

# Predicting postsecondary attendance by family income in the United States using multilevel regression with poststratification

Benjamin T. Skinner\*  
University of Florida

William R. Doyle  
Vanderbilt University

April 5, 2022

## Abstract

Despite tens of billions of dollars in yearly public spending to fund grants for higher education for youth from low-income families, no government agency tracks how many young people from low-income families enroll in higher education by state. Proxy measures like the number of college students who receive Pell grants address the number of already enrolled students who come from low-income families rather than tracking the rate of enrollment among the overall number of young people who are from low-income families. Estimates of postsecondary enrollment among low-income young people from U.S. Census surveys likely overestimate enrollment among this population due to their design and administration. In this paper we use multilevel regression with poststratification to estimate postsecondary attendance rates by family income in 50 states and the District of Columbia for a recent cohort of young people. Our application of Bayesian techniques for estimation and inference allow us to make appropriate statements of uncertainty regarding our estimates, including the probability that low-income young people will attend higher education in a given state.

**Keywords:** college access, low income, multilevel regression with poststratification

---

\**Corresponding author:* btskinner@coe.ufl.edu. We thank Richard Blissett, Taylor Burtch, Josh Clinton, Chris Fonnesebeck, Heather Jack, Daniel Klasik, and participants at both 2021 ASHE and 2016 AEFPP conferences for their helpful comments on various iterations of this paper. All errors and mistakes in interpretation remain our own.

In the United States, the vast majority of federal, state, and institutional financial aid money goes to individuals from low-income families in order to increase their college participation rates (Abraham & Clark, 2006; Baum & Payea, 2013). Thanks to large-scale longitudinal surveys of high school-aged young people regularly conducted since the late 1970s, we have reliable national and regional level estimates of postsecondary attendance by family income over time (Adelman et al., 2003; Berkner & Chavez, 1997; R. Bozick & Lauff, 2007). Despite decades of financial aid directed toward low-income youth, estimates from these surveys have consistently shown large gaps in postsecondary attendance by family income and wealth (Lovenheim & Reynolds, 2013). According to the most recent federal longitudinal survey of students transitioning from high school to young adulthood, 49 percent of low-income young people attend some form of postsecondary education compared to 74 percent of middle-income young people (Duprey et al., 2020).

In order to evaluate the efficacy of state-level financial aid policies, it would be beneficial to have state-specific estimates of college attendance by family income. By virtue of their design, however, nationally representative federal surveys like those noted above cannot be used to make inferences at the state level. Looking to other data sources, states with well-structured administrative data systems may be able to track the inverse, the proportion of college students from low-income families currently enrolled in college,  $P(\text{low income} \mid \text{college})$ . If their systems span the P-20 pipeline, they may even be able to provide estimates of the proportion of interest,  $P(\text{college} \mid \text{low income})$ , among those students who remain within the state system. Should these latter estimates exist, however, they are not widely available. Across the country, we find no systematic measures of postsecondary enrollment by family income at the state level.

Other means of estimating low-income youth enrollment in college have their own problems. State-specific estimates of low-income youth college enrollment from U.S. Census Bureau surveys cannot be trusted due to their household-based sampling designs. To determine income-based eligibility, financial aid programs typically rely on the Free Application for

Federal Student Aid (FAFSA), which asks questions about parental income in addition to student income (Dynarski & Scott-Clayton, 2013). Dependent students, even those living apart from the parents or guardians who financially support them, have their need categorized by their family’s financial situation via the FAFSA. The household-based sampling designs of Census surveys only capture incomes of people in a single household, which is not the same as family or parental income for many young people. Because Census data often misclassify financially dependent young people from moderate to high-income families as low income (U.S. Census Bureau, 2009), estimates from these sources likely overestimate rates of college attendance among low-income young people.

Based on these difficulties, it may be tempting to use estimates of  $P(\text{low income} \mid \text{college})$  taken from institutional sources such as the Integrated Postsecondary Education Data System (IPEDS) as proxies for estimates of  $P(\text{college} \mid \text{low income})$ . Nevertheless, there is no necessary correlation between the proportion of low-income college students in a state with the proportion of youth from low-income families who enroll in higher education:  $P(\text{low income} \mid \text{college}) \neq P(\text{college} \mid \text{low income})$ . It is entirely possible that states with high proportions of low-income college students still have large populations of low-income young people who fail to enroll. Conversely, it is also possible that many low-income young people who enroll are not identified as such because they do not complete the FAFSA and therefore do not benefit from financial supports they might otherwise receive (Kofoed, 2017). A lack of clear estimates of enrollment by family income at the state level means that neither federal and state policymakers nor college administrators have clear evidence that their policies are effective in increasing college participation among low-income youth.

To fill this informational gap, we use multilevel regression with poststratification, MRP (Park et al., 2004), to estimate postsecondary attendance rates by family income in each of the 50 states and the District of Columbia. In this multi-step procedure, we first estimate the probability of postsecondary enrollment by family income among strata of a recent cohort of students surveyed in the High School Longitudinal Study of 2009 (Duprey et al., 2020).

We next weight these predictions with matched state-level population characteristics in order to estimate college attendance rates by family income within each state. To validate our procedure, we apply the MRP procedure to simulated data with known properties as well as compare MRP estimates to representative estimates that can be computed using survey-provided sampling weights. We conclude by comparing our MRP estimates to Census- and IPEDS-derived estimates of low-income youth enrollment.

Our contribution to the literature is to provide the best possible state-specific estimates of the probability of college attendance among low-income youth. In addition to finding that low-income young persons are less likely than their higher-income peers to attend college overall, we find substantial variation in both the difference in attendance rates between these groups and their respective attendance rates across the states and the District of Columbia. Our estimates differ substantively from those taken directly from the U.S. Census as well as proxy measures such as the proportion of currently enrolled students who are Pell eligible and those who fall within the lowest family income band. While we believe our principled MRP procedure offers robust estimates of college enrollment among low-income youth across the United States, our ultimate hope is that our findings are obviated by better data collection and reporting in the future.

## Background

It is well-established that young people from low-income families are less likely to attend higher education than their peers (Chen et al., 2017), even though higher education could serve as a means of increasing their opportunity (Corak, 2013). Many researchers have documented the struggles of young people from low-income families to attend higher education (Bell et al., 2009; Dowd, 2004; McDonough, 1997; Perna & Jones, 2013). In addition to poorer access to advising resources at their schools and often less accurate beliefs about the potential costs and benefits of college (Perna, 2006), one particular reason for comparatively lower

rates of college attendance may be that young people from low-income families are more price responsive than their peers (Bartik et al., 2021; Carneiro & Heckman, 2002; Deming & Dynarski, 2009; Dynarski, 2002; Dynarski et al., 2021; Heller, 1997; T. J. Kane, 2006; Leslie & Brinkman, 1987), substantiating the need for income-based financial aid (T. Kane, 1999).

Given that young people from low-income families face the biggest barriers to attending college, the nation should keep track of attendance rates in higher education by family income (Deming & Dynarski, 2009). Furthermore, these data should be available at the state level, as states are the key players in setting higher education policy (Callan & Finney, 1997; Denning, 2019; Gurantz, 2022; Richardson et al., 1999; Scott-Clayton & Schudde, 2020; Zumeta et al., 2012). State leaders decide the amount of higher education that will be supplied, its price, and play a large role in determining how much financial aid will be provided to students (Zumeta et al., 2012). Yet state leaders make all of these decisions without clearly knowing how enrollment levels across different family income groups will be impacted. Federal data sets from the U.S. Census and Education Department cannot provide accurate estimates of college attendance rates by family income at the state level due to their design. If the states themselves are tracking enrollment by income, we have not been able to find any reports of their findings. A comprehensive review of state policy documents yielded no studies that track the proportion of individuals who attend college by family income.

Calculating postsecondary attendance rates by family income at the national level using federal surveys is a relatively simple task often done in the context of a larger project. Reports generated by the National Center for Education Statistics (NCES) find substantial gaps in postsecondary attendance by family income, even after taking into account the academic preparation of high school graduates (Berkner & Chavez, 1997; R. Bozick & Lauff, 2007). Duncan & Murnane (2011) find that gaps in postsecondary attendance by income did not close over the two decades prior to their study. NCES longitudinal surveys are designed to be generalizable to a national and often regional population of high school students (e.g., Ingels et al., 2014). In the case of the High School Longitudinal Study of 2009 (HSLSO9),

they are representative in a small subset of states (Duprey et al., 2020). These surveys are not designed, however, to support estimates at the state level across all states. Because their complex sampling procedures do not include stratification at the state level or cluster sampling within states, sampled students are not representative of the state in which they attended high school, even when restricted-use data files that indicate students’ states of enrollment are used (Ingels et al., 2014).

The U.S. Census Bureau’s American Community Survey (ACS) and Current Population Survey (CPS) represent another pair of data sources that one might consider using to estimate college attendance by family income level. Nevertheless, these surveys are plagued by a common problem that is particularly severe in the case of young people such as college students. Both Census surveys utilize a sampling procedure based on the households where people live (U.S. Census Bureau, 2012a, 2013). This design does not generally present a problem for estimating overall family income levels for older respondents, even at the state level. Estimating family income levels for young adults, however, represents a challenge. Survey questionnaires ask for information on all persons who are residents of the household at the time the forms are received. Residency follows the “two-month rule,” which says that

[i]n general, people who are away from the sample unit for two months or less are considered to be current residents, even though they are not staying there when the interview is conducted, while people who have been or will be away for more than two months are considered not to be current residents (U.S. Census Bureau, 2009, pp. 6–2).

While specific allowances are made for child dependents under 18 years of age who are away for boarding school, none are made for college-age dependents who are away for longer than the two month period.

The sampling design of Census surveys such as ACS means that when researchers limit analyses to observations of students enrolled in college, they are likely to miss a significant proportion of students whose families were interviewed while they were living apart for school.

In addition, researchers are likely to receive limited family information on young people living in group housing. For this latter group, data limitations are due to the fact that when a household is surveyed, members are asked to report their own income, which is later combined to create an estimate of total household income. While incomes of close relatives within a household are separately combined to form a family income, the income of family members outside of the household, which includes parents of a still dependent survey respondent, is not included in this estimate. Thus, if a surveyed young person lives in a household other than that of their parents or guardians for longer than the two-month window—in a household made up of other college students, for example—their reported family income will reflect only their individual earnings and will likely be lower than expected considering their true financial dependency. This would be the case even among students who came from high-income families.

The end result is that family income levels reported by young people between 18- and 24-years-old are unreliable (U.S. Census Bureau, 2012b, 2012a, 2013). Because estimates of college attendance by family income taken from household-based Census data are likely to overstate the proportion of young people from low-income families in higher education, they should not be used to estimate college attendance rates by family income. Otherwise, such estimates, were they to be used in the evaluation of student aid policies, could lead researchers and policymakers alike to be more sanguine about the efficacy of family income-based access policies than is warranted.

What we do have at the state and institutional level are proxy measures of the proportion of students *already in college* who are from low-income families,  $P(\text{low income} \mid \text{college})$ . Many states and institutions, for instance, report the proportion or number of students who are Pell eligible (Tebbs & Turner, 2005). Yet indirect measures like the percentage of students who are Pell eligible do not directly measure the college-going rates of young people from low-income families. These measures ignore young persons who are not enrolled in higher education or, in the case of Pell grants, do not apply for federal student aid.

As a result, researchers and policymakers do not currently know what percentage of young people from different family income levels transition from high school to college in each state.

Understanding the impact of state and federal policy on postsecondary attendance rates requires an inverse of the above estimate, that is, an estimate of the probability that a young person will enroll in college conditional on being from a low-income family,  $P(\text{college} \mid \text{low income})$ . We have established that such estimates cannot be produced by a straightforward analysis of currently available federal longitudinal or cross-sectional data sets. Instead, we must utilize both types of data in order to come up with estimates of the likely level of college attendance by income in each state. Our approach, which we describe in more detail in the next section, begins with estimating the probability of enrollment in postsecondary education using data that covers the period from when students are first enrolled in high school to when they are 18-19 years old, graduated from high school, and eligible to enroll in college. From this federal longitudinal survey, we use characteristics of high school students aged 14-15 years-old that include their family income, race/ethnicity, and gender to predict the probability that a student will attend postsecondary education in the year after high school graduation. We then turn to state-level Census data from the base year of the longitudinal survey to get counts of the number of 14 and 15 year olds in each state with the same combinations of characteristics as those in our sample of high school students. Using these counts to weight our estimates, we both bypass issues with determining family income for young people aged 18 and 19 years old and are able to produce representative state-level estimates of college attendance among low-income young persons.

## Methodology

To recover state-level estimates of enrollment among low-income young persons in college, we use multilevel regression with poststratification (MRP), a statistical technique that has been widely used in the political science literature to estimate public opinion (Gao et al.,



2019; Gelman et al., 2010; Gelman & Little, 1997; Howe et al., 2015; Jonathan P. Kastellec et al., 2019; Kennedy & Gelman, 2019; Lax & Phillips, 2009; Lei et al., 2017; Lipps & Schraff, 2019; Little, 1993; Pacheco, 2011; Park et al., 2004; Wang et al., 2015; Warshaw & Rodden, 2012). Researchers in other disciplines such as public health (Downes et al., 2018; Eke et al., 2016; Zhang et al., 2014) and education policy (Ortagus et al., 2021) have also used MRP to produce representative estimates using non-representative data.

MRP works using two data sets and two primary analysis steps. First, a multilevel model of the form

$$P(y_i = 1) = \text{logit}^{-1}(\beta_0 + \sum_{k=1}^K \alpha_{j[i]}^k + Z_i \gamma) \quad (1)$$

is fit to a binary outcome of interest,  $y_i$ , using observations,  $i$ , from the first data set that contains non-representative survey responses. The outcome could be voting for a particular candidate, supporting a policy position, or, in our case, enrolling in college. The right-hand side parameters in the model include a grand mean,  $\beta_0$ , and a suite of random intercepts,  $\alpha^k$ , indicating demographic categories and geographic areas that separate each observation into a limited number of population cells,  $j$ . In addition, the right-hand side of the model includes second-level covariates and parameters,  $Z_i \gamma$ , that are associated with the area to which one wishes to poststratify (e.g., state). Once fit, predicted probabilities,  $\pi_j$ , are computed for each population cell in the data set. As one example, one would predict the likelihood that a low-income ( $\alpha^1$ ) white ( $\alpha^2$ ) male ( $\alpha^3$ ) high school graduate ( $\alpha^4$ ) from Kentucky ( $\alpha^5$ ) enrolls in college. The total number of population cells for which predictions are computed would equal the cross of all categories in  $\alpha^k$ .

In the second step, values of  $\pi_j$  are aggregated to the area of interest,  $\theta_S$ , using

$$\theta_S = \frac{\sum_{j \in S} N_j \pi_j}{\sum_{j \in S} N_j} \quad (2)$$

which reweights each demographic cell's predicted probability using corresponding population counts,  $N_j$ , from the second data set. Population cell counts in the second poststratification

data set are often constructed using Census data. In general, the area of interest,  $\theta_S$ , could be any geographic or institutional (see Ortagus et al., 2021) level for which population cell counts can be computed. The geographic area of interest in our study is the state.

Demographic cells indicated in the first data set and multilevel regression must match those available in the second poststratification matrix.<sup>1</sup> This often limits the unit-level information that can be used to predict the outcome. For example, high school GPA, which would be positively correlated with college enrollment, cannot be included in the multilevel model because high school GPA is not included in Census data.<sup>2</sup> For this reason, second-level covariates are important for improving fit and producing better area-specific poststratified predictions (Park et al., 2004).

With this approach, we use nationally-representative data to first estimate the probability of college attendance by individual characteristics, including race/ethnicity, gender and age. We then use state-level estimates of the numbers of individuals in each of those categories to predict enrollment rates by income at the state level. Before estimating unknown enrollment rates among low-income young persons using real data, we begin with a validation of MRP.

## Validation of MRP

To demonstrate MRP’s validity, we offer two validation exercises. We begin with a simulation exercise in which we follow the MRP procedure on data with known properties. With this simulation, we follow the example of Park et al. (2004), but with modifications that make it more pertinent to our use case. Next, after describing our analysis data and fitting our model, we validate our empirical MRP estimates against known estimates. Specifically, we

---

<sup>1</sup>Jonathan P. Kestellec et al. (2015) propose a method for using a poststratification matrix both comprised of non-Census variables (in their case, partisanship) and that incorporates uncertainty in the poststratification matrix. It remains true, however, that demographic cells,  $j$ , in the primary data set and poststratification matrix must correspond.

<sup>2</sup>This was not the case for Ortagus et al. (2021), who used administrative data and poststratified to colleges in the original sample frame rather than a geographic area. Nevertheless it is true in our study, in which we use publicly available Census data to construct the poststratification matrix.

provide national-level estimates of low-income college enrollment supported by HSLS09 as well as state-level estimates in the ten states for which it is representative—California, Florida, Georgia, Michigan, North Carolina, Ohio, Pennsylvania, Tennessee, Texas, and Washington. Being able to recover correct coverage of the ground truth for otherwise known values will provide one level of support for our application of MRP to the problem of estimating state-specific college enrollment among low-income young persons across the United States.

## Simulation

Simulated data used in this exercise represent a population of students across the 50 states plus the District of Columbia and has been constructed to resemble the HSLS09 data set we use in our primary analyses of low-income youth enrollment in college.<sup>3</sup> The total population size in these simulated data is  $N = 1,000,027$  and is spread out across the states in rough correspondence to the relative population size of each state. Every unit,  $i$ , has three observed characteristics: a continuous value of  $x$  between 10,000 and 50,000, a categorical value of  $v \in \{1, 2, 3, 4, 5, 6\}$ , and a categorical value of  $w \in \{0, 1\}$ . For each observation, we construct a dummy variable,  $D$ , in which  $D = 1$  when  $x < 22,000$  and  $D = 0$  otherwise.  $D = 1$  for approximately 23% of the population. We generate a binary outcome measure,  $y$ , using a random draw from a Bernoulli distribution in which the latent probability is a linear combination of various individual- and state-level characteristics plus noise. In our analyses, we are interested in obtaining an unbiased estimate of  $P(y = 1)$ , which we call  $\Theta$ , across the states and by values of  $D$ .

When creating the simulated population the following rules apply:

1. Observations for whom  $D = 1$  are less likely to have  $y = 1$  than those for whom  $D = 0$ , meaning that  $\Theta_{D=1} < \Theta_{D=0}$ .
2. Values of  $\Theta$  vary across observed characteristics,  $w$  and  $v$ , both through main effects and interaction with  $D$ .

---

<sup>3</sup>Data are simulated using the R package, `fabricatr` (Blair et al., 2021).

3. Values of  $\Theta_{state}$  also vary across states as a function of three state-level covariates,  $Z$ , region, and random noise.
4. The effect of  $D$  on  $\Theta_{state}$  varies across states due to different distributions of individual characteristics,  $w$  and  $v$ , across the states as well as interactions between  $D$  and random state-level noise.

### Samples from the simulated population

For our analyses, we draw four samples from the population. Two are samples of  $N = 10,000$  (1% of the population) and two are smaller samples of  $N = 1,000$  (0.1% of the population). Within each sample size, we draw a simple random sample where all population units have an equal probability of selection. For the other two samples, we use weights such that observations with the following characteristics are more likely to be selected: (1)  $D = 1$ ; (2) small  $v$  groups (those representing less than 10% of population); and (3) from the Northeastern region of the (simulated) United States. Samples are drawn from the full population of observations with no sub-sampling within region or state. Throughout the next section, we distinguish between true population values ( $\Theta$ ), observed sample values ( $\theta$ ), and estimates ( $\hat{\theta}$ ) of  $P(y = 1)$ .

### Comparison of samples and estimates to population values

Figure 1 compares population differences in  $\Theta$  by  $D = 0$  (left panel) and  $D = 1$  (right panel) to estimates using each sample. Population values are presented to the left of the black line within each facet. Except for the population values all points represent observed values,  $\theta$ . Though there is some bias in that estimates  $\theta_{D=0}$  are a little high whereas estimates  $\theta_{D=1}$  are low, they are within  $\pm 5$  percentage points (p.p.) of the true value. As with NCES data sets like HSLS, national estimates provide a reasonable approximation of the truth.<sup>4</sup>

State-level population differences in  $\Theta_{state}$  are presented in the top panel of Figure 2.

---

<sup>4</sup>When including standard errors and inverse sampling weights for the weighted sample, the true population value is contained in the 95% confidence intervals in 10/12 estimates (using weighted samples both with and without weighting adjustment).

This panel shows that  $\Theta_{state|D=1} < \Theta_{state|D=0}$  across all states, with variation in both values of  $\Theta_{state}$  across the states. Observed values of  $\theta_{state}$  across the four samples are presented in the four bottom panels of Figure 2. As can be seen in the figure, state-level values of  $\theta_{state}$  are generally highly biased. In some states,  $\theta_{state} = 100$ , meaning that values of  $D$  are perfectly collinear with outcomes  $y$ . In other states, only one level of  $D$  is represented, that is, some conditions are unobserved in that state. Because none of the four samples were stratified by state, it is only within some of the larger states that an estimate  $\hat{\theta}_{state}$  (with accompanying estimates of uncertainty, not presented) somewhat approximates the true value,  $\Theta_{state}$ . Using inverse weights with the weighted samples does not improve state-level estimates, which is expected since the weights apply to a national-level sampling procedure. Figure 2 demonstrates the problem with using national longitudinal surveys to estimate state-level characteristics: given that the sampling procedure is not designed to provide state-level estimates, those estimates will be highly biased in many states.

Figure 3 directly compares observed state-level values of  $\theta_{state}$  and the true population values,  $\Theta_{state}$ . Each facet represents a different sample. On the x-axis is the sample value of  $\theta_{state}$  and on the y-axis is the population (true) value of  $\Theta_{state}$ . States are plotted by their abbreviation, with red values representing when  $D = 1$  and blue green values when  $D = 0$ . The 45-degree line in each facet represents equality between sample estimate and the true value. Across the four samples, values of  $\theta$  are on average lower when  $D = 1$ , shown by their comparatively lower position on the 45-degree line. In the larger 1% samples, estimates of  $\theta$  when  $D = 0$  are closer to the line, likely due to their representing about 77% of the simulated population. In all state-specific samples, values of  $\theta$  among  $D = 1$  observations are more biased away from the population truth, particularly for some of the smaller states. In some cases, there are estimates of 0 or 100, suggesting very few observed values in the state. Some states do not have estimates in a particular sample. For example, no observations from Alaska were sampled in the 0.1% weighted sample so an estimate for the state does not show up at all in bottom right facet. Across states that were sampled, bias is even greater among

both  $D$  conditions in the smaller 0.1% samples, with even more extreme values of  $\theta_{state}$ .

Though these data are simulated, the samples from the population behave as we might expect national surveys designed like HSL09 to behave. While we can recover reasonable national-level estimates from sample data sets, small area estimates at the state level are highly biased and cannot be recovered with supplied survey weights. The simulation study at this point reflects the reality of the situation for analysts. Using standard techniques, we cannot obtain reliable estimates of the parameter of interest—the probability that a young person from a low-income family will enroll in college. We next turn to our proposed solution to see how well it can provide estimates with credible intervals that cover the known simulated population parameters.

### Using MRP on simulated data

To recover state-level estimates, we fit the following multilevel regression to each sample type:

$$P(y_i = 1) = \text{logit}^{-1}(\beta_0 + \beta_1 * D_i + \alpha_{s[i]}^{state} + \alpha_{s[i]}^{state.D} + Z_i\gamma). \quad (3)$$

In addition to a grand mean parameter,  $\beta_0$ , each state is given a random intercept,  $\alpha^{state}$ , and the key covariate of interest,  $D$ , is estimated as both a main effect,  $\beta_1$ , and in interaction with each state,  $\alpha^{state.D}$ . We also include state-level covariates,  $Z$ , which we consider to be known from the population data set. Even though we have further observation level covariates,  $v$  and  $w$ , that we know are part of the data generating process, we fit a simplified model that only separates observations within each state into two demographic categories:  $D = 0$  and  $D = 1$ .

We use a Bayesian approach for estimation and inference. While a full explanation of Bayesian statistics is beyond the scope of this paper, the primary difference between Bayesian and frequentist inferential approaches is that a Bayesian approach treats the unknown parameter as a random variable while frequentist statistics treat the unknown parameter as

a fixed value (Gelman et al., 2013). Bayesian approaches estimate the distribution of the unknown parameter,  $\theta$ , using

$$P(\theta|X) \propto P(X|\theta) \times P(\theta),$$

where the posterior distribution,  $P(\theta|X)$ , is proportional to the likelihood of the data given the parameter,  $P(X|\theta)$ , times a prior,  $P(\theta)$ , which represents the initial state of the analyst’s beliefs about the distribution of the population parameter. This posterior distribution represents an updated estimate of the likely distribution of the unknown parameter after having taken into account the likelihood of the data. Our discussion will focus on summary measures—mostly quantiles—of the posterior distribution of various parameters. The summary measures in our work will discuss the probability that the population parameter is in a certain range—the credible interval. In contrast, frequentist statistics focus on the probability of observing values in an infinite number of repeated samples under assumptions about the sampling distribution (null hypothesis), which is generally not an estimate of interest in most policy applications (Gelman et al., 1995).

After fitting equation 3 and producing cell-specific predicted probabilities,  $\hat{\pi}_j$ , we poststratify our estimates using population counts from the original population data set.<sup>5</sup> These counts come from collapsing the population data set ( $N = 1,000,027$ ) by summing matching demographic cells ( $D \in \{0, 1\}$ ) within each state. This is the only way that we use the population data set in our MRP procedure for the simulation. Figure 4 compares true values of  $\Theta_{state}$  to the median of our poststratified estimates,  $\hat{\theta}_{state}$ . As with Figure 3, each facet represents a unique sample. Across all samples, the poststratified estimates are much closer to the true population values. This is most apparent in the  $D = 1$  values, which are shrunk toward the 45 degree line (signaling less bias). The  $D = 0$  values are even more tightly bunched on the 45 degree line and no state-specific estimate sits at an extreme of 0 or

---

<sup>5</sup>We use the Stan probabilistic programming language (Stan Development Team, 2021) in conjunction with the R programming language (R Core Team, 2021) to fit all Bayesian models in our paper. For equation 3, we assign all regression parameters weakly informative normal priors:  $\alpha, \beta, \gamma \sim N(0, \sigma)$ ;  $\sigma \sim N_+(0, 1)$ .

100. Another important result is that we are able to provide estimates for states which were not included in the weighted sample, such as Alaska.

To make the comparison clearer, Figure 5 plots the population truth ( $\Theta_{state}$ , red circle), observed sample mean ( $\theta_{state}$ , green triangle), and poststratified estimate ( $\hat{\theta}_{state}$ , blue square as median of the poststratified posterior distribution with 95% credible intervals) across both 0.1% samples across all states. Except in a few cases, the 95% CIs of the poststratified estimates contain the true value across the states, even when the state sample average is either very different from the ground truth or unobserved because the state was not included in the sampling—again, see Alaska in both weighted sample panels, which does not contain an observed sample mean estimate. This is true in both the simple random sample as well as the weighted random sample, the latter of which oversampled units based on values of  $D$  as well as characteristics we did not include in the model.

We conclude from this simulation exercise that MRP is a viable solution for recovering state-level estimates of college enrollment among low-income young persons. With data generated to replicate conditions very similar to the ones that apply in our actual data, these simulation results show that MRP generates posterior distributions that overwhelmingly include the true value of the population parameter even when (1) we use a simplified model of the data generating process and (2) the data for a given state is minimal or nonexistent. These results do not mechanically follow from the construction of the sample: while each sampling procedure produces naive estimates that are quite different from the actual values, MRP estimates, which rely on the types of data available to us in our empirical approach, reliably include known true values. In the next section, we describe the model and data we use to estimate state-level college enrollment among low-income young persons.



# Estimating college enrollment among low-income youth

To estimate the proportion of low-income youth who attend college, we use the same MRP procedure outlined in the simulation. In the first step, we use student-level data from the most recent NCES survey, HSLS09, to estimate the likelihood of enrolling in college. Our multilevel logistic regression takes the form,

$$P(y_i = 1) = \text{logit}^{-1}(\beta_0 + \beta^{\text{lowinc}} * \text{lowinc}_i + \beta^{\text{female}} * \text{female}_i + \alpha_{re[i]}^{\text{race/ethnicity}} + \alpha_{r[i]}^{\text{region}} + \alpha_{s[i]}^{\text{state}} + \alpha_{sl[i]}^{\text{state.lowinc}} + Z_i\gamma), \quad (4)$$

in which we use random effects parameters for student characteristics that include gender, race/ethnicity, and income status as well as state and regional indicators. We fully interact income status, our binary covariate of interest, with state indicators to account for potential differences in low-income college enrollment across the states. Second-level covariates in  $Z$  include state-level percentages of adults with a Bachelor’s degree or higher, the proportion of college students who attend public two-year institutions, and the average tuition at four-year public institutions for the years aligning with the year of on-time college enrollment for students surveyed in HSLS09. We give all regression parameters weakly informative normal priors:  $\alpha, \beta, \gamma \sim N(0, \sigma); \sigma \sim N_+(0, 1)$ .

Our second level estimates come from the literature on predictors of low-income college enrollment. Based on national surveys, analysts have shown that students who live in states with higher educational attainment are more likely to attend postsecondary education, as are students who live in states with more community colleges (Robert Bozick, 2009; Doyle & Skinner, 2016). Similarly, the “sticker price” of four-year colleges has been shown to be a reliable predictor of on-time college enrollment, as colleges with lower tuition are perceived to be more affordable, regardless of actual net price (Hemelt & Marcotte, 2011).

After fitting equation 4, we create predicted probabilities of enrollment for each demographic cell in our data set. Once we have calculated enrollment probabilities for each

demographic cell,  $N_J = 1,224$  unique cells, we compute the weighted probability of enrollment for each state,  $\theta_S$ , with equation 2 using population counts as weights. Within each state, we estimate two values of  $\theta_S$ ,  $\theta_S^{hi}$  and  $\theta_S^{lo}$ , representing the likelihood of college enrollment among middle/high-income and low-income youth, respectively.

## Data

Data for our study come from two primary sources. Unit-level data come from a student longitudinal survey conducted by the National Center for Education Statistics, the High School Longitudinal Study of 2009 (HSL09, Duprey et al., 2020). As with prior NCES longitudinal surveys, HSL09 tracked a nationally-representative sample of high schools ninth-graders starting in 2009 as they moved through high school and either enrolled in postsecondary institutions or entered the workforce. Unlike prior NCES surveys, HSL09 is representative in 10 states—California, Florida, Georgia, Michigan, North Carolina, Ohio, Pennsylvania, Tennessee, Texas, and Washington—though it cannot be used to construct state-level estimates across the entire United States.

We use information from the first and fourth wave of the survey. Data on the state of residence for all students in the base year of the survey come from the restricted-use data files. We use the variable `x4fb16enrstat` to determine those students who enrolled within one-year of high school graduation. We use the base-year family income variable, `x1famincome`, which discretizes incomes into 13 bins of non-equal size, to construct a binary variable of income status. We assign low-income status for all students with family incomes below \$35,000, which equals between 150% and 185% of the federal poverty line for a family of four in 2009. Our primary unit-level data set has  $N = 13,020$  observations, 26% of whom are coded as low-income.

To construct our poststratification matrix of demographic cell sizes, we use population data from the American Community Survey. We select data from the period that

best aligns with the year the HSLS09 sample, which are single-year population estimates for 14 to 15 year olds in 2009. Second-level state characteristics are taken from various NCES Digest of Education statistics tables covering 2013, the year of on-time college enrollment for students in the HSLS09 sample.

## Results

We begin with a second validation of MRP, this time producing estimates with real data that we are supported by HSLS09 survey weights. The top panel of Figure 6 compares national estimates of college enrollment by income. In both panels, the left (red) estimate is from a simple weighted mean, using the most appropriate longitudinal weight supplied by HSLS09, `w4w1stu`, in which the center point is the mean and the vertical lines the 95% confidence intervals. The right estimate (blue green) are our MRP estimates, with the center dots representing the median value ( $\hat{\theta}_{q50}$ ) of the full posterior distribution and the vertical lines the 95% credible intervals. Though the weighted and MRP estimates in the each facet (and the state-level comparisons in the bottom panel) are not strictly commensurate due to their different statistical frameworks (frequentist vs. Bayesian), they represent the most likely comparison in terms of common estimation practice.

Using survey weights, the direct national estimate of low-income student college enrollment is 49.2% (46.1,52.2) compared to a MRP estimate of 50.5% [48.7,52.2]. The direct estimate of college enrollment among middle-high income students is 74.1% (72.8,75.4) compared to a MRP estimate of 77.7% [76.1,78.6]. In both situations, the MRP estimate is somewhat higher than the weighted estimate, though the weighted mean estimate for low-income student enrollment—the primary estimate of interest—is contained within the credible interval of the MRP estimate. In the bottom panel of Figure 6, we compare estimates across the ten representative states in HSLS09. Among these estimates, 16 of 20 weighted survey mean estimates are contained within the 95% credible intervals of their corresponding

MRP estimates. All but two weighted survey estimates are within 5 p.p. of the poststratified posterior median value and half are within 3 p.p. This second validation exercise using real data shows that while MRP estimates do not perfectly align with those produced using survey weights, they are commensurate more often than not and when different, only seldomly provide estimates that are meaningfully different.

We now turn to our primary results in Figure 7, which shows the probability of college enrollment across the 50 states and the District of Columbia for both low-income (red) and middle-to-high-income (blue green) high school graduates. The figure is ordered by low-income enrollment, median values ( $\hat{\theta}_{50}^{state}$ ) which range from 34% [21.9,48.3] of recent high school graduates in Utah to 67.8% [56.8,78.6] in Mississippi, a range of 33.8 p.p. with a median of 50.4% [45.7,56] in Michigan. The within-state difference between low- and middle-to-high-income enrollment ranges between 19.9 p.p. [12.2,26.4] in Mississippi and 30.3 p.p. [26.4,36.8] in Arizona with a median of 27.4 p.p. [23,31.2] in California.

While the size of the 95% credible intervals demonstrate some uncertainty in our estimates, which at the extreme ranges 33.2 p.p. [28.8,62] for low-income enrollment in the District of Columbia, we make two notes. First, that the credible intervals for low- and middle-to-high-income estimates within 49 of 50 states and D.C. do not overlap provides strong evidence that the on-time enrollment rates of low-income students are lower than that of their higher income peers across the country. Second, we are able to provide estimates in low population states and, in the case of the District of Columbia, a location entirely unsampled by HSL09.

Finally, we compare our estimates of low-income student enrollment with potentially cognate estimates that are publicly available. We consider three measures: (1) a direct measure of low-income student enrollment, (2) the percentage of students enrolled in college who receive Pell grant funding, and (3) the percentage of students enrolled in college whose families earn less than \$30,000 a year.

The first cognate measure, which comes from IPUMS microcensus data for the

American Community Survey, is the state-level percentages of 18 and 19 year olds in 2013 who have earned a high school credential (diploma or GED) but not a postsecondary credential and who are currently enrolled in college as undergraduates. We further filter this group to those young people with family incomes less than \$37,800 per year, which equals the \$35,000 per year cut-off we use for the HSL09 sample adjusted for inflation.<sup>6</sup> The percentage point difference between these estimates of low-income college enrollment and our own are shown in the top panel of Figure 8. States are ordered from least to most difference. Values above zero, shown by the dashed line, indicate that direct Census estimates are higher than those obtained using MRP. As expected based on the way Census constructs family income, that is, not accounting for financial dependency across households, Census estimates skew higher than MRP estimates, with a median difference of 23.1 p.p. [-2.6,45.2]. In practical terms, this means that using direct estimates from the Census are likely to overstate the participation on low-income young persons in college by a large margin in many states.

Panel B of Figure 8 compares MRP estimates on the  $y$ -axis to the percentage of first-time in college undergraduates who participated in the Pell grant program in the 2013-2014 academic year on the  $x$ -axis. Vertical lines indicate the central 95 credible interval for each of the state-level estimates. The diagonal dashed line represents the point at which the two estimates align. The percentage of students who are Pell eligible is computed using institution-level data from IPEDS, aggregated to the state level using undergraduate full-time equivalent (FTE) counts as weights. Unlike the first cognate measure, which represents the probability of enrolling in college conditional on being low income,  $P(\text{college} \mid \text{low income})$ , the relation of interest and what we estimate using MRP, the Pell estimate measures the probability of being low-income conditional on being enrolled in college,  $P(\text{low income} \mid \text{college})$ . Nevertheless, the percentage of students who either use Pell or are Pell eligible is often used as a proxy for low-income enrollment. The figure shows little or no relationship between our measure and the proportion of students who are Pell eligible.

---

<sup>6</sup>We use an 8% adjustment, which applies to the period between September 2009 and September 2013 according to the Consumer Price Index inflation calculator at [https://www.bls.gov/data/inflation\\_calculator.htm](https://www.bls.gov/data/inflation_calculator.htm).

In panel C of Figure 8, we compare MRP estimates to the percentage of first-time in college undergraduates with family incomes less than \$30,000 in the 2013-2014 academic year, which comes from the student financial aid component of IPEDS. This income category represents the closest approximation to the low-income category we use for our primary MRP estimates. We are limited to reporting only this grouping of low-income students given IPEDS' reporting standards for enrollment by income. As with the Pell measure, this financial aid cohort measure is averaged from the institution level to the state level using undergraduate FTE enrollments as weights. It also represents the inverse probability of the MRP estimates:  $P(\text{low income} \mid \text{college})$  rather than  $P(\text{college} \mid \text{low income})$ . Panel C shows that as with the Pell-eligible measures, the proportion of the financial aid cohort that is low-income does not correlate with our MRP estimates of the state-level proportions of low-income young persons who enroll in higher education.

Together, the comparisons presented in Figure 8 demonstrate that three common cognate measures of low-income youth college enrollment are not strongly aligned with more principled MRP estimates. Using U.S. Census-reported information on enrollment by household income substantially over-reports the proportion of low-income young people who enroll in higher education. As we find no observable relationship between our MRP estimates and either measure provided by IPEDS, we provide evidence that using a measure of  $P(\text{low income} \mid \text{college})$  is a poor substitute for the actual measure of interest,  $P(\text{college} \mid \text{low income})$ . This latter finding should not be unexpected since measures based on the probability of being low-income given that a person is enrolled in college have no necessary reason to be correlated with the more policy relevant measure of the probability of enrolling in college given that a person is low-income, despite an understandable desire to substitute the former for the latter as a matter of convenience. Based on these results, we conclude that while existing measures of low-income youth enrollment in college do not capture the policy-relevant outcome, our MRP-based estimates do.

## Conclusion

The significance of this study is two-fold. First, we offer state-level estimates of low-income college enrollment for a recent cohort of young persons. These estimates provide better evidence of the efficacy of policies meant to support low-income youth enrollment in college at the state level, where many aid policies are set and funded. Second, we generate our estimates using a statistical procedure, MRP, that we believe can be usefully applied to other education policy questions for which data at the proper level of inference is otherwise limited (e.g., Ortagus et al., 2021).

The goal of federal grant aid in the form of the Pell Grant program has been to increase attendance rates in postsecondary education (T. Kane, 1999). The primary focus of efforts to increase participation in the Pell Grant program has been young people (Deming & Dynarski, 2009). Similarly, the primary goal of many state policymakers has been to increase attendance rates in college among low-income young people through lower tuition and financial aid programs (Zumeta et al., 2012). Based on our findings here, considerable work remains to accomplish these goals.

Our results highlight several salient facts regarding attendance in higher education by income across states. First, fewer low-income young people attend college in every state than do their middle-income or high-income peers. The median gap between low-income postsecondary attendance and high-income attendance credibly ranges 33.8 p.p. across states. This means that despite 50 years of sustained effort to close the college attendance gap by income, low-income young persons still face considerable barriers to attending college, even in the states with the lowest net prices. Second, the probability of attending college for low-income young people varies dramatically across states. In the lowest performing states, about 34 percent of low-income young people attend higher education, while in the highest performing states, about 67.8 percent of low-income young people attend higher education. This means that a low-income young person’s chance for attending college depends crucially on their state of residence.

In future research, these estimated postsecondary attendance rates could be related to state policy to help find patterns regarding the effectiveness of various policies. While it is well-established that lowering the price of higher education increases college attendance rates, particularly among young people, more research can be done to establish the effectiveness of other policies to encourage more young people to enroll. Technical issues aside, the failure of current data sources to provide reliable estimates on low-income postsecondary enrollment by state has substantial implications for education policy and cannot be overstated. Even with a combined \$40 billion in spending from public sources on grant aid for low-income students, we know next to nothing about college attendance patterns by income at the state level (Baum & Payea, 2013). It is our hope that the methods proposed here will no longer be necessary in the years to come, as this data will be collected and reported.



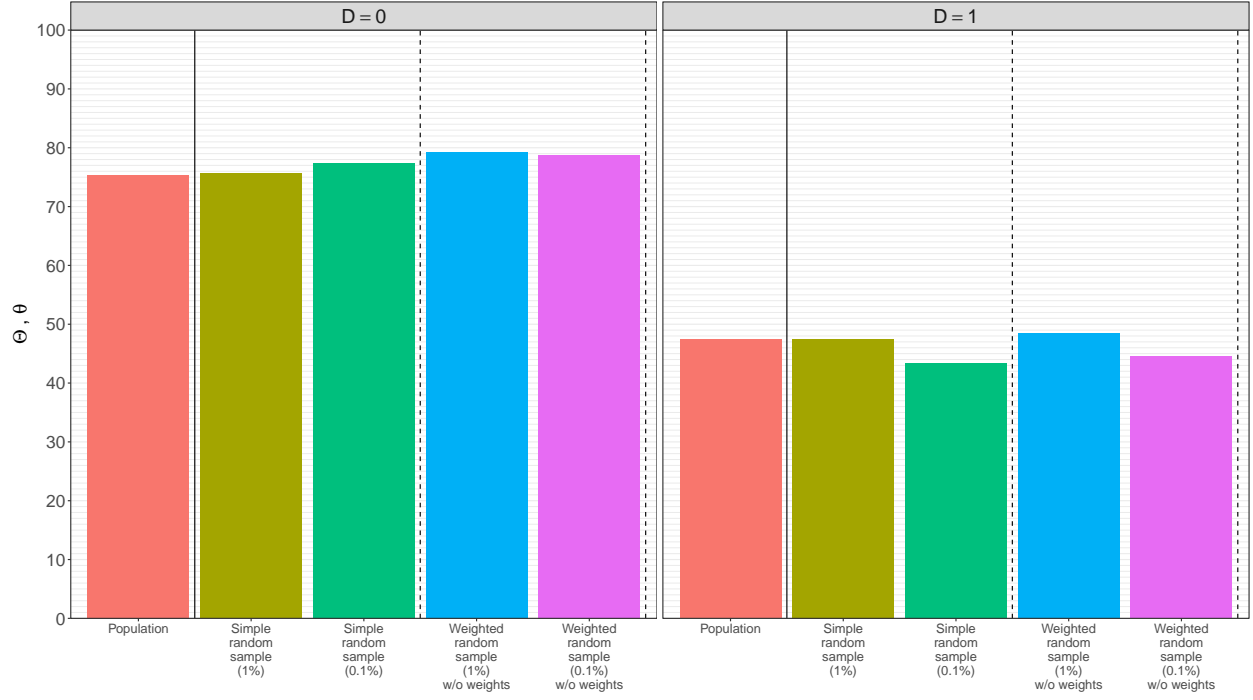
## References

- Abraham, K. G., & Clark, M. A. (2006). Financial aid and students' college decisions: Evidence from the district of columbia tuition assistance grant program. *Journal of Human Resources*, 41(3), 578–610. <https://doi.org/10.3368/jhr.XLI.3.578>
- Adelman, C., Daniel, B., & Berkovits, I. (2003). *Postsecondary attainment, attendance, curriculum, and performance: Selected results from the NELS:88/2000 postsecondary education transcript study (PETS), 2000* (NCES No. 2003-394). National Center for Education Statistics. <http://nces.ed.gov/pubs2003/2003394.pdf>
- Bartik, T. J., Hershbein, B., & Lachowska, M. (2021). The effects of the Kalamazoo promise scholarship on college enrollment and completion. *Journal of Human Resources*, 56(1), 269–310. <https://doi.org/10.3368/jhr.56.1.0416-7824R4>
- Baum, S., & Payea, K. (2013). *Trends in student aid 2013* [Trends in Higher Education Series]. The College Board.
- Bell, A. D., Rowan-Kenyon, H. T., & Perna, L. W. (2009). College knowledge of 9th and 11th grade students: Variation by school and state context. *The Journal of Higher Education*, 80(6), 663–685. <https://doi.org/10.1353/jhe.0.0074>
- Berkner, L., & Chavez, L. (1997). *Access to postsecondary education for the 1992 high school graduate* (No. 98-105). National Center for Education Statistics. <http://nces.ed.gov/pubs/98/98105.pdf>
- Blair, G., Cooper, J., Coppock, A., Humphreys, M., Rudkin, A., & Fultz, N. (2021). *Fabricatr: Imagine your data before you collect it*. <https://CRAN.R-project.org/package=fabricatr>
- Bozick, Robert. (2009). Job opportunities, economic resources, and the postsecondary destinations of American youth. *Demography*, 46(3), 493–512.
- Bozick, R., & Lauff, E. (2007). *Education longitudinal study of 2002 (ELS:2002): A first look at the initial postsecondary experiences of the sophomore class of 2002* (NCES No. 2008-308). National Center for Education Statistics. <http://nces.ed.gov/pubs2008/2008308.pdf>
- Callan, P., & Finney, J. (1997). *Public and private financing of higher education*. ACE/Oryx Press.
- Carneiro, P., & Heckman, J. J. (2002). The evidence on credit constraints in post-secondary schooling. *The Economic Journal*, 112(482), 705–734. <https://doi.org/10.1111/1468-0297.00075>
- Chen, X., Lauff, E., Arbeit, C. A., Henke, R., Skomsvold, P., & Hufford, J. (2017). *Early millennials: The sophomore class of 2002 a decade later* (NCES No. 2017-437). U.S. Department of Education. National Center for Education Statistics.
- Corak, M. (2013). Income inequality, equality of opportunity, and intergenerational mobility. *The Journal of Economic Perspectives*, 27(3), 79–102.
- Deming, D., & Dynarski, S. (2009). *Into college, out of poverty? Policies to increase the postsecondary attainment of the poor* (NBER Working Paper No. 15387). National Bureau of Economic Research. <http://www.nber.org/papers/w15387>
- Denning, J. T. (2019). Born under a lucky star: Financial aid, college completion, labor supply, and credit constraints. *Journal of Human Resources*, 54(3), 760–784.
- Dowd, A. C. (2004). Income and financial aid effects on persistence and degree attainment in public colleges. *Education Policy Analysis Archives*, 12(21).

- Downes, M., Gurrin, L. C., English, D. R., Pirkis, J., Currier, D., Spittal, M. J., & Carlin, J. B. (2018). Multilevel regression and poststratification: A modelling approach to estimating population quantities from highly selected survey samples. *American Journal of Epidemiology*, 187(8), 1780–1790.
- Doyle, W. R., & Skinner, B. T. (2016). Estimating the education-earnings equation using geographic variation. *Economics of Education Review*, 53, 254–267.
- Duncan, G. J., & Murnane, R. J. (2011). *Whither opportunity?: Rising inequality, schools, and children's life chances*. Russell Sage Foundation.
- Duprey, P., M. D., Wilson, D. H., Jewell, D. M., Brown, D. S., Caves, L. R., Kinney, S. K., Mattox, T. L., Smith Ritchie, N., Rogers, J. E., Spagnardi, C. M., & Wescott, J. D. (2020). *High school longitudinal study of 2009 (HSL:09) postsecondary education transcript study and student financial aid records collection data file documentation* (NCES No. 2020-004). U.S. Department of Education. Washington, DC: National Center for Education Statistics, Institute of Education Sciences. <https://nces.ed.gov/pubsearch/pubinfo.asp?pubid=2020004>
- Dynarski, S. (2002). The behavioral and distributional implications of aid for college. *The American Economic Review*, 92(2), 279–285. <http://www.jstor.org/stable/3083417>
- Dynarski, S., Libassi, C., Micheltore, K., & Owen, S. (2021). Closing the gap: The effect of reducing complexity and uncertainty in college pricing on the choices of low-income students. *American Economic Review*, 111(6), 1721–1756. <https://doi.org/10.1257/aer.20200451>
- Dynarski, S., & Scott-Clayton, J. (2013). *Financial aid policy: Lessons from research* (Working Paper No. 18710). National Bureau of Economic Research. <https://doi.org/10.3386/w18710>
- Eke, P. I., Zhang, X., Lu, H., Wei, L., Thornton-Evans, G., Greenlund, K. J., Holt, J. B., & Croft, J. B. (2016). Predicting periodontitis at state and local levels in the United States. *Journal of Dental Research*, 95(5), 515–522.
- Gao, Y., Kennedy, L., Simpson, D., & Gelman, A. (2019). Improving multilevel regression and poststratification with structured priors. *arXiv Preprint arXiv:1908.06716*.
- Gelman, A., Carlin, J. B., Stern, H. S., Dunson, D. B., Vehtari, A., & Rubin, D. B. (2013). *Bayesian data analysis* (3rd ed.). Chapman; Hall/CRC.
- Gelman, A., Carlin, J. B., Stern, H. S., & Rubin, D. B. (1995). *Bayesian data analysis*. Chapman; Hall/CRC.
- Gelman, A., Lee, D., & Ghitza, Y. (2010). Public opinion on health care reform. *The Forum*, 8(1).
- Gelman, A., & Little, T. C. (1997). Poststratification into many categories using hierarchical logistic regression. *Survey Methodology*, 23, 127–135.
- Gurantz, O. (2022). Impacts of state aid for nontraditional students on educational and labor market outcomes. *Journal of Human Resources*, 57(1), 1–32.
- Heller, D. E. (1997). Student price response in higher education: An update to Leslie and Brinkman. *Journal of Higher Education*, 68(6), 624–659.
- Hemelt, S. W., & Marcotte, D. E. (2011). The impact of tuition increases on enrollment at public colleges and universities. *Educational Evaluation and Policy Analysis*, 33(4), 435–457.
- Howe, P. D., Mildemberger, M., Marlon, J. R., & Leiserowitz, A. (2015). Geographic variation

- in opinions on climate change at state and local scales in the USA. *Nature Climate Change*, 5, 596–603.
- Ingels, S. J., Pratt, D. J., Alexander, C. P., Jewell, D. M., Lauff, E., Mattox, T. L., & Wilson, D. (2014). *Education longitudinal study of 2002 third follow-up data file documentation* (NCES No. 2014-364). National Center for Education Statistics.
- Kane, T. (1999). Has financial aid policy succeeded in ensuring access to college? In *The price of admission*. Brookings.
- Kane, T. J. (2006). Public intervention in post-secondary education. *Handbook of the Economics of Education*, 2, 1369–1401.
- Kastellec, Jonathan P., Lax, J. R., Malecki, M., & Phillips, J. H. (2015). Polarizing the electoral connection: Partisan representation in Supreme Court confirmation politics. *The Journal of Politics*, 77(3), 787–804. <https://doi.org/0.1086/681261>
- Kastellec, Jonathan P., Lax, J. R., & Phillips, J. H. (2019). *Estimating state public opinion with multi-level regression and poststratification using R*.
- Kennedy, L., & Gelman, A. (2019). *Know your population and know your model: Using model-based regression and poststratification to generalize findings beyond the observed sample* (arXiv Preprint No. 1906.11323v1).
- Kofoed, M. S. (2017). To apply or not to apply: FAFSA completion and financial aid gaps. *Research in Higher Education*, 58(1), 1–39. <https://doi.org/10.1007/s11162-016-9418-y>
- Lax, J. R., & Phillips, J. H. (2009). How should we estimate public opinion in the states? *American Journal of Political Science*, 53(1), 107–121.
- Lei, R., Gelman, A., & Ghitza, Y. (2017). The 2008 election: A preregistered replication analysis. *Statistics and Public Policy*, 4(1), 1–8.
- Leslie, L. L., & Brinkman, P. T. (1987). Student price response in higher education: The student demand studies. *Journal of Higher Education*, 58(2), 181–204.
- Lipps, J., & Schraff, D. (2019). Estimating subnational preferences across the European Union. *Political Science Research and Methods*, 1–9.
- Little, R. J. A. (1993). Post-stratification: A modeler’s perspective. *Journal of the American Statistical Association*, 88(423), 1001–1012. <https://doi.org/10.1080/01621459.1993.10476368>
- Lovenheim, M. F., & Reynolds, C. L. (2013). The effect of housing wealth on college choice: Evidence from the housing boom. *Journal of Human Resources*, 48(1), 1–35. <http://jhr.uwpress.org/content/48/1/1>
- McDonough, P. M. (1997). *Choosing colleges: How social class and schools structure opportunity*. SUNY Press.
- Ortagus, J. C., Skinner, B. T., & Tanner, M. J. (2021). Investigating why academically successful community college students leave college without a degree. *AERA Open*, 7(1), 1–17. <https://doi.org/10.1177/23328584211065724>
- Pacheco, J. (2011). Using national surveys to measure dynamic U.S. State public opinion. *State Politics & Policy Quarterly*, 11(4), 415–439.
- Park, D. K., Gelman, A., & Bafumi, J. (2004). Bayesian multilevel estimation with poststratification: State-level estimates from national polls. *Political Analysis*, 12(4), 375–385.
- Perna, L. W. (2006). Studying college access and choice: A proposed conceptual model. In J. C. Smart (Ed.), *Higher education: Handbook of theory and research* (Vol. 21, pp. 99–157). Springer Netherlands. [https://doi.org/10.1007/1-4020-4512-3\\_3](https://doi.org/10.1007/1-4020-4512-3_3)

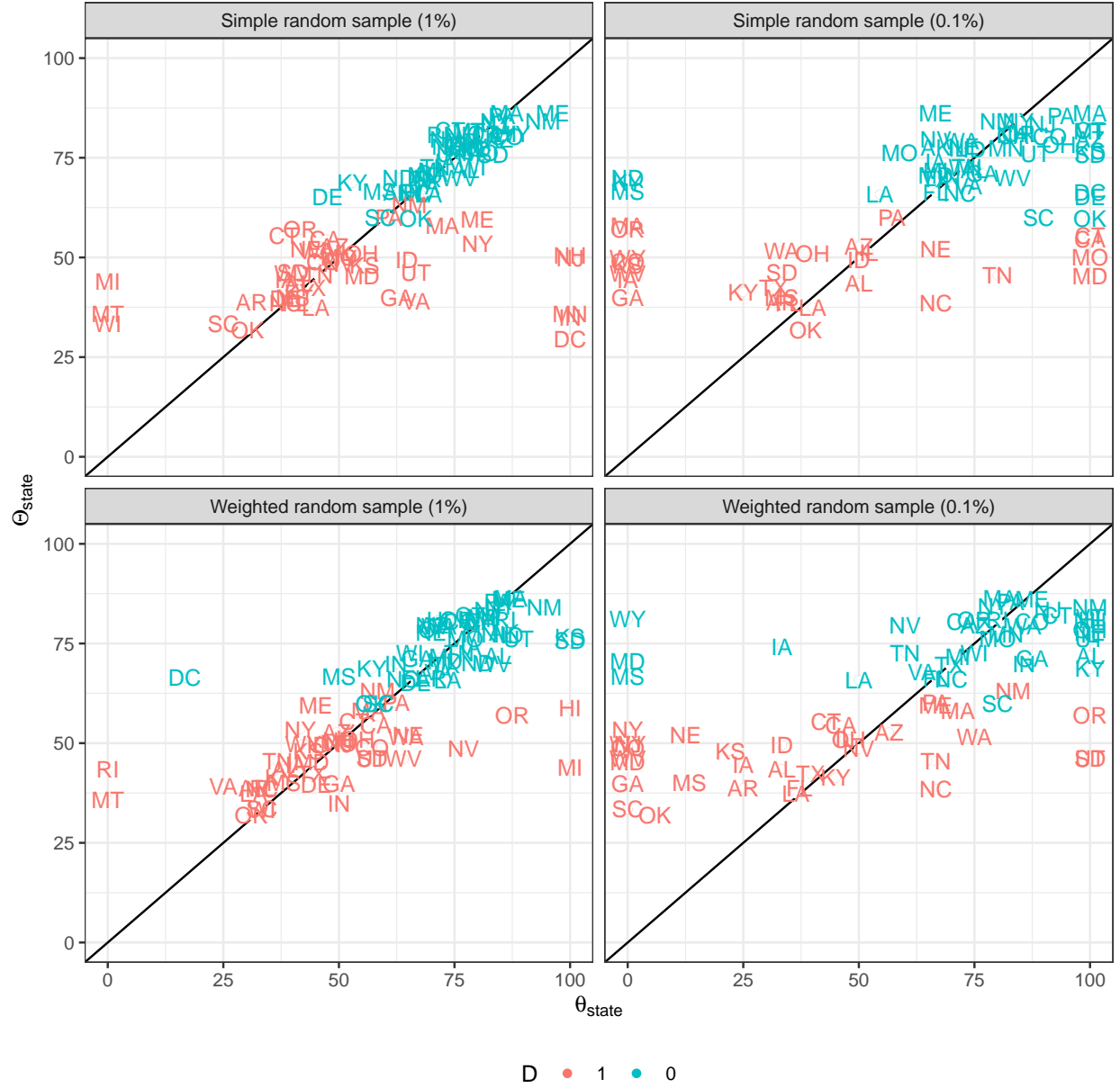
- Perna, L. W., & Jones, A. (2013). *The state of college access and completion: Improving college success for students from underrepresented groups*. Routledge.
- R Core Team. (2021). *R: A language and environment for statistical computing*. R Foundation for Statistical Computing. <https://www.R-project.org/>
- Richardson, R., Bracco, K. R., Callan, P., & Finney, J. (1999). *Designing state higher education systems for a new century*. ACE/Oryx Press.
- Scott-Clayton, J., & Schudde, L. (2020). The consequences of performance standards in need-based aid: Evidence from community colleges. *Journal of Human Resources*, 55(4), N.PAG.
- Stan Development Team. (2021). *Stan modeling language users guide and reference manual, 2.28*. <https://mc-stan.org>
- Tebbs, J., & Turner, S. (2005). Low-income students a caution about using data on pell grant recipients. *Change: The Magazine of Higher Learning*, 37(4), 34–43.
- U.S. Census Bureau. (2009). *Design and methodology American Community Survey [ACS-DM1]*. United States Census Bureau.
- U.S. Census Bureau. (2012a). *American community survey and Puerto Rico community survey 2012 subject definitions*. United States Census Bureau. [https://www2.census.gov/programs-surveys/acs/tech\\_docs/subject\\_definitions/2012\\_ACSSubjectDefinitions.pdf](https://www2.census.gov/programs-surveys/acs/tech_docs/subject_definitions/2012_ACSSubjectDefinitions.pdf)
- U.S. Census Bureau. (2012b). *Current population survey, October 2012 school enrollment and internet use supplement file [CPS-12]*. United States Census Bureau. <http://www.census.gov/prod/techdoc/cps/cpsoct12.pdf>
- U.S. Census Bureau. (2013). *Instructions for applying statistical testing to the 2010-2012 ACS 3-year data and the 2008-2012 ACS 5-year data*. United States Census Bureau.
- Wang, W., Rothschild, D., Goel, S., & Gelman, A. (2015). Forecasting elections with non-representative polls. *International Journal of Forecasting*, 31(3), 980–991.
- Warshaw, C., & Rodden, J. (2012). How should we measure district-level public opinion on individual issues? *The Journal of Politics*, 74(1), 203–219.
- Zhang, X., Holt, J. B., Lu, H., Wheaton, A. G., Ford, E. S., Greenlund, K. J., & Croft, J. B. (2014). Multilevel regression and poststratification for small-area estimation of population health outcomes: A case study of chronic obstructive pulmonary disease prevalence using the behavioral risk factor surveillance system. *American Journal of Epidemiology*, 179(8), 1025–1033.
- Zumeta, W., Breneman, D. W., Callan, P. M., & Finney, J. E. (2012). *Financing American higher education in the era of globalization*. Harvard Education Press.



**Figure 1:** From simulated data, national values of  $\Theta$  compared to observed values of  $\theta$  by  $D \in \{0, 1\}$ . Samples of 1% and 0.1% have  $N = 10,000$  and  $N = 1,000$  observations, respectively. Weighted random samples oversample some subpopulations and estimates are presented without sampling weights. All estimates were computed as simple mean statistics.



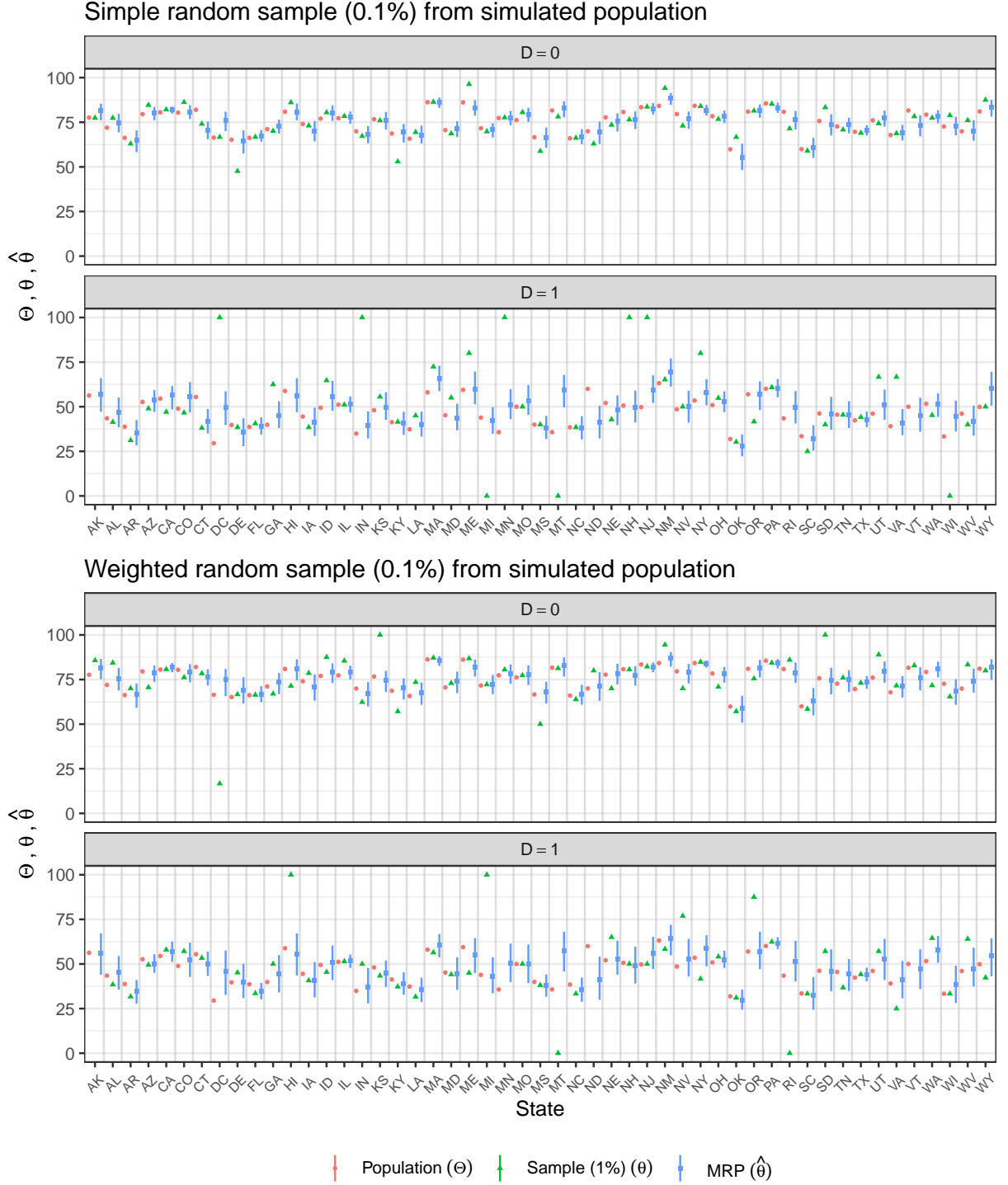
**Figure 2:** From simulated data, state-specific values of  $\Theta_{state}$  (Population) compared to observed values of  $\theta_{state}$  (Samples) by  $D \in \{0, 1\}$ . Samples of 1% and 0.1% have  $N = 10,000$  and  $N = 1,000$  observations, respectively. Weighted random samples are presented without weighting adjustment.



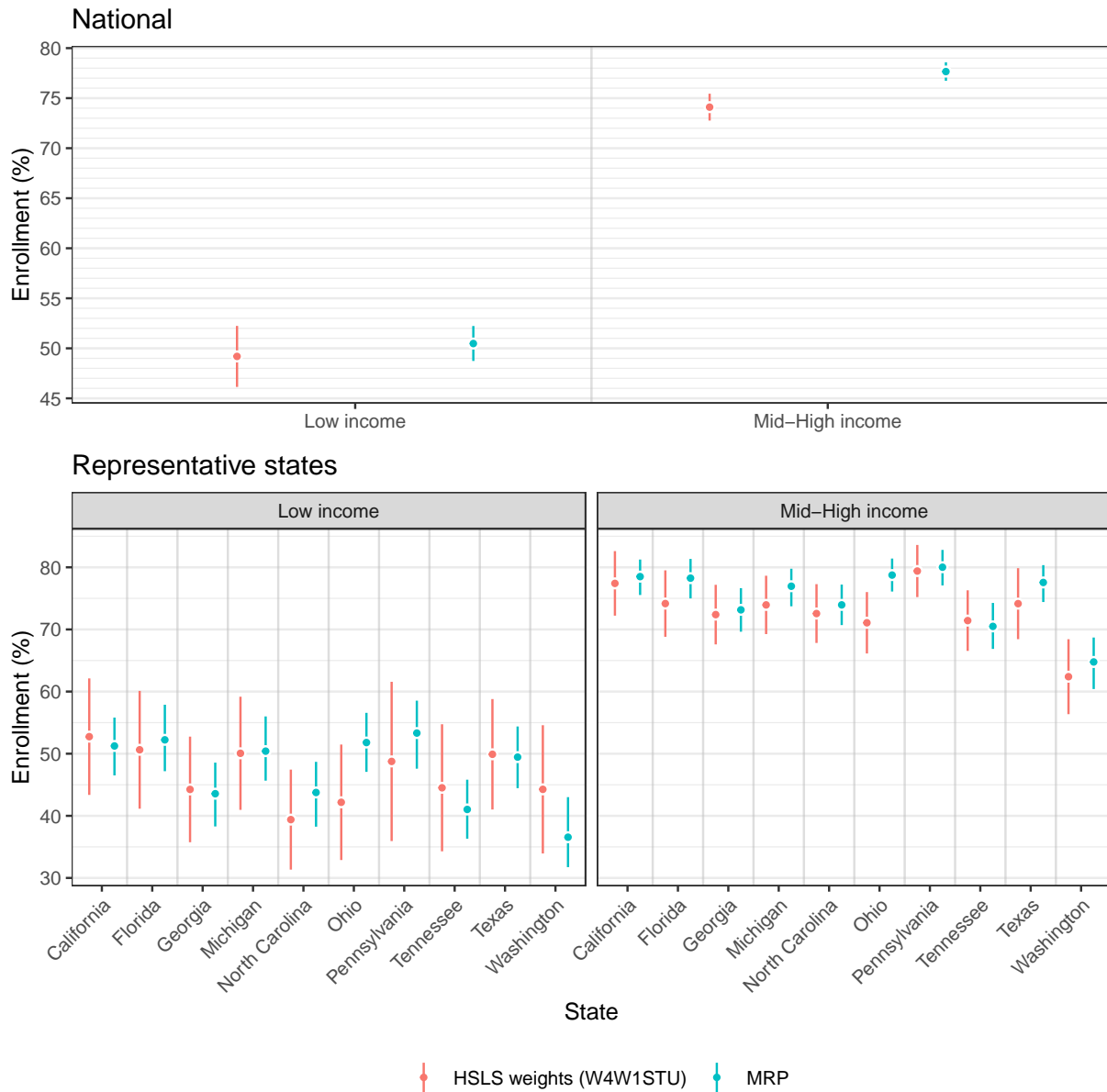
**Figure 3:** From simulated data, a comparison of observed values of  $\theta_{state}$  across samples to true values of  $\Theta_{state}$ . States are plotted twice in each facet, once for each value of  $D$ , except in the case when the sample does not have observations for a particular state by  $D$  combination. The 45 degree line on each plot represents equality between unobserved true ( $\Theta_{state}$ ) and observed ( $\theta_{state}$ ) values. Weighted samples are presented without weighted adjustments.



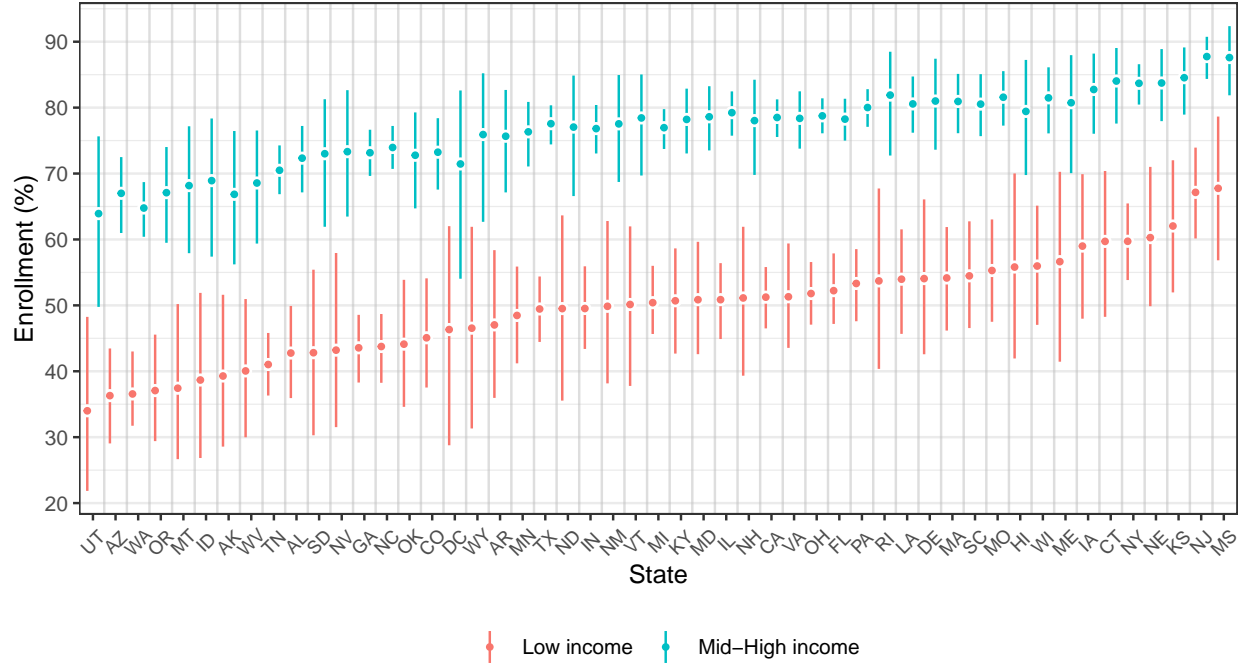




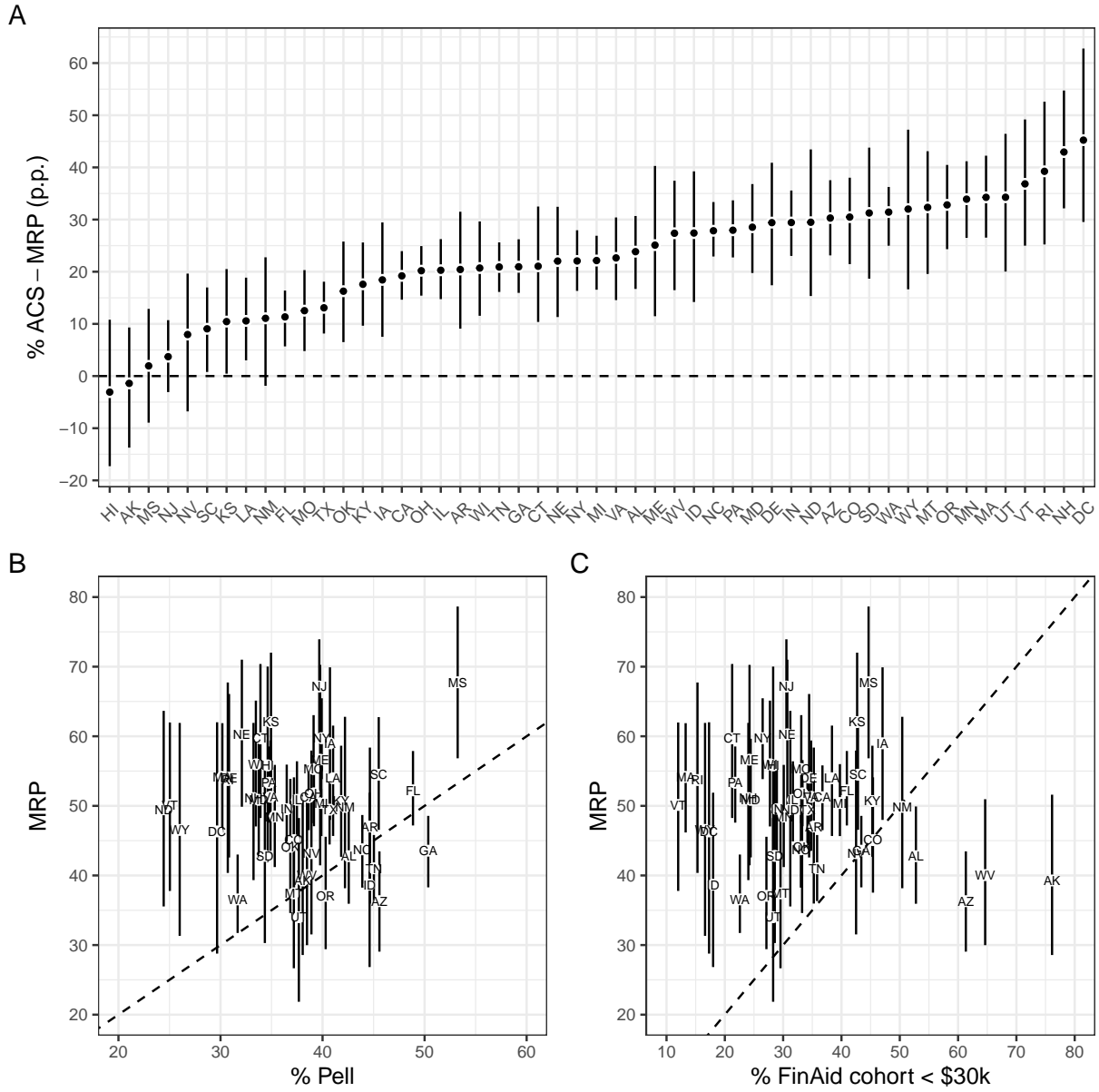
**Figure 5:** Comparison of true ( $\Theta_{state}$ ), observed ( $\theta_{state}$ ), and poststratified estimates ( $\hat{\theta}_{state}$ ) for the 0.1% simple random sample (top panels) and 0.1% weighted random sample (bottom panels). Lines on the poststratified estimates represent 95% credible intervals. The top and bottom facets for each sample group represent values when  $D$  is equal to 0 and 1, respectively.



**Figure 6:** Validation of MRP estimates. In the top panel, MRP estimates of national enrollment rates by income group are compared to estimates computed using HSLs09 survey weights (W4W1STU). HSLs09 is also representative of 10 states. In the bottom two panels, state-level MRP estimates in these states are compared to those produced using the same HSLs09 survey weights.



**Figure 7:** State-level MRP estimates of college enrollment by income group. Each line represents the 95% credible interval estimate of  $\theta_{attend}$ , with the center dot representing the median value ( $\theta_{q50}$ ). States are ordered by low-income student enrollment rates.



**Figure 8:** Comparison of MRP estimates to proxy measures of low-income student enrollment. Top panel estimates come from the American Community Survey (ACS) and both institution-level measures come from the Integrated Postsecondary Education Data System (IPEDS); all are aggregated to the state level. The top panel subtracts ACS measures of low-income student enrollment in 2013 from the MRP estimates of the same measure. The left bottom panel compares average student population that receives Pell grant funding and the right bottom panel compares the percentage of the financial aid cohort in the < \$30,000 category out of the full time, first time in college cohort to MRP estimates, respectively.