

# Prediction of the 2018 Midterm Elections

Jiafeng Chen\*      Joon Yang†

October 19, 2018

## 1 Introduction

## 2 Statistical model and prediction function

Let

$$Y_i = \frac{\text{Republican}\%}{\text{Republican}\% + \text{Democrat}\%}$$

be the proportion of Republican vote share between the major parties in district  $i$ .<sup>1</sup> Let  $x_i$  be a list of features for the district. Assume the (misspecified) probability model

$$Y_i \sim \mathcal{N}(\mu_{i0}, \sigma_0^2).$$

where  $\mu_{i0} = x_i^T \beta_0$ . The model is misspecified since  $Y_i \in [0, 1]$  but Normal is supported on  $\mathbb{R}$ . We justify the misspecification by noting that  $\text{Beta}(a, b) \xrightarrow{d} \mathcal{N}(\mu, \sigma^2)$  where  $a, b \rightarrow \infty$  in such a way that expectation and variance are fixed at  $\mu, \sigma$ .<sup>2</sup> Let

$$\hat{\beta}_0 = \arg \min_{\beta} \sum_{i=1}^N (y_i - x_i^T \beta)^2 + \lambda (\alpha \|\beta\|_1 + (1 - \alpha) \|\beta\|_2)$$

be fitted with an elastic net regularizer over the training data, where  $\lambda$  is chosen via  $K$ -fold cross validation and  $\alpha$  is some fixed constant, say 0.9. Let  $\hat{\sigma}_0^2 = \sum_{i=1}^N (y_i - x_i^T \hat{\beta}_0)^2$  be fitted as the variance of the residuals. Let  $\hat{\mu}_{i0} = x_i^T \hat{\beta}_0$ .

---

\*Harvard College, [jiafengchen@college.harvard.edu](mailto:jiafengchen@college.harvard.edu)

†Harvard College, [joonhyukyang@college.harvard.edu](mailto:joonhyukyang@college.harvard.edu)

<sup>1</sup>Assuming districts have no time dimension: i.e. Alabama-01 is represented by different  $i, j$ 's across two different years.

<sup>2</sup>Alternatively, we could fit some generalized linear model with link function implied by a Beta distribution, but Normal linear model provides a lot of computational ease.

For a district that corresponds to an upcoming election, we form a prior

$$Y_i \sim \mathcal{N}(\hat{\mu}_{i0}, \hat{\sigma}_0^2).$$

Note that in such a formulation, we ignore the sampling variance of  $\hat{\mu}_{i0}$  and  $\hat{\sigma}_0^2$ ,<sup>3</sup> instead forming a plug-in estimate, appealing to the law of large numbers.<sup>4</sup> To obtain a more accurate and timely prediction for the district, we update the prior in two steps.

First, to take into account the “blue wave,” or the general prediction for the Democratic party to win, we update our prior via the generic congressional ballot.<sup>5</sup> Formally, we model generic congressional poll as  $Z_G \mid Y_i \sim \mathcal{N}(Y_i, \sigma_G^2)$ , where  $\sigma_G^2$  is estimated from the 90% confidence interval provided by FiveThirtyEight. We update our prior to form an intermediate posterior:

$$Y_i \mid Z_G \sim \mathcal{N}\left(\frac{\sigma_0^2}{\sigma_G^2 + \sigma_0^2} Z_G + \frac{\sigma_G^2}{\sigma_0^2 + \sigma_G^2} \mu_{i0}, (\sigma_G^{-2} + \sigma_0^{-2})^{-1}\right) =: \mathcal{N}(\mu_{1i}, \sigma_{1i}^2).$$

From now on, we drop the conditioning on  $Z_G$ .

Second, for some districts, we observe a number of district-specific polls  $Z_{i1}, \dots, Z_{iJ}$ . The variance across polls is much higher than implied by a simple Beta-Binomial model, where one assumes that each poll is an independent  $\text{Bin}(n_j, p)$  where  $p$  is sampled from a Beta (again, approximately Normal) prior. As such, we hesitate from using a simple Beta-Binomial updating procedure and opt for the following model: We assume that poll  $j$  has an independent bias  $\epsilon_j \sim \mathcal{N}(0, \sigma_p^2)$ , where  $\sigma_p$  is estimated as the empirical variance of poll outcomes in a district:

$$\begin{aligned} Y_i &\sim \mathcal{N}(\mu_{1i}, \sigma_1^2) \quad \epsilon_j \sim \mathcal{N}(0, \sigma_p^2) \\ Z_{ij} \mid Y_i &\sim n^{-1} \cdot \text{Bin}(n, Y_i + \epsilon_j). \end{aligned}$$

It’s somewhat difficult to compute the posterior  $Y_i \mid Z_{i1}, \dots, Z_{iJ}$ . Instead we may assume the misspecified model

$$\begin{aligned} Y_i &\sim \mathcal{N}(\mu_{1i}, \sigma_1^2) \\ Z_{ij} \mid Y_i &\sim \mathcal{N}\left(Y_i, \frac{1}{4n_j} + \sigma_p^2\right), \end{aligned}$$

so as to (a) use a Normal model and (b) ignore the dependence of  $\mathbb{V}(Z_{ij})$  on  $Y_i$ . We now

---

<sup>3</sup>The elastic net regularizer in the fitting method for  $\beta$  makes the sampling variance of  $\hat{\beta}_0$  difficult to compute.

<sup>4</sup>From this point on, we drop the hat on  $\mu_{i0}, \sigma_0$ .

<sup>5</sup><https://projects.fivethirtyeight.com/congress-generic-ballot-polls/>

have the posterior by, say, sequentially Bayesian updating:

$$Y_i \mid Z_{i1}, \dots, Z_{iJ} \sim \mathcal{N}(\mu_{2i}, \sigma_{2i}^2).$$

We predict  $\hat{Y}_i = \mu_{2i}$ , and, naturally

$$\widehat{\text{Winner}}_i = \begin{cases} \text{Republican} & \mu_{2i} > .5 \\ \text{Democrat} & \mu_{2i} < .5 \end{cases}$$

### **3 Fitting Method**

### **4 Results**

### **5 Uncertainty and robustness**

### **6 Conclusion**