

# An Analysis of Daily Predictions for the 2008 United States Presidential Election

**Steven E. Rigdon**

*Southern Illinois University Edwardsville*

**Edward C. Sewell**

*Southern Illinois University Edwardsville*

**Sheldon H. Jacobson**

*University of Illinois*

**Wendy K. T. Cho**

*University of Illinois*

**Christopher J. Rigdon**

*Southern Illinois University Edwardsville*

*We describe a Bayesian model for analyses of the 2008 presidential election polls that incorporates third party candidates as well as undecided voters. The state-by-state results were used in a recursive formula for the total electoral votes for Barack Obama and for John McCain. A web site was updated almost daily with the state-by-state projections and the posterior distribution for the number of electoral votes for each candidate. The probability of Obama winning was nearly one late in the summer before dropping below one half during mid-September, and finally recovering to essentially one five weeks before the election. The presentation is accessible to readers with an intermediate level of statistics.*

*Key Words: Electoral College, Bayesian inference, posterior distribution, recursive algorithm, Dirichlet distribution.*

## Introduction

A candidate for President of the United States wins election by either (1) winning 270 or more electoral votes, or (2) being elected by the House of Representatives if no candidate receives at least 270 electoral votes. Thus, predicting the winner of the presidential election should involve state-by-state polls rather than national polls.

During the run-up to the 2008 presidential election, we used a Bayesian model to estimate the probability that

each candidate would win each state. Then, using the recursive algorithm of Kaplan and Barnett (2003), we obtained the posterior distribution for the number of electoral votes for each candidate. The results were posted (almost) daily at an internet site, <http://election08.cs.uiuc.edu/>. The details of the Bayesian method require an intermediate to advanced background in statistics, but the results require only a basic understanding of statistics.

Christensen and Florence (2008) proposed a model with similar purposes. Our model differs from theirs in a number of respects. First, our model accounts for undecided voters and third-party candidates, whereas Christensen and Florence assumed only two possible responses. Second, we studied five different swing scenarios among the undecideds. Third, we used the Kaplan-Barnett recursive algorithm to get the posterior distribution for the number of electoral votes. Christensen and Florence used simulation to get the posterior distribution, but simulation is not necessary since the exact posterior can be obtained.

Another popular web site was [fivethirtyeight.com](http://fivethirtyeight.com), run by Nate Silver, the noted baseball scholar. Silver studied the state by state polls, and created models for demographically similar states and how the changes in one state affect the other state. For example, a poll released in South Dakota would have implications for North Dakota and Montana, even on weeks when no polls were released in North Dakota or Montana. Our model treats states independently, using data from only one state at a time. Like Christensen and Florence (2008), Silver used simulation to estimate the number of electoral votes for each candidate.

We present our Bayesian model in the next section and the recursive algorithm for the posterior distribution is given in the section after that. The fourth section discusses the daily predictions that we made and posted on the web page, including a dynamic graph showing how the posterior changed over time. The last two sections discuss some of the most likely scenarios and an interesting tie scenario.

### Bayesian Model

Polls taken in each state asked the subjects who they would vote for if the election were held today. Subjects then responded with Obama, McCain, a third party candidate, or undecided (or refuse to answer). We thus considered a multinomial model with these four categories. We have combined all third-party candidates together into one category. If there were a third-party candidate with a reasonable chance of winning a state, we would have had to give that candidate an extra category, but in 2008, there were no third party candidates polling more than about 4% in any state. We called the fourth category “Undecided” even though some subjects may have made up their minds, but refused to answer, and we called the members of this category “undecideds.”

The observed polling data in each state therefore consists of a vector  $\mathbf{X}$  that has a multinomial distribution, which

we write  $MULT(n, p_1, p_2, p_3, p_4)$ ; here,  $p_i$  is the probability that a subject responds with category  $i$ , and  $n$  is the sample size. The probabilities must sum to one; that is  $p_1 + p_2 + p_3 + p_4 = 1$ . The conjugate prior for the multinomial is the Dirichlet, which for four dimensions has probability density function

$f(p_1, p_2, p_3, p_4) = c p_1^{b_1-1} p_2^{b_2-1} p_3^{b_3-1} p_4^{b_4-1}$ ,  $p_i \geq 0$ ,  $p_1 + p_2 + p_3 + p_4 = 1$ , where  $c = \Gamma(\sum b_i) / \prod \Gamma(b_i)$ . Marginal distributions of two or more  $p_i$ 's are Dirichlet and the marginal of a single  $p_i$  is a beta distribution with parameters  $b_i$  and  $S - b_i$ , where  $S = \sum b_k$ .

For the Dirichlet distribution, the expected value and variance are

$$E(p_i) = \frac{b_i}{S} \quad \text{and} \quad V(p_i) = \frac{\frac{b_i}{S} \left(1 - \frac{b_i}{S}\right)}{S+1}.$$

Fixing  $E(p_i)$  for  $i=1,2,3,4$  is not sufficient for determining the prior distribution. For example,  $b_1=b_2=4$ ,  $b_3=b_4=1$  gives the same prior means as  $b_1=b_2=40$ ,  $b_3=b_4=10$ . Larger values of  $S = \sum b_k$  lead to smaller prior variances. For the selection of prior parameters for the Democratic and Republican candidates, we used the concept of the *normal vote* (Converse, 1966; Nardulli, 2005, Rigdon et al., 2009).

The normal vote is a measure of the underlying partisan attachments in each state, devoid of short-term forces. For example, the normal vote in Massachusetts, one of the “bluest” (i.e. Democratic) states is

Democratic = 62%  
 Republican = 37%  
 Third Party = 1%

For Texas, one of the “reddest” (i.e. Republican) states, the normal vote is

Democratic = 38%  
 Republican = 61%  
 Third Party = 1%

We then took the proportion of votes for third-party candidates away from the major party candidates in equal amounts, and we allocated 3% for undecided voters. In a process of trial and error, we computed the prior probabilities of Obama winning Texas, and McCain winning Massachusetts. (Texas and Massachusetts were chosen because they have been so one-sided in past elections; Texas for Republicans and Massachusetts for Democrats.) We wanted to choose the parameters so that these probabilities were small, but not minuscule. We settled on  $S=40$ , which is equivalent, in the amount of information it contains, to a preliminary sample of 40 likely voters from each state. With these choices, we have

$$b_1 = 40(1 - 0.03)(NV_1 - \frac{1}{2}C_3) = 38.8(NV_1 - \frac{1}{2}C_3)$$

$$b_2 = 40(1 - 0.03)(NV_2 - \frac{1}{2}C_3) = 38.8(NV_2 - \frac{1}{2}C_3)$$

$$b_3 = 40(1 - 0.03)C_3 = 38.8C_3$$

$$b_4 = 40 \times 0.03 = 1.2,$$

where  $NV_i$  is the normal vote for candidate  $i$  in the state under consideration, and  $C_3$  is the proportion of votes going to all third-party candidates combined. For states in which polls were frequently taken, there was sufficient information so that the prior distribution had little effect on the posterior. For states where polls were infrequently taken, such as Hawaii and North Dakota, the prior exerted much greater influence on the posterior. See Rigdon et al. (2009) for a further discussion of the choice of the prior parameters.

Then, given the observed multinomial data  $\mathbf{X}$  from the state poll, the posterior distribution is

$$p|\mathbf{X} \sim \text{DIRICHLET}(x_1 + b_1, x_2 + b_2, x_3 + b_3, x_4 + b_4).$$

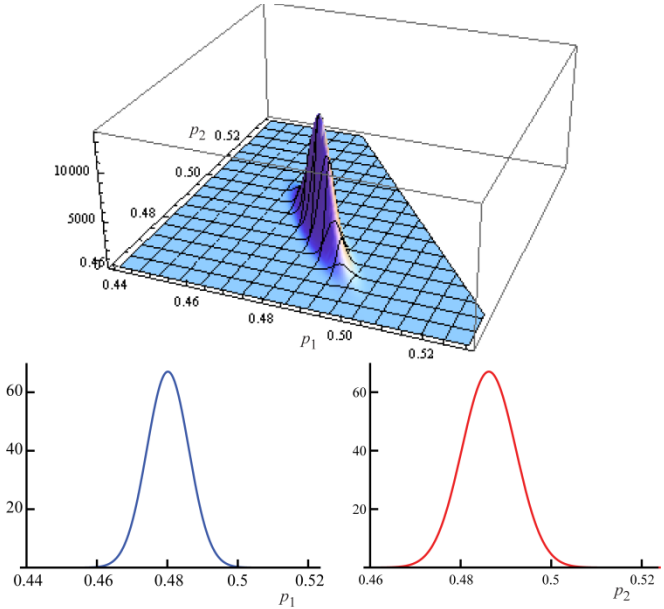
Subjects who responded “undecided” play a key role in the determination of the probability of winning a state. Since the behavior of undecideds can turn an election one way or the other, we looked at five different scenarios about how they would cast their votes. The three numbers below give the split of the undecideds to Obama/McCain/Third Party. These scenarios provide gaps of 0, 5, and 10 percentage points between the Democratic and Republican split of the undecideds.

1. Strong Democratic swing (undecideds split 53.2/43.2/3.6, 10 point gap)
2. Weak Democratic swing (undecideds split 50.6/45.6/3.8, 5 point gap)
3. Neutral scenario (undecideds split 48/48/4, 0 point gap)
4. Weak Republican swing (undecideds split 45.6/50.6/3.8, 5 point gap)
5. Strong Republican swing (undecideds split 43.2/53.2/3.6, 10 point gap).

The posterior probability that Obama would win a state under the neutral scenario is obtained by integrating the marginal posterior of  $(p_1, p_2)$  over the region where  $p_1 > p_2$ . For the other scenarios, the split of undecideds becomes a factor and we must integrate over the appropriate region in the  $(p_1, p_2, p_4)$  space. For example, under the weak Democratic swing, the probability that Obama wins is obtained by integrating over the region where  $p_1 + 0.506p_4 > p_2 + 0.456p_4$ . Details of the required integrations and of the Bayesian analysis in general, can be found in Rigdon et al. (2009).

To illustrate the posterior, consider Missouri, the state where the polls were the closest. The joint distribution for  $(p_1, p_2)$ , and the marginals for  $p_1$  (Obama) and  $p_2$

(McCain), are shown in Figure 1. Since McCain was slightly ahead in Missouri, the posterior for McCain is just a bit to the right of the posterior for Obama. Doing the integration described above gave 0.302 for the probability that Obama would win Missouri.



**Figure 1.** Joint posterior for  $(p_1, p_2)$  (top) and marginal posteriors for  $p_1$  (bottom left) and  $p_2$  (bottom right).

With the results, we can compute the probability that each candidate will win each state under each swing scenario. Table 1 gives the posterior probabilities for all 51 states (really, the 50 states and the District of Columbia) under each scenario. These probabilities are conditioned on the data as of November 3, 2008, the day before election day.

Once all of the state-by-state probabilities are determined, we can use the recursive formula of Kaplan and Barnett (2003), described in the next section, to determine the posterior distribution of the number of electoral votes.

### Recursive Algorithm for Electoral Vote Posterior

We treat the electoral votes in each state as a winner-take-all proposition. Two states, Maine and Nebraska, allow splitting their electoral votes, with one going to the winner of each congressional district and two to the overall winner of the state. This splitting of electoral votes had not happened until the 2008 election when one district in Nebraska went for Obama, while the other districts and the state as a whole went for McCain. We assumed that such splitting would not occur, although with more detailed polling information, such as polls from each

**Table 1.** Posterior probabilities that Obama wins each state under the five swing scenarios. The one state that our model missed, Indiana, is shown in bold.

State	Electoral Votes	Strong Democratic Swing	Weak Democratic Swing	No Swing	Weak Republican Swing	Strong Republican Swing
Alabama	9	0	0	0	0	0
Alaska	3	0	0	0	0	0
Arizona	10	0.023	0.015	0.010	0.007	0.004
Arkansas	6	0.003	0.002	0.002	0.002	0.001
California	55	1	1	1	1	1
Colorado	9	1	1	1	1	1
Connecticut	7	1	1	1	1	1
Delaware	3	1	1	1	1	1
Florida	27	0.990	0.984	0.974	0.960	0.940
Georgia	15	0.049	0.039	0.031	0.025	0.020
Hawaii	4	1	1	1	1	1
Idaho	4	0	0	0	0	0
Illinois	21	1	1	1	1	1
<b>Indiana</b>	<b>11</b>	<b>0.433</b>	<b>0.350</b>	<b>0.273</b>	<b>0.206</b>	<b>0.151</b>
Iowa	7	1	1	1	1	1
Kansas	6	0.001	0	0	0	0
Kentucky	8	0	0	0	0	0
Louisiana	9	0	0	0	0	0
Maine	4	0.999	0.999	0.999	0.999	0.998
Maryland	10	0.994	0.993	0.992	0.991	0.990
Massachusetts	12	1	1	1	1	1
Michigan	17	1	1	1	1	1
Minnesota	10	1	1	1	1	1
Mississippi	6	0.001	0.001	0	0	0
Missouri	11	0.410	0.355	0.302	0.254	0.210
Montana	3	0.387	0.318	0.256	0.200	0.157
Nebraska	5	0	0	0	0	0
Nevada	5	1	1	1	1	1
New Hampshire	4	1	1	1	1	1
New Jersey	15	1	1	1	1	1
New Mexico	5	1	1	1	1	1
New York	31	1	1	1	1	1
North Carolina	15	0.795	0.755	0.708	0.657	0.607
North Dakota	3	0.305	0.274	0.245	0.218	0.206
Ohio	20	1	1	1	1	0.997
Oklahoma	7	0	0	0	0	0
Oregon	7	1	1	1	1	1
Pennsylvania	21	1	1	1	1	1
Rhode Island	4	1	1	1	1	1
South Carolina	8	0.001	0.001	0.001	0.001	0
South Dakota	3	0.001	0.001	0.001	0	0
Tennessee	11	0	0	0	0	0
Texas	34	0	0	0	0	0
Utah	5	0	0	0	0	0
Vermont	3	1	1	1	1	1
Virginia	13	1	1	0.999	0.999	0.999
Washington	11	1	1	1	1	1
Washington DC	3	1	1	1	1	1
West Virginia	5	0	0	0	0	0
Wisconsin	10	1	1	1	1	1
Wyoming	3	0	0	0	0	0

congressional district, this could be accounted for. We also assumed that there were no faithless electors (electors who cast their vote for a candidate who did not win the state, as happened in 1988 when one elector from West Virginia voted for Lloyd Bentson, in 2000 when one elector from the District of Columbia submitted a blank ballot, and most recently (2004), when one elector from Minnesota voted for John Edwards [sic]).

Under these assumptions, we can find the posterior distribution for the number of electoral votes for a candidate using a recursive formula developed by Kaplan and Barnett (2003). Our description of the algorithm follows the reasoning in their paper.

Consider a particular candidate, say Barack Obama. Suppose that state  $i$  has  $v_i$  electoral votes. We will let  $V_i$  denote the random variable that is equal to the number of electoral votes that Obama wins in state  $i$ , and we let  $p_i$  denote the posterior probability that Obama wins state  $i$ . Then by the winner-take-all nature of the electoral college, the possible values for  $V_i$  are 0 and  $v_i$ , with probabilities

$$P(V_i = v) = \begin{cases} p_i, & v = v_i \\ 1 - p_i, & v = 0. \end{cases}$$

Now, let  $T_k$  be the number of electoral votes for Obama in states 1 through  $k$ ; that is

$$T_k = \sum_{i=1}^k V_i. \tag{1}$$

The total number of electoral votes for Obama is then  $T_{51}$  (since we count DC as a state as far as the Electoral College goes). The posterior mean and variance can be computed immediately from equation (1):

$$E(T_{51}) = \sum_{i=1}^{51} E(V_i) = \sum_{i=1}^{51} p_i v_i$$

and

$$Var(T_{51}) = \sum_{i=1}^{51} Var(V_i) = \sum_{i=1}^{51} p_i(1-p_i)v_i.$$

Now, since  $T_{k+1} = T_k + V_{k+1}$  we can write

$$P(T_{k+1} = t) = (1 - p_{k+1})P(T_k = t) + p_{k+1}P(T_k = t - v_{k+1}) \tag{2}$$

for  $t = 0, 1, 2, \dots, 538$  and  $k = 1, 2, \dots, 51$ , where  $T_0 = 0$ . The formula in (2) is the recursive algorithm of Kaplan and Barnett (2003).

To illustrate the algorithm, suppose the posterior probabilities for Obama winning Missouri, Indiana, and North Carolina are as shown in Table 2. (These were three of the closest states in the 2008 election.). If we treat these states in this order (Missouri, Indiana, and North Carolina), then after the first state, Missouri, which has 11 electoral votes,

$$P(T_1 = t) = (1 - p_1)P(T_0 = t) + p_1P(T_0 = t - 11).$$

These probabilities on the right are zero unless  $t=0$  or  $t=11$ . This gives

$$P(T_1 = 0) = 1 - p_1 = 0.6$$

$$P(T_1 = 11) = p_1 = 0.4$$

$$P(T_1 = t) = 0, \quad t \neq 0, 11$$

**Table 2.** Hypothetical probabilities that Obama wins Missouri, Indiana, North Carolina.

$i$	State	Posterior Probability of Winning the State $p_i$	Number of Electoral Votes for the State $v_i$
1	Missouri	0.4	11
2	Indiana	0.5	11
3	North Carolina	0.6	15

Next, consider state 2, Indiana, which also has 11 electoral votes:

$$P(T_2 = t) = (1 - p_2)P(T_0 = t) + p_2P(T_1 = t - 11).$$

The possible values for  $T_1$  are 0, 11, and 22, with probabilities

$$P(T_2 = 0) = (1 - p_2)P(T_1 = 0) + p_2P(T_1 = 0 - 11) = 0.5 \times 0.6 + 0.5 \times 0 = 0.3$$

$$\begin{aligned} P(T_2 = 11) &= (1 - p_2)P(T_1 = 11) + p_2P(T_1 = 11 - 11) \\ &= (1 - p_2)p_1 + p_2(1 - p_1) = 0.5 \times 0.4 + 0.5 \times 0.6 = 0.5 \end{aligned}$$

$$\begin{aligned} P(T_2 = 22) &= (1 - p_2)P(T_1 = 22) + p_2P(T_1 = 22 - 11) \\ &= 0.5 \times 0 + 0.5 \times 0.4 = 0.2. \end{aligned}$$

For  $i=3$  (North Carolina, which has 15 electoral votes) we have

$$P(T_3 = t) = (1 - p_3)P(T_1 = 22) + p_3P(T_2 = t - 15)$$

The possible values for  $T_3$  are 0, 11, 15, 22, 26, and 37, with probabilities

$T$	0	11	15	22	26	37
$P(T_3 = t)$	0.12	0.20	0.18	0.08	0.30	0.12

This process continues until all 51 states have been accounted for. The computing effort for this algorithm is proportional to the product of the number of states and the total number of electoral votes. This is much better than considering all  $2^{51} = 2,251,799,813,685,248$  possible outcomes.

### Daily Predictions

During the two months prior to the November 4, 2008 presidential election, we kept track of all state polls (from [www.realclearpolitics.com](http://www.realclearpolitics.com)) and (almost) daily we updated the state-by-state probabilities. We then found the posterior distribution of each  $p_i, i = 1, 2, \dots, 51$ . We ran the Kaplan and Barnett (2003) algorithm to find the posterior distribution for electoral votes.

Christensen and Florence (2008) discuss the choice of weighting schemes for the most recent polls. A number of plans can be constructed. The simplest is to use the most recent poll and ignore all others. Another is to use all polls within the last  $n$  days. Christensen and Florence propose the following weight functions, where  $t$  is the age of the poll in days,

$$w_1(t) = \begin{cases} 1-t/70, & t \leq 56 \\ 0.2, & t > 56 \end{cases}$$

and

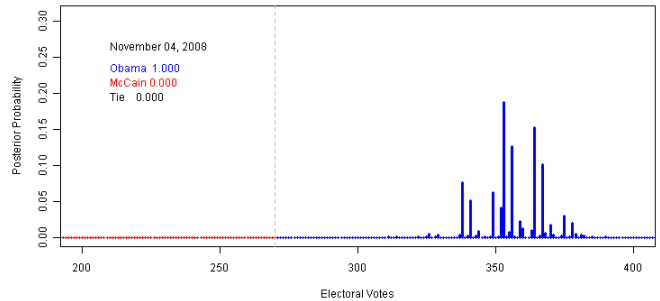
$$w_2(t) = \begin{cases} 1-t/14, & t \leq 13 \\ 0.05, & t > 13. \end{cases}$$

The first gives more weight to older samples than the second; conversely,  $w_2(t)$  gives heavy weight to the most recent polls, and very little weight to older polls.

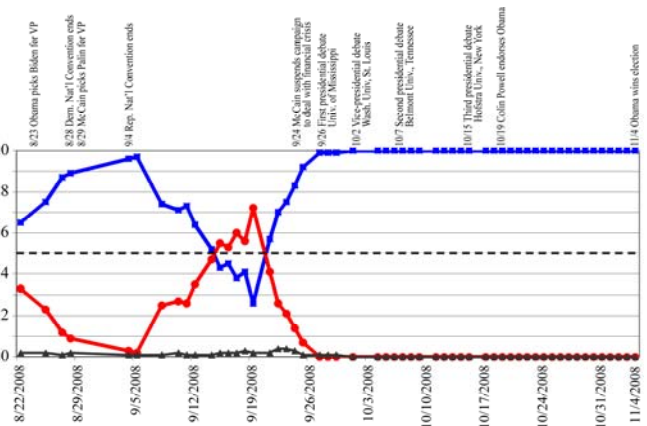
We considered several weighting schemes. One was to weight the sample sizes by  $w(t) = 2^{-t}$ . For example, if one poll with 500 voters was released today ( $t=0$ ) and the only other poll was released three days ago ( $t=3$ ) and had a sample of 800 voters, then the weights would be  $w(0)=1$  and  $2(3) = 2^{-3} = 1/8$ , with effective sample sizes of 500 and  $100 = 800 \times 1/8$ . Thus the combined sample would consist of 600 likely voters, and we would scale the numbers for each candidate. For example, if 54% (270) of the 500 likely voters were for Obama in the most recent poll and if 60% (480) of the 800 were for Obama in the poll that was three days old, then the number for Obama in the combined poll would be  $270 + 480 \times 1/8 = 330$ . Thus, we would treat it as if the sample of size 600 had 330 likely voters who prefer Obama, yielding a percentage of 55%. This weighting scheme places heavy weights on the most recent polls and some, but small, weight on previous polls. After a poll is a week or so old, the effect is negligible.

We also considered weighting schemes that gave a weight of 1 if the poll was less than  $K$  days old, and 0 otherwise. In the end, we used this scheme. We combined all polls so long as they were less than seven days old; otherwise we discarded them.

Figure 2 shows the posterior distribution for the number of electoral votes that we posted on the morning of November 4, the day of the election. Clearly, the probability of Obama winning was nearly one at that time. The probability of Obama winning did, however fluctuate over time. Figure 3 shows the probability of each candidate winning the election from about mid-August until the day of the election. At the top of Figure 3 we show some of the important milestones in the election. McCain's rise in the polls followed the selection of Sarah Palin (8/29) and the end of the Republican National Convention (9/3). For a brief time in mid-September, McCain had the higher probability of winning. This lead quickly evaporated, and Obama's probability of winning went back to nearly one. Very few noteworthy events happened in mid-September, so it is difficult to determine the reason for McCain's fall.



**Figure 2.** Posterior probability histogram for the number of electoral votes for Barack Obama on November 4, 2008.



**Figure 3.** Posterior probability of winning the 2008 presidential election. Obama (Squares), McCain (Circles), and Tie (Triangles).

A dynamic graphic is available (at the journal's web site, along with the spreadsheets containing the data) that shows the posterior distribution changing across time,

from early August through election day. Early on, most of the posterior distribution (for Obama) lies to the right of 270. This gradually shifts downward until September 19, when the probability of winning the election was approximately 0.26 for Obama and 0.72 for McCain (and a small probability for a tie). From late September through early October, the posterior shifted back in Obama's favor until nearly all of the probability for Obama was above 270.

### Most Likely Scenarios

It is rather remarkable that so few states were in play in the 2008 election. In 41 of the 51 states (including DC), the leading candidate had a probability of 0.999 or higher of winning. Thus, among the  $2^{51}$  possible outcomes, most had practically zero chance of occurring. However, among all these possible outcomes, the most likely one, given the prior information and all of the poll results, is the one where the leading candidate wins each state. The most likely outcome is shown in Table 3. Christensen and Florence (2008) made a case for independence among the states. Under the assumption of independence, the most likely scenario has probability

$$P = P(M \text{ wins AL} \wedge M \text{ wins AK} \wedge M \text{ wins AR} \wedge O \text{ wins AR} \wedge \dots \wedge O \text{ wins WI} \wedge M \text{ wins WY})$$

$$= P(M \text{ wins AL})P(M \text{ wins AK})P(M \text{ wins AR})P(O \text{ wins AR}) \times \dots \times P(O \text{ wins WI})P(M \text{ wins WY}).$$

where "M" indicates McCain and "O" indicates Obama. At the time of the last polls, McCain had the greater chance of winning Alabama (AL), Alaska (AK), Arizona (AZ), ..., and Wyoming (WY), while Obama had the greater chance of winning Arkansas (AK), ..., and Wisconsin (WI). Although this is the most likely scenario, it has probability of only 0.186 given the prior and the data. The next most likely scenario, is that the closest state (in terms of the posterior probabilities of winning) goes to the candidate with the smaller posterior probability, in this case, Missouri, and all other states go to the leading candidate. The most likely scenarios continue with all states but one going to the candidate with the higher probability of winning through the top five states, measured in the closeness of the posterior probabilities. Then come various combinations of two states. Table 3 shows the 26 most likely scenarios.

In the 2008 election, the leading candidate won all states except Indiana, where McCain was slightly ahead in the polls, yet Obama won the state. This is the fourth most likely scenario in Table 4. There was also one congressional district in Nebraska that went for Obama, so Nebraska split its electoral vote with four for McCain and

**Table 3** The most likely scenario is that the leading candidate wins in each state. This table gives the leading candidate on the day of the election.

Obama wins: California, Colorado, Connecticut, Delaware, District of Columbia, Florida, Hawaii, Illinois, Iowa, Maine, Maryland, Massachusetts, Michigan, Minnesota, Nevada, New Hampshire, New Jersey, New Mexico, New York, North Carolina, Ohio, Oregon, Pennsylvania, Rhode Island, Vermont, Virginia, Washington, Wisconsin
McCain wins: 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, Wyoming

**Table 4.** Most likely scenarios, given prior information and all poll results. The fourth most likely scenario is what actually occurred, with the exception of the split of Nebraska's electoral votes.

Rank	Scenario	Posterior Probability
1	Leading Candidate Wins Each State	0.186
2	All but MO	0.080
3	All but NC	0.076
4	<b>All but IN</b>	<b>0.069</b>
5	All but MT	0.065
6	All but ND	0.062
7	All but MO, NC	0.033
8	All but MO, IN	0.029
9	All but NC, IN	0.028
10	All but MO, MT	0.028
11	All but NC, MT	0.027
12	All but MO, ND	0.027
13	All but NC, ND	0.025
14	All but IN, MT	0.024
15	All but IN, ND	0.023
16	All but MT, ND	0.022
17	All but MO, IN, NC	0.012
18	All but MO, NC, MT	0.011
19	All but MO, NC, ND	0.011
20	All but MO, IN, MT	0.010
21	All but NC, IN, MT	0.010
22	All but MO, IN, ND	0.010
23	All but NC, IN, ND	0.009
24	All but MO, MT, ND	0.009
25	All but NC, MT, ND	0.009
26	All but IN, MT, ND	0.008
	Sum of Above Probabilities	0.902

one for Obama. In our final prediction, the posterior expected number of electoral votes for Obama was 356.3 and for McCain, it was 181.7. In the end, Obama won 365 and McCain won 173 electoral votes. In terms of states, we predicted correctly all states except Indiana and one Congressional District in Nebraska. Thus, counting electoral votes in the various states where each

candidate was ahead, we predicted 353 for Obama and 185 for McCain. These numbers are off by 12 each: 11 for the electoral votes in Indiana, and 1 for the Second Congressional District in Nebraska.

### Tie Scenarios

The probability of a 269-269 tie was small over most of the period before the election. The greatest probability, approximately 0.04, occurred on September 22 and 23. The key state to a tie was New Hampshire. On September 23, if every state except New Hampshire had gone for the candidate ahead in the polls and if New Hampshire had gone for McCain (who was trailing slightly at the time) then the Electoral College would have been a 269-269 tie. If this had occurred, the election of President would have gone to the House of Representatives of the 111<sup>th</sup> Congress, that is the Congress that took office in January 2009. In this election, each state (not each representative) has one vote. This would have almost certainly led to Obama winning the presidency. There were 34 states where the number of Democratic representatives exceeded the number of Republicans, whereas there were only 14 states where Republicans dominated. Two states (Georgia 6-6 and Idaho 1-1) had a tie. There were four other states where the margin was just one and the presidential candidate of the other party won that state. In Delaware, there is just one representative, a Republican, and Obama won the state. In a case like this, the one representative might be pressured to vote for Obama, since Obama carried Delaware handily. Conversely, North Dakota has just one representative, a Democrat and McCain won the state. Similarly, Tennessee and West Virginia have delegations where the Democrats have a one seat lead, and McCain carried both states. Even if all six of these states (Georgia, Idaho, Delaware, North Dakota, Tennessee, and West Virginia) had gone for McCain, it would not have been enough; Obama would have carried 31 states to McCain's 19.

### Summary and Conclusion

A Bayesian model using four categories (Democratic candidate, Republican candidate, all third-party candidates combined, and undecided) can be used to model the voting behavior in each state. Our analysis includes five swing scenarios for the undecided voters. Once the probabilities for each candidate winning a state are computed, we apply the recursive formula of Kaplan and Barnett (2003) for obtaining the posterior distribution of the number of electoral votes. From the posterior distribution for electoral votes, we can compute the probability that each candidate wins the election. This process was conducted almost daily throughout the two months lead-

ing up to the 2008 presidential election. Our model missed only one state (Indiana) and one congressional district in Nebraska.

### REFERENCES

- Christensen, W. F. and Florence, L. W. 2008. Predicting Presidential and Other Multistage Election Outcomes Using State-Level Pre-Election Polls. *The American Statistician* 26: 1-10.
- Converse, P. 1966. The Concept of a Normal Vote. In A. Campbell, P. Converse, W. Miller, and D. Stokes (Eds.). *Elections and the Political Order*, New York: Wiley & Sons, 9-39.
- Kaplan, E. H. and Barnett, A. 2003. A New Approach to Estimating the Probability of Winning the Presidency. *Operations Research* 51: 32-40.
- Nardulli, P. F. 2005. *Popular Efficacy in the Democratic Era: A Re-examination of Electoral Accountability in the United States, 1828-2000*. Princeton: Princeton University Press.
- Real Clear Politics (2009) [realclearpolitics.com](http://realclearpolitics.com)
- Rigdon, S. E., Jacobson, S. H., Cho, W. K. T., Sewell, E. C. and Rigdon, C. J. 2009. A Bayesian Prediction Model for the U.S. Presidential Election. *American Politics Research* 37: 700--724.
- Silver, N. (2009) [fivethirtyeight.com](http://fivethirtyeight.com)

**Acknowledgements:** The authors would like to thank the reviewer for helpful comments that have certainly led to a clearer presentation of this case study.

This research has been supported in part by the National Science Foundation (IIS 08-27540 SGER). The web site, <http://election08.cs.uiuc.edu>, was hosted by the Department of Computer Science at the University of Illinois. Special thanks to Eric Johnson, William Kormos, Aukrit Unahalekhaka, and Andrew Keeney for their efforts in developing and maintaining the web site.

Correspondence: [srigdon@siue.edu](mailto:srigdon@siue.edu)