Long-Range State-Level 2024 Presidential Election Forecast: How Can You Forecast an Election When You Don't Know Who the Candidates Are Yet?

Jay A. DeSart, Utah Valley University, USA

ABSTRACT This model generates projections of the national popular vote and Electoral College votes a year in advance of the U.S. Presidential Election, before each party's nominees are known. It forecasts the Democratic two-party popular vote in each state and the District of Columbia. It uses four independent variables: national head-to-head polling data 13 months prior to the election, the states' prior election result, a party-adjusted home state advantage dummy variable, and a party adjusted variable simply counting the number of consecutive terms the current incumbent party has occupied the White House. New to this year's model is a polling average approach that encompasses all possible candidate matchups for whom data is available. This year's forecast suggests a distinct possibility of an Electoral College misfire benefitting the Republicans.

he developments of July 21, 2024, brought one of the fundamental challenges of election forecasting into sharp focus. Forecast models that generate predictions based on the individual head-to-head matchups between candidates were dealt a significant blow when President Joe Biden announced that he was withdrawing from the 2024 presidential election. This was particularly the case for my own model, which generates a forecast one year prior to Election Day. I had used it to generate long-range forecasts for both the 2016 and 2020 elections and had presented its preliminary forecast for 2024 at last year's American Political Science Association (APSA) meeting in Los Angeles (DeSart 2023). In what now seems in hindsight to be a prescient statement, the focus of that presentation was on the challenges of generating an election forecast well before the nominees are known.

The unique contribution of this model to the forecasting literature is that it pushes the lead-time envelope by producing both national popular-vote and Electoral College forecasts a year in advance of the election, long before the nominees of each party are known. This typically has necessitated generating a matrix of conditional forecasts representing the various potential matchups between each Republican candidate against each Democratic candidate.

It is a laborious process because even though it can provide a glimpse of how each potential pairing might end up, it typically means that data must be collected for each matchup to generate several different point estimates. For example, in October 2015, there were six Republican candidates (i.e., Jeb Bush, Ben Carson, Ted Cruz, Carli Fiorina, Marco Rubio, and Donald Trump) and two Democratic candidates (i.e., Hillary Clinton and Bernie Sanders) for whom there was available polling data necessary for the model to generate a prediction. The first presentation of this model at the 2015 Iowa Conference on Presidential Politics (DeSart 2015) produced 11 separate forecasts (no polling data were available for the matchup between Jeb Bush and Bernie Sanders). The preliminary forecast I presented at the 2019 APSA Annual Meeting (DeSart 2019) featured a one-by-five matrix showing conditional forecasts pitting the incumbent Donald Trump against five potential challengers (i.e., Joe Biden, Bernie Sanders, Elizabeth Warren, Kamala Harris, and Pete Buttigieg).

Jay A. DeSart ^[D] *is chair of the history and political science department at Utah Valley University. He can be reached at jdesart@uvu.edu.*

Biden's surprise withdrawal late in the campaign season this year highlights the necessity of taking this broad approach if we wanted to cover the multitude of possibilities that developed over the next several months. However, this problem was not unique to the unusual circumstances of the current campaign. Michael Bloomberg's surge in the polls in December 2019—one month It is not surprising that the lagged dependent variable as a predictor dominates the model because the bivariate correlation is quite strong (i.e., Pearson's r=0.90). States' relative positions to one another in terms of the partisan distribution of their election results typically do not shift much from one election to the next, as demonstrated in figure 1.

The developments of July 21, 2024, brought one of the fundamental challenges of election forecasting into sharp focus.

after I announced my long-range forecast for 2020—threatened to render the forecast meaningless because pollsters had not even begun treating Bloomberg as a potential nominee until December. Despite the dramatic difference in the outcomes from 2016 to 2020, the Democratic share of the two-party popular vote (D2PPV) was remarkably stable over that time, with a Pearson's r of 0.99.

The unique contribution of this model to the forecasting literature is that it pushes the lead-time envelope by producing both national popular-vote and Electoral College forecasts a year in advance of the election, long before the nominees of each party are known.

It was this particular challenge—not knowing for sure who the nominees would be that far in advance—that prompted a different approach than I proposed a year ago: that is, averaging the available polling data across the potential matchups to generate a single-point estimate.

THE LONG-RANGE STATE-LEVEL FORECAST MODEL

The long-range state-level model (DeSart 2015, 2019, 2021) generates state-level popular-vote forecasts in each of the 50 states and the District of Columbia. These predictions then can be extrapolated to national-level forecasts by (1) awarding each state's electoral votes to the predicted popular-vote winner; and (2) calculating a turnout-weighted average of each state's popularvote forecast to generate a national popular-vote projection.

The state-level forecasts are the prediction of the Democratic share of the two-party popular vote in each state as a function of the following four variables:

- *Prior Result*: The share of the two-party popular vote won by the Democratic candidate in each state in the previous presidential election.
- *Polls*: The average Democratic two-party share of national headto-head polls taken 13 months in advance of the election, in October of the year prior to the election.
- Home State: A dummy variable for each state indicating its status as a home state for each candidate, signed according to party (i.e., positive for a Democrat and negative for a Republican). Due to the inclusion of the lagged dependent variable as a predictor variable, it also is necessary to have an opposite signed value for the party-adjusted home-state dummy variable from the previous election to account for the removal of previous candidates' advantages.
- *Consecutive Terms*: A simple party-adjusted variable (i.e., positive for Democrats and negative for Republicans) that captures the number of terms a party has consecutively occupied the White House going into the election.

The line in figure 1 represents the pattern if the 2016 results had been a perfect predictor of the 2020 results. Although it is clear that the correlation is strong, there nevertheless is a fairly systemic shift from 2016 to 2020. It is that systemic shift that the model ultimately is trying to capture.

Of course, the challenge is the availability of suitable predictor variables so far in advance of the election that can explain the shift. The main predictor designed to accomplish that is the Terms variable. The two-term penalty is a now well-documented phenomenon in presidential elections. Norpoth (1995) pointed out the cyclical pattern in the popular-vote outcomes of presidential elections across time. Abramowitz (1988) used a two-term penalty term in his Time-for-Change model. The consecutive-term-count variable in this model attempts to capture the cyclical shift from one election to the next.

THE CHALLENGE OF PRE-NOMINATION GENERAL ELECTION FORECASTS

The main challenge in attempting to generate a forecast so far in advance of the election was that we did not yet know who the nominees would be. Given that two key variables in the model, Polls and Home State, are dependent on the specific matchup between candidates, not knowing which candidates would face off against one another a year later complicated the scenario. The process of generating this forecast involved calculating the polling averages for each matchup; generating the state-level popular-vote predictions (while accounting for the home-state advantages of each candidate); extrapolating the state-level predictions to national-level outcome by calculating a turnout-weighted average for a national popular-vote projection; awarding Electoral College votes to the candidates based on the state popular-vote forecast; and finally running Monte Carlo simulations to calculate state and national win probabilities (DeSart 2024).

The larger the field, the more time-consuming this process becomes. There are two potential solutions to this problem. One is simply to make a judgment on which candidates are the two most



likely to win their party's nomination. That might be easier in some years than in others. When an incumbent president is running for reelection, it typically narrows the field on one side of the ballot; however, this still can produce a long list of potential challengers. Even so, the developments of 2024 show that even when an incumbent is running for reelection, there is no guarantee that he or she eventually will be the nominee.

Choosing among a field of potential candidates who are considered the "most likely" to win the nomination is a challenge in and of itself so far in advance of an election. In November 1991, the clear frontrunner in polling for the Democratic nomination was California Governor Jerry Brown. The eventual nominee, Arkansas Governor Bill Clinton, was still trailing behind Brown, Iowa Senator Tom Harkin, and Nebraska Senator Bob Kerrey at that time. Generally speaking, the media outlets and polling organizations have done a fairly good job of assessing the field of candidates early on when deciding which candidates are "top-tier" and which ones likely will be considered "also rans." Up to this point, it has been relatively easy to find the necessary head-to-head polling data between the two eventual nominees 13 months in advance of the election in order to produce a forecast. Nevertheless, it is easy to think of a scenario in which the list of declared candidates 13 months ahead of the election might not include the eventual nominee.

Given that the main goal of this model is to *accurately* project an outcome of the election a year in advance and not leave out potential candidates, we either must generate a full matrix of all possible matchups for which there are available data or come up with a way to produce a single-point estimate that shows the likelihood of one party's nominee winning over the other. In 2016 and 2020, I opted for the former approach, with all of the work and challenges that it entailed. Fortunately, those matrices included the eventual matchup. Although there is variation in the projections generated across the different matchups, they generally tend to point in the same direction, with a few exceptions. Matchups with lesser-known candidates (who often eventually drop out) tend to have much less available polling data and to have closer margins and higher proportions of undecided respondents. The variation across projections provides a glimpse of the potential impact that candidate characteristics—at least those that are known that far in advance of the election—may have on the outcome. However, it may be of greater interest to focus on the broader partisan context that the candidates will face on Election Day and whether we can capture that context so far in advance of an election.

A DATA-AVERAGING APPROACH

If the intent is to capture the overall context rather than the specific matchups, then it may be beneficial (and less work) to average the available data instead of trying to generate a forecast for each possible matchup. In addition to the amount of work involved in producing multiple forecasts across the multiple candidate pairings, there is an issue regarding the reliability of those forecasts that rely on a relatively small number of polling datapoints. Fortunately, there were numerous polls conducted 13 months before the general election that asked respondents to choose between Hillary Clinton and Donald Trump in 2015 and between Joe Biden and Donald Trump in 2019. However, for example, if the 2016 matchup had turned into a race between Bernie Sanders and Ted Cruz or between Hillary Clinton and Mike Huckabee, the forecasts for those contests would be based on a single survey and therefore highly dependent on any sampling issues present in that one poll. As it was, there were 10 polls conducted in October 2015 that pitted Hillary Clinton against Donald Trump. Therefore, the sampling errors in any one of those polls potentially could be ameliorated by averaging it with the other polls.

Poll aggregation is a widespread practice among those who report poll results. RealClearPolitics, FiveThirtyEight, and the now-defunct Pollster.com all use an averaging technique to estimate the "true" population parameter of various survey questions, not only election polls. Several forecast models use the strategy of averaging polling data as one of the predictor variables (DeSart and Holbrook 2003; Graefe 2018; Graefe et al. 2014; Holbrook 2008, 2012). The central-limit theorem in probability theory suggests that the mean of a sampling distribution should be equal to the true population mean, assuming that said distribution consists of sample statistics from unbiased probability samples (Billingsley 1995).

The approach I used for the 2024 forecast extended that principle to averaging the polls not only within each matchup but also across *all possible matchups*. The result of such an approach should yield a measure of the overall partisan context going into the election by not only balancing out the random sampling errors—as suggested by the central-limit theorem—but also muting the candidate-specific effects that each particular matchup brings to each specific poll result. However, given that polling organizations generally tend to ask more frequently about top-tier candidates than the also-rans, greater weight will be given to the poll results featuring candidates who are most likely the actual candidates and who will face one another in the general election. The presence of the polls featuring lesser-known candidates likely will moderate the overall mean to capture something closer to the general partisan context of the election at large.

The inclusion of polling data as an independent variable introduces a potential source of prediction error into the model. Given the apparent polling "misfires" in 2016, we might question the efficacy of adding polls as a predictor. This contribution to the prediction can be mitigated by polling error, especially if there is systemic error in the polls that underestimate or overestimate the support for a particular candidate.

That concern should be alleviated by two factors. First, whereas polling error is clearly a concern, the direction and extent of that error varies from one election to the next, and the average polling error generally has declined over time. Although the average error in 2020 was slightly higher than in 2016, it nevertheless was on par with polling errors since the 1960s (Clinton et al. 2021). Furthermore, the polling misfires in 2016 were mostly present at the state level whereas the national-level polls performed quite well (Kennedy et al. 2017)—and this model relies on national polls. Second, the Polls variable performs quite well in the model, achieving statistical significance in the hypothesized direction. Ultimately, despite any potential issue for polling error in any given election, the polls provide a useful contribution to the explanation of election results over time.

ALLOCATING HOME-STATE ADVANTAGE

One problem created by averaging the polling data across matchups is that we lose the specific nature of determining the home-state advantage in each candidate pairing. When there is a head-to-head forecast, it is simply a matter of assigning the partyadjusted home-state dummy variable to each candidate's home state. If all polls across all matchups are averaged together, we lose the ability to assign that dummy variable to any specific state.

To circumvent that issue, I chose to allocate the home-state advantage proportionally across the field of candidates from each party. Instead of assigning a value of o or 1, I allocated a value to the home-state variable equal to the proportion of time that each candidate appeared in all of the polls from that candidate's party. For example, in all 123 of the polling matchups conducted in October 2019, Bernie Sanders was listed as the Democratic candidate in 20 matchups. Therefore, under the proposed approach, a value of 0.162 (20/123) would be given to the home-state variable for Sanders' home state of Vermont. Conversely, a value of 0.407 would be given to Delaware to account for the fact that Joe Biden was the Democratic candidate in 50 of those matchups.¹

Table 1 shows the impact that using this approach would have had on the model's *a priori* forecasts in 2016 and 2020. In each instance, the national popular-vote projection in each forecast was improved over the projection that used only the polls featuring the two eventual nominees in each election. This suggests that a dataaveraging approach would tend to improve the overall performance of the model moving forward.

THE MODEL FOR 2024

With those modifications, along with incorporating the data from 2020 to update the model, the coefficients I used to generate the forecast for 2024 are listed in table 2. All four variables remain statistically significant with the inclusion of the observations from 2020. In addition, the values of their coefficients remained fairly stable compared to those used in previous elections.

Using these coefficients and the new approach that this article proposes, I generated a forecast of the 2024 presidential election in November 2023. Using the RealClearPolitics Election Polls archive, I obtained a total of 85 general election polling matchups from October 2023. Joe Biden was listed as the Democratic candidate in 80 of the polls and Donald Trump was listed as the Republican candidate in 68 polls. Biden and Trump were featured as the matchup in 65 of the 85 polls. The remainder featured various matchups between either of these candidates with potential opponents, including Bernie Sanders, Kamala Harris, Nikki Haley, Ron DeSantis, and Mitt Romney. The distribution of the matchups and the resultant impact on calculating the potential home-state-advantage variable are presented in table 3.

Table 1 Performance of Polling Averages

			Matchup			Averaged Across Matchups			
Year	Result (National D2PPV%)	Poll Average	Forecast	Forecast Error	Poll Average	Forecast	Forecast Error		
2016	51.1	50.7	50.3	-0.8	51.0	50.5	-0.6		
2020	52.3	54.9	54.8	+2.3	52.8	51.6	-0.7		

Table 2 Updated Long-Range State-Level Forecast Model

Independent	Unstandardized	Standard	
Variable	Regression Coefficient	Error	
Prior Result	1.017	0.018	
Previous October Polls	0.524	0.051	
Home-State Advantage	2.544	0.747	
Number of Terms	-0.988	0.136	
Constant	-27.423	2.785	
R ² =0.90			
S.E. _{y x} =3.08			
N=350			

In-Sample Model Performance over Time

	1996	2000	2004	2008	2012	2016	2020	OVERALL	
States Correctly Predicted	88%	84%	94%	92%	100%	90%	96%	92%	
Mean Absolute Error	2.74	2.57	2.01	2.75	1.88	2.56	1.47	2.28	
National-Level Predictions (Excluding DC)									
National Popular Vote	53.6	48.6	49.7	52.7	52.5	49.8	54.1		
Error	-1.1	-1.6	+1.0	-0.9	+0.6	-1.2	+1.8	1.2 †	
Electoral College	399	231	228	333	329	269	347		
Error	+23	-33	-21	-28	0	+39	+44	27 †	

Note: † Mean absolute error.

Table 3 Distribution of Candidate Appearances in Polls and Resultant Home-State-Advantage Variable for 2024

Party	Candidate	Home State	Number of Polls	Home-State-Advantage Dummy
Republican	Donald Trump	Florida	68	-0.953
	Ron DeSantis	Florida	12	
	Nikki Haley	South Carolina	4	-0.047
	Mitt Romney	Utah	1	-0.011
Democrat	Joe Biden	Delaware	80	-0.059
	Kamala Harris	California	4	0.047
	Bernie Sanders	Vermont	1	0.011

Given the way that the home-state-advantage variable is coded, along with the lagged-result variable as a predictor, the advantage already was present in the data for the incumbent president running for reelection. This would result in a value of o for Delaware in 2024. However, since a few polls did not feature Biden as a candidate, a slight adjustment was made to account for that. In 2020, Trump still was considered a New York resident, so the Florida advantage needed to be accounted for in the 2024 forecast. In addition, the impact of the presumed New York advantage for Trump in 2020 was accounted for by coding New York with a value of 1 in the home-state-advantage variable. This would represent the return to "normal" for New York in 2024.

Given the new coefficients in table 2, these values for the home-state variable in table 3, and polling data from October 2023, the model generated a forecast of the 2024 election suggesting that an Electoral College misfire was a distinct possibility. The model projected that the Democratic candidate would win the national two-party popular vote 50.7% to 49.3%. However, when I used the forecast's state-level point estimates and simply awarded each state's electoral votes to the candidate forecasted to win a majority of the two-party popular vote, the projected Electoral College had the Republican candidate winning a majority 226 to 312.

Table 4 lists the results of the Monte Carlo simulations in which 100,000 elections were generated with the model's statelevel predictions, while allowing them to randomly vary in a normal distribution around that point estimate using the model's standard error of the estimate (3.08). The mean of the distribution of these simulated election results represents the model's forecasts of both the national popular-vote and Electoral College outcomes.

Table 4 Monte Carlo Simulation Results (2024 Forecast)

	Democratic Candidate	Republican Candidate
Mean Projected Share of National Two-Party Popular Vote	50.7%	49.3%
95% Confidence Interval	49.4%–51.9%	48.1%–50.6%
Mean Projected Electoral College Vote Total	256	282
95% Confidence Interval	218–306	232–320

Table 5 presents the distribution of outcomes in these simulated elections based on which candidate wins a majority of the national two-party popular vote and which candidate wins an Electoral College majority. These results demonstrate that the model suggested that another Electoral College misfire was a distinct possibility. The Democratic candidate won the national popular vote in 86% of the 100,000 simulated outcomes. However, the Democratic candidate won an Electoral College majority in 25.2% of the simulated elections. This means that these projections suggested a 61% chance of a repeat of the elections of 2000 and 2016, wherein the Republican lost the popular vote but won a

Table 5 National Popular Vote and Electoral College Vote Outcomes in 100,000 Simulated 2024 Presidential Elections

		National Popula	National Popular-Vote Results		
		Republican Wins	Democrat Wins		
Electoral College Result	Republican Wins	12.8%	60.9%		
	Tie	0.1%	1.0%		
	Democrat Wins	1.0%	24.2%		

Admittedly, the rather wide confidence interval on the Electoral College projection presented in table 4 does not instill much comfort in the reliability of the model. The conclusion that could be drawn is that the model suggests that neither candidate had an insignificant chance of winning. That being the case, it is reasonable to ask whether this model is worthy of attention at all.

Therein lies the downside of a long-range forecast. In that regard, it is not unlike hurricane forecast models that have an everwidening "cone of uncertainty" the longer the time frame of the

Overall, these results suggest that the Republicans had an approximate 74% chance of regaining the White House as a result of this year's election.

majority in the Electoral College. Overall, these results suggest that the Republicans had an approximate 74% chance of regaining the White House as a result of this year's election.

CONCLUDING THOUGHTS

At the very least, it should be clear that the model suggested that the election would be very close, perhaps closer than it was in 2020. However, we must recognize that we were dealing with an unprecedented set of circumstances this year. Biden's late departure only solidified the need to move beyond a model that was tied specifically to individual matchups. Trump's felony convictions earlier in 2024 also injected a degree of potential uncertainty. This degree of uncertainty is what makes election forecasting challengingespecially a year in advance of the election. Even so, this model performed reasonably well in capturing the systematic shifts from one election to the next, even when the specific candidates who eventually appeared on the ballot was in doubt. It is entirely possible that in any given election, the candidates appearing on the ballot in November were not even under consideration in the polls 13 months prior to the election. Therefore, a data-averaging approach so far in advance should mitigate some of those elements of uncertainty. Although it could work to mute the candidatespecific factors that may affect the outcome, it allows us to gauge the overall partisan context underlying the dynamics of the campaign-especially when the specific candidates themselves are somewhat in doubt.

prediction. Table 1 demonstrates that projections tend to center around the actual result with a reasonable amount of variability and that they generally tend to point in the "correct" direction. Even so, the fact that the model's projections for this year's election left so much in doubt is testimony to the level of uncertainty that voters may have had in October 2023. Biden's withdrawal in July appears to have removed a significant source of that uncertainty, and more recent polling data suggested a much more favorable context for the Democrats than this model's projections suggested. Ultimately, the results of the election issued a verdict on this model's efficacy and the utility of using such a long leadtime in predicting the outcome of elections.

DATA AVAILABILITY STATEMENT

Research documentation and data that support the findings of this study are openly available at the *PS: Political Science & Politics* Harvard Dataverse at https://doi.org/10.7910/DVN/HZGLY9.

CONFLICTS OF INTEREST

The author declares that there are no ethical issues or conflicts of interest in this research.

NOTE

¹ We might question whether this approach of averaging the home-state advantage is completely necessary because it often could result in a relatively negligible value

for the variable for those candidates who appear in few polling matchups. Nevertheless, to hold true to the spirit of the data-averaging approach, I deemed it necessary to keep the model internally consistent.

REFERENCES

- Abramowitz, Alan. 1988. "An Improved Model for Predicting Presidential Election Outcomes." *PS: Political Science & Politics* 21 (4): 843–47.
- Billingsley, Patrick. 1995. *Probability and Measure*. 3rd Edition. Hoboken, NJ: John Wiley & Sons.
- Clinton, Josh, Jennifer Agiesta, Megan Brenan, Camille Burge, Marjorie Connelly, Ariel Edwards-Levy, et al. 2021. "*Task Force on 2020 Pre-Election Polling: An Evaluation of the 2020 General Election Poll.*" Alexandria, VA: American Association for Public Opinion Research.
- DeSart, Jay. 2015. "State Electoral Histories, Regime Age, and Long-Range Presidential Election Forecasts: Predicting the 2016 Presidential Election." Paper presented at the Iowa Conference on Presidential Politics. Sauk Centre, MN.
- DeSart, Jay. 2019. "A Long-Range State-Level Presidential Election Forecast Model in 2016 and 2020." *Paper presented at the Annual Meeting of the American Political Science Association*. Washington, DC.
- DeSart, Jay. 2021. "A Long-Range State-Level Forecast of the 2020 Presidential Election." *PS: Political Science & Politics* 54 (1): 73–76.

- DeSart, Jay. 2023. "But You Don't Even Know Who the Nominees Are Yet!" Preliminary 2024 Forecast." Paper presented at the Annual Meeting of the American Political Science Association. Los Angeles, CA.
- DeSart, Jay. 2024. "Replication Data for 'Long-Range State-Level 2024 Presidential Election Forecast: How Can You Forecast an Election When You Don't Know Who the Candidates Are Yet?" *PS: Political Science & Politics*. DOI:10.7910/DVN/ HZGLY9.
- DeSart, Jay, and Thomas Holbrook. 2003. "Statewide Trial Heat Polls and the 2000 Presidential Election: A Forecast Model." *Social Science Quarterly* 84 (3): 561–73.
- Graefe, Andreas. 2018. "Predicting Elections: Experts, Polls, and Fundamentals." Judgement and Decision Making 13 (4): 334–44. DOI:10.1017/S1930297500009219.
- Graefe, Andreas, J. Scott Armstrong, Randall J. Jones, and Alfred G. Cuzán. 2014. "Combining Forecasts: An Application to Elections." *International Journal of Forecasting* 30 (1): 43–54.
- Holbrook, Thomas M. 2008. "Incumbency, National Conditions, and the 2008 Presidential Election." *PS: Political Science & Politics* 41 (4): 709–12.
- Holbrook, Thomas M. 2012. "Incumbency, National Conditions, and the 2012 Presidential Election." *PS: Political Science & Politics* 45 (4): 640–43.
- Kennedy, Courtney, Mark Blumenthal, Scott Clement, Joshua Clinton, Claire Durand, Charles Franklin, et al. 2017. *An Evaluation of 2016 Election Polls in the United States*. Alexandria, VA: American Association for Political Opinion Research.
- Norpoth, Helmut. 1995. "Is Clinton Doomed? An Early Forecast for 1996." PS: Political Science & Politics 28 (2): 201–7.