

# Improved Bayesian Ethnorace Prediction

Bertrand Wilden

Last updated on 03 May, 2021

## Abstract

Quantitative social science research on race and ethnicity can be constrained by the lack of available individual-level data with these markers. In response to this issue, methods have been developed to predict individuals' race and ethnicity using Bayes' rule. I expand upon existing methods by incorporating new information into the algorithm. Along with some other adjustments, these improvements lead to substantial gains in predictive performance when validated against official state voter files. The largest of these predictive gains are found for African Americans—a group which existing methods find difficult to predict accurately. Furthermore, I apply my algorithm to published research using the old methods and demonstrate the substantive implications of switching to more accurate predictions.

## Introduction

Research on racial and ethnic disparities in the United States can be constrained by the lack of available individual-level data with these markers (Trounstein 2020; Kuk, Hajnal, and Lajevardi 2020; Burch 2013). This is particularly true when racial geography is a key aspect of the research design. In these contexts, researchers are often forced to use aggregate-level data on group proportions, such as county statistics, to draw inferences. But this method is susceptible to the modifiable areal unit problem, whereby statistical bias is introduced due to arbitrary and unequal geographic unit sizes (Fotheringham and Wong 1991). Furthermore, some research designs all but require individual-level data. Geographic regression discontinuity—an increasingly common causal inference technique—typically relies on individually geocoded addresses (Keele and Titiunik 2018; Velez and Newman 2019; Cantoni 2020). Therefore if race or ethnicity are central elements of the research question, individual-level identifiers for these categories are likely necessary.

To overcome some of these challenges, I describe a method for imputing individual ethnoraces directly. This method uses Bayes’ rule to make predictions by combining information from nationwide distributions of six ethnorace categories over other characteristics, such as names, geographies, political party identification, and others. Intuitively, an individual with a name that is highly unique to a particular ethnorace, and who lives in an area with many other people of the same ethnorace, will be given a high predicted probability of belonging to that group.

My implementation builds off existing methods which use a similar prediction algorithm (Elliott et al. 2009; Imai and Khanna 2016; Voicu 2018). The method described by Imai and Khanna (2016) in particular, and the associated R package `wru`, has become popular in recent studies on race and ethnicity. Its application has been used in work on racial protests and voting patterns (Enos, Kaufman, and Sands 2019); disparities in campaign financing (Grumbach and Sahn 2020; Grumbach, Sahn, and Staszak 2020), evictions (Hepburn, Louis, and Desmond 2020), and voter turnout (Fraga 2018); the impact of electoral institutions on local representation (Abott and Magazinnik 2020); and public health issues such as suicide rates (Studdert et al. 2020).

My ethnorace prediction method improves upon Imai and Khanna (2016) in several ways. Whereas their method only takes as inputs distributions over surnames, geolocation, political party, age, and gender, my implementation adds information from a nationwide list of first names (Tzioumis 2018) as well as address characteristics. Additionally, I incorporate insights from the machine learning literature to further improve predictive performance. The result of these modifications is a substantial increase in predictive power compared to Imai and Khanna (2016) in validation tests. These predictive gains are particularly strong in regards to correctly classifying African Americans—especially in contexts where fine-grain geolocation data are unavailable. Under certain data constraints, the probability that an individual predicted to be Black self-identifies as Black is nearly twice as high. My method is available in an easy-to-use R package `bper`.<sup>1</sup>

In the next section I provide some background and specifics regarding the inputs and outputs of my ethnorace prediction method. Then I explain the methodology and compare my implementation with previous versions. Next, I demonstrate the predictive performance of my method when validated against the combined North Carolina and Florida voter file ( $n = 21,164,503$ ). The inclusion of self-reported ethnorace identifiers in these data allows me to compare my predictions

---

<sup>1</sup>`bper`: Bayesian Prediction for Ethnicity and Race. <https://github.com/bwilden/bper>

against a ground truth. And lastly, I replicate Grumbach and Sahn’s (2020) descriptive findings in their study on racial disparities in campaign finance. I compare the empirical results using my prediction method against those from Imai and Khanna (2016). This demonstrates some of the substantive implications of using improved predictions.

## Data

### Outputs

Before discussing the methodology further, I want to first define what “predicting ethnicity or race” means. These are categories which, although relatively immutable compared to other identities, do not have universally accepted delineations and meanings (Omi and Winant 2014). I follow the convention from previous ethnorace prediction methods by using the US Census Bureau categorizations (Elliott et al. 2009; Imai and Khanna 2016; Voicu 2018). In this framework, individuals can be classified as non-Hispanic White, non-Hispanic Black or African American, non-Hispanic Asian and Pacific Islander, non-Hispanic American Indian and Alaska Native, Hispanic or Latino alone, and non-Hispanic Other Race.<sup>2</sup> Because Hispanic identity is defined by the Census, and understood commonly, as an ethnicity, rather than a race, I use the term “ethnorace” in this paper to refer to any of the previously-mentioned groups.

There are several benefits to using the Census ethnorace categorization. These definitions capture a common understanding of race and ethnicity in the US, and correspond to groups studied frequently in social science research. The data sources of these groups’ distributions that serve as inputs to the prediction formula also rely on the Census categorization. There is currently no feasible alternative data source that I know of that offers the same level of detail as the Census. This also facilitates comparison of my method against previous ethnorace prediction methods.<sup>3</sup> One downside to using the Census categories, however, is that it obscures substantial heterogeneity that may exist within each group. Within Asian Americans and Latinos, for example, there is

---

<sup>2</sup>Non-Hispanic Other Race includes individuals who identify as belonging to two or more race/ethnicities, as well as those who may not identify with the other Census categories.

<sup>3</sup>Unlike my method, Imai and Khanna (2016) do not include a separate category for American Indian and Alaska Native. In order to create similar comparison groups, I recode all predicted American Indian and Alaska Native individuals as Other Race during the validation exercises. If desired, however, the `bper` package will produce predicted probabilities for the American Indian and Alaska Native category.

considerable variation in terms of national ancestry. Furthermore, the unfortunate necessity of a catch-all Other Race category, which also includes those who identify as multi-racial, ensures that important sources of diversity are lost.<sup>4</sup> Researchers interested in using this ethnorace prediction method should be aware of these shortcomings. Nonetheless, because the output groups correlate strongly with many important outcomes in social science research, there are still many applications where this prediction method will be useful.

## Inputs

### First Names

The first names list I use comes from Tzioumis (2017). It is drawn from mortgage applications and contains ethnorace counts in each of the six groups across 4,250 first names. Relative to Census data, this list of first names may be unrepresentative of the larger US population. Mortgage applicants are wealthier than the average American, and are more likely to be employed. To the extent that first name distributions differ by ethnorace given these unobserved characteristics, this may be a concern. But the predictive benefits from using first name data, as I will demonstrate, likely overcome these worries.

### Last Names

For my last names data, I draw from the 2010 Census Surnames List.<sup>5</sup> This list comes from the 2010 decennial Census and contains over 160,000 common US last names (those occurring 100 or more times in the population). I also append any last names from the 2000 Census Surnames List that are not found in the 2010 list to increase the probability the algorithm will match an individual's last name. Like the first names list, these data include counts of individuals in each of the six ethnorace categories across each last name. Once the 2020 Census Surnames List is released, these data will be readily incorporated as well.<sup>6</sup>

---

<sup>4</sup>This is hinted at empirically by the method's poor predictive performance for the Other Race category.

<sup>5</sup>[https://www.census.gov/topics/population/genealogy/data/2010\\_\\_surnames.html](https://www.census.gov/topics/population/genealogy/data/2010__surnames.html)

<sup>6</sup>Because most of the input data sources are from the 2010 decennial Census, the predictions are likely most accurate when applied to data sets from around that time period. The algorithm could be improved in the future by allowing flexibility in which years to draw input data from.

## Geolocations

My ethnorace distributions by geolocations come from the 2010 decennial Census, accessed via IPUMS NHGIS.<sup>7</sup> In decreasing order of mean population, these geographies include *State*, *County*, *Census Place*, *ZIP Code*, and *Census Block*. Predictions tend to improve with more precise levels of geography. With this in mind, my implementation automatically matches each individual to the most fine-grain level of geography available. As with the last names input data, once the 2020 Census data is released the algorithm will include the updated geolocation distributions.

## Party Identification

My party identification data come from a 2012 Gallup poll.<sup>8</sup> The three categories of political party I include are Republican, Democrat, and Other (including Independents and “don’t knows”). The Gallup report tells me both the probability that an individual with a given ethnorace belongs to a particular political party, and the probability that an individual with a given political party identifies with a particular ethnorace. Party ID by ethnorace has likely fluctuated since 2012, yet the inclusion of these data still modestly improve predictions for recent data sets.

## Age and Gender

Like my geolocation data, age and gender distributions come from the 2010 decennial Census, accessed via IPUMS NHGIS. These variables do not contain much predictive power in terms of ethnoracial classification, but nevertheless, I find that their inclusion in the algorithm helps slightly.

## Multi-unit Address

These data refer to ethnorace distributions over multi-unit housing occupancy. Individuals are matched to these probabilities if their address contains “Apt”, “Unit”, “#”, or other such identifier. Unfortunately, I was not able to find these distributions for the 2010 decennial Census, so instead I use those from the year 2000. Again, hopefully I will be able to incorporate data from the 2020 decennial Census here when, or if, it becomes available.

---

<sup>7</sup>Steven Manson, Jonathan Schroeder, David Van Riper, Tracy Kugler, and Steven Ruggles. IPUMS National Historical Geographic Information System: Version 15.0 [dataset]. Minneapolis, MN: IPUMS. 2020. <http://doi.org/10.18128/D050.V15.0>

<sup>8</sup><https://news.gallup.com/poll/160373/democrats-racially-diverse-republicans-mostly-white.aspx>

## Data Structure

The raw data sources I describe above, with the exception of party ID,<sup>9</sup> all contain counts of individuals with a particular attribute (i.e. the first name JOHN, or the ZIP Code “92092”) per ethnorace category. Taking proportions by cell across a given attribute tells us  $Pr(Ethnorace|Attribute)$ , and taking proportions by cell across a given ethnorace group tells us  $Pr(Attribute|Ethnorace)$ . These two conditional probabilities form the building blocks of the classification algorithm described later.

If any cell in the input data is empty (i.e. if there are no individuals of a particular ethnorace with some attribute), then the conditional probabilities  $Pr(Ethnorace|Attribute)$  and  $Pr(Attribute|Ethnorace)$  will be zero. As will become clear in the Methodology section, if either of those two probabilities for an individual equal zero for a given ethnorace, the algorithm will predict a zero percent probability that the individual belongs to that ethnorace. This will occur even if some other attributes about that individual predict a high probability of belonging to the ethnorace. For example, an individual could have first and last names that are highly predictive of being Hispanic, but reside in a Census block which had zero Hispanic occupants at the time of the 2010 decennial Census. Blocks typically contain only around 400 individuals—so this is a real possibility. For this hypothetical person the input data claims that  $Pr(Hispanic|CensusBlock) = 0$ , which yields  $Pr(Hispanic) = 0$  due to the structure of the prediction algorithm. To resolve this issue, I apply a technique from the machine learning literature known as Laplace smoothing to my input data. This works by adding some constant, or psuedocount, to number of individuals in every cell in the input data, then calculating the conditional probabilities  $Pr(Ethnorace|Attribute)$  and  $Pr(Attribute|Ethnorace)$ .<sup>10</sup> Through validation tests, I found that Laplace smoothing led to significant gains in predictive performance for Asian individuals in particular. I also conjecture that, absent this smoothing technique, the algorithm’s predictions are too beholden to the specifics of the 2010 decennial Census. The predictions will generalize better to other time periods without the rigid assumptions of zero conditional probability for some attribute/ethnorace combinations.

---

<sup>9</sup>Party ID percentages by ethnorace are directly available in the Gallup report.

<sup>10</sup>Missing counts are only a problem in the first names, last names, and geolocation data so psuedocounts are only added for those inputs. The exact value for the Laplace smoothing psuedocount could be any number greater than zero, but through out-of-sample validation I found 5 to be the optimal value.

## Methodology

In general terms, the method computes predicted probabilities for each of the six aforementioned ethnorace categories for each individual. Then each individual is classified into the category corresponding to the highest predicted probability. These predicted probabilities can be stated more formally as the conditional probability of identifying as a particular ethnorace for an individual with a particular profile of first name, last name, geolocation, party ID, age, gender, and address type. Bayes' rule provides a template for how to answer this sort of conditional probability problem.

$$Pr(R = r|X) = \frac{Pr(X|R = r)Pr(R = r)}{Pr(X)} \quad (1)$$

Where  $R$  is an individual's true ethnorace,  $r$  is one of six possible ethnorace categories (White, Black, Asian, Native American, Hispanic, or Other race), and  $X$  is the joint probability of an individual having a particular profile of attributes (first name, last name, geolocation, party ID, age, gender, and address type). Unfortunately, the joint probability  $X$  in Equation (1) is intractable due to both data constraints and the astronomically large number of combinations of possible attribute profiles. If however, we assume conditional independence of ethnorace among each attribute in  $X$ , we can rewrite Equation (1) in terms of less complex conditional probabilities:

$$Pr(R = r|X) = \frac{Pr(R = r|x') \prod_{j=1}^6 Pr(x_j|R = r)}{\sum_{i=1}^6 Pr(R = r_i|x') \prod_{j=1}^6 Pr(x_j|R = r_i)} \quad (2)$$

Where  $x$  is the vector of individual attributes indexed by  $j$ . The particular attribute  $x'$  comes from using the chain rule to decompose the joint probability  $Pr(R = r, X)$ . The choice of which attribute to use for  $x'$  is atheoretical, but all previous prediction methods have used last names (Elliott et al. 2009; Imai and Khanna 2016; Voicu 2018). During my validation exercises, I found that the choice of  $x'$  has potentially large consequences for predictive performance. This is because the expression  $Pr(R = r|x')$  is typically much greater than any  $Pr(x_j|R = r)$ , and therefore contributes more weight to the final posterior probability.<sup>11</sup> For example, using last names for  $x'$

---

<sup>11</sup>To see why, imagine  $x = \textit{Smith}$  and  $r = \textit{Black}$ . The probability that an individual is Black given they have the last name Smith,  $Pr(R = r|x)$ , is 0.23. But the probability that an individual is named Smith given that they are Black,  $Pr(x|R = r)$ , is 0.0000029.

appears to help predictions of Whites—but to the detriment of non-Whites. In light of these trade-offs, my method cycles through every attribute as the choice of  $x'$  and computes  $Pr(R = r|X)$  for each. These posterior probabilities are then averaged within each ethnoracial category to generate final predicted probabilities that an individual belongs to a particular ethnorace. The result of this smoothing is more balanced predictions across each ethnorace.

The conditional independence assumption necessary for transforming equation (1) to (2) says that knowing both a particular attribute of an individual, and that individual’s ethnorace, should give us no extra knowledge of any other attribute for that individual. Stated formally,  $Pr(x_j|R = r, X) = Pr(x_j|R = r)$  for all  $x_j$ . This assumption is almost certainly violated in the present context. One example that has been demonstrated empirically is that last name distributions by race vary across regions in the US (Crabtree and Chykina 2018). Violations of the conditional independence assumption are commonplace in most applications of similar “naive” Bayes classification algorithms. Nevertheless, these prediction methods perform well in many contexts (Lewis 1998; Domingos and Pazzani 1997; Rish 2001). This is likely because of the decision rule governing the final classifications—the posterior probabilities of the true class do not have to necessarily be statistically valid, they only need to be higher than those of every other class to be accurately classified.

Fortunately, as I will demonstrate, the predicted probabilities produced by my method are actually very well calibrated. This means that, unlike typical naive Bayes predicted probabilities which cluster close to 0 and 1 (Zadrozny and Elkan, n.d.), the output probabilities from my algorithm accurately reflect the uncertainty in the predictions. Therefore researchers can feel confident in incorporating these measures of uncertainty into their analyses when using the method if they so choose.

## Validation

To test the performance of the model, I apply the predictions to the combined North Carolina and Florida State voter file.<sup>12</sup> These files contain snapshots of the registered voters in their respective states and provide individual-level data for first names, last names, address, political party, age, gender, and crucially self-identified ethnorace. After combining the two voter files, I then geocoded

---

<sup>12</sup>Retrieved May 2020



each unique address in the sample. This allowed me to match individual observations to Census places and blocks, and ZIP Codes. Then I applied the prediction algorithm described above using the `bper` package and calculated each individuals' predicted ethnorace. In order to compare my method against an existing benchmark, I also used the `wru` package (Imai and Khanna 2016) to calculate ethnorace predictions for the same individuals.

Unlike typical machine learning techniques, my prediction method does not fit the model on some subsample, or training set, from these data and then compare predictions against a held-out test set. Instead, the conditional probabilities for each attribute and ethnorace described in the Inputs section are merged into voter file from the input data sources. This allows the entire voter file to be used for validation. And because the input data come from nationally representative samples, the risk of overfitting due to any peculiarities of the North Carolina/Florida electorate are minimized. Together, these two states represent 21,164,503 individuals. Compared to nationwide percentages, Florida has a higher proportion of Hispanics and North Carolina has a higher proportion of African Americans. When combined they form a reasonably ethnoracially diverse population—2% Asian, 16% Black, 11% Hispanic, 6% Other Race, 65% White.

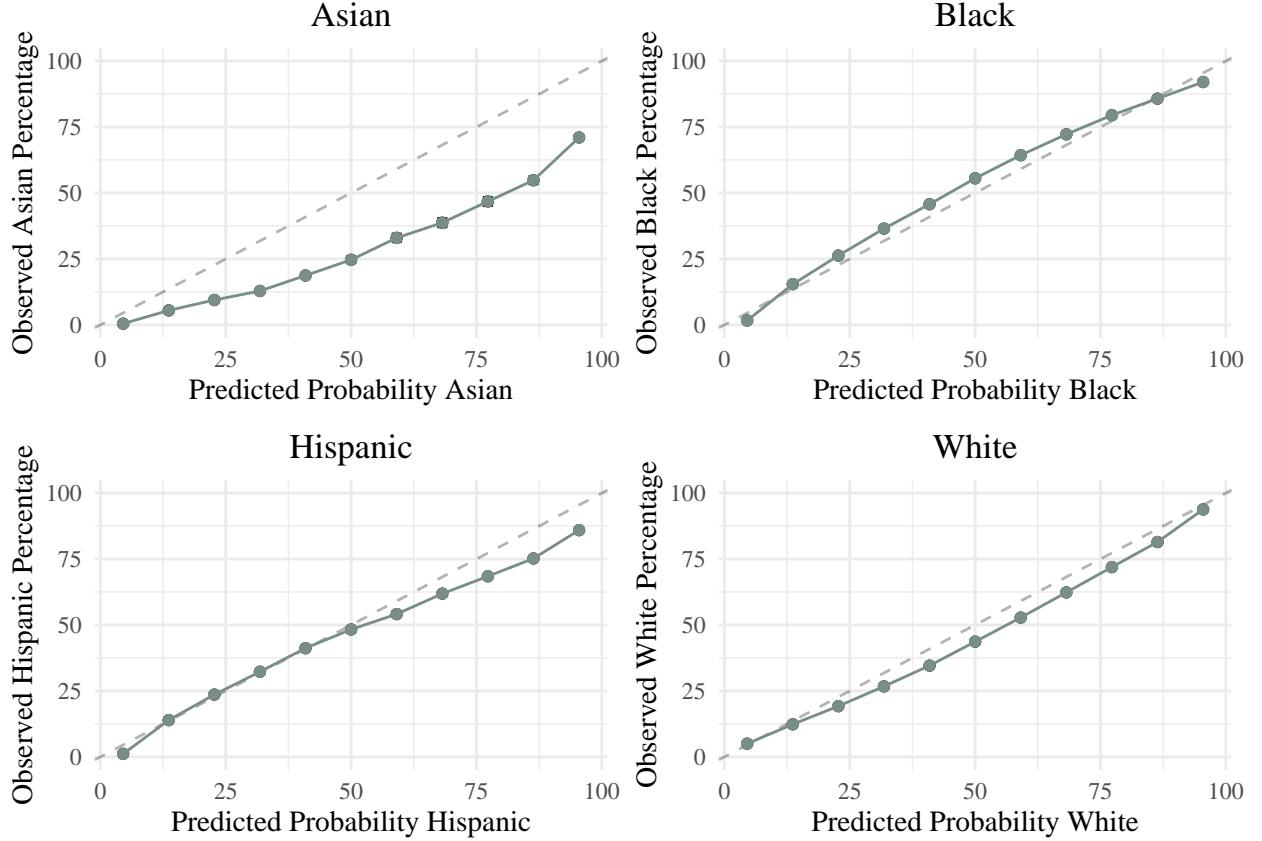
## Calibration

As mentioned previously, the predicted probabilities produced by `bper` are remarkably well calibrated. This means that they accurately reflect the true probability that an individual belongs to a particular ethnorace. We can test this empirically by dividing the predicted probabilities for each ethnorace, for each individual in the sample, into evenly spaced bins (in this case 11). Then we compute the percentage of observed individuals of a particular ethnorace who fall in each bin. If the predicted probabilities are well calibrated, we would expect to see a one-to-one relationship between these two groups. For example, under good calibration we should observe roughly 25% of individuals whose predicted probability of being White is in the neighborhood of 0.25 to actually be White.

Figure 1 displays the probability calibration plots for Asian, Black, Hispanic, and White predictions in the Florida/North Carolina sample. As shown by their proximity to the dashed 45 degree line, the predictions for Black, Hispanic, and White individuals appear to be well calibrated. For each bin of predicted probabilities, we observe essentially the corresponding percentage of individ-

uals of that ethnorace in the sample. Calibration for Asian predictions, on the other hand, appear somewhat worse. Across the range of predicted probabilities, we observe between 5 to 20% fewer Asian individuals in each bin. This means that the algorithm has a tendency to overestimate Asian predicted probabilities.

Figure 1: Probability Calibration Plots by Ethnorace



Rather than eye-balling plots like those in Figure 1, we can obtain objective measures of calibration by calculating Brier scores for each ethnorace. A Brier score is defined by

$$BS = \frac{1}{N} \sum_{t=1}^N (f_t - o_t)^2$$

Where  $f_t$  is the predicted probability an individual  $t$  belongs to a particular ethnorace and  $o_t$  is the binary observed ethnorace value (1 if the individual belongs to a particular ethnorace, 0 otherwise). Brier scores range from 0 (good calibration) to 1 (bad calibration) and are analogous to a mean squared error. The Brier scores for the Florida/North Carolina sample are: Asian 0.012; Black 0.056; Hispanic 0.034; White 0.097. These are all very low—indicating that the average

predicted probability accurately reflects the uncertainty in ethnorace classification. The fact that the Brier score for Asian is lower than every other group might be puzzling given the calibration plots in Figure 1. However, this low score mostly reflects the abundance of low Asian predicted probabilities for individuals who are not Asian. The rarity of Asians in the Florida/North Carolina sample, therefore, produces a low Brier score for this group.

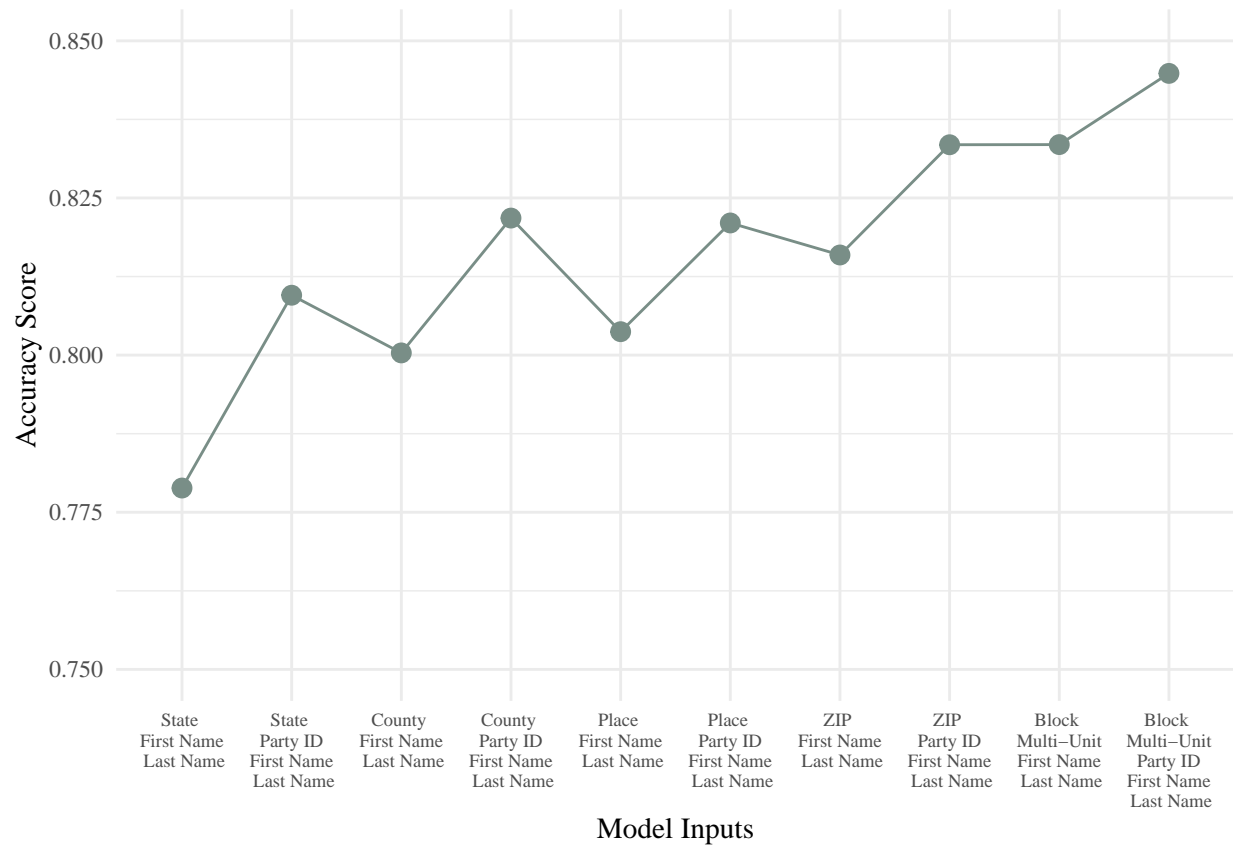
Given the discovery of over-inflated predicted probabilities for Asian Americans, I subsequently adjusted the algorithm to down-weight these values by 50%. Reducing the initial Asian predicted probabilities *greatly* increased overall predictive performance for this group without influencing the results for other ethnoraces. While the down-weighting for Asian predictions can be defended on the basis of the calibration plots in Figure 1, even better predictive results could be obtained through a more structured approach to adjusting predicted probabilities post-hoc. Platt scaling and isotonic regression are two methods developed to improve the calibration, and the subsequent predictive performance, of similar algorithms (Zadrozny and Elkan, n.d.).

## Predictive Performance

While investigating the calibration of the predictions is important for assessing the reliability of the predicted probabilities, calibration does not tell us how well the method ranks these probabilities when generating final ethnorace predictions. Instead, the a common metric for assessing individual-level predictive performance is the Accuracy score, or Overall Error Rate. This number is the proportion of correctly classified individuals in the sample. I ran models separately for different combinations of input variables to mimic data availability constraints in real-world applications, and then calculated the Accuracy score for each. Figure 2 displays a summary of the results of these different models.

As expected, the model with the greatest number of input data sources, and at the most precise geographic level, on the far right of the figure performs the best in terms of overall Accuracy. Using Census blocks, multi-unit address, party ID, first names, and last names, this model correctly identifies the ethnorace of 84.5% of individuals in the sample. The upward trend in model accuracy seen in the figure corresponds to shrinking the size of the geolocation variable used. Moving from state to county, from county to place, from place to ZIP code, and from ZIP code to block all improve the overall predictive performance in terms of overall Accuracy. I include the input for

Figure 2: Accuracy Scores by Input Data



multi-unit occupancy only for the Census block models because this reflects the practical contexts where **bper** might be used. If a researcher has access to individual-level addresses, they should be able to geocode these to find the matching Census blocks and should also be able to parse the residency type (multi-unit or stand-alone). But researchers relying on more aggregate geographies likely do not have access to individual addresses, and hence the residency characteristics, for their sample.<sup>13</sup> In Figure 2 I pair each geolocation variable with a model missing the party ID input. The predictive gains from using party ID appear to diminish slightly with smaller levels of geography.

Accuracy scores, however, are an incomplete metric for assessing predictive performance. In contexts where the true distribution of classes is highly imbalanced, Accuracy can provide overly-optimistic results. For example, if we were to simply classify every individual as White in the North Carolina/Florida voter file, we would achieve 65% Accuracy without even trying! We can evaluate the models in a more rigorous way by looking at each ethnorace category separately.

Two better metrics to assess group-level predictions are Precision and Recall. Precision is the percentage of correctly classified individuals among all individuals predicted to belong to a specific ethnorace. It answers the question of how likely an individual’s predicted ethnorace in our sample matches their true, or self-identified, ethnorace. Recall, also known as Sensitivity or the True Positive Rate, is the percentage of all individuals who belong to a given ethnorace group which the model correctly classifies.

Precision and Recall reflect substantively important concerns for real-world applications of the method, and the inherent trade-offs between optimizing for either metric provide a balanced assessment of the predictive performance. On the one hand, Precision rewards very conservative classification procedures. We could, for example, only classify individuals as White if their predicted probability of being White was greater than 99%. This would ensure a very high Precision score for Whites because we are only capturing the low-hanging fruit. A conservative classification procedure like this, however, would likely result in extremely low Recall for Whites. If we only capture the low-hanging fruit, a greater share of White individuals will be mis-classified as non-White. Likewise, optimizing the algorithm for perfect Recall for Whites is trivial. By classifying every individual as White, we ensure 100% of Whites are correctly classified. Of course this procedure would result

---

<sup>13</sup>In the event that data is available, adding multi-unit occupancy inputs to aggregate geographies, such as ZIP codes, places, counties, or states *greatly* enhances the predictive accuracy of the model.

in extremely low Precision for Whites because every non-White individual would be classified as White as well. Achieving both high Precision and high Recall, therefore, is a difficult task.

Figure 3: Precision/Recall Scores by Ethnorace and Input Data

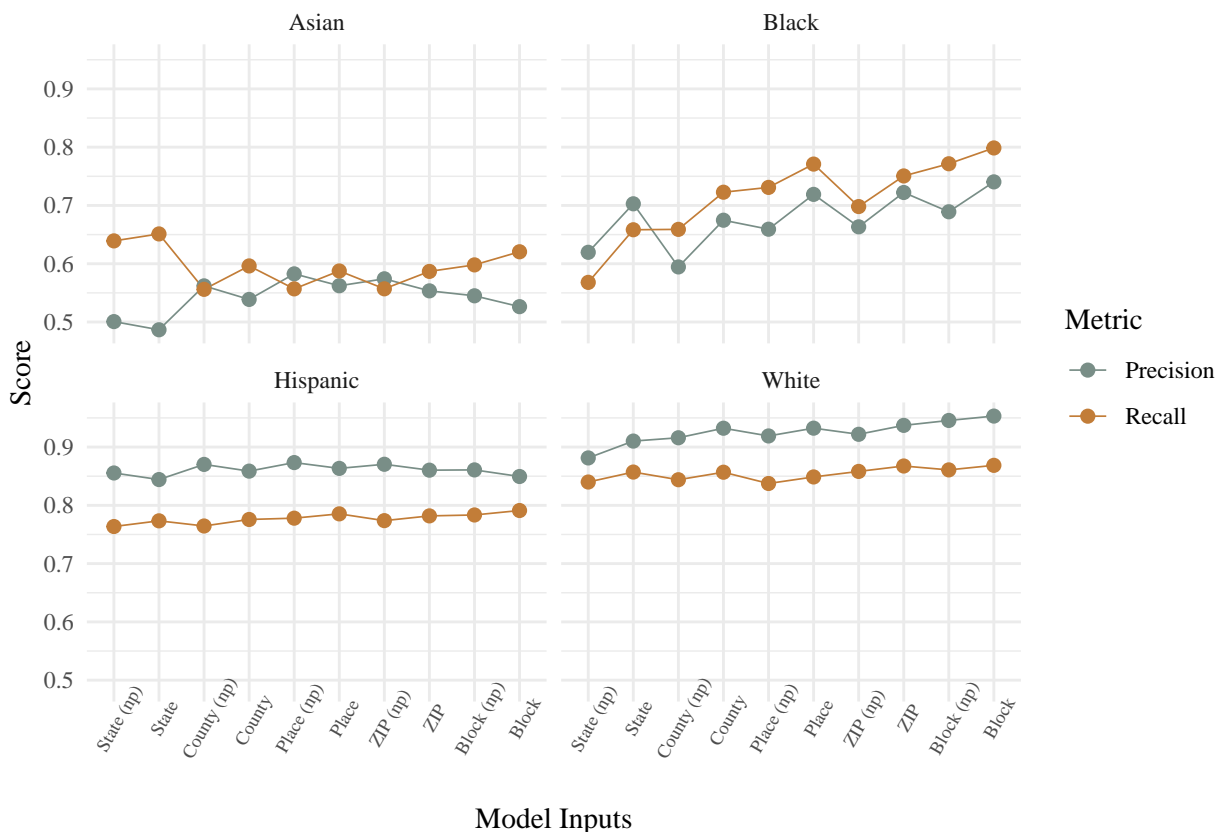


Figure 3 displays the Precision and Recall scores broken down by ethnorace for the same models used in Figure 2.<sup>14</sup> For both Whites and Hispanics, Precision and Recall remain high across all models. Using the most input data available (the models on the far right of the figure), White Precision and Recall are 0.953 and 0.869, respectively. Hispanic Precision and Recall in this model are 0.849 and 0.791, respectively. Hispanic predictive performance across all models likely remains high due to the distinctiveness of Spanish surnames. These metrics are uniformly lower for Asians—likely reflecting the relative rarity of Asian individuals in the sample and the heterogeneous qualities of this particular group.<sup>15</sup> When using Census blocks as the unit of geography, predictions for African Americans are quite strong (Precision: 0.74, Recall: 0.799). But predictive performance

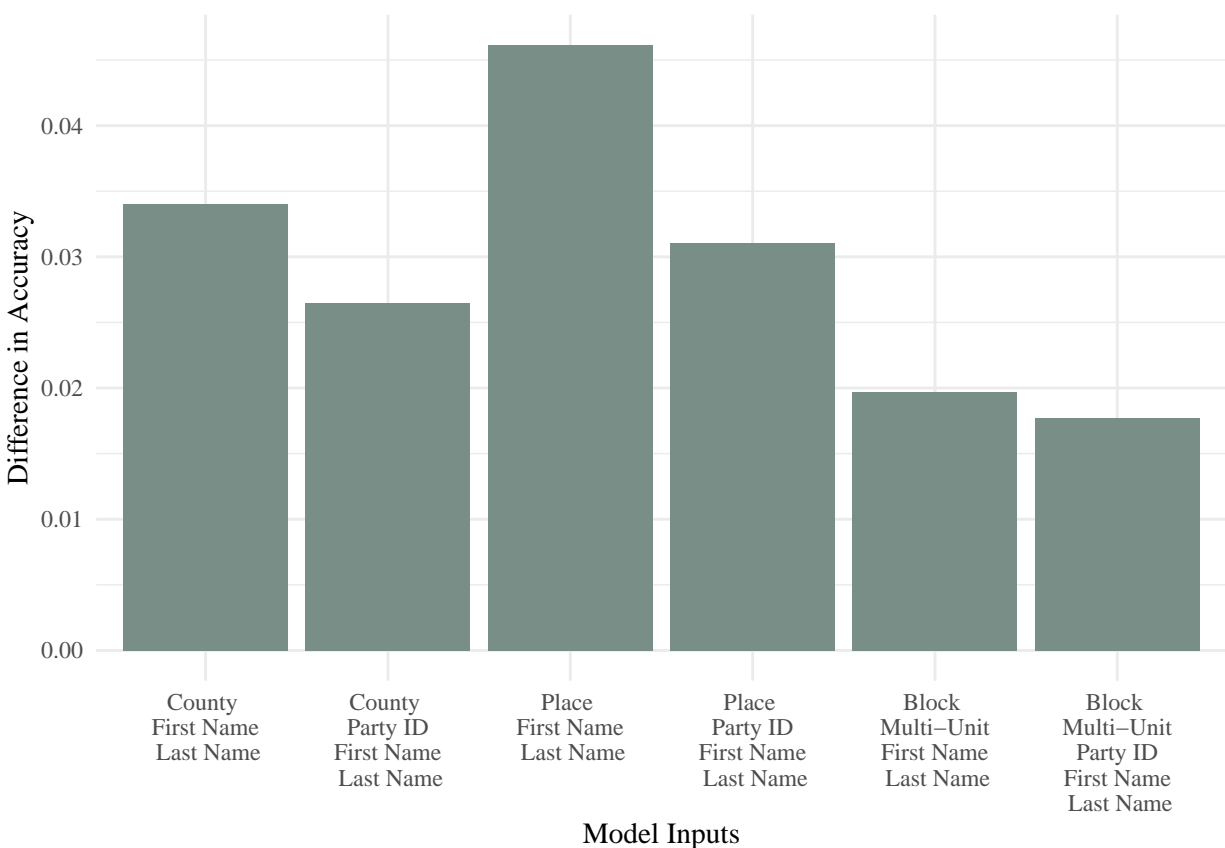
<sup>14</sup>The model names are abbreviated by the level of geography, and all use first name and last name inputs. Models with (np) do not use the party ID inputs.

<sup>15</sup>I investigated which ethnoraces Asian Americans were wrongly predicted to be and found a wide spread: 65% White, 18% Hispanic, 14% Black, 1% Other

falls with broader geolocations. At the state level, without party ID, Black Precision is 0.62 and Recall is 0.568. This difference in predictive performance between geographies could be explained by the legacy of Black segregation in at least the North Carolina/Florida sample.

The classification metrics detailed above provide some information about the predictive shortcomings of my method. These limitations notwithstanding, however, the ethnorace predictions made by **bper** are nearly uniformly better than those from **wru** (Imai and Khanna 2016). Figure 4 displays the difference in Accuracy between **bper** and **wru** using the same input data discussed previously.<sup>16</sup> The baseline of zero in the figure represents the Accuracy scores from each model using **wru**, and the height of the bars display the change in Accuracy using **bper**. Across all model types, **bper** scores between 1.8 and 4.6 percentage points higher on this metric. Regardless of input data, **bper** classifies a higher proportion of individuals correctly in the North Carolina/Florida voter file.

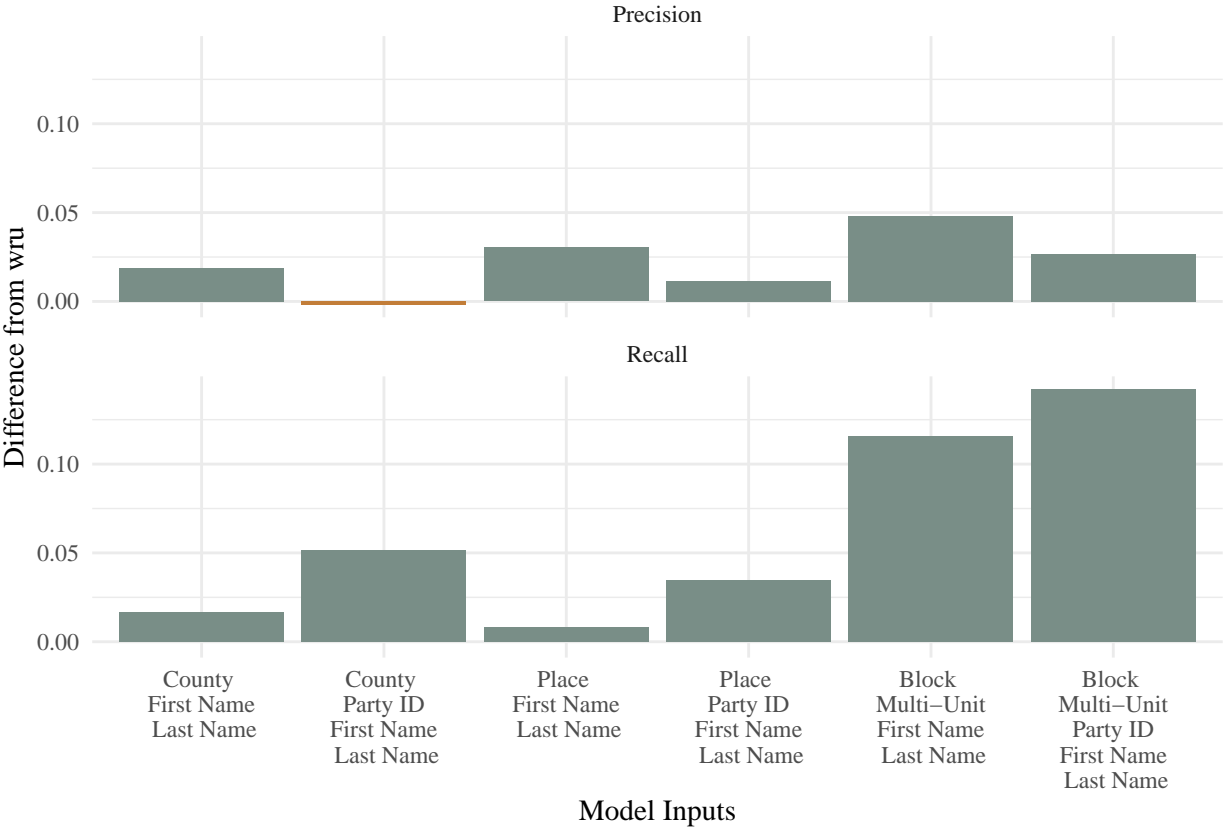
Figure 4: Accuracy Score Comparison to **wru**



<sup>16</sup>The **wru** package does not provide ZIP code or state level predictions so those models are excluded from the comparison.

Figures 5 through 8 show the same comparisons between **wru** and **bper** for Precision and Recall across different ethnoraces. For Asians in Figure 5, the comparison strongly favors **bper**. Across almost all sources of input data, **bper** out-performs **wru** in terms of both Precision and Recall. In particular, using block geolocations leads to between a 10 and 12 percentage point increase in Recall. A greater proportion of self-identified Asian individuals are correctly classified. Using county geolocations along with party ID (the model second from the left) we see a very slight decrease in Precision, but this is compensated for a 5 percentage point increase in Recall under these conditions.

Figure 5: Precision and Recall Score Comparison to **wru**: Asian



The comparison among predictions for Black individuals in Figure 6 are striking. My method shows dramatic gains in Precision for all models while also modestly increasing Recall. Without party ID, **bper** improves upon **wru** by over 25 percentage points in Precision using place or county geolocations. This means that individuals predicted to be Black by **bper** are over 25 percentage points more likely to self-identify as Black relative to the predictions generated by **wru**. The modest gains in Recall across all models signifies that **bper** is correctly predicting a greater share of self-



identified African Americans in the sample as well.

Figure 6: Precision and Recall Score Comparison to **wru**: Black

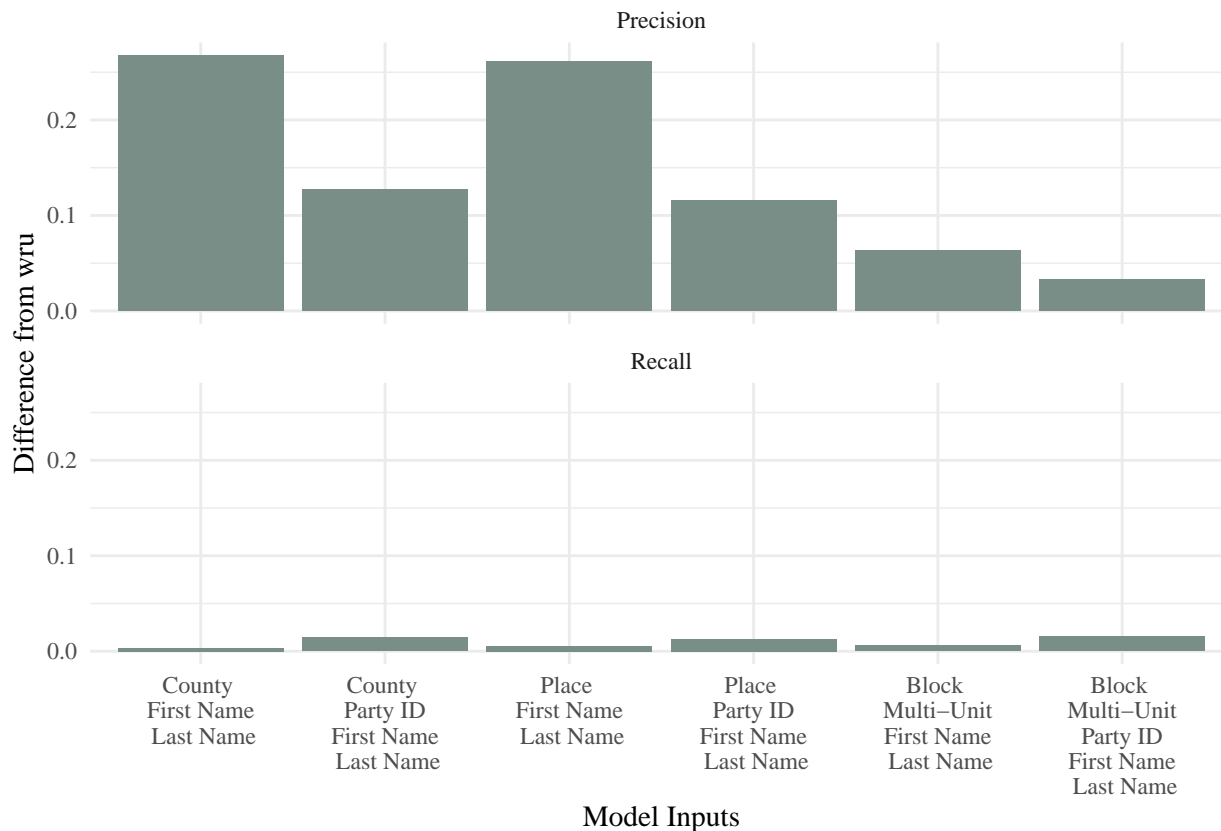


Figure 7 shows the comparison in predictions among Hispanics. Again, here **bper** out-performs **wru** in both Precision and Recall for every model. The magnitudes of these differences, however, are smaller than those for Black or Asian predictions. This is likely due to distinct Spanish last names, which already provide a lot of predictive information for Hispanics. Therefore, the additional predictive improvements in **bper**, such as first name information and smoothing, have less to contribute.

Lastly, Figure 8 displays the results for the White predictions. In all models except county and place geolocations (without party ID), **bper** out-performs **wru** on both Precision and Recall for White individuals. In the two exceptions, however, losses to Precision are compensated with gains in Recall. This means that, while **bper** is doing a better job classifying a higher proportion of Whites, it is mis-classifying more non-Whites as White compared to **wru** for county and place geolocations without party ID.

In sum, across almost all models of varying input data **bper** outperforms **wru** on both Preci-

Figure 7: Precision and Recall Score Comparison to `wru`: Hispanic

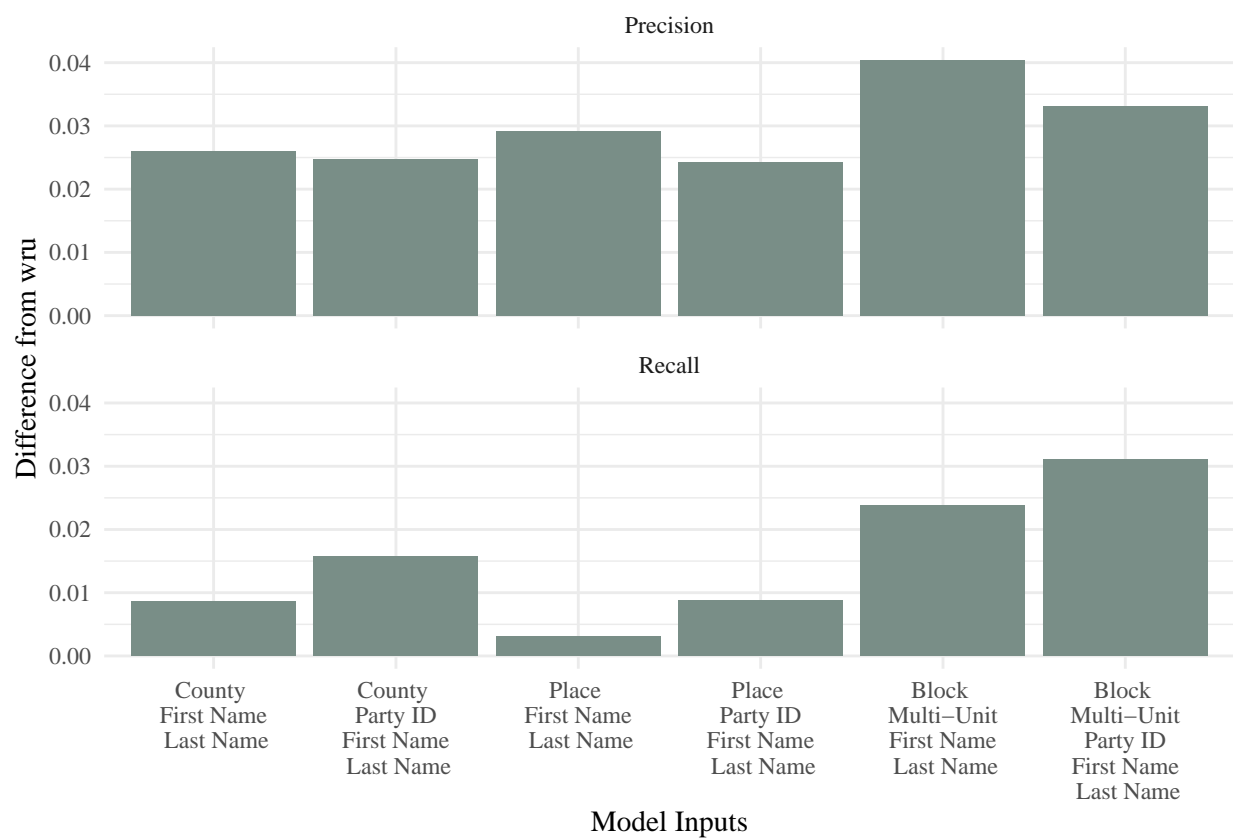
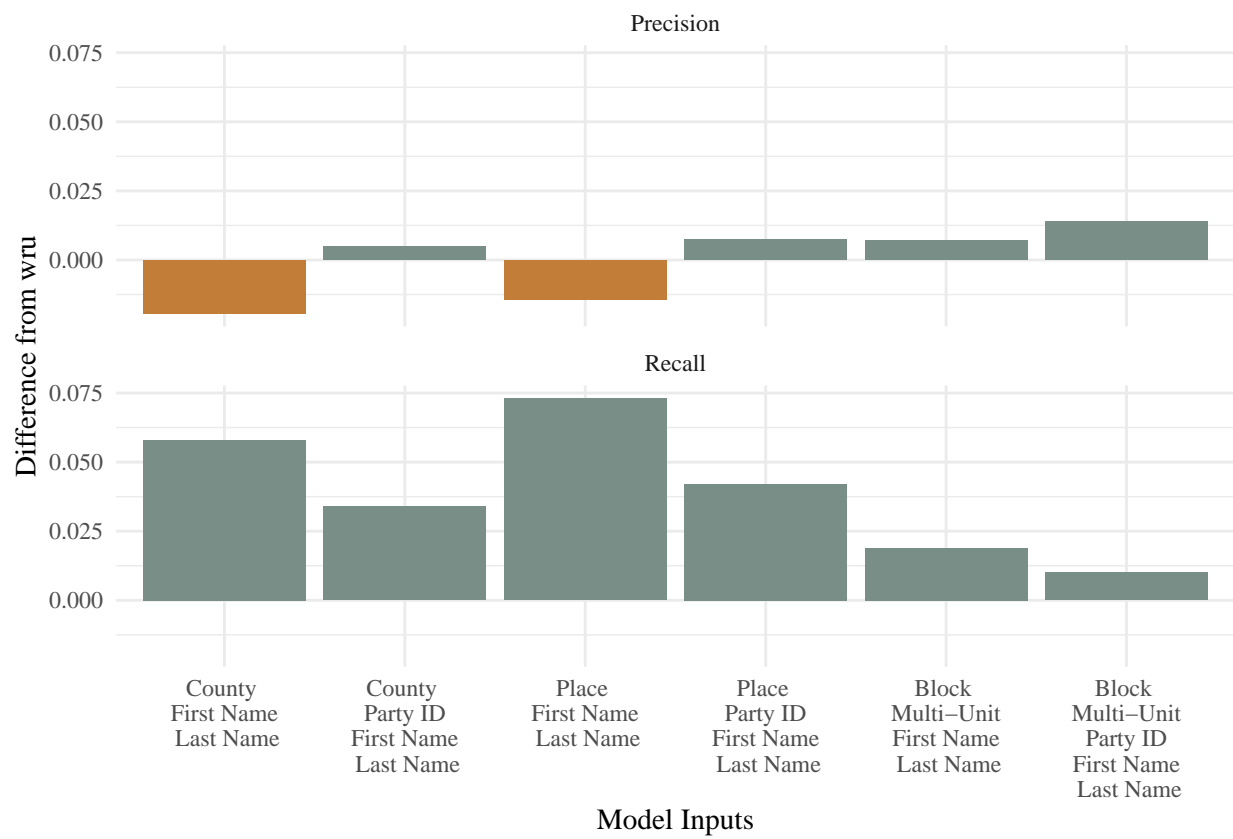


Figure 8: Precision and Recall Score Comparison to `wru`: White



sion and Recall for each group. When we transition to more aggregate geographies, **bper** shows dramatic gains in Precision for Asians and African Americans. Without party ID, there is some loss to Precision for Whites in aggregate geographies, but gains in Recall are more than triple the magnitude for these models.

## Replication Study

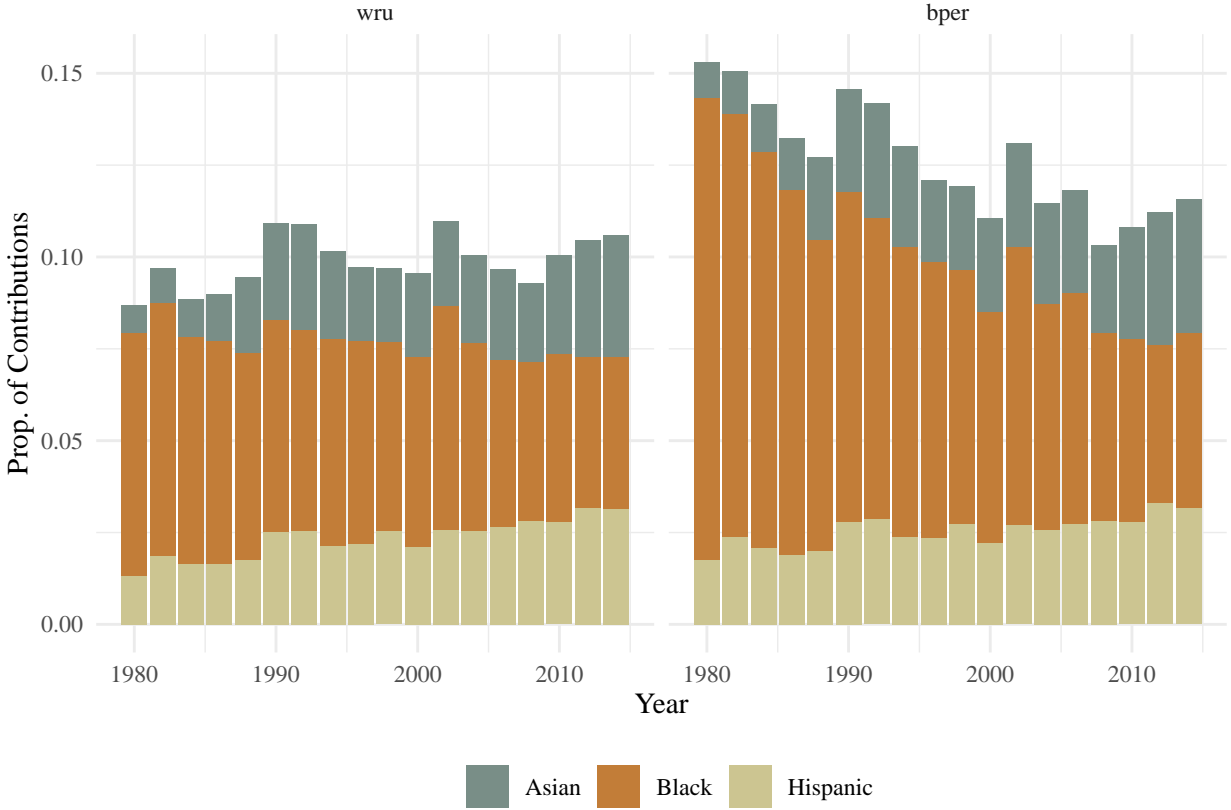
In “Race and Representation in Campaign Finance” (2020), Jacob Grumbach and Alexander Sahn investigate racial inequality in campaign contributions. Descriptively, they find that Latinos and African Americans are underrepresented among individuals who donate to US House campaigns. They also use differences-in-differences and regression discontinuity analyses to show that when candidates of color run, the contributor class becomes more ethnoracially representative because these candidates garner more co-ethnic donations.

Grumbach and Sahn rely primarily on the Dataset on Ideology and Money in Elections (DIME) from Bonica (2013). These data contain information—such as first names, last names, genders, addresses (for contributors), and political parties (for recipients)—for both the contributors and recipients of campaign contributions from 1980 to 2014. Crucially, however, the DIME data do not contain individuals’ self-reported race or ethnicity. Grumbach and Sahn therefore use **wru** to predict the ethnorace of both contributors and recipients. By replacing the ethnorace predictions used in their study with those from **bper** I am able to demonstrate the substantive implications of switching to more accurate predictions.

The replication files for their study, however, do not contain the code which implements the ethnorace prediction steps. It has proved impossible to replicate their study exactly, so instead I reproduced their descriptive analysis to the best of my ability using both **wru** and **bper**. The DIME data contain ZIP codes for contributors, so I use that level of geography along with first name, last name, and gender as input data in **bper** for predicting contributor ethnoraces. Candidate data is limited to first names, last names, party ID, and states. For the **wru** predictions I use last names and counties for contributors because **wru** does not allow for ZIP code geographies. And for candidate **wru** predictions I use last names and party ID.

In their original paper, Grumbach and Sahn appear to match candidates to counties from

Figure 9: Ethnoracial Composition of the Contributor Class (1980 - 2014)



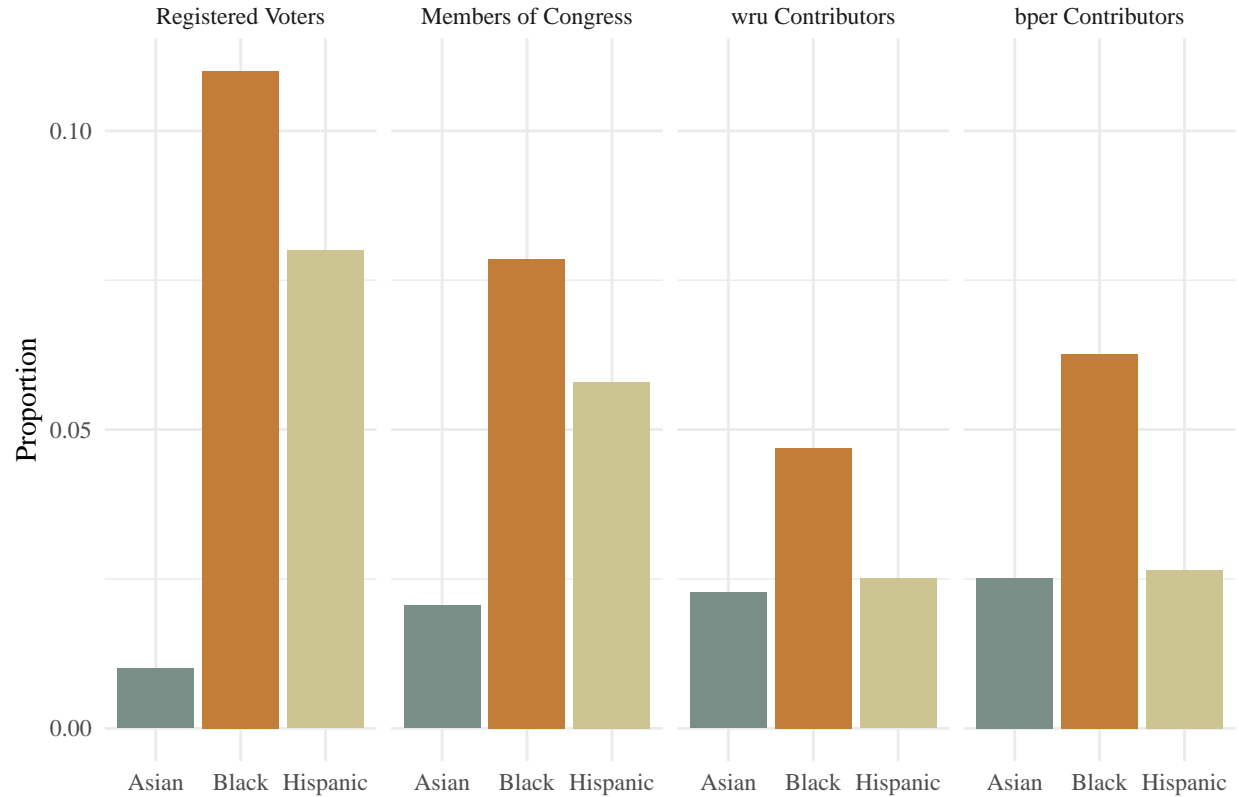
which their donors reside during their race prediction coding. There is no guarantee, however, that candidates and donors reside in the same county—or that county demographics of donors match those of the candidates. Furthermore, because the DIME data is at the level of individual contribution, this procedure classifies candidates as potentially belonging to different ethnoraces over time. For example, Rep. Jim Bacchus appears in Grumbach and Sahn’s replication data as Black, White, and Other in various cycles. Roughly 1,360 candidates have multiple predicted ethnoraces in their data but these errors are obscured in subsequent aggregation procedures. For this reason I calculate my own ethnorace probabilities with **wru** separately for this replication study. However, this means that the results I report using **wru** differ slightly from those reported in the published article.

Figure 9 (Figure 2 in Grumbach and Sahn 2020) shows the proportion of total campaign contributions by ethnorace by year.<sup>17</sup> The left-hand panel shows these results using **wru**’s predictions and matches closely the figure from the published article. Compared to **wru**, the results from **bper** show

<sup>17</sup>As in the original figure, contributions from Whites are omitted to display Asian, Black, and Hispanic contributions better.

substantially larger shares of contributions from individuals of color. Additionally, there appears to be a dramatic decline over time in contributions from African Americans using **bper** predictions. African Americans accounted for 12% of the contributions in 1980 compared to only 5% in 2014. This pattern exists, but is less clear in the **wru** panel (6% Black contributions in 1980 compared to 4% in 2014). Switching to **bper** predictions thereby reveals a contributor class which is more diverse than the one described by Grumbach and Sahn. However, the **bper** predictions show a worrying decline in this diversity—the explanation for which could be the subject of future research.

Figure 10: Ethnoracial Composition of the Contributor Class Versus Electorate



Whereas Figure 9 showed the proportions of *contributions* by ethnorace, Figure 10 (Figure 3 in original) displays the proportions of *contributors* by ethnorace relative to other forms of political participation. Both **wru** and **bper** predictions show that Black and Hispanic Americans are underrepresented in the contributor class. This finding, along with the increased representation of Asian Americans relative to their proportion among registered voters, replicates Grumbach and Sahn’s descriptive analysis from the published article. However, **bper** predictions show a substantially higher share of Black contributors compared to the predictions from **wru**. We would need to

disaggregate these results by election cycle to make a comparison to Figure 9, but it appears that African Americans may be more likely to contribute relative to Asian and Hispanic Americans, but make smaller average contributions.

Figure 11: Average Total Contributions by Ethnorace

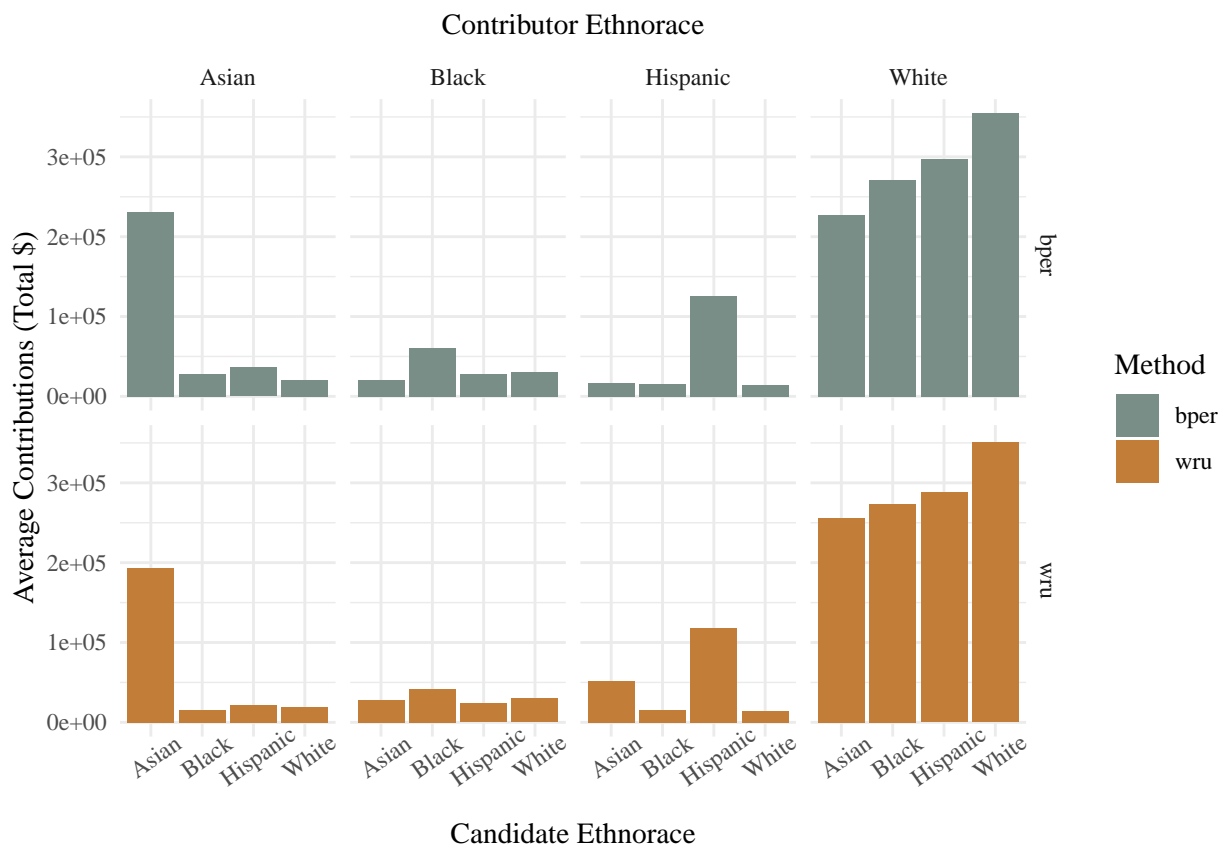


Figure 11 (Figure 4 in original) uses both the ethnorace of contributors (columns) and the ethnorace of candidates (x-axis). It displays the average total contributions candidates receive from contributors of a given ethnorace. Similar to the original study, I find strong evidence of co-ethnic contributing. Among each ethnorace's contributors, either a plurality or a majority of contributions go to candidates who share the same ethnorace. This finding is supported more by the **bper** predictions than the **wru** predictions. Contributions from Black donors appear to be more heavily concentrated among Black candidates when using **bper**. The same pattern of increased co-ethnic concentration appears for Hispanic contributors and candidates. The descriptive findings from Figure 11 form the inspiration for the subsequent causal analysis Grumbach and Sahn pursue in their study. Given some of the patterns of increased co-ethnic contributing from the **bper**

predictions, therefore, we may expect to see stronger effects in replications of their differences-in-differences and regression discontinuity results.<sup>18</sup>

## Conclusion

This paper describes a powerful method for predicting individuals' race or ethnicity based on other known attributes. Through validation tests, I have shown that my method produces better predictions than those used in recent social science research. These improvements are not trivial. If researchers only have access to individuals' name and county or place, my method nearly doubles the probability that someone predicted to be Black self-identifies as Black relative to predictions generated from other methods (Imai and Khanna 2016). These improved predictions reveal new empirical findings as demonstrated by my replication of Grumbach and Sahn (2020).

---

<sup>18</sup>These replications could not be completed in time for the current draft.



## References

- Abott, Carolyn, and Asya Magazinnik. 2020. "At-Large Elections and Minority Representation in Local Government." *American Journal of Political Science* 64 (3): 717–33. <https://doi.org/10.1111/ajps.12512>.
- Burch, Traci. 2013. *Trading Democracy for Justice: Criminal Convictions and the Decline of Neighborhood Political Participation*.
- Cantoni, Enrico. 2020. "A Precinct Too Far: Turnout and Voting Costs." *American Economic Journal: Applied Economics* 12 (1): 61–85. <https://doi.org/10.1257/app.20180306>.
- Crabtree, Charles, and Volha Chykina. 2018. "Last Name Selection in Audit Studies." *Sociological Science* 5: 21–28. <https://doi.org/10.15195/v5.a2>.
- Domingos, Pedro, and Michael Pazzani. 1997. "Beyond Independence: Conditions for the Optimality of the Simple Bayesian Classifier." *Machine Learning* 29: 103–30.
- Elliott, Marc N., Peter A. Morrison, Allen Fremont, Daniel F. McCaffrey, Philip Pantoja, and Nicole Lurie. 2009. "Using the Census Bureau's Surname List to Improve Estimates of Race/Ethnicity and Associated Disparities." *Health Services and Outcomes Research Methodology* 9 (2): 69–83. <https://doi.org/10.1007/s10742-009-0047-1>.
- Enos, Ryan D., Aaron R. Kaufman, and Melissa L. Sands. 2019. "Can Violent Protest Change Local Policy Support? Evidence from the Aftermath of the 1992 Los Angeles Riot." *American Political Science Review* 113 (4): 1012–28. <https://doi.org/10.1017/S0003055419000340>.
- Fotheringham, A S, and D W S Wong. 1991. "The Modifiable Areal Unit Problem in Multivariate Statistical Analysis." *Environment and Planning A: Economy and Space* 23 (7): 1025–44. <https://doi.org/10.1068/a231025>.
- Fraga, Bernard L. 2018. *The Turnout Gap. Race, Ethnicity, and Political Inequality in a Diversifying America*. Cambridge University Press.
- Grumbach, Jacob M., and Alexander Sahn. 2020. "Race and Representation in Campaign Finance." *American Political Science Review* 114 (1): 206–21. <https://doi.org/10.1017/S0003055419000637>.
- Grumbach, Jacob M., Alexander Sahn, and Sarah Staszak. 2020. "Gender, Race, and Intersectionality in Campaign Finance." *Political Behavior*, June. <https://doi.org/10.1007/s11109-020->

09619-0.

- Hepburn, Peter, Renee Louis, and Matthew Desmond. 2020. "Racial and Gender Disparities Among Evicted Americans." *Sociological Science* 7: 649–62. <https://doi.org/10.15195/v7.a27>.
- Imai, Kosuke, and Kabir Khanna. 2016. "Improving Ecological Inference by Predicting Individual Ethnicity from Voter Registration Records." *Political Analysis* 24 (2): 263–72. <https://doi.org/10.1093/pan/mpw001>.
- Keele, Luke, and Rocío Titunik. 2018. "Geographic Natural Experiments with Interference: The Effect of All-Mail Voting on Turnout in Colorado." *CESifo Economic Studies* 64 (2): 127–49. <https://doi.org/10.1093/cesifo/ify004>.
- Kuk, John, Zoltan Hajnal, and Nazita Lajevardi. 2020. "A Disproportionate Burden: Strict Voter Identification Laws and Minority Turnout." *Politics, Groups, and Identities*, June, 1–9. <https://doi.org/10.1080/21565503.2020.1773280>.
- Lewis, David D. 1998. "Naive (Bayes) at Forty: The Independence Assumption in Information Retrieval." In *Machine Learning: ECML-98*, edited by Claire Nédellec and Céline Rouveirol, 1398:4–15. Berlin, Heidelberg: Springer Berlin Heidelberg. <https://doi.org/10.1007/BFb0026666>.
- Omi, Michael, and Howard Winant. 2014. *Racial Formation in the United States*. Routledge.
- Rish, Irina. 2001. "An Empirical Study of the Naive Bayes Classifier." *IJCAI 2001 Workshop on Empirical Methods in Artificial Intelligence* 3 (22): 41–46.
- Studdert, David M., Yifan Zhang, Sonja A. Swanson, Lea Prince, Jonathan A. Rodden, Erin E. Holsinger, Matthew J. Spittal, Garen J. Wintemute, and Matthew Miller. 2020. "Handgun Ownership and Suicide in California." *New England Journal of Medicine* 382 (23): 2220–9. <https://doi.org/10.1056/NEJMsa1916744>.
- Trounstine, Jessica. 2020. "The Geography of Inequality: How Land Use Regulation Produces Segregation." *American Political Science Review* 114 (2): 443–55. <https://doi.org/10.1017/S0003055419000844>.
- Tzioumis, Konstantinos. 2018. "Demographic Aspects of First Names." *Scientific Data* 5 (1): 180025. <https://doi.org/10.1038/sdata.2018.25>.
- Velez, Yamil Ricardo, and Benjamin J. Newman. 2019. "Tuning in, Not Turning Out: Evaluating the Impact of Ethnic Television on Political Participation." *American Journal of Political*

*Science* 63 (4): 808–23. <https://doi.org/10.1111/ajps.12427>.

Voicu, Ioan. 2018. “Using First Name Information to Improve Race and Ethnicity Classification.”

*Statistics and Public Policy* 5 (1): 1–13. <https://doi.org/10.1080/2330443X.2018.1427012>.

Zadrozny, Bianca, and Charles Elkan. n.d. “Obtaining Calibrated Probability Estimates from Decision Trees and Naive Bayesian Classifiers,” 8.