

October, 2015. Incomplete Draft. Please do not cite without permission.

Hierarchical Modeling of Non-Representative Polls

Wei Wang

Columbia University

ww2243@columbia.edu

In this chapter, I will discuss an application of hierarchical modeling to non-representative survey sampling. As it is mentioned in the last chapter, there is a dichotomy in modern survey research, the camp of *describers* and the camp of *modelers*. However, at the heart of modern opinion polling, for both *describers* and *modelers*, is representative sampling, built around the goal that every individual in a particular target population (e.g., registered or likely U.S. voters) has the same probability of being sampled. Non-representative sampling has fallen out of favor among pollsters as a result of its inherent bias. I will show that, using an example of a highly-biased poll on US presidential election conducted on Xbox gaming platform, that hierarchical sampling can be used to remedy the bias and help extract useful information from non-representative polls. This is a joint work with David Rothschild, Sharad Goel, and Andrew Gelman, and is published in (Wang et al. 2014).

A BRIEF HISTORY OF REPRESENTATIVE SAMPLING VS NON-REPRESENTATIVE SAMPLING

The wide-scale adoption of representative polling can largely be traced to a pivotal polling mishap in the 1936 U.S. presidential election campaign. During that campaign, the popular magazine *Literary Digest* conducted a mail-in survey that attracted over two million responses, a huge sample even by modern standards. The magazine, however, incorrectly predicted a landslide victory for Republican candidate Alf Landon over the incumbent Franklin Roosevelt. Roosevelt, in fact, decisively won the election, carrying every state except for Maine and Vermont. As pollsters and academics have since pointed out, the magazine's pool of respondents was highly biased: it consisted mostly of auto and telephone owners as well as the magazine's own subscribers, which underrepresented Roosevelt's core constituencies (Squire 1988). During that same campaign, pioneering pollsters, including George Gallup, Archibald Crossley, and Elmo Roper, used considerably smaller but representative samples to predict the election outcome with reasonable (Gosnell 1937). Accordingly, non-representative or "convenience sampling" rapidly fell out of favor with polling experts. Methods used for sampling have evolved over time, from address-based, in-home interview sampling in the 1930s to random digit dialing after the growth of landlines and cellphones; nevertheless, leading polling organizations continue to put immense effort into obtaining representative samples.

Two recent trends spur the interest for non-representative polls. First, representative sampling is not nearly as representative as its name suggests, and it is becoming less so. Random digit dialing (RDD), the standard method in modern representative polling, has suffered increasingly high non-response rates, both due to the general public's growing reluctance to answer phone surveys, and expanding technical means to screen unsolicited calls (Keeter et al. 2006). By one measure, RDD response rates have decreased from 36% in 1997 to 9% in 2012 (Kohut et al. 2012). With such low response rates, even if the initial pool of targets is representative, those who ultimately answer the phone and elect to respond are almost certainly not, calling into question the statistical benefits of such an approach. Related to dropping response rates is a corresponding increase in cost, in both time and money, as one needs to contact more and more potential respondents to find one willing to participate. The second trend driving this research is that with recent technological innovations, it is increasingly convenient and cost-effective to collect large numbers of highly non-representative samples

via online surveys. What took several months for the *Literary Digest* editors to collect in 1936 can now take only a few days and can cost just pennies per response. The challenge, of course, is to extract meaningful signal from these unconventional samples.

It is worth noting that the so-called “Big Data” is more often than not a convenient sample, with potentially huge selection bias. Without adequately addressing this issue first, any conclusion drawn from big data analysis might be misleading.

Xbox Data

The analysis is based on an opt-in poll continuously available on the Xbox gaming platform during the 45 days preceding the 2012 U.S. presidential election. Each day, three to five questions were posted, one of which gauged voter intention with the standard query, “If the election were held today, who would you vote for?”. Respondents were allowed to answer at most once per day. The first time they participated in an Xbox poll, respondents were additionally asked to provide basic demographic information about themselves, including their sex, race, age, education, state, party ID, political ideology, and for whom they voted in the 2008 presidential election. In total, 750,148 interviews were conducted with 345,858 unique respondents—over 30,000 of whom completed five or more polls—making this one of the largest ever election panel studies.

Despite the large sample size, the pool of Xbox respondents is far from representative of the voting population. Figure~1 compares the demographic composition of the Xbox participants to that of the general electorate, as estimated via the 2012 national exit poll. For ease of interpretation, in Figure~1 states are grouped into 4 categories: (1) battleground states (Colorado, Florida, Iowa, New Hampshire, Ohio, and Virginia), the five states with the highest amounts of TV spending plus New Hampshire, which had the highest per-capita spending; (2) quasi-battleground states (Michigan, Minnesota, North Carolina, Nevada, New Mexico, Pennsylvania, and Wisconsin), which round out the states where the campaigns and their affiliates made major TV buys; (3) solid Obama states (California, Connecticut, District of Columbia, Delaware, Hawaii, Illinois, Maine, Maryland, Massachusetts, New Jersey, New York, Oregon, Rhode Island, Vermont, and Washington); and (4) solid Romney states (Alabama, Alaska, Arizona, Arkansas, Georgia, Idaho, Indiana, Kansas, Kentucky, Louisiana, Mississippi, Missouri, Montana, Nebraska, North Dakota, Oklahoma, South Carolina, South Dakota, Tennessee, Texas, Utah, West Virginia, and Wyoming). The most striking differ-

ences are for age and sex. As one might expect, young men dominate the Xbox population: 18-to-29-year-olds comprise 65% of the Xbox dataset, compared to 19% in the exit poll; and men make up 93% of the Xbox sample but only 47% of the electorate. Political scientists have long observed that both age and sex are strongly correlated with voting preferences (Kaufmann and Petrocik 1999), and indeed these discrepancies are apparent in the unadjusted time-series of Xbox voter intent shown in Figure 2. In contrast to estimates based on traditional, representative polls (indicated by the dotted blue line in Figure 2), the uncorrected Xbox sample suggests a landslide victory for Mitt Romney, reminiscent of the infamous *Literary Digest* error.

ESTIMATING VOTER INTENT WITH MULTILEVEL REGRESSION AND POSTSTRATIFICATION

Multilevel regression and poststratification

To transform the raw Xbox data into accurate estimates of voter intent in the general electorate, I make use of the rich demographic information that respondents provide. In particular I *poststratify* the raw Xbox responses to mimic a representative sample of likely voters. Poststratification is a popular method for correcting for known differences between sample and target populations (Little 1993). The core idea is to partition the population into cells (e.g., based on combinations of various demographic attributes), use the sample to estimate the response variable within each cell, and finally to aggregate the cell-level estimates up to a population-level estimate by weighting each cell by its relative proportion in the population. Using y to indicate the outcome of interest, the poststratification estimate is defined by,

$$\hat{y}^{\text{PS}} = \frac{\sum_{j=1}^J N_j \hat{y}_j}{\sum_{j=1}^J N_j}$$

where \hat{y}_j is the estimate of y in cell j , and N_j is the size of the j -th cell in the population. An estimate of y can be analogously derived at any subpopulation level s (e.g., voter intent in a particular state) by

$$\hat{y}_s^{\text{PS}} = \frac{\sum_{j \in J_s} N_j \hat{y}_j}{\sum_{j \in J_s} N_j}$$

where J_s is the set of all cells that comprise s . As is readily apparent from the form of the poststratification estimator, the key is to obtain accurate cell-level estimates, as well as estimates for the cell sizes.

One of the most common ways to generate cell-level estimates is to simply average sample responses within each cell. If it is assumed that within a cell the sample is drawn at random from the larger population, this yields an unbiased estimate. However, this assumption of cell-level simple random sampling is only reasonable when the partition is sufficiently fine; on the other hand, as the partition becomes finer, the cells become sparse, and the empirical sample averages become unstable. I address these issues by instead generating cell-level estimates via a regularized regression model, namely multilevel regression.

This combined model-based poststratification strategy, known as multilevel regression and poststratification (MRP), has been used to obtain accurate small-area subgroup estimates, such as for public opinion and voter turnout in individual states and demographic subgroups [Park, Gelman, and Bafumi (2004); Lax and Phillips (2009); Ghitza and Gelman (2013)].

More formally, applying MRP in this setting comprises two steps. First a Bayesian hierarchical model is fit to obtain estimates for sparse poststratification cells; second, one averages over the cells, weighting by a measure of forecasted voter turnout, to get state and national-level estimates. Specifically, I generate the cells by considering all possible combinations of sex (2 categories), race (4 categories), age (4 categories), education (4 categories), state (51 categories), party ID (3 categories), ideology (3 categories) and 2008 vote (3 categories), which partition the data into 176,256 cells. {All demographic variables are collected prior to respondents' first poll, alleviating concerns that respondents may adjust their demographic responses to be inline with their voter intention (e.g., a new Obama supporter switching his or her party ID from Republican to Democrat). I fit two, nested multilevel logistic regressions to estimate candidate support in each cell. The first of the two models predicts whether a respondent supports a major-party candidate (i.e., Obama or Romney), and the second predicts support for Obama given that the respondent supports a major-party candidate. Following the notation of (Gelman and Hill 2007), the first model is given by

$$\begin{aligned} \Pr(Y_i \in \{\text{Obama, Romney}\}) = & \\ & \text{logit}^{-1}(\alpha_0 + \alpha_1(\text{state last vote share}) \\ & + a_{j[i]}^{\text{state}} + a_{j[i]}^{\text{edu}} + a_{j[i]}^{\text{sex}} + a_{j[i]}^{\text{age}} + a_{j[i]}^{\text{race}} + a_{j[i]}^{\text{party ID}} + b_{j[i]}^{\text{ideology}} + b_{j[i]}^{\text{last vote}}) \end{aligned} \quad (1)$$

where α_0 is the fixed baseline intercept, and α_1 is the fixed slope for Obama's fraction of two-party vote share in the respondent's state in the last presidential election. The terms $a_{j[i]}^{\text{state}}$, $a_{j[i]}^{\text{edu}}$, $a_{j[i]}^{\text{sex}}$ and so on—which in general is denote by $a_{j[i]}^{\text{var}}$ —correspond to varying coefficients associated with each categorical variable. Here the subscript $j[i]$ indicates the cell to which the i -th respondent belongs. For example, $a_{j[i]}^{\text{age}}$ takes values from $\{a_{18-29}^{\text{age}}, a_{30-44}^{\text{age}}, a_{45-64}^{\text{age}}, a_{65+}^{\text{age}}\}$ depending on the cell membership of the i -th respondent. The varying coefficients $a_{j[i]}^{\text{var}}$ are given independent prior distributions

$$a_{j[i]}^{\text{var}} \sim N(0, \sigma_{\text{var}}^2).$$

To complete the full Bayesian specification, the variance parameters are assigned a hyperprior distribution

$$\sigma_{\text{var}}^2 \sim \text{inv-}\chi^2(\nu, \sigma_0^2),$$

with a weak prior specification for the remaining parameters, ν and σ_0 . The benefit of using a multilevel model is that estimates for relatively sparse cells can be improved through “borrowing strength” from demographically similar cells that have richer data. Similarly, the second model is defined by

$$\begin{aligned} \Pr(Y_i = \text{Obama} \mid Y_i \in \{\text{Obama, Romney}\}) = & \\ & \text{logit}^{-1}(\beta_0 + \beta_1(\text{state last vote share}) \\ & + b_{j[i]}^{\text{state}} + b_{j[i]}^{\text{edu}} + b_{j[i]}^{\text{sex}} + b_{j[i]}^{\text{age}} + b_{j[i]}^{\text{race}} + b_{j[i]}^{\text{party ID}} + b_{j[i]}^{\text{ideology}} + b_{j[i]}^{\text{last vote}}) \end{aligned} \quad (2)$$

and

$$\begin{aligned} b_{j[i]}^{\text{var}} &\sim N(0, \eta_{\text{var}}^2), \\ \eta_{\text{var}}^2 &\sim \text{inv-}\chi^2(\mu, \eta_0^2). \end{aligned}$$

Jointly, Eqs.~(1) and (2) define a Bayesian model that describes the data. Ideally, a fully Bayesian analysis would be performed to obtain the posterior distribution of the parameters. However, for computational convenience, I use the approximate marginal maximum likelihood estimates obtained from the `glmer()` function in the R package `lme4` (Bates, Maechler, and Bolker 2013).

Having detailed the multilevel regression step, I now turn to poststratification, where cell-level estimates are weighted by the proportion of the electorate in each cell and aggregated to the appropriate level (e.g., state or national). To compute cell weights, cross-tabulated population data is needed. One commonly used source for such data is the Current Population Survey (CPS); however, the CPS does not include some key poststratification variables, such as party identification. I thus instead use exit poll data from the 2008 presidential election. Exit polls are conducted on election day outside voting stations to record the choices of exiting voters, and they are generally used by researchers and news media to analyze the demographic breakdown of the vote (after a post-election adjustment that aligns the weighted responses to the reported state-by-state election results). In total, 101,638 respondents were surveyed in the state and national exit polls. I use the exit poll from 2008, not 2012, because this means that in theory the method as described here could have been used to generate real-time predictions during the 2012 election campaign. Admittedly, this approach puts my prediction at a disadvantage since the demographic shifts of the intervening four years cannot be captured. While combining exit poll and CPS data would arguably yield improved results, for simplicity and transparency I exclusively use the 2008 exit poll summaries for poststratification.

National and State Voter Intent

Figure~3 shows the adjusted two-party Obama support for the last 45 days of the election. The daily voter intents for two-party Obama support at the national level are illustrated in Figure 3. Compared with the uncorrected estimates in Figure 2, the MRP-adjusted estimates yield a much more reasonable timeline of Obama's standing over the course of the final weeks of the campaign. With a clear advantage at the beginning, Obama's support slipped rapidly after the first presidential debate—though never falling below 50%—and gradually recovered, building up a decisive lead in the final days.

On the day before the election, the estimate of voter intent is off by a mere 0.6 percentage points from the actual outcome (indicated by the dotted horizontal line). Voter intent in the weeks prior to the election does not directly equate to

an estimate of vote share on election day—a point I return to in Section~???. As such, it is difficult to evaluate the accuracy of the full time-series of estimates. Nonetheless, note that the estimates are not only intuitively reasonable, but that they are also inline with prevailing estimates based on traditional, representative polls. In particular, the estimates roughly track—and are even arguably better than—those from Pollster.com, one of the leading poll aggregators during the 2012 campaign.

National vote share receives considerable media attention, but state-level estimates are particularly relevant for many stakeholders given the role of the Electoral College in selecting the winner (Rothschild 2013). Forecasting the joint probability of victory for each candidate in state-by-state races is a challenging problem due to the interdependencies in state outcomes, %and the joint electoral votes has not yet become the standard forecast the logistical difficulties of measuring state-level vote preference, and the effort required to combine information from various sources (Lock and Gelman 2010). The MRP framework, however, provides a straightforward methodology for generating state-level results. Namely, I use the same cell-level estimates employed in the national estimate, as generated via the multilevel model in Eqs. (1) and (2), and I then poststratify to each state’s demographic composition. In this manner, the Xbox responses can be used to construct estimates of voter intent over the last 45 days of the campaign for all 51 Electoral College races.

Figure 4 shows two-party Obama support for the 12 states with the most electoral votes. The state timelines share similar trends (e.g., support for Obama dropping after the first debate), but also have their own idiosyncratic movements, an indication of a reasonable blend of national and state-level signals. To demonstrate the accuracy of the MRP-adjusted estimates, I plot, in dotted blue lines in Figure~4, the estimates generated by Pollster.com, which are broadly consistent with the state-level MRP estimates. Moreover, across the 51 Electoral College races, the mean and median absolute errors of the estimates on the day before the election are just 2.5 and 1.8 percentage points, respectively.

Voter intent for demographic subgroups

Apart from Electoral College races, election forecasting often focuses on candidate preference among demographic subpopulations. Such forecasts are of significant importance in modern political campaigns, which often employ targeted campaign strategies (Hillygus and Shields 2009). In the highly non-representative Xbox survey, certain subpopulations are heavily underrepresented and plausibly

suffer from strong self-selection problems. This begs the question, how accurate the estimates for older women based on a platform that caters to mostly young men?

It is straightforward in MRP to estimate voter intent among any collection of demographic cells: I again use the same cell-level estimates as in the national and state settings, but poststratify to the desired target population. For example, to estimate voter intent among women, the poststratification weights are based on the relative number of women in each demographic cell. To illustrate this approach, I compute Xbox estimates of Obama support for each level of the categorical variables (e.g., males, females, Whites, Blacks, etc.) on the day before the election, and compare those with the actual voting behavior of those same groups as estimated by the 2012 national exit poll. As seen in Figure~5, the Xbox estimates are remarkably accurate, with a median absolute difference of 1.5 percentage points between the Xbox and the exit poll numbers. Note that Respondents' 2008 vote was not asked on the 2012 exit poll, so I exclude that comparison from Figure~5.

Not only do the Xbox data facilitate accurate estimation of voter intent across these single-dimensional demographic categories, but they also do surprisingly well at estimating two-way interactions (e.g., candidate support among 18–29 year-old Hispanics, and liberal college graduates). Figure~6 shows this result, plotting the Xbox estimates against those derived from the exit polling data for each of the 149 two-dimensional demographic subgroups. Note that state contestedness is excluded from the two-way interaction groups since the 2012 state exit polls are not yet available, and the 2012 national exit poll does not have enough data to reliably estimate state interactions; 2008 vote is also excluded, as it was not asked in the 2012 exit poll. The “other” race category was also dropped as it was not consistently defined across the Xbox and exit poll datasets. Most points lie close to the diagonal, indicating that the Xbox and exit poll estimates are in agreement. Specifically, for women who are 65 and older—a group whose preferences one might a priori believe are hard to estimate from the Xbox data—the difference between Xbox and the exit poll is a mere one percentage point (49.5% and 48.5%, respectively). Across all the two-way interaction groups, the median absolute difference is just 2.4 percentage points. As indicated by the size of the points in Figure~6, the largest differences occur for relatively small demographic subgroups (e.g., liberal Republicans), for which both the Xbox and exit poll estimates are less reliable. For the 30 largest demographic subgroups, Figure~7 lists the differences between Xbox and exit poll estimates. Among these largest subgroups, the median absolute difference drops to just 1.9 percentage

points.

FORECASTING ELECTION DAY OUTCOME

Converting Voter Intent to Forecasts

As mentioned above, daily estimates of voter intent do not directly correspond to estimates of vote share on election day. There are two key factors for this deviation. First, opinion polls (both representative and non-representative ones) only gauge voter preference on the particular day when the poll is conducted, with the question typically phrased as, “if the election were held today.” Political scientists and pollsters have long observed that such stated preferences are prone to several biases, including the anti-incumbency bias, in which the incumbent’s polling numbers tend to be lower than the ultimate outcome (Campbell 2008), and the fading early lead bias, in which a big lead early in the campaign tends to diminish as the election gets closer (Erikson and Wlezien 2008). Moreover, voters’ attitudes are affected by information revealed over the course of the campaign, so preferences weeks or months before election day are at best a noisy indicator of one’s eventual vote. Second, estimates of vote share require a model of likely voters. That is, opinion polls measure preferences among a hypothetical voter pool, and are thus accurate only to the extent that this pool captures those who actually turn out to vote on election day. Both of these factors introduce significant complications in forecasting election day outcomes.

To convert daily estimates of voter intent to election day predictions—which I hereafter refer to as (???) voter intent—I compare daily voter intent in previous elections to the ultimate outcomes in those elections. Specifically, I collected historical data from three previous U.S. presidential elections, in 2000, 2004, and 2008. For each year, I obtained top-line (i.e., not individual-level) national and state estimates of voter intent from all available polls conducted in those elections. The polling data are obtained from Pollster.com and RealClearPolitics.com. From this collection of polling data, I then constructed daily estimates of voter intent by taking a moving average of the poll numbers, in a similar manner to the major poll aggregators. Note that I rely on traditional, representative polls to reconstruct historical voter intent; in principle, however, I could have started with non-representative polls if such data were available in previous election cycles.

I next infer a mapping from voter intent to election outcomes by regressing election day vote share on the historical time-series of voter intent. The key difference between the approach in this chapter and previous related work (Erikson and

Wlezien 2008; Rothschild 2009) is that I explicitly model state-level correlations, via nested national and state models and correlated error terms. Specifically, I first fit a national model given by

$$y_e^{\text{US}} = a_0 + a_1 x_{t,e}^{\text{US}} + a_2 |x_{t,e}^{\text{US}}| x_{t,e}^{\text{US}} + a_3 t x_{t,e}^{\text{US}} + \eta(t, e)$$

where y_e^{US} is the national election day vote share of the incumbent party candidate in election year e , $x_{t,e}^{\text{US}}$ is the national voter intent of the incumbent party candidate at t days before the election in year e , and $\eta \sim N(0, \sigma^2)$ is the error term. Both y_e^{US} and $x_{t,e}^{\text{US}}$ are offset by 0.5, so the values run from $-\$0.5$ to 0.5 rather than 0 to 1. The term involving the absolute value of voter intent pulls the vote share prediction toward 50%, capturing the diminishing early lead effect. I do not include a main effect for time since it seems unlikely that the number of days until the election itself contributes to the final vote share directly, but rather time contributes through its interaction with the voter intent (which it is include in the model).

Similarly, the state model is given by

$$y_{s,e}^{\text{ST}} = b_0 + b_1 x_{s,t,e}^{\text{ST}} + b_2 |x_{s,t,e}^{\text{ST}}| x_{s,t,e}^{\text{ST}} + b_3 t x_{s,t,e}^{\text{ST}} + \varepsilon(s, t, e)$$

where $y_{s,e}^{\text{ST}}$ is the election day state vote share of the state's incumbent party candidate at day t , $x_{s,t,e}^{\text{ST}}$ is the state voter intent at day t , and ε is the error term. The outcome $y_{s,e}^{\text{ST}}$ is offset by the national projected vote share on that day as fit with the national calibration model, and $x_{s,t,e}^{\text{ST}}$ is offset by that day's national voter intent. Furthermore, I impose two restrictions on the magnitude and correlation structure of the error term $\varepsilon(s, t, e)$. First, since the uncertainty naturally decreases as the election gets closer (as t becomes smaller), I apply the heteroscedastic structure $\text{Var}(\varepsilon(s, t, e)) = (t + a)^2$, where a is a constant to be estimated from the data. Second, the state-specific movements within each election year are allowed to be correlated. For simplicity, and as in (Chen, Ingersoll, and Kaplan 2008), I assume these correlations are uniform (i.e., all pairwise correlations are the same), which creates one more parameter to be estimated from the data. I fit the full calibration model with the `glms()` function in the R package `nlme` (Pinheiro et al. 2012).

In summary, the procedure for generating election day forecasts proceeds in three steps:

1. Estimate the joint distribution of state and national voter intent by applying MRP to the Xbox data, as described in Section ??.
2. Fit the nested calibration model described above on historical data to obtain point estimates for the parameters, including estimates for the error terms.
3. Convert the distribution of voter intent to election day forecasts via the fitted calibration model.

National and state election day forecasts

Figure 8 plots the projected vote shares and pointwise 95% confidence bands over time for the 12 states with the most electoral votes. Though these time-series look quite reasonable, it is difficult to assess their accuracy as there are no ground truth estimates to compare with in the weeks prior to the election. As a starting point, I compare the state-level estimates to those generated by prediction markets, which are widely considered to be among the most accurate sources for political predictions~(Rothschild 2013; Wolfers and Zitzewitz 2004). For each state, prediction markets produce daily probabilities of victory. Though Figure~8 plots the forecasts in terms of expected vote share, this estimation procedure in fact yields the full distribution of outcomes, and so I can likewise convert my estimates to probabilistic forecasts. Figure~9 shows this comparison, where the prediction market estimate is derived by averaging the two largest election markets, Betfair and Intrade. My probabilistic estimates are largely consistent with the prediction market probabilities. In fact, for races with little uncertainty (e.g., Texas and Massachusetts), the Xbox estimates do not seem to suffer from the long-shot bias common to prediction markets (Rothschild 2009), and instead yield probabilities closer to 0 or 1. For tighter races, the Xbox estimates—although still highly correlated with the prediction market probabilities—look more volatile, especially in the early part of the 45-day period. Since the ground truth is not clearly defined, it is difficult to evaluate which method—Xbox or prediction markets—yields better results. From a Bayesian perspective, if one believes the stability shown by prediction markets, this could be incorporated into the structure of the Xbox calibration model.

With the full state-level outcome distribution, I can also estimate the distribution of Electoral College votes. Figure~10 plots the median projected electoral votes for Obama over the last 45-days of the election, together with the 95% confidence band. In particular, on the day before the election, my model estimates Obama had an 88% chance of victory, in line with estimates based on traditional

polling data. For example, Simon Jackman predicted Obama had a 91% chance of victory, using a method built from (Jackman 2005). Zooming in on the day before the election, Figure~11 shows the full predicted distribution of electoral votes for Obama. Compared to the actual 332 votes that Obama captured, I estimate a median of 312 votes, with the most likely outcome being 303. Though this distribution of Electoral College outcomes seems reasonable, it does appear to have higher variance than one might expect. In particular, the extreme outcomes seem to have unrealistically high likelihood of occurring, which is likely a byproduct of the calibration model not fully capturing the state-level correlation structure. Nonetheless, given that my forecasts are based on a highly biased convenience sample of respondents, the model predictions are remarkably good.

CONCLUSION

Forecasts not only need to be accurate, but also relevant, timely, and cost-effective. In this chapter, I construct election forecasts satisfying all of these requirements using extremely non-representative data. Though the data were collected on a proprietary polling platform, in principle one can aggregate such non-representative samples at a fraction of the cost of conventional survey designs. Moreover, the data produce forecasts that are both relevant and timely, as they can be updated faster and more regularly than standard election polls. Thus, the key question—and one of the main contributions of this chapter—is to assess the extent to which one can generate accurate predictions from non-representative samples. Since there is limited ground truth for election forecasts, definitely establishing the accuracy of my predictions is difficult. Nevertheless, I show that the MRP-adjusted and calibrated Xbox estimates are both intuitively reasonable, and are also quite similar to those generated by more traditional means.

The greatest impact of non-representative polling will likely not be for presidential elections, but rather for smaller, local elections and specialized survey settings, where it is impractical to deploy traditional methods due to cost and time constraints. For example, non-representative polls could be used in Congressional elections, where there are currently only sparse polling data. Non-representative polls could also supplement traditional surveys (e.g., the General Social Survey) by offering preliminary results at shorter intervals. General Social Survey, which is . Finally, when there is a need to identify and track pivotal events that affect public opinion, non-representative polling offers the possibility of cost-effective continuous data collection. Standard representative polling will certainly con-

tinue to be an invaluable tool for the foreseeable future. However, 75 years after the *Literary Digest* failure, non-representative polling (followed by appropriate post-data adjustment) is due for further exploration, for election forecasting and in social research more generally.

BIBLIOGRAPHY

Bates, Douglas, Martin Maechler, and Ben Bolker. 2013. *Lme4: Linear Mixed-Effects Models Using S4 Classes*. <http://CRAN.R-project.org/package=lme4>.

Campbell, James E. 2008. *The American Campaign: US Presidential Campaigns and the National Vote*. Texas A&M University Press.

Chen, M Keith, Jonathan E Ingersoll, and Edward H Kaplan. 2008. "Modeling a Presidential Prediction Market." *Management Science* 54(8): 1381–94.

Erikson, Robert S, and Christopher Wlezien. 2008. "Are Political Markets Really Superior to Polls as Election Predictors?" *Public Opinion Quarterly* 72(2): 190–215.

Gelman, Andrew, and Jennifer Hill. 2007. *Data Analysis Using Regression and Multilevel/hierarchical Models*. Cambridge University Press.

Ghitza, Yair, and Andrew Gelman. 2013. "Deep Interactions with MRP: Election Turnout and Voting Patterns Among Small Electoral Subgroups." *American Journal of Political Science* 57(3): 762–76.

Gosnell, Harold F. 1937. "How Accurate Were the Polls?" *Public Opinion Quarterly* 1(1): 97–105.

Hillygus, D Sunshine, and Todd G Shields. 2009. *The Persuadable Voter: Wedge Issues in Presidential Campaigns*. Princeton University Press.

Jackman, Simon. 2005. "Pooling the Polls over an Election Campaign." *Australian Journal of Political Science* 40(4): 499–517.

Kaufmann, Karen M, and John R Petrocik. 1999. "The Changing Politics of American Men: Understanding the Sources of the Gender Gap." *American Journal of Political Science* 43(3): 864–87.

Keeter, Scott, Courtney Kennedy, Michael Dimock, Jonathan Best, and Peyton Craighill. 2006. "Gauging the Impact of Growing Nonresponse on Estimates from a National RDD Telephone Survey." *Public Opinion Quarterly* 70(5): 759–79.

Kohut, Andrew, Scott Keeter, Carroll Doherty, Michael Dimock, and Leah Christian. 2012. "Assessing the Representativeness of Public Opinion Surveys." *Pew Research Center for The People & The Press* 15(May): 2012.

Lax, Jeffrey R, and Justin H Phillips. 2009. "How Should We Estimate Public Opinion in the States?" *American Journal of Political Science* 53(1): 107–21.

Little, Roderick JA. 1993. "Post-Stratification: A Modeler's Perspective." *Journal of the American Statistical Association* 88(423): 1001–12.

Lock, Kari, and Andrew Gelman. 2010. "Bayesian Combination of State Polls and Election Forecasts." *Political Analysis* 18(3): 337–48.

Park, David K, Andrew Gelman, and Joseph Bafumi. 2004. "Bayesian Multilevel Estimation with Poststratification: State-Level Estimates from National Polls." *Political Analysis* 12(4): 375–85.

Pinheiro, Jose, Douglas Bates, Saikat DebRoy, Deepayan Sarkar, and R Core Team. 2012. *Nlme: Linear and Nonlinear Mixed Effects Models*.

Rothschild, David. 2009. "Forecasting Elections Comparing Prediction Markets, Polls, and Their Biases." *Public Opinion Quarterly* 73(5): 895–916.

———. 2013. "Combining Forecasts: Accurate, Relevant, and Timely."

Squire, Peverill. 1988. "Why the 1936 Literary Digest Poll Failed." *Public Opinion Quarterly* 52(1): 125–33.

Wang, Wei, David Rothschild, Sharad Goel, and Andrew Gelman. 2014. "Forecasting Elections with Non-Representative Polls." *International Journal of Forecasting*.

Wolfers, Justin, and Eric Zitzewitz. 2004. *Prediction Markets*. National Bureau of Economic Research.

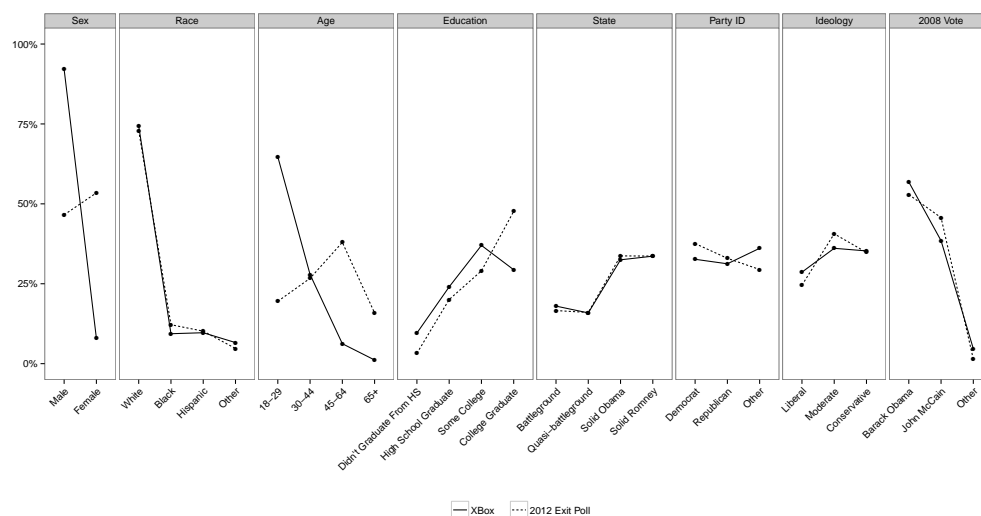


Figure 1: A comparison of the demographic, partisan, and 2008 vote distribution in the Xbox dataset and the 2012 electorate (as measured by adjusted exit polls). The sex and age distributions, as one might expect, exhibit considerable differences.

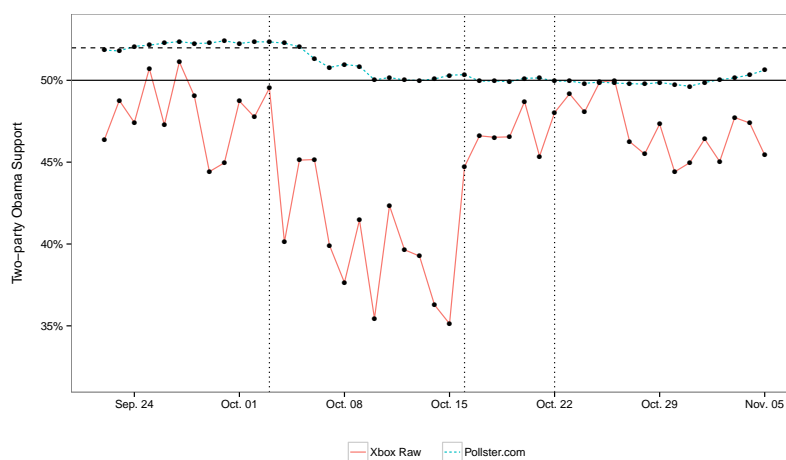


Figure 2: Daily (unadjusted) Xbox estimates of two-party Obama support during the 45 days leading up to the 2012 presidential election, which suggest a landslide victory for Mitt Romney. The dotted blue line indicates a consensus average of traditional polls (the daily aggregated polling results from Pollster.com), the horizontal dashed line at 52% indicates the actual two-party vote share obtained by Barack Obama, and the vertical dotted lines give the dates of the three presidential debates.

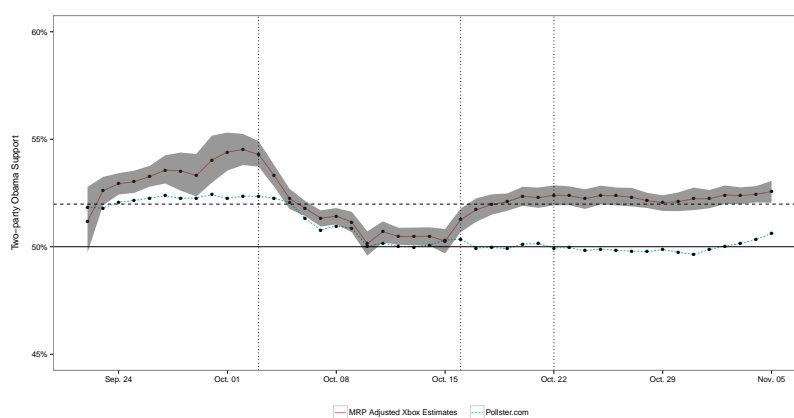


Figure 3: National MRP-adjusted voter intent of two-party Obama support over the 45-day period and the associated 95% confidence bands. The horizontal dashed line indicates the actual two-party Obama vote share. The three vertical dotted lines indicate the presidential debates. Compared with the raw responses in Figure 2, the MRP-adjusted voter intent is much more reasonable, and voter intent in the last few days is very close to the actual outcome. For comparison, the daily aggregated polling results from Pollster.com, shown as the blue dotted line, are further away from the actual vote share than the estimates generated from the Xbox data in the last few days.

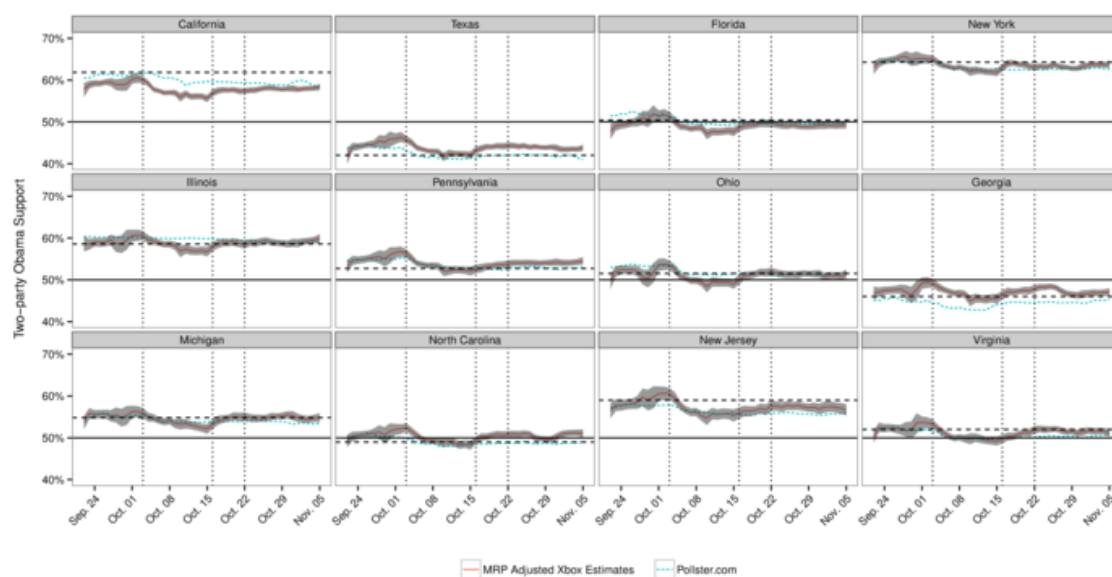


Figure 4: MRP-adjusted daily voter intent for the 12 states with the most electoral votes, and the associated 95% confidence bands. The horizontal dashed lines in each panel give the actual two-party Obama vote shares in that state. The mean and median absolute errors of the last day voter intent across the 51 Electoral College races are 2.5 and 1.8 percentage points, respectively. The state-by-state daily aggregated polling results from Pollster.com, given in the dotted blue lines, are broadly consistent with the estimates from the Xbox data.

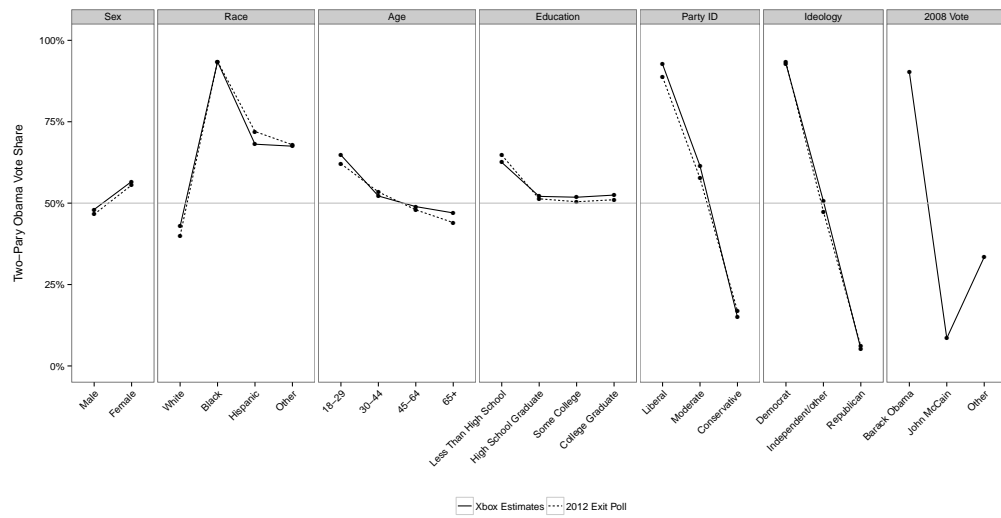


Figure 5: Comparison of two-party Obama vote share for various demographic subgroups, as estimated from the 2012 national exit poll and from the Xbox data on the day before the election.

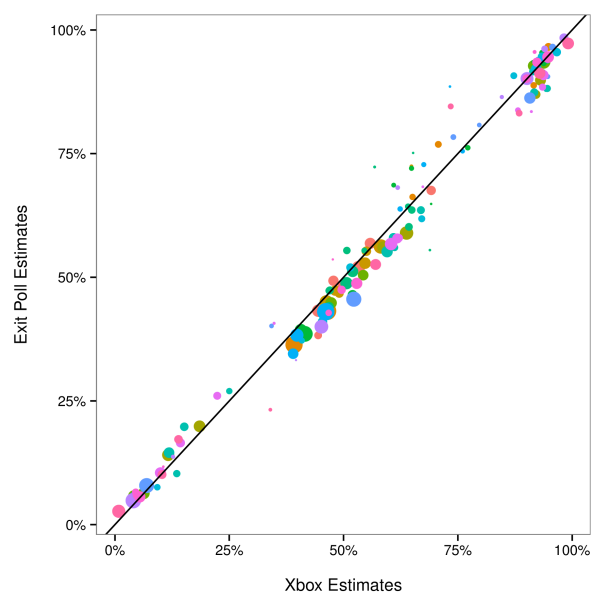


Figure 6: Two-party Obama support as estimated from the 2012 national exit poll and from the Xbox data on the day before the election, for various two-way interaction demographic subgroups (e.g., 65+ year-old women). The sizes of the dots are proportional to the population sizes of the corresponding subgroups. Subgroups within the same two-way interaction category (e.g., age by sex) have the same color.

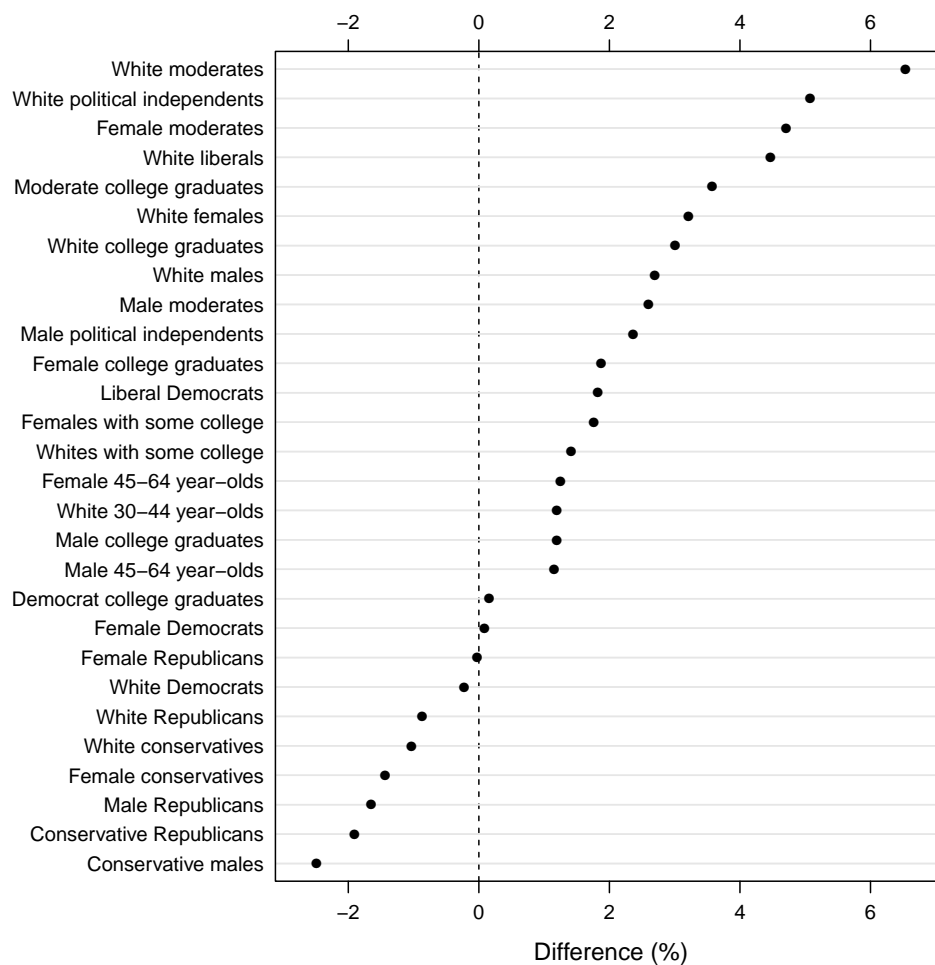


Figure 7: Differences between the Xbox MRP-adjusted estimates and the exit poll estimates for the 30 largest two-dimensional demographic subgroups, ordered by the difference. Positive values indicate the Xbox estimate is larger than the corresponding exit poll estimate. Among these 30 subgroups, the median and mean absolute differences are 1.9 and 2.2 percentage points, respectively.

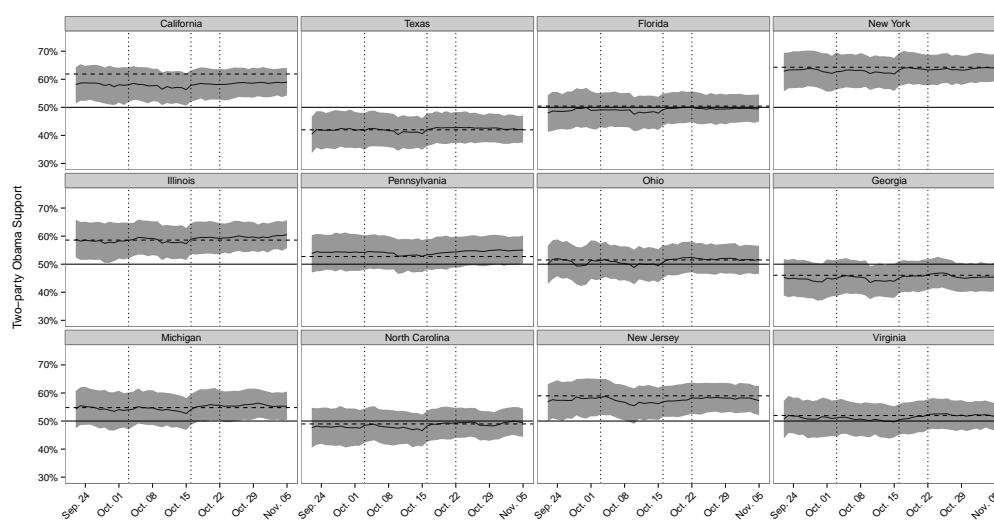


Figure 8: Projected Obama share of the two-party vote on election day for each of the 12 states with the most electoral votes, and associated 95% confidence bands. Compared to the MRP-adjusted voter intent in Figure 4, the projected two-party Obama support is more stable, and the North Carolina race switches direction after applying the calibration model. Additionally, the confidence bands become much wider and give more reasonable state-by-state probabilities of Obama victories.

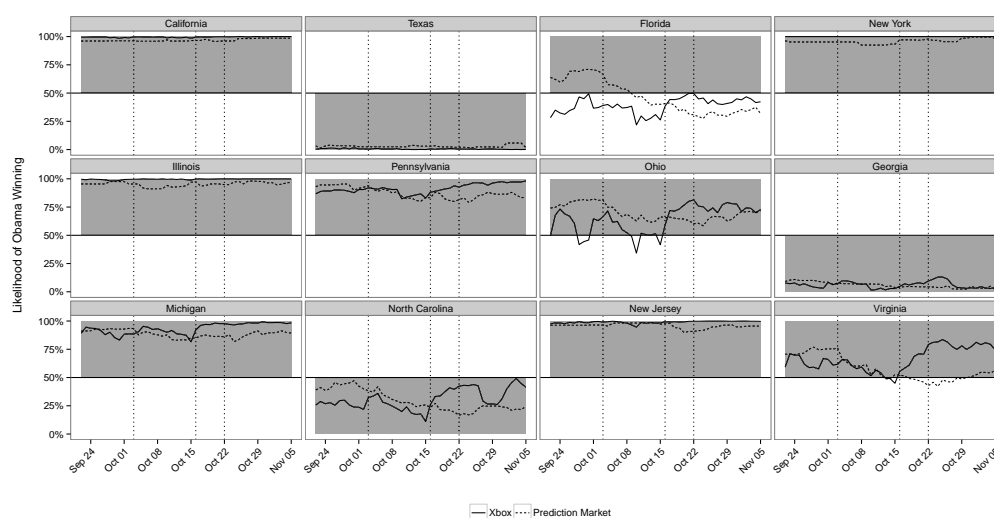


Figure 9: Comparison between the probability of Obama winning the 12 largest Electoral College races based on Xbox data and on prediction market data. The prediction market data are the average of the raw Betfair and Intrade prices from winner-take-all markets. The three vertical lines represent the dates of three presidential debates. The shaded halves indicate the direction that race went.

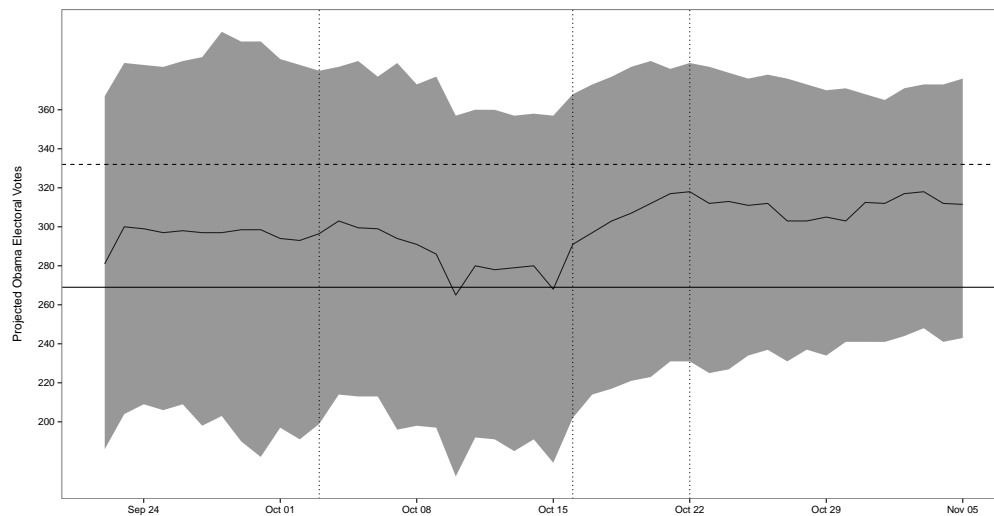


Figure 10: Daily projections of Obama electoral votes in the 45-day period leading up to the 2012 election and associated 95% confidence bands. The solid line represents the median of the daily distribution. The horizontal dashed line represents the actual electoral votes, 332, that Obama captured in 2012 election. Three vertical dotted lines indicate the dates of three presidential debates.

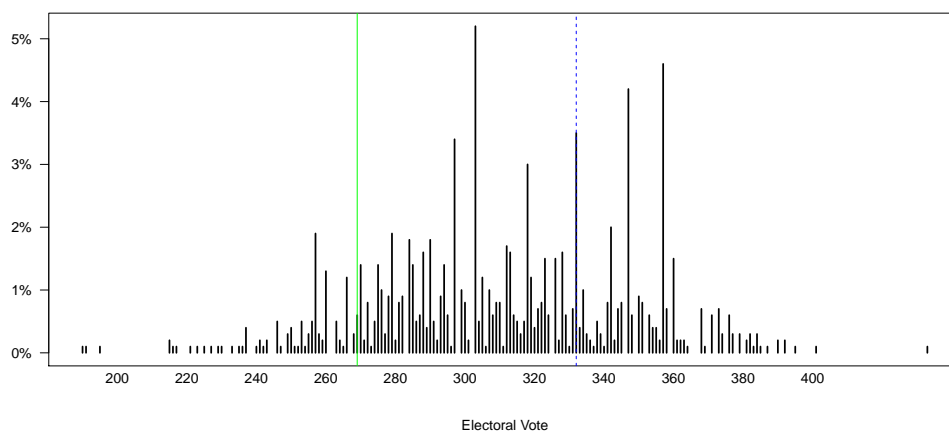


Figure 11: Projected distribution of electoral votes for Obama one day before the election. The green vertical dotted line represents 269, the minimum number of electoral votes that Obama needs for a tie. The blue vertical dashed line gives 332, the actual number of electoral votes captured by Obama. The estimated likelihood of Obama winning the electoral vote is 88%.