We believe that for each of these three predictors, a change in the predictor produces a proportional change in Trump's projected vote share. It further assumes that the residuals are independent and not autocorrelated, and homoscedastic—the variance of these errors is constant across all levels of predictors.

We also assumes the normality of residuals. No perfect multicollinearity among the predictors is taken into consideration; otherwise, high correlations may hide the individual effects of the predictors and, consequently, the estimates are not that reliable. Finally, our model assumes that all the relevant variables determining the response have been included in order to avoid any omitted variable bias. Nevertheless, minor deviations will be noted, and future improvements may consider other models in case significant violations of assumptions are found.

## 4.2 Model Summary

We fit the linear model to the data, which results in the following summary:

```
                          OLS Regression Results
==============================================================================
Dep. Variable:        scaled_trump_pct   R-squared:                       0.053
Model:                             OLS   Adj. R-squared:                  0.050
Method:                  Least Squares   F-statistic:                     17.75
Date:                 Sun, 03 Nov 2024   Prob (F-statistic):           3.22e-11
Time:                         21:38:30   Log-Likelihood:                 1842.7
No. Observations:                  965   AIC:                            -3677.
Df Residuals:                      961   BIC:                            -3658.
Df Model:                            3
Covariance Type:             nonrobust
==============================================================================
                 coef    std err          t      P>|t|      [0.025      0.975]
------------------------------------------------------------------------------
const          0.5324      0.007     71.847      0.000       0.518       0.547
numeric_grade -0.0160      0.004     -4.129      0.000      -0.024      -0.008
sample_size -2.28e-06   4.76e-07     -4.788      0.000   -3.21e-06   -1.35e-06
pollscore     -0.0064      0.004     -1.746      0.081      -0.014       0.001
==============================================================================
Omnibus:                      300.765   Durbin-Watson:                   1.768
Prob(Omnibus):                  0.000   Jarque-Bera (JB):             2247.037
Skew:                          -1.217   Prob(JB):                         0.00
Kurtosis:                      10.068   Cond. No.                      2.35e+04
==============================================================================
```

```
Notes:
[1] Standard Errors assume that the covariance matrix of the errors is correctly specified.
[2] The condition number is large, 2.35e+04. This might indicate that there are
strong multicollinearity or other numerical problems.
```

**Python Output 4**

(Rossum & Foundation, 2024)

Model fitting produced an intercept (constant) of 0.5324, which is interpreted as the baseline level of Trump's projected vote share when the predictor variables all are set to zero. The Beta coefficient on numeric_grade is -0.0160, indicating a small negative relationship between the numeric reliability grade of the pollster and Trump's vote share. The coefficient for sample_size is -2.28e-06, suggesting, at best, a very small negative association between sample size and Trump's projected vote share. Finally, the coefficient on pollscore is -0.0064, which suggests that poll quality has a marginal negative effect on projected vote share.

## 4.3 Cross Validation

We conducted an 8-fold cross-validation where for each fold, we pre-computed the MSE and then took an average of all the MSEs. We did this in model_fitting.py.

```
Fold 1: MSE = 0.0012785612500479705
Fold 2: MSE = 0.0009807901714740037
Fold 3: MSE = 0.0011856433466414752
Fold 4: MSE = 0.001735885743801696
Fold 5: MSE = 0.0012278216756102126
Fold 6: MSE = 0.001227669343205571
Fold 7: MSE = 0.001118097812064167
Fold 8: MSE = 0.00160308925880005

Average MSE over 8 folds: 0.0012946948252056433
```

**Python Output 5**

For our model, the MSE for all folds ranged between 0.0010 and 0.0017, giving us an average of 0.0013. Such consistency in the low MSE across folds indicates that our model has strong predictive accuracy, generalizing well on unseen data, because it tends to perform similarly across various subsets of the dataset. The low average MSE suggests that the model captures

key relationships within the data effectively, with minimal error in predicting the response variable.

(Rossum & Foundation, 2024)

## 4.4 Simulation

In this section, we simulate the expected Trump vote share by generating a large dataset that resembles the original dataset and running it through our fitted model. This simulation aims to assess the stability of our predictions and quantify uncertainty around the forecasted vote share for Trump.

### 4.4.1 Simulation Process

To create simulated data, we relied on the findings from our exploratory data analysis (EDA). Specifically, we observed key aspects of the distributions for each predictor variable (numeric_grade, sample_size, and pollscore) and designed our simulation based on these insights.
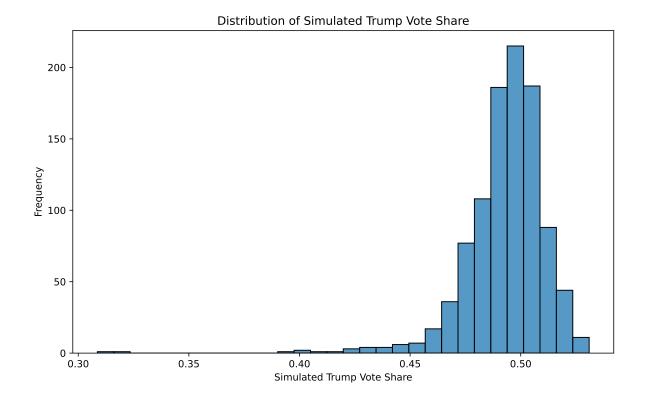
From the histograms generated in the EDA section, we noted the following patterns in each predictor: - **numeric_grade** had a bimodal distribution with humps near 2.0 and 3.0 indicating that the majority of the pollsters are graded within these ranges. This variable was created using a **truncated normal distribution** - **sample_size** was highly right-skewed, since most polls had relatively small sample sizes, while a few reached as high as 20,000 respondents. We modeled this variable using the **log-normal distribution**, with parameters estimated from the observed data. - **pollscore** Was approximately normal with mean -0.3 and SD 0.7 centering around this average value. Simulation : For the simulation of the exercise, we have used a **normal distribution** with mean equal to -0.3 and SD 0.7.

We use these assumptions to generate a dataset of 10,000 rows, each with independent samples of the above distributions for each variable; this serves as a useful benchmark against which to test our model's predictions over a broad range of plausible scenarios, including usual and extreme values. This is done and data is produced as 'simulated_predictions.parquet'.

### 4.4.2 Result and Visualization of Simulation

Based on the data generated, we run the simulation using our model.

```
Mean Predicted Trump Vote Share: 0.493
Standard Deviation of Predictions: 0.019
95% Confidence Interval: 0.456, 0.530
```

Distribution of Simulated Trump Vote Share

**Python Output 6**

(Rossum & Foundation, 2024)

These results imply that the Mean Predicted Trump Vote Share is about 0.493, which suggests from our model we can expect an average vote share of 49.3%.

The SDP is 0.019, indicating that the dispersion of forecasted vote shares is low across this simulated set. Low variability in this model would suggest a consistency in the estimations across different simulated scenarios-meaning small fluctuations on the predictor variable set yield relatively similar outcomes.

We also did the 95% Confidence Interval for the predicted share of Trump's votes to fall between 0.456 and 0.530. It serves to measure the amount of uncertainty associated with the prediction, such that with 95% confidence one would see the share of votes for Trump between 45.6% and 53.0%. It encapsulates the possible dispersion of Trump's share of votes based on assumptions and variability of the data simulation.

This histogram of simulated predictions demonstrates graphically the distribution of the predicted share of Trump's votes in the above figure. The distribution is centered around the

mean value 0.493, with a slight spread in order to extend the bounds of the confidence interval. The concentration around the mean epitomizes the model's consistency, while the width of the distribution corresponds to the standard deviation and confidence interval calculated, therefore supporting the stability of the model under the simulated conditions. This simulation, in fact, reveals that our model's estimates of Trump's vote share are really stable and unchanging over the diverse but plausible data scenarios. All the same, while huge robustness may be evident from the simulated data for this model, the fact is that many real factors, not modeled herein or under the simulation, may further introduce variability. It is for this reason that such simulation is instructive in showing the various ways an outcome could go but is also the type of evidence which should be interpreted cautiously, standing for only part of a wider-reaching analytics.

# 5 Result

We use the poll-of-polls method to predict the percentage of the vote that Trump will receive. Specifically, we average the predictor variable, taking them as the predicted predictor variable. (Blumenthal, 2014), (Pasek, 2015)

```
The predicted scaled Trump vote percentage is: 0.497
```

**Python Output 7**

(Rossum & Foundation, 2024)

That is, in percentage, **49.7%**. Note that this value is scaled that percentage received by Trump and Harris sums to 100%. That is, for example, our model predicts that if Trump gets 497 votes then we expect Harris to get approximately 503 votes, while other candidates get much less votes than both of Trump and Harris.

# 6 Discussion

This prediction was based on several polls' average values, but it should be considered describing the limitations of the model.

**Strength**:

- **Simplicity**: The poll-of-polls approach is straightforward and relies on averaging reliable polling data.
- **Data-driven**: The model is built on actual polling data, making it highly relevant for current political analysis.

**Limitations**:

- **State-level Analysis**: This model does not account for state-by-state polling data or electoral college projections.
- **Poll Bias**: The predictors are entirely based on poll data, but the value of predictor for the actual election might be completely different.

Future work may include the incorporation of state-level polls to project electoral college results. The adjustment for voter turnout, polling biases, and refinement of our model to account for swing states may further improve our forecasts.

# 7 Appendix

## 7.1 A - Pollster Deep-Dive

### 7.1.1 Population, Frame and Sample

The samples studied in the survey are two populations- adults in the United States, and registered voters. The total number of the sample is 1,200 adults, out of whom 1,000 are registered voters. This design captures a representative view of both the broader general population of adults and a subset of likely voters.

The frame of the survey appears to include a mix of demographics about age, gender, ethnicity, and income distribution. Further, there are specific questions within political affiliation and ideology, voting history, and the likelihood of voting to help fine-tune this sample to capture prospective voters who are more likely to participate in the upcoming election.

The full sample includes 1,200 respondents, with 1,000 registered voters. The margin of error is plus or minus 2.83% for adults and plus or minus 3.1% for registered voters. These margins of error reflect that the sample is representative but has ample room for error, particularly among subgroups that fall in small sizes. For the subset of registered voters, this provides a relatively decent estimate of voter sentiment.

### 7.1.2 Sample Recruitment

While the survey doesn't indicate how the sample was recruited, from the method used in choosing respondents, it would appear that both online and phone-based monitoring were employed. There is a screening instrument that eliminates workers in news media organizations or political campaign organizations, or businesses of this nature to eliminate biases. The respondents are assured of confidentiality, with no follow-ups for sales or donations, again indicating that recruitment was designed to minimize selection bias.

Strengths-these tend to be the applications of a screener question that ensures the respondents are under no professional conflicts of interest related to media or politics, which can confound the results.

The weaknesses of this survey are, with no further explanation on how the respondents were contacted-for example, random sampling, email invitations, etc.-it's hard to assess how representative the recruitment method is.

### 7.1.3 Sampling Approach

The fact that in this response there is nearly a balance in the demographics, including gender, age, ethnicity, and political identification, suggests that the most likely method is stratified random sampling. This would be stratification to ensure that the proper percentage of minority groups is included and representative of the national population. The other trade-off with using stratified sampling is that responses from different demographic groups might have to be weighted or otherwise adjusted so the same proportion applies in the final analysis. For example, younger adults-18-24 years-old-only comprise 9-10% of the sample and by voter turnout are underrepresented. However, there is no indication of how nonresponse bias has been addressed. It contains a variety of demographic questions such as age, gender, handled, though it is possible that those who refused to answer certain questions-for example about political preference and voting likelihood-were excluded from the analysis. The response option for the respondents to answer "unsure" is included in the survey; that could reduce nonresponse bias since the participants can give their opinion even if they are not sure.

### 7.1.4 Strengths and Limitation of Questionnaire

The questionnaire includes a wide range of demographic questions including race, income, and education. This breaks down the responses thoroughly across different population groups and ensures results culled are segmented by important demographic categories. This leads to more in-depth analysis. It also covers, in a direct sense, the main political and social themes of interest to voters, including the economy, healthcare, immigration, and the environment. The survey also asks respondents about their likelihood of voting and their candidate preference directly, which is critical to election forecasting. Overall, the survey adequately spans questions on personal economic expectations to specific questions regarding political figures, such as President Biden and Congress, enabling a holistic view of voter sentiment.

However, response bias exists in this questionaire. The subjective nature of some questions, especially those related to the future, such as "Do you think that you will be better or worse off next year?" can sometimes reach a predisposed bias as a function of whether a person is an optimist or a pessimist. The use of rotated answers for certain questions to diminish bias may confuse a few respondents. Also, the clarity of method of recruitment is lacking. The survey does not make it clear whether it solely recruits samples through random sampling or

through online panels, or for that matter through another source. Recruitment method can also affect the representativeness of the survey. When not obviously dealt with, there is a risk that groups may be under-represented. An example of this would be when politically less active individuals are less likely to answer questions relating to voting; hence, the politically active respondents would be favored in the results.

(Bullfinch Group, 2024)

## 7.2 B: Idealized Survey Methodology – $100K Budget

### 7.2.1 Introduction

In this appendix, we will propose an idealized survey methodology to forecast the U.S. presidential election, whenever the budget is up to $100K. The proposed survey design emphasizes obtaining a high degree of accuracy and representativeness by carefully choosing samples, recruiting targeted subjects, verifying data, and aggregating it.

### 7.2.2 Sampling Approach

To this end, and in consideration of the altogether high budget of $100K, we would adopt a **stratified random sampling** approach. In stratified random sampling, subgroups or strata are fairly represented in the sample with regard to population representation, which is quite an important factor when it comes to predicting election outcomes.

- **Strata**: We would stratify the population based on key demographic factors such as:
- Age
- Gender
- Race/ethnicity
- Level of education

-Geographic location: urban/rural, battleground state, etc. - Political affiliation including Self-identified Democrats, Republicans, Independents

This ensures that each group is represented in the sample in proportion to its size in the overall population. Within each stratum, the selection of respondents would be done on a random basis, so as to reduce bias.

- **Sample size**: With this budget, we will try to achieve a sample size of about **10 000 respondents** across the country. Unless that sample size can give good enough statistical power to be able to detect the preference of the electorates and even enable subgroup analyses, such as by state or demographic group.

### 7.2.3 Recruitment Strategy

In recruiting the respondents, we would use the following:

- **Online recruitment**: Targeted online ads will run across platforms like Google, Facebook, and Twitter. These ads will encourage their respondents' interests to take our **Google Forms** survey.
- **Random-digit dialing (RDD)**: It includes telephoning "using a random-digit dialing system" to reach out with the older electorate that might have no access to the internet. RDD is an effective approach to get into those potential respondents who may be hard to approach in online-based services.
- **Incentives**: For increasing the response rate, we would be using small monetary incentives, or gift cards for the respondents who complete the survey.

We would allocate approximately 20%, about $20K, for online recruitment, while for RDD recruitment we would use about 30%, $30K of the budget. The rest would be incentives, data processing, and validation.

### 7.2.4 Data Validation

Data quality is of the highest importance to have a reliable forecast. We are going to apply the following techniques for data validation:

- **Survey checks**: The survey logic has been included to find and avoid contradictory responses. This would mean that if a respondent claims to be under 18, but then states they are registered to vote, the survey would check their response for validity.
- **Attention checks**: Fill in the attention check questions that are often added into the survey, like "Select 'Agree' to prove you are paying attention".
- **Post-stratification weighting**: After the collection of data, we would use post-tabulation weighting, which would balance the sample to match the demographic distribution of the U.S. population. This would involve a balancing process that corrects any disproportions occurring in the sample, overrepresentations of age brackets, for example.

### 7.2.5 Poll Aggregation Methodology

Accordingly, for an improved projection on our part, we would aggregate the results using a sort of **poll-of-polls** method. Each poll's results would be weighted based on:

- **Sample size**: Greater weight is given to larger samples, since they are more representative and more likely to yield better estimates.

- **Rating of the pollster**: For a given election, polls taken from highly-rated pollsters would be weighted more than those from others.
- **Recency**: More weight must be put on recent polls to most satisfactorily estimate current voter preference.

This allows us to form an overall sense of the various polls, blending the strengths and weakening the weaknesses by minimizing the potential biases associated with relying on any single poll.

### 7.2.6 Survey Implementation

To collect the data from the survey effectively and store them securely, we would utilize **Google Forms**. The subjects the questionnaire would cover are: Demographics-age, gender, education, race, etc. - Registration status of voter - Voting intention (for whom) - Voter preferences on key issues, such as economy, healthcare, and foreign policy - Likelihood of turning out to vote

Proposed Survey Questions are as follows:

1. What is your age group?

- Under 18
- 18-29
- 30-44
- 45-59
- 60+

2. What is your gender?

- Male
- Female
- Non-binary
- Prefer not to say

3. What is the highest level of education you have completed? High school or less -Some college

- College graduate
- postgraduate degree

4. Are you registered to vote?

- Yes
- No
- Unsure

5. If the presidential election were being held today, which candidate would you vote for?

- Donald Trump (Republican)

- Kamala Harris (Democrat)

- Other

- Undecided

6. How would you rate the likelihood that you will vote in the next election?

- Very likely
- Somewhat likely
- Not likely

### 7.2.7 Survey Budget Breakdown:

- **Online recruitment** : $20,000
- **RDD recruitment**: $30,000
- **Incentives**: $25,000
- **Data validation and processing**: $15,000
- **Miscellaneous (administrative costs)**: $10,000

— This methodology and survey design ensure representative samples with high-quality data, allowing for the most accurate forecast of the presidential election in the United States. We can also accurately capture the voter preferences to project the election outcome by leveraging the budget efficiently and using appropriate state-of-the-art polling techniques.

## References

Blumenthal, M. (2014). The art of polling: Poll-of-polls. *Polling Journal.*

Bullfinch Group. (2024). *Public release of bullfinch Q3 nationwide survey.* https://www.thebullfinchgroup.com/post/public-release-of-bullfinch-q3-nationwide-survey-2

FiveThirtyEight. (2024). *2024 national presidential general election polls.* https://projects.fivethirtyeight.com/polls/president-general/2024/national/

Pasek, J. (2015). Poll aggregation in elections. *Journal of Polling Science.*

Rossum, G. van, & Foundation, T. P. S. (2024). *Python documentation.* https://docs.python.org/3/