

ENDOGENOUS STRATIFICATION IN RANDOMIZED EXPERIMENTS

Alberto Abadie, Matthew M. Chingos, and Martin R. West*

Abstract—Policymakers are often interested in estimating how policy interventions affect the outcomes of those most in need of help. This concern has motivated the practice of disaggregating experimental results by groups constructed on the basis of an index of baseline characteristics that predicts the values of individual outcomes without the treatment. This paper shows that substantial biases may arise in practice if the index is estimated by regressing the outcome variable on baseline characteristics for the full sample of experimental controls. We propose alternative methods that correct this bias and show that they behave well in realistic scenarios.

I. Introduction

RECENT years have seen rapid growth in the use of randomized experiments in the social sciences. This resulted in part from the “credibility revolution” in empirical research, which increased scrutiny of the validity of conditions that allow credible estimation of treatment effects (Angrist & Pischke, 2010; Murnane & Willett, 2011). The main advantage of a large and well-executed randomized experiment is that the researcher can confidently rule out the possibility that unobserved differences between the treatment and control groups could explain the study’s results.

In addition to allowing estimation of average treatment effects, experiments also make it possible to obtain unbiased estimates of treatment effects for subgroups. Subgroup treatment effects are of particular interest to policymakers seeking to target policies on those most likely to benefit.¹ As a general rule, subgroups must be created based on characteristics that are either immutable (e.g., race) or observed before randomization (e.g., on a baseline survey) so that they could not possibly have been affected by the treatment.

Received for publication February 19, 2016. Revision accepted for publication November 3, 2017. Editor: Bryan S. Graham.

* Abadie: MIT; Chingos: Urban Institute; West: Harvard University.

We thank Beth Akers, Josh Angrist, Matias Cattaneo, Gary Chamberlain, David Deming, Sara Goldrick-Rab, Josh Goodman, Jerry Hausman, Guido Imbens, Max Kasy, Larry Katz, Amanda Pallais, Paul Peterson, Russ Whitehurst, participants in many seminars, and the editor and referees for helpful comments and discussions, and Jeremy Ferwerda for developing *estrat* (available at SSC), a Stata package that calculates the leave-one-out and repeated split-sample endogenous stratification estimators considered in this study.

A supplemental appendix is available online at http://www.mitpressjournals.org/doi/suppl/10.1162/rest_a_00732.

¹ Bitler, Gelbach, and Hoynes (2013), Crump et al. (2008, 2009), and Djebbari and Smith (2008) study estimation and testing of subgroup treatment effects.

However, many researchers and policymakers are interested in estimating how treatments affect those most in need of help, that is, those who would attain unfavorable outcomes in the absence of the treatment. Treatment parameters of this nature depend on the joint distribution of potential outcomes with and without treatment, which is not identified by randomization (see Heckman, Smith, & Clements, 1997). A practical solution to this problem is to combine baseline characteristics into a single index that reflects each participant’s predicted outcome without treatment and conduct separate analysis for subgroups of participants defined in terms of intervals of the predicted outcome without treatment.

A well-known implementation of this idea is the use of data on out-of-sample untreated units to estimate a prediction model for the outcome variable, which can then be applied to predict outcomes without treatment for the experimental units. This approach is common in medical research, where validated risk models are available to stratify experimental subjects based on their predicted probability of certain health outcomes (Kent & Hayward, 2007).

However, experimental studies in the social sciences often lack externally validated models that can be employed to predict the outcomes that experimental units would attain in the absence of the treatment. A potential approach to this problem that is gaining popularity among empirical researchers is to use in-sample information on the relationship between the outcome of interest and covariates for the experimental controls to estimate potential outcomes without treatment for all experimental units. We call this practice endogenous stratification because it uses in-sample data on the outcome variable to stratify the sample.

Endogenous stratification is typically implemented in practice by first regressing the outcome variable on baseline characteristics using the full sample of experimental controls, and then using the coefficients from this regression to generate predicted potential outcomes without treatment for all sample units. Unfortunately, as we will show, this procedure generates estimators of treatment effects that are substantially biased, and the bias follows a predictable pattern: results are biased upward for individuals with low predicted outcomes and biased downward for individuals with high predicted outcomes.

This bias pattern matches the results of several recent experimental studies that use this procedure and estimate strong positive effects for individuals with low predicted

outcomes and, in some cases, negative effects for individuals with high predicted outcomes. For example, a working paper by Goldrick-Rab et al. (2011) reports that a Wisconsin need-based financial aid program for postsecondary education had no overall impacts on college enrollment or college persistence among eligible students as a whole. Looking separately at subgroups based on predicted persistence, however, the study finds large positive effects on enrollment after three years for students in the bottom third of predicted persistence and almost equally large negative effects for students in the top third of predicted persistence.² A working paper by Dynarski, Hyman, and Schanzenbach (2011) analyzing long-term impacts of the Project STAR experiment similarly finds that assignment to a small class in grades K–3 increased college enrollment rates among the quintile of students with the lowest ex ante probability to enroll by 11 percentage points but had no impact on students in the top four quintiles. Pane et al. (2014) report experimental estimates of the effects of a technology-based algebra curriculum on the test scores of middle and high school students disaggregated by quintiles of predicted test scores. For middle school students exposed to the program in the first year of its implementation, they find “potentially moderately large positive treatment effects in the lowest quintile and small negative effects of treatment in the highest two quintiles.” Hemelt, Roth, and Eaton (2012) find no significant average impacts in an experimental evaluation of the effects of two elementary school interventions on college enrollment or degree receipt. They report, however, significant positive impacts on two-year college enrollment for both interventions and on associate’s degree completion for one of the interventions when they restrict the sample to students in the bottom quartile of the in-sample predicted probability of college attendance. Rodriguez-Planas (2012) reports that a mentoring program for adolescents reduced risky behavior and improved educational attainment for students in the top half of the risk distribution but increased risky behavior in the bottom half.³

²Goldrick-Rab et al. (2011) report that for students in the bottom-third group of predicted persistence, grant receipt was associated with an increase of 17 percentage points in enrollment three years after they started college. Conversely, for students in the top-third group of predicted persistence, grant receipt was associated with a decrease of 15 percentage points in enrollment three years after the start of college. These findings were characterized by the authors as “exploratory” but received widespread media coverage, including articles in the *Chronicle of Higher Education*, *Inside Higher Education*, and *Education Week*. In a related paper on the design of randomized experiments, Harris and Goldrick-Rab (2012) discuss potential explanations for the unexpected heterogeneity in their impact estimates based on full-sample endogenous stratification.

³We should note that because of recent concerns about the properties of endogenous stratification estimators raised in part by previous versions of this paper, endogenous stratification estimates do not appear in the published versions of two of the studies described here; see Dynarski, Hyman, and Schanzenbach (2013) and Hemelt, Roth, and Eaton (2013), or in a subsequent working paper on the grant program evaluated in Goldrick-Rab et al. (2011) by the same authors; see Goldrick-Rab et al. (2012). Rodriguez-Planas (2014) uses the analysis and estimators of this paper to update Rodriguez-Planas (2012) correcting for endogenous stratification biases.

Endogenous stratification also plays a supporting role in Angrist and Lavy’s (2009) experimental evaluation of a cash incentive program aimed at increasing matriculation certification rates for Israeli high school students. In order to test whether the program was most effective for girls on the certification margin, the researchers first group students by baseline test scores. They also, however, report results for students grouped by ex ante certification probability based on a broader set of background characteristics as “a check on the notion that high lagged scores identify students who have a shot at classification” (p. 1396).

The possibility of bias arising from endogenous stratification has been previously acknowledged in the evaluation literature (see, e.g., Peck, 2003), in statistics (Hansen, 2008), and in economics (Sanbonmatsu et al., 2006, and Giné, Goldberg, and Yang, 2012), but the size and significance of the bias in realistic evaluation settings are not well understood.⁴ A deceptively comforting property of the bias is that it vanishes as sample size increases, under weak regularity conditions. However, as we demonstrate in this paper using data from the National JTPA Study and the Project STAR experiment, biases resulting from endogenous stratification can completely alter the quantitative and qualitative conclusions of empirical studies.

In the remainder of this paper, we first describe in more detail the increasingly popular practice of stratifying experimental data by groups constructed on the basis of the predicted values from a regression of the outcome on baseline covariates for the full sample of experimental controls. We next explain why this method generates biases and describe the direction of those biases. We then describe leave-one-out and repeated split sample procedures that generate consistent estimators and show that the biases of these estimators are substantially lower than the bias of the full sample estimator in two realistic scenarios. We use data from the National JTPA Study and the Project STAR experiment to demonstrate the performance of endogenous stratification estimators and the magnitude of their biases. We restrict our attention to randomized experiments because this is the setting where endogenous stratification is typically used. However, similarly large biases may arise from endogenous stratification in observational studies.

II. Using Control Group Data to Create Predicted Outcomes

We begin by describing in detail the endogenous stratification method already outlined, which aims to classify study participants into groups based on their predicted value of the outcome variable in the absence of the treatment. Suppose

⁴Hausman and Wise (1977) and Hausman and Wise (1981), from which we borrow the term *endogenous stratification*, study the related problem of biased sampling in randomized experiments. Altonji and Segal (1996) study biases that arise in the context of efficient generalized methods of moments estimation for reasons that are related to those that explain the bias of the full sample endogenous stratification estimator.

that the sample consists of N observations of the triple (y, w, \mathbf{x}) , where y is an outcome variable, w is the treatment, and \mathbf{x} is a vector of baseline characteristics. When the object of interest is the average treatment effect, which in a randomized experiment is equal to $\tau = E[y|w = 1] - E[y|w = 0]$, researchers typically compare sample average outcomes for the treated and the control groups:

$$\widehat{\tau} = \frac{\sum_{i=1}^N y_i w_i}{\sum_{i=1}^N w_i} - \frac{\sum_{i=1}^N y_i (1 - w_i)}{\sum_{i=1}^N (1 - w_i)}.$$

As discussed, researchers sometimes aim to compare treated and nontreated after stratifying on a predictor of the outcome in the absence of the treatment. To our knowledge, most studies that use endogenous stratification implement it roughly as follows:

1. Regress the outcome variable on a set of baseline characteristics using the control group only. The regression coefficients are

$$\widehat{\boldsymbol{\beta}} = \left(\sum_{i=1}^N \mathbf{x}_i (1 - w_i) \mathbf{x}_i' \right)^{-1} \sum_{i=1}^N \mathbf{x}_i (1 - w_i) y_i.$$

2. Use the estimated coefficients to generate predicted outcome values for all participants (both treatment and control groups), $\mathbf{x}_i' \widehat{\boldsymbol{\beta}}$.
3. Divide participants into groups based on their predicted outcomes. Typically unit i is assigned to group k if $\mathbf{x}_i' \widehat{\boldsymbol{\beta}}$ falls in some interval delimited by c_{k-1} and c_k . The interval limits may be fixed or could be quantiles of the empirical distribution of $\mathbf{x}_i' \widehat{\boldsymbol{\beta}}$. Many authors use a three-bin classification scheme of low, medium, and high predicted outcomes.
4. Estimate treatment effects for each of the subgroups,

$$\widehat{\tau}_k = \frac{\sum_{i=1}^N y_i I_{[w_i=1, c_{k-1} < \mathbf{x}_i' \widehat{\boldsymbol{\beta}} \leq c_k]}}{\sum_{i=1}^N I_{[w_i=1, c_{k-1} < \mathbf{x}_i' \widehat{\boldsymbol{\beta}} \leq c_k]}} - \frac{\sum_{i=1}^N y_i I_{[w_i=0, c_{k-1} < \mathbf{x}_i' \widehat{\boldsymbol{\beta}} \leq c_k]}}{\sum_{i=1}^N I_{[w_i=0, c_{k-1} < \mathbf{x}_i' \widehat{\boldsymbol{\beta}} \leq c_k]}}$$

where I_A is the indicator function that takes values 1 if event A is realized and value 0 otherwise. Alternatively, treatment effect estimates could be computed after controlling for a set of covariates using regression.

For example, in their study of the impact of a need-based grant, Goldrick-Rab et al. (2011) regress college persistence on baseline characteristics using only observations from the control group, generate predicted probabilities of college persistence for all students, classify students into three equal-sized groups based on their ex ante predicted probability, and then estimate treatment effects for each of the three groups.

This is a simple and direct approach to stratification, which has great intuitive appeal. Moreover, it is easy to show that under usual regularity conditions, $\widehat{\tau}_k$ converges to

$$\tau_k = E[y|w = 1, c_{k-1} < \mathbf{x}'\boldsymbol{\beta} \leq c_k] - E[y|w = 0, c_{k-1} < \mathbf{x}'\boldsymbol{\beta} \leq c_k].$$

As we demonstrate in this paper, however, $\widehat{\tau}_k$ is biased in finite samples, and the bias follows a predictable pattern.

Here we provide an intuitive explanation of the bias. To simplify the exposition, suppose that predicted outcomes are divided into three groups (low, medium, high). Let $\boldsymbol{\beta} = (E[\mathbf{x}\mathbf{x}'|w = 0])^{-1} E[\mathbf{x}y|w = 0]$ be the population counterpart of $\widehat{\boldsymbol{\beta}}$, and let $e_i = y_i - \mathbf{x}_i'\boldsymbol{\beta}$ be the regression error. In a finite sample, untreated observations with large negative values for e_i tend to be overfitted, so we expect $\mathbf{x}_i'\widehat{\boldsymbol{\beta}} < \mathbf{x}_i'\boldsymbol{\beta}$, which pushes these observations toward the lower interval of predicted outcomes.⁵ This creates a negative bias in the average outcome among control observations that fall into the lower interval for $\mathbf{x}_i'\boldsymbol{\beta}$ and, therefore, a positive bias in the average treatment effect estimated for that group. Analogously, average treatment effect estimators for the upper intervals of predicted outcomes are biased downward. Endogenous stratification results in a predictable pattern: average treatment effect estimators are biased upward for individuals with low predicted outcomes and biased downward for individuals with high predicted outcomes. As we will demonstrate, because the finite sample bias of the endogenous stratification estimator is created by overfitting, this bias tends to be more pronounced when the number of observations is small and the dimensionality of \mathbf{x}_i is large.

A natural solution to the overfitting issue is provided by leave-one-out estimators. This is the approach followed in Sanbonmatsu et al. (2006). Harvill, Peck, and Bell (2013) propose a variant of this approach based on 10-fold cross-validation. Let

$$\widehat{\boldsymbol{\beta}}_{(-i)} = \left(\sum_{j \neq i} \mathbf{x}_j (1 - w_j) \mathbf{x}_j' \right)^{-1} \sum_{j \neq i} \mathbf{x}_j (1 - w_j) y_j,$$

be the regression coefficients estimators that discard observation i . Overfitting is precluded by not allowing the outcome, y_i , of each observation to contribute to the estimation of its own predicted value, $\mathbf{x}_i'\widehat{\boldsymbol{\beta}}_{(-i)}$. Because only untreated observations are employed in the estimation of $\widehat{\boldsymbol{\beta}}_{(-i)}$ and $\widehat{\boldsymbol{\beta}}$, if i is a treated observation, then $\widehat{\boldsymbol{\beta}}_{(-i)} = \widehat{\boldsymbol{\beta}}$. We consider the following leave-one-out estimator of τ_k :

$$\widehat{\tau}_k^{LOO} = \frac{\sum_{i=1}^N y_i I_{[w_i=1, c_{k-1} < \mathbf{x}_i' \widehat{\boldsymbol{\beta}}_{(-i)} \leq c_k]}}{\sum_{i=1}^N I_{[w_i=1, c_{k-1} < \mathbf{x}_i' \widehat{\boldsymbol{\beta}}_{(-i)} \leq c_k]}} - \frac{\sum_{i=1}^N y_i I_{[w_i=0, c_{k-1} < \mathbf{x}_i' \widehat{\boldsymbol{\beta}}_{(-i)} \leq c_k]}}{\sum_{i=1}^N I_{[w_i=0, c_{k-1} < \mathbf{x}_i' \widehat{\boldsymbol{\beta}}_{(-i)} \leq c_k]}}.$$

⁵ Overfitting arises when the sample regression values, $\mathbf{x}_i'\widehat{\boldsymbol{\beta}}$, are closer to the outcome values, y_i , than the population regression values, $\mathbf{x}_i'\boldsymbol{\beta}$. The appendix provides a formal explanation of overfitting in a regression model.

Under weak assumptions, it can be seen that both $\widehat{\tau}_k$ and $\widehat{\tau}_k^{LOO}$ are consistent estimators of τ_k . Moreover, $\widehat{\tau}_k$ and $\widehat{\tau}_k^{LOO}$ have the same large sample distribution.⁶ However, we show in section IV that $\widehat{\tau}_k$ is substantially biased in two realistic scenarios, while $\widehat{\tau}_k^{LOO}$ is not. A separate issue in the estimation of τ_k is that first-step estimation of β affects the large-sample distribution of the estimator (see the online appendix for a derivation of the large-sample distribution of $\widehat{\tau}_k$ and $\widehat{\tau}_k^{LOO}$). The contribution of the estimation of β to the variance of $\widehat{\tau}_k$ has been ignored in empirical practice.

Another way to avoid overfitting is sample splitting. We consider a repeated split sample estimator. In each repetition, m , the untreated sample is randomly divided into two groups, which we will call the prediction and the estimation groups. Let $v_{im} = 0$ if untreated observation i is assigned the prediction group in repetition m and $v_{im} = 1$ if it is assigned to the estimation group. In each repetition, m , we estimate β using only the observations in the prediction group:

$$\widehat{\beta}_m = \left(\sum_{i=1}^N \mathbf{x}_i (1 - w_i) (1 - v_{im}) \mathbf{x}_i' \right)^{-1} \times \sum_{i=1}^N \mathbf{x}_i (1 - w_i) (1 - v_{im}) y_i.$$

For each repetition, m , the split sample estimator of τ_k is

$$\widehat{\tau}_{km}^{SS} = \frac{\sum_{i=1}^N y_i I_{[w_i=1, c_{k-1} < \mathbf{x}_i' \widehat{\beta}_m \leq c_k]}}{\sum_{i=1}^N I_{[w_i=1, c_{k-1} < \mathbf{x}_i' \widehat{\beta}_m \leq c_k]}} - \frac{\sum_{i=1}^N y_i I_{[w_i=0, v_{im}=1, c_{k-1} < \mathbf{x}_i' \widehat{\beta}_m \leq c_k]}}{\sum_{i=1}^N I_{[w_i=0, v_{im}=1, c_{k-1} < \mathbf{x}_i' \widehat{\beta}_m \leq c_k]}}.$$

We then average $\widehat{\tau}_{km}^{SS}$ over M repetitions to obtain the repeated split sample estimator:

$$\widehat{\tau}_k^{RSS} = \frac{1}{M} \sum_{m=1}^M \widehat{\tau}_{km}^{SS}.$$

The repeated split sample estimator is asymptotically unbiased and Normal but, unlike the leave-one-out estimator, its large sample distribution does not coincide with the large sample distribution of the full-sample endogenous stratification estimator. For large M , however, the difference between the large-sample distribution of the repeated split sample estimator and the large-sample distribution of the full-sample and leave-one-out estimators is small.⁷

Given that $\widehat{\tau}_k^{LOO}$ and $\widehat{\tau}_k^{RSS}$ require higher computational effort than $\widehat{\tau}_k$, computational aspects of the estimation deserve some comment. First, we note that leave-one-out

versions of $\mathbf{x}_i' \widehat{\beta}$ can be efficiently computed through the well-known formula

$$\mathbf{x}_i' \widehat{\beta}_{(-i)} = \mathbf{x}_i' \widehat{\beta} - \frac{h_{Ni}}{1 - h_{Ni}} (y_i - \mathbf{x}_i' \widehat{\beta}),$$

where h_{Ni} is the leverage of observation i .⁸ This implies that all values of $\mathbf{x}_i' \widehat{\beta}_{(-i)}$ for the experimental controls can be computed with a single regression of the outcome on the covariates for the entire sample of experimental controls. While the repeated split sample estimator, $\widehat{\tau}_k^{RSS}$, is more computationally demanding than $\widehat{\tau}_k$ or $\widehat{\tau}_k^{LOO}$ (especially for large values of M), the calculations of the components, $\widehat{\tau}_{km}^{SS}$, are independent tasks that can easily be divided in batches and parallelized.

In the next section, we apply the estimators we have described to the analysis of data from two well-known experimental studies: the National JTPA Study and the Tennessee Project STAR experiment.

III. Evidence of Large Biases in Two Actual Applications

To demonstrate the performance of the estimators described in the previous section and the magnitude of their biases in realistic scenarios, we use data from two randomized evaluations: the National JTPA Study, an evaluation of an employment and training program in the United States, and the kindergarten cohort of the Tennessee Project STAR class-size experiment.

A. The National JTPA Study

We first examine data from the National JTPA Study, a large experimental evaluation of an employment and training program commissioned by the U.S. Department of Labor in the late 1980s. The study data have been extensively analyzed by Orr et al. (1996), Bloom et al. (1997), and many others.⁹ The study randomized access to JTPA services to applicants in sixteen service delivery areas (SDAs), across the United States. Randomized assignment was done after applicants were deemed eligible for the program and recommended to one of three possible JTPA service strategies: on-the-job training/job search assistance, classroom training, and other services. Individuals in the treatment group were provided with access to JTPA services, and individuals in the control group were excluded from JTPA services for an eighteen-month period after randomization. We use data for the sample of adult males recommended to the on-the-job training/job search assistance service strategy and discard three SDAs with few observations. Our sample consists of 1,681 treated observations and 849 untreated observations, for a total of 2,530 observations in thirteen SDAs.¹⁰ In this

⁸ See the online appendix for a precise definition of h_{Ni} .

⁹ Moreover, the historical recounting in Peck (2013) suggests that the practice of endogenous stratification may have originated within the JTPA evaluation.

¹⁰ See the online appendix for detailed information on sample selection and estimation methods.

⁶ Proofs of these and other formal statements made in this paper are provided in the online appendix.

⁷ See the online appendix for a proof.

TABLE 1.—JTPA ESTIMATION RESULTS

| A. Average Treatment Effect | | | | | | |
|--|--------------------------|-----------------------|--------------------------|-------------------------|-----------------------|--------------------------|
| | Unadjusted | | | Adjusted | | |
| $\hat{\tau}$ | 1,516.49* (807.27) | | | 1,207.22 (763.54) | | |
| B. Average Treatment Effect by Predicted Outcome Group | | | | | | |
| | Unadjusted | | | Adjusted | | |
| | Low | Medium | High | Low | Medium | High |
| $\hat{\tau}_k$ | 2,379.65** (1,151.07) | -719.38 (1,474.81) | 2,397.26 (1,672.62) | 2,011.70* (1,150.68) | -554.65 (1,482.32) | 1,769.03 (1,639.06) |
| $\hat{\tau}_k^{LOO}$ | 573.74 (1,201.33) | 35.31 (1,509.30) | 3,646.53** (1,727.08) | 173.45 (1,213.25) | 172.28 (1,513.70) | 3,118.17* (1,679.62) |
| $\hat{\tau}_k^{RSS}$ | 788.75 (1,027.47) | 254.25 (1,092.85) | 3,569.41** (1,496.73) | 412.01 (1,042.17) | 181.81 (1,087.51) | 2,942.69** (1,454.16) |
| $\hat{\tau}_k^{PREV}$ | 1,278.88 (1,221.96) | -67.95 (1,284.77) | 3,972.21** (1,497.47) | 822.05 (1,235.13) | -150.89 (1,274.45) | 3,146.85** (1,430.37) |

Bootstrap standard errors, based on 1,000 bootstrap repetitions, are reported in parentheses. $\hat{\tau}_k^{RSS}$ uses 100 repetitions, with 425 untreated observations in the prediction group and 424 in the estimation group. The “unadjusted” estimates are differences in mean outcomes between treated and nontreated. The “adjusted” estimates are regression coefficients on the treatment variable in a linear regression that includes the list of covariates detailed in section III. Statistically significant at *0.10 and **0.05.

example, w_i is an indicator of a randomized offer of JTPA services, y_i is nominal thirty-month earnings in U.S. dollars after randomization, and x_i includes age, age squared, marital status, previous earnings, indicators for having worked fewer than thirteen weeks during the year previous to randomization, having a high school diploma, being African American, and being Hispanic, as well as SDA indicators.

Table 1 reports estimates for the JTPA sample. The first row reports two treatment effect estimates. The “unadjusted” estimate is the difference in outcome means between treated and controls; the “adjusted” estimate is the coefficient on the treatment indicator in a linear regression of the outcome variable, y_i , on the treatment indicator, w_i , and the covariates, x_i , listed above. The unadjusted estimate suggests a \$1,516 effect on thirty-month earnings. This estimate is significant at the 10% level. Regression adjustment reduces the point estimate to \$1,207, which becomes marginally nonsignificant at the 10% level. The rest of table 1 reports average treatment effects by predicted-outcome group. The first set of estimates corresponds to $\hat{\tau}_k$, the full-sample endogenous stratification estimator. This estimator produces a large and significant effect for the low predicted-outcome group. The unadjusted estimate is \$2,380 and significant at the 5% level. This represents a 12.6% effect on thirty-month earnings once we divide it by the average value of thirty-month earnings among the experimental controls. It also represents an effect that is 57% higher than the corresponding unadjusted estimate for the average treatment effect in the first row of the table. The adjusted estimate is \$2,012, similarly large, and significant at the 10% level. For the high predicted-outcome group, the estimates are also large but not statistically significant at conventional test levels. For the middle predicted-outcome group, the estimates are negative but of moderate magnitude and not statistically significant. All in all, the full-sample endogenous stratification estimates provide a much more favorable picture of JTPA effectiveness relative to the average treatment effects reported on the first row. The bulk of the effect seems to be concentrated on the

low predicted-outcome group, precisely the one in most need of help, with more diffuse effects estimated for the middle and high predicted-outcome groups.

The next two sets of estimates reported in table 1 correspond to the leave-one-out estimator, $\hat{\tau}_k^{LOO}$, and the repeated split sample estimator $\hat{\tau}_k^{RSS}$, with number of repetitions, M , equal to 100. These two estimators, which avoid overfitting bias arising from the estimation of β , produce results that are substantially different from those obtained with the full-sample endogenous stratification estimator, $\hat{\tau}_k$. Relative to the $\hat{\tau}_k$ estimates, the $\hat{\tau}_k^{LOO}$ and $\hat{\tau}_k^{RSS}$ estimates are substantially smaller for the low predicted-outcome group and substantially larger for the high predicted-outcome group. For the high predicted-outcome group, we obtain unadjusted estimates of \$3,647 (leave-one-out) and \$3,569 (split sample), both significant at the 5% level, and adjusted estimates of \$3,118 (leave-one-out) and \$2,943 (split sample) significant at the 10% and 5% levels, respectively. The estimates for the low and middle predicted-outcome groups are small in magnitude and not statistically significant. These results place the bulk of the treatment effect on the high predicted-outcome group and do not provide substantial statistical evidence of beneficial effects for the low and middle predicted-outcome groups.¹¹ The comparison of estimates produced with the full sample endogenous stratification estimator and the leave-one-out and split sample estimators suggest that the overfitting bias of the full sample endogenous stratification estimator is of substantial magnitude and dramatically changes the qualitative and quantitative interpretations of the results.

As a further check on the magnitude of endogenous stratification biases in the analysis of the National JTPA Study data, table 1 reports a last set of treatment effects estimates,

¹¹ This is loosely consistent with the findings in Abadie, Angrist, and Imbens (2002), who report large JTPA effects at the upper tail of the distribution of earnings for male trainees and no discernible effects at the middle or lower parts of the distribution.

TABLE 2.—STAR ESTIMATION RESULTS

| A: Average Treatment Effect | | Unadjusted | | | Adjusted | | |
|--|----------------------|----------------------|-----------------------|----------------------|----------------------|-----------------------|--|
| $\hat{\tau}$ | | 0.1659** (0.0329) | | | 0.1892** (0.0294) | | |
| B: Average Treatment Effect by Predicted Outcome Group | | | | | | | |
| | Unadjusted | | | Adjusted | | | |
| | Low | Medium | High | Low | Medium | High | |
| $\hat{\tau}_k$ | 0.3705** (0.0521) | 0.2688** (0.0655) | -0.1330** (0.0636) | 0.3908** (0.0509) | 0.3023** (0.0678) | -0.1242** (0.0614) | |
| $\hat{\tau}_k^{LOO}$ | 0.3277** (0.0547) | 0.2499** (0.0670) | -0.0486 (0.0654) | 0.3440** (0.0519) | 0.2730** (0.0696) | -0.0660 (0.0634) | |
| $\hat{\tau}_k^{RSS}$ | 0.3152** (0.0467) | 0.2617** (0.0505) | -0.0520 (0.0567) | 0.3130** (0.0459) | 0.3005** (0.0526) | -0.0374 (0.0552) | |

Bootstrap standard errors, based on 1,000 bootstrap repetitions, are reported in parentheses. $\hat{\tau}_k^{RSS}$ uses 100 repetitions, with 1,009 untreated observations in the prediction group and 1,008 in the estimation group. The “unadjusted” estimates are differences in mean outcomes between treated and nontreated. The “adjusted” estimates are regression coefficients on the treatment variable in a linear regression that includes the list of covariates detailed in section III. Statistically significant at *0.10 and **0.05.

which are stratified using data on earnings before randomization. The National JTPA Study data include individual earnings during the twelve months before randomization. We use the sorting of the experimental subjects in terms of prerandomization earnings to approximate how the experimental subjects would have been sorted in terms of earnings in the absence of the treatment. We construct the estimator $\hat{\tau}_k^{PREV}$ in the same way as $\hat{\tau}_k$ but using previous earnings, instead of predicted earnings, to divide the individuals into three groups of approximately equal size. Notice that because previous earnings is a baseline characteristic, $\hat{\tau}_k^{PREV}$ is not affected by overfitting bias. As shown on the bottom of table 1, stratification on previous earnings produces results similar to those obtained with $\hat{\tau}_k^{LOO}$ and $\hat{\tau}_k^{RSS}$: large and significant effects for the high predicted-outcome group and smaller and nonsignificant effects for the middle and low predicted-outcome groups.

B. The Project STAR Experiment

Our second example uses data from the Project STAR class-size study. In Project STAR, students in 79 Tennessee schools were randomly assigned to small, regular-size, and regular-size classes with a teacher’s aide. Krueger (1999) analyzes the STAR data set and provides detailed explanations of the STAR experiment. For our analysis, we use the 3,764 students who entered the study in kindergarten and were assigned to small classes or to regular-size classes (without a teacher’s aide). Our outcome variable is standardized end-of-the-year kindergarten math test scores.¹² The covariates are indicators for African-American, female, eligibility for the free lunch program, and school attended. We discard observations with missing values in any of these variables.

Results for the STAR experiment data are reported in table 2. The adjusted and unadjusted estimators of the

average treatment effect on the first row of table 2 show positive and significant effects. Using a simple difference in means, the effect of small classes is estimated as 0.1659 of the regular class standard deviation in math test scores and 0.1892 of the same standard deviation when we use a regression-adjusted estimator.¹³ In both cases, the estimates are significant at the 5% level. For the low and middle predicted-outcomes groups, the full-sample endogenous stratification estimator, $\hat{\tau}_k$, produces estimates that are positive and roughly double the average treatment effects estimates on the first row of the table. Counterintuitively, however, the full sample endogenous stratification estimates for the high predicted-outcome group are negative and significant. They seem to suggest that being assigned to small classes was detrimental for students predicted to obtain high math scores if all students had remained in regular-size classes. We deem this result counterintuitive because it implies that reductions in the student/teacher ratio have detrimental effects on average for a large group of students. Notice that the magnitudes of the negative effects estimated for the high predicted-outcome group are substantial: smaller, but not far from the positive average treatment effects reported in the first row of the table. We will see that the large and significant negative effect for the high predicted-outcome group disappears when the leave-one-out or the repeated split sample procedures are used for estimation. Indeed, the leave-one-out and repeated split sample estimates on the two bottom rows of table 2 suggest positive, significant, and large effects on the low and middle predicted-outcome groups and effects of small magnitude and not reaching statistical significance at conventional test levels for the high predicted-outcome group. As for the JTPA, the qualitative and quantitative interpretations of

¹² Standardized test scores are computed dividing raw test scores by the standard deviation of the distribution of the scores in regular-size classes.

¹³ To be consistent with much of the previous literature on the STAR experiment, we report both regression-adjusted and unadjusted estimates. Because the probability of assignment to a small class varied by school in the STAR experiment, the regression-adjusted estimator is most relevant in this setting. As in Krueger (1999), however, covariate regression adjustment does not substantially change our estimates.

the STAR experiment results change dramatically when the leave-one-out or the repeated split sample estimators are used instead of the full-sample endogenous stratification estimator.

In this section, we have used data from two well-known and influential experimental studies to investigate the magnitude of the distortion that overfitting may induce on endogenous stratification estimators. In the next section, we use Monte Carlo simulations to assess the magnitude of the biases of the different estimators considered in section II. To keep the exercise as realistic as possible, in two of our simulations we choose data-generating processes that mimic the features of the JTPA and STAR data sets.

IV. Simulation Evidence on the Behavior of Endogenous Stratification Estimators

This section reports simulation evidence on the finite sample behavior of endogenous stratification estimators. We run Monte Carlo simulations in three settings. In the first two Monte Carlo simulations, we make use of the JTPA and STAR data sets to assess the magnitudes of biases and other distortions to inference in realistic scenarios. The third and fourth Monte Carlo simulations use computer-generated data to investigate how the bias of endogenous stratification estimators changes when the sample size or the number of covariates changes.

In the JTPA-based simulation, we first use the JTPA control units to estimate a two-part model for the distribution of earnings conditional on the covariates of the adjusted estimates in table 1. The two-part model consists of a logit specification for the probability of 0 earnings and a Box-Cox model for positive earnings.¹⁴ In each Monte Carlo iteration, we draw 2,530 observations—the same number of observations as in the JTPA sample—from the empirical distribution of the covariates in the JTPA sample. Next, we use the estimated two-part model to generate earnings data for each observation in the Monte Carlo sample. Then we randomly assign 1,681 observations to the treatment group and 849 observations to the control group, to match the numbers of treated and control units in the original JTPA sample. Finally, in each Monte Carlo iteration, we compute the full-sample, leave-one-out, and repeated split sample endogenous stratification estimates. For the repeated split sample estimator, we use 100 repetitions ($M = 100$), which provide a reasonable balance between precision of the estimators and computational time. We also compute the value of the unfeasible estimator, $\hat{\tau}_k^{UNF}$, obtained by stratification on the population regression function (which can be calculated from the estimated parameters of the two-part model by simulation). We conduct 10,000 Monte Carlo iterations.

Figure 1 reports the Monte Carlo distributions of the endogenous stratification estimators that divide the

experimental sample into three categories of predicted earnings of roughly equal size (bottom third, middle third, and top third). To economize space, this figure shows only the distribution of the unadjusted estimators.¹⁵ Because assignment to the treatment and control groups is randomized in our simulation and because the process that generates earnings data is the same for treated and controls, it follows that the average effect of the treatment in the simulations is equal to 0 unconditionally as well as conditional on the covariates. As a result, unbiased estimators should generate Monte Carlo distributions centered at 0. The first plot of figure 1 shows the Monte Carlo distribution of the full-sample endogenous stratification estimator of average treatment effects conditional on predicted earnings group. The pattern of the distribution of the average treatment effect estimator for the bottom, middle, and top third predicted earnings groups matches the directions of the biases discussed in section II. That is, $\hat{\tau}_k$ is biased upward for the low predicted-earnings group and downward for the high predicted-earnings group. The remaining three plots of figure 1 do not provide evidence of substantial biases for the leave-one-out, repeated split sample, or unfeasible estimators. These three estimators produce Monte Carlo distributions that are centered at 0 for each predicted-earnings category.

Table 3 reports biases, 1 minus the coverage rates of nominal 0.95 confidence intervals based on the Normal approximation, and root mean square error (root-MSE) values for endogenous stratification estimators in the JTPA-based Monte Carlo simulation. In addition to the estimators considered in figure 1, we compute a single split sample estimator, $\hat{\tau}_k^{SSS}$, which is defined like the repeated split sample estimator but with $M = 1$ (see Peck, 2003). The full-sample endogenous stratification estimator is subject to substantial distortions for the low and high predicted-earnings group. The magnitude of the bias in each these two groups is more than \$1,000, which is substantial compared to the \$1,516 and \$1,207 unadjusted and adjusted average effect estimates in the JTPA data. As reflected in figure 1, the bias is positive for the low predicted-earnings group and negative for the high predicted-earnings group. Biases are uniformly small for the leave-one-out, repeated split sample, and unfeasible estimators, but the leave-one-out estimator has higher biases than the repeated split sample and the unfeasible estimator. Similar results emerge for coverage rates and mean square errors. The full-sample endogenous stratification estimator produces substantially higher than nominal coverage rates and substantially higher root-MSE than the leave-one-out and repeated split sample estimators for the low and high predicted-income categories. The repeated split sample estimator dominates in terms of root-MSE. The single split sample estimators produce small biases and close to nominal coverage rates, but they have root-MSE values consistently higher than the full-sample endogenous stratification estimator.

¹⁴ Additional details about the simulation models can be found in the online appendix.

¹⁵ Simulation results for unadjusted and adjusted estimators are very similar, as reflected in tables 3 to 6.

FIGURE 1.—DISTRIBUTIONS OF THE ESTIMATORS IN THE JTPA SIMULATION

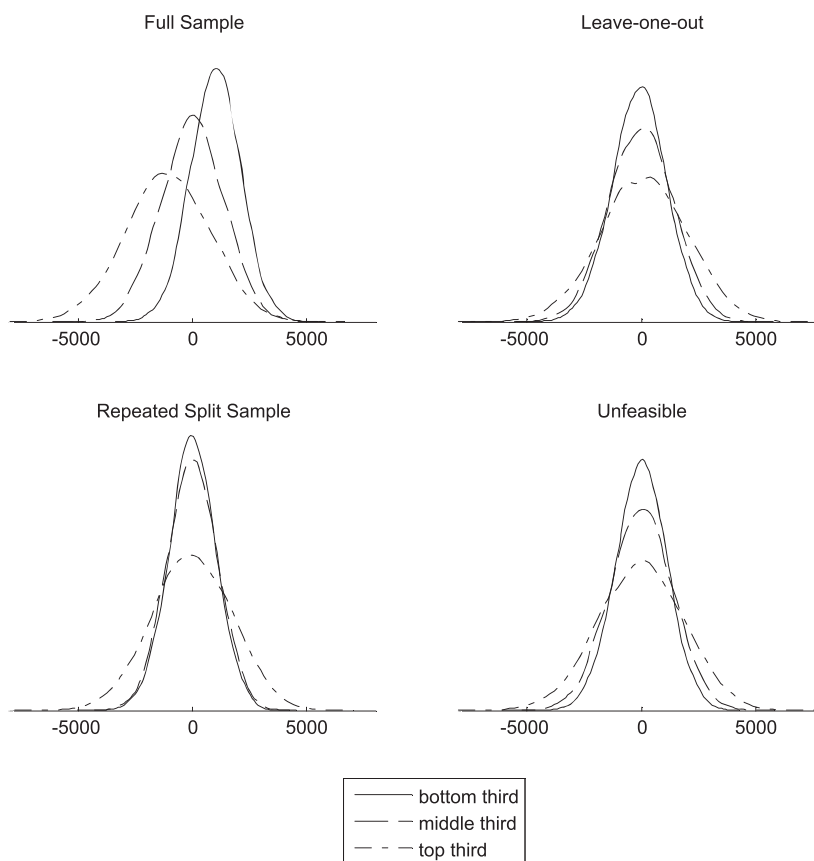


TABLE 3.—JTPA SIMULATION RESULTS

| A. Bias | | | | | | |
|---|------------|----------|-----------|----------|----------|-----------|
| | Unadjusted | | | Adjusted | | |
| | Low | Medium | High | Low | Medium | High |
| $\hat{\tau}_k$ | 1,017.51 | -4.81 | -1,082.42 | 1,017.60 | -0.98 | -1,062.39 |
| $\hat{\tau}_k^{LOO}$ | -88.23 | -23.28 | 96.57 | -59.08 | -54.01 | 42.86 |
| $\hat{\tau}_k^{RSS}$ | -2.74 | -2.34 | -20.75 | -3.30 | -5.96 | -17.56 |
| $\hat{\tau}_k^{SSS}$ | 5.62 | -9.88 | 6.62 | 1.34 | -11.24 | 1.39 |
| $\hat{\tau}_k^{UNF}$ | -1.50 | -8.56 | -16.85 | -2.50 | -9.04 | -11.67 |
| B. One Minus Coverage Rates for Nominal 0.95 CI | | | | | | |
| | Unadjusted | | | Adjusted | | |
| | Low | Medium | High | Low | Medium | High |
| $\hat{\tau}_k$ | 0.152 | 0.049 | 0.089 | 0.154 | 0.050 | 0.089 |
| $\hat{\tau}_k^{LOO}$ | 0.051 | 0.048 | 0.050 | 0.051 | 0.049 | 0.051 |
| $\hat{\tau}_k^{RSS}$ | 0.051 | 0.048 | 0.049 | 0.052 | 0.048 | 0.050 |
| $\hat{\tau}_k^{SSS}$ | 0.050 | 0.048 | 0.051 | 0.052 | 0.049 | 0.050 |
| $\hat{\tau}_k^{UNF}$ | 0.053 | 0.050 | 0.050 | 0.053 | 0.051 | 0.051 |
| C. Root-MSE | | | | | | |
| | Unadjusted | | | Adjusted | | |
| | Low | Medium | High | Low | Medium | High |
| $\hat{\tau}_k$ | 1,492.26 | 1,364.27 | 2,145.78 | 1,489.89 | 1,375.04 | 2,065.87 |
| $\hat{\tau}_k^{LOO}$ | 1,192.35 | 1,399.74 | 1,895.93 | 1,180.13 | 1,398.76 | 1,800.86 |
| $\hat{\tau}_k^{RSS}$ | 1,031.61 | 1,101.14 | 1,751.43 | 1,022.53 | 1,103.53 | 1,660.50 |
| $\hat{\tau}_k^{SSS}$ | 1,500.97 | 1,797.74 | 2,383.17 | 1,493.40 | 1,792.14 | 2,271.90 |
| $\hat{\tau}_k^{UNF}$ | 1,119.34 | 1,372.67 | 1,867.51 | 1,118.52 | 1,383.76 | 1,792.25 |

Averages over 10,000 simulations. See section IV and the online appendix for details.

FIGURE 2.—DISTRIBUTIONS OF THE ESTIMATORS IN THE STAR SIMULATION

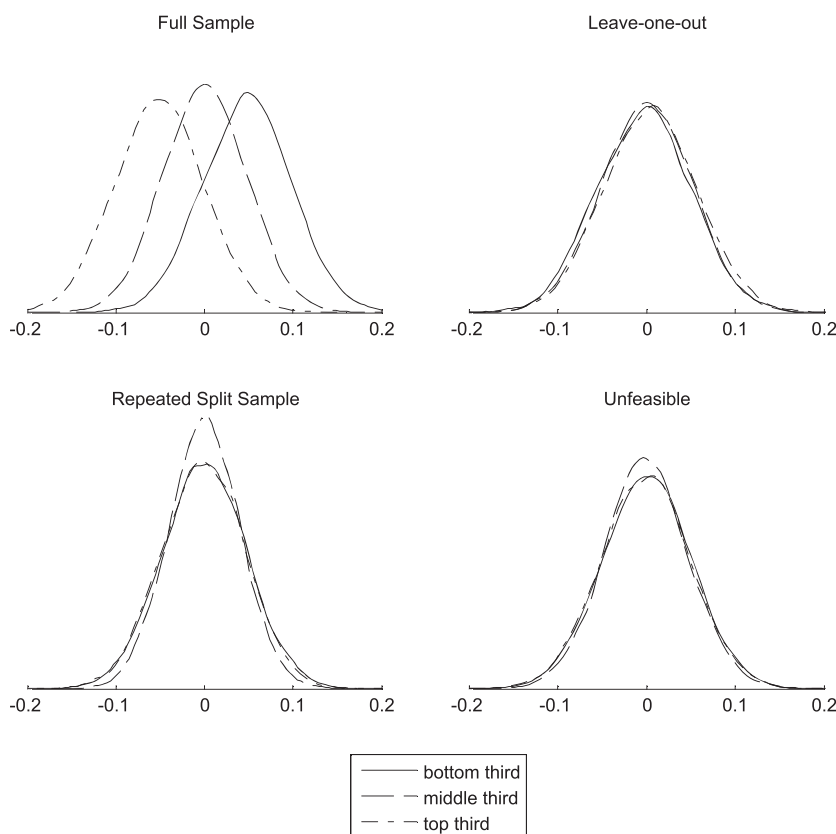


Figure 2 and table 4 report simulation results for the STAR-based Monte Carlo simulation. For this simulation, the data-generating process is based on a linear model with Normal errors. The model is estimated using data for STAR students in regular size classes. The results are qualitatively identical to those obtained in the JTPA-based simulation. The biases of $\hat{\tau}_k$ are around 0.05 and -0.05 for the low and medium predicted test score groups, respectively. These are sizable magnitudes compared to the STAR effect estimates in table 2. Also, as in the JTPA-based simulation, for the low and high predicted-outcome groups, coverage rates of the full sample endogenous stratification estimator are heavily distorted, and root-MSE values are larger than for the leave-one-out and the repeated split sample estimators. The repeated split sample estimator has the lowest root-MSE, and single sample splits produce root-MSE values larger than any other estimator with the exception of the full-sample endogenous stratification estimator.

The analysis of how average treatment effects covary with predicted outcomes without the treatment can also be based on a regression equation with interaction terms, such as

$$y_i = \alpha_0 + (\mathbf{x}'_i\boldsymbol{\beta})\alpha_1 + w_i\alpha_2 + w_i(\mathbf{x}'_i\boldsymbol{\beta})\alpha_3 + u_i,$$

where u_i is a regression error orthogonal to the included regressors. A negative sign of α_3 would indicate average

treatment effects inversely related to $\mathbf{x}'\boldsymbol{\beta}$. Under the data-generating processes employed in our simulations, α_3 is equal to 0. Figure 3 reports Monte Carlo distributions of estimators of α_3 for the JTPA-based and STAR-based simulations. The full sample and leave-one-out estimators use $\hat{\boldsymbol{\beta}}$ and $\hat{\boldsymbol{\beta}}_{(-i)}$, respectively, instead of $\boldsymbol{\beta}$, and estimate the regression equation by ordinary least squares. The unfeasible estimator uses the true value, $\boldsymbol{\beta}$. For $m = 1, \dots, M$, the repeated split sample estimator uses $\hat{\boldsymbol{\beta}}_m$ and averages the resulting estimates of α_3 over the M repetitions. Finally, we also report the distribution of the estimator of α_3 given by one-step nonlinear least squares estimation of the regression equation above. The one-step nature of the nonlinear least squares estimator implies that predicted outcomes are fitted to all experimental units, and not only to the units in the control group.¹⁶ The results in figure 3 are consistent with our previous evidence on the performance of estimators that stratify on subgroups of predicted values. The Monte Carlo distributions of the leave-one-out, repeated split sample, nonlinear least squares, and unfeasible estimators are all centered at 0. In contrast, the full-sample endogenous stratification estimator of α_3 is negatively biased.

¹⁶We thank Gary Chamberlain for suggesting this estimator. Nonlinear least squares estimation of the regression equation above uses the normalization $\alpha_0 = 0$ and $\alpha_1 = 1$ to ensure that the regression parameters are properly defined.

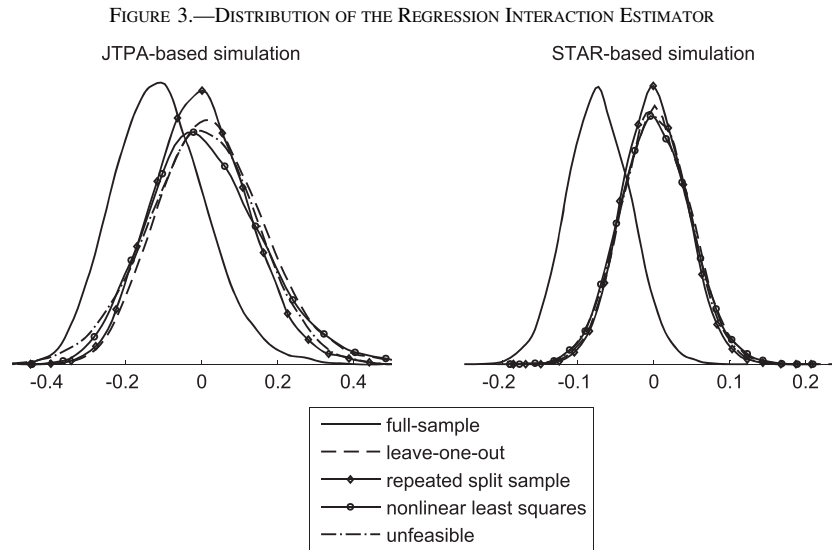
TABLE 4.—STAR SIMULATION RESULTS

| A. Bias | | | | | | |
|----------------------|------------|---------|---------|----------|---------|---------|
| | Unadjusted | | | Adjusted | | |
| | Low | Medium | High | Low | Medium | High |
| $\hat{\tau}_k$ | 0.0483 | 0.0006 | -0.0511 | 0.0487 | 0.0010 | -0.0506 |
| $\hat{\tau}_k^{LOO}$ | -0.0025 | 0.0005 | 0.0046 | 0.0028 | -0.0025 | -0.0075 |
| $\hat{\tau}_k^{RSS}$ | 0.0001 | -0.0000 | -0.0012 | 0.0002 | 0.0001 | -0.0010 |
| $\hat{\tau}_k^{SSS}$ | -0.0005 | 0.0004 | -0.0017 | -0.0002 | 0.0003 | -0.0015 |
| $\hat{\tau}_k^{UNF}$ | 0.0004 | -0.0003 | -0.0009 | 0.0002 | -0.0002 | -0.0006 |

| B. One Minus Coverage Rates for Nominal 0.95 CI | | | | | | |
|---|------------|--------|-------|----------|--------|-------|
| | Unadjusted | | | Adjusted | | |
| | Low | Medium | High | Low | Medium | High |
| $\hat{\tau}_k$ | 0.161 | 0.051 | 0.178 | 0.178 | 0.049 | 0.191 |
| $\hat{\tau}_k^{LOO}$ | 0.048 | 0.050 | 0.051 | 0.050 | 0.049 | 0.056 |
| $\hat{\tau}_k^{RSS}$ | 0.053 | 0.051 | 0.048 | 0.052 | 0.051 | 0.049 |
| $\hat{\tau}_k^{SSS}$ | 0.051 | 0.050 | 0.052 | 0.052 | 0.052 | 0.050 |
| $\hat{\tau}_k^{UNF}$ | 0.051 | 0.052 | 0.049 | 0.051 | 0.050 | 0.050 |

| C: Root-MSE | | | | | | |
|----------------------|------------|--------|--------|----------|--------|--------|
| | Unadjusted | | | Adjusted | | |
| | Low | Medium | High | Low | Medium | High |
| $\hat{\tau}_k$ | 0.0695 | 0.0472 | 0.0716 | 0.0677 | 0.0471 | 0.0691 |
| $\hat{\tau}_k^{LOO}$ | 0.0526 | 0.0509 | 0.0530 | 0.0492 | 0.0507 | 0.0494 |
| $\hat{\tau}_k^{RSS}$ | 0.0473 | 0.0402 | 0.0470 | 0.0444 | 0.0399 | 0.0439 |
| $\hat{\tau}_k^{SSS}$ | 0.0617 | 0.0589 | 0.0615 | 0.0577 | 0.0583 | 0.0571 |
| $\hat{\tau}_k^{UNF}$ | 0.0501 | 0.0469 | 0.0503 | 0.0480 | 0.0475 | 0.0476 |

Averages over 10,000 simulations. See section IV and the online appendix for details.



The third and fourth Monte Carlo simulations use computer-generated data only. The purpose of these simulations is to demonstrate how the bias of endogenous stratification estimators changes with changes in the sample size and the number of covariates. The data-generating model for the third simulation is

$$y_i = 1 + \sum_{l=1}^{40} z_{li} + v_i$$

for $i = 1, \dots, N$, where the variables z_{li} have independent Standard Normal distributions, and the variable v_i has an independent Normal distribution with variance equal to 60. As a result, the unconditional variance of y_i is equal to 100. In each Monte Carlo simulation, the sample is divided at random into two equally sized treated and control groups. Predicted outcomes are computed using data for the control group to estimate

$$y_i = \alpha + \mathbf{x}'_{Ki} \boldsymbol{\beta}_K + u_{Ki}$$

by least squares, where \mathbf{x}_{Ki} is the $(K \times 1)$ -vector $(z_{1i}, \dots, z_{Ki})'$, for $K \leq 40$. That is, \mathbf{x}_{Ki} contains the values of the first K regressors in z_{1i}, \dots, z_{40i} . The data-generating process implies that α is equal to 1, $\boldsymbol{\beta}_K$ is a $(K \times 1)$ -vector of 1s, $u_{Ki} = z_{K+1i} + \dots + z_{40i} + v_i$ if $K < 40$, and $u_{40i} = v_i$. We run Monte Carlo simulations for sample sizes $N = 200$, $N = 1,000$, and $N = 5,000$, and numbers of included regressors $K = 10$, $K = 20$, and $K = 40$.

The results are reported in table 5. To economize space, we omit results on the single split sample estimator and report bias results only. Coverage rate and root-MSE results are available on request. The magnitudes of the biases in table 5 are easily understood when compared to the standard deviation of the outcome, which is equal to 10. As expected, the bias of the full-sample endogenous stratification estimator is particularly severe when the sample size is small or the number of included regressors is large, because in both cases, significant overfitting may occur. The increase in bias resulting from increasing the number of regressors is particularly severe when the sample size is small, $N = 200$. The biases of the leave-one-out, repeated split sample, and unfeasible estimators are negligible in most cases and consistently smaller than the bias of the full-sample endogenous stratification estimators, although the leave-one-out estimator tends to produce larger biases than the repeated split sample and unfeasible estimators.

The bias of the full-sample endogenous stratification estimator increases with K in spite of the fact that as K increases, each additional included regressor has the same explanatory power as each of the regressors included in simulations with smaller K . Our final simulation studies a setting where each additional included regressor has lower explanatory power than the previously included ones. Consider

$$y_i = 1 + \sum_{l=1}^{40} \rho^{l-1} z_{li} + \tilde{v}_i,$$

where the variables z_{li} have independent Standard Normal distributions as before, and the variable \tilde{v}_i has an independent Normal distribution with a variance such that the variance of y_i is equal to 100. Table 6 reports the biases of the endogenous stratification estimators across Monte Carlo simulations under the new data-generating process (with $\rho = 0.80$). The biases of the full sample endogenous stratification estimator are larger than in the previous simulation. Their magnitudes increase faster than in the previous simulation when the number of included covariates increases and decrease more slowly than in the previous simulation when the number of observations increases. Relative to the full sample estimator, the biases of the leave-one-out, repeated split sample, and unfeasible estimators are smaller and less sensitive to changes in the number of included covariates and sample size.

TABLE 5.—BIAS IN SIMULATIONS USING ARTIFICIAL DATA (CONSTANT REGRESSION COEFFICIENTS)

| | K = 10 | | | | | | | | | K = 20 | | | | | | | | | K = 40 | | | | | | | | | | | |
|----------------------|------------|--------|-------|----------|--------|-------|------------|--------|-------|----------|--------|-------|------------|--------|-------|----------|--------|-------|------------|--------|-------|----------|--------|-------|-------|-------|-------|-------|-------|-------|
| | Unadjusted | | | Adjusted | | | Unadjusted | | | Adjusted | | | Unadjusted | | | Adjusted | | | Unadjusted | | | Adjusted | | | | | | | | |
| | Low | Medium | High | Low | Medium | High | Low | Medium | High | Low | Medium | High | Low | Medium | High | Low | Medium | High | Low | Medium | High | Low | Medium | High | | | | | | |
| $N = 200$ | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | |
| $\hat{\tau}_k$ | 2.24 | -0.05 | -2.28 | 2.23 | -0.06 | -2.27 | 3.15 | -0.00 | -3.13 | 3.12 | -0.01 | -3.08 | 4.09 | 0.02 | -4.06 | 3.89 | 0.04 | -3.88 | 4.09 | 0.02 | -4.06 | 3.89 | 0.04 | -3.88 | 4.09 | 0.02 | -4.06 | 3.89 | 0.04 | -3.88 |
| $\hat{\tau}_k^{LOO}$ | -0.32 | -0.03 | 0.25 | -0.07 | -0.04 | 0.01 | -0.18 | 0.00 | 0.17 | 0.30 | 0.01 | -0.26 | -0.09 | 0.00 | 0.13 | 0.89 | 0.02 | -0.88 | -0.09 | 0.00 | 0.13 | 0.89 | 0.02 | -0.88 | -0.09 | 0.00 | 0.13 | 0.89 | 0.02 | -0.88 |
| $\hat{\tau}_k^{RSS}$ | -0.03 | -0.04 | -0.04 | -0.02 | -0.04 | -0.05 | 0.01 | -0.00 | -0.01 | 0.03 | 0.01 | 0.01 | 0.03 | 0.01 | -0.00 | 0.04 | 0.02 | -0.01 | 0.03 | 0.01 | -0.00 | 0.04 | 0.02 | -0.01 | 0.03 | 0.01 | -0.00 | 0.04 | 0.02 | -0.01 |
| $\hat{\tau}_k^{UNF}$ | -0.04 | -0.04 | -0.03 | -0.02 | -0.03 | -0.04 | 0.03 | 0.00 | 0.00 | 0.02 | -0.01 | 0.03 | 0.02 | -0.01 | 0.01 | -0.01 | 0.05 | 0.04 | 0.02 | -0.01 | 0.01 | -0.01 | 0.05 | 0.04 | 0.02 | -0.01 | 0.01 | -0.01 | 0.05 | 0.04 |
| $N = 1,000$ | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | |
| $\hat{\tau}_k$ | 0.55 | 0.01 | -0.54 | 0.55 | 0.01 | -0.54 | 0.71 | -0.01 | -0.71 | 0.70 | -0.00 | -0.70 | 0.82 | 0.01 | -0.82 | 0.81 | 0.01 | -0.82 | 0.82 | 0.01 | -0.82 | 0.81 | 0.01 | -0.82 | 0.82 | 0.01 | -0.82 | 0.81 | 0.01 | -0.82 |
| $\hat{\tau}_k^{LOO}$ | -0.05 | 0.01 | 0.07 | -0.05 | 0.01 | 0.06 | -0.05 | -0.01 | 0.05 | -0.05 | -0.00 | 0.05 | -0.01 | -0.00 | 0.01 | -0.00 | -0.01 | -0.01 | -0.01 | -0.00 | 0.01 | -0.00 | -0.01 | -0.01 | -0.00 | 0.01 | -0.00 | -0.01 | -0.01 | -0.00 |
| $\hat{\tau}_k^{RSS}$ | 0.02 | 0.00 | -0.00 | 0.02 | 0.01 | -0.00 | -0.01 | -0.01 | 0.01 | -0.01 | -0.00 | 0.01 | -0.01 | -0.00 | -0.00 | 0.01 | 0.01 | -0.01 | -0.00 | 0.00 | 0.01 | -0.00 | 0.01 | -0.00 | -0.00 | 0.00 | -0.01 | 0.01 | -0.00 | -0.00 |
| $\hat{\tau}_k^{UNF}$ | 0.01 | 0.01 | 0.00 | 0.00 | 0.01 | 0.00 | -0.00 | -0.01 | -0.01 | -0.00 | -0.00 | -0.01 | -0.00 | -0.00 | 0.00 | -0.00 | -0.00 | -0.01 | -0.00 | 0.00 | 0.00 | -0.00 | 0.01 | -0.00 | -0.00 | 0.00 | -0.01 | 0.01 | -0.00 | -0.00 |
| $N = 5,000$ | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | |
| $\hat{\tau}_k$ | 0.12 | -0.00 | -0.12 | 0.12 | -0.00 | -0.12 | 0.15 | 0.00 | -0.16 | 0.15 | 0.00 | -0.15 | 0.16 | -0.00 | -0.16 | 0.16 | -0.00 | -0.16 | 0.16 | -0.00 | -0.16 | 0.16 | -0.00 | -0.16 | 0.16 | -0.00 | -0.16 | 0.16 | -0.00 | -0.16 |
| $\hat{\tau}_k^{LOO}$ | -0.01 | -0.00 | 0.01 | -0.01 | -0.00 | 0.01 | -0.01 | 0.00 | 0.00 | -0.00 | 0.00 | 0.00 | -0.01 | -0.00 | 0.01 | -0.01 | -0.00 | -0.00 | -0.00 | -0.00 | 0.00 | -0.00 | -0.00 | -0.00 | -0.00 | 0.00 | -0.00 | -0.00 | -0.00 | -0.00 |
| $\hat{\tau}_k^{RSS}$ | 0.01 | -0.00 | -0.01 | 0.01 | -0.00 | -0.01 | 0.00 | -0.00 | -0.01 | 0.01 | -0.00 | -0.01 | 0.00 | -0.00 | -0.01 | 0.00 | -0.00 | -0.00 | -0.00 | -0.00 | 0.00 | -0.00 | -0.00 | -0.00 | -0.00 | 0.00 | -0.00 | -0.00 | -0.00 | -0.00 |
| $\hat{\tau}_k^{UNF}$ | 0.01 | -0.00 | -0.01 | 0.01 | -0.00 | -0.01 | 0.01 | -0.01 | -0.01 | 0.01 | -0.00 | -0.01 | 0.01 | -0.00 | -0.01 | 0.01 | -0.00 | -0.00 | -0.00 | -0.00 | 0.00 | -0.00 | -0.00 | -0.00 | -0.00 | 0.00 | -0.00 | -0.00 | -0.00 | -0.00 |

Averages over 10,000 simulations. K is the number of included regressors. See section IV and the online appendix for details.

TABLE 6.—BIAS IN SIMULATIONS USING ARTIFICIAL DATA (DECAYING REGRESSION COEFFICIENTS)

| | K = 10 | | | | | | K = 20 | | | | | | K = 40 | | | | | |
|----------------------|------------|-------|----------|-------|------------|-------|----------|-------|------------|-------|----------|-------|------------|-------|----------|-------|--|--|
| | Unadjusted | | Adjusted | | Unadjusted | | Adjusted | | Unadjusted | | Adjusted | | Unadjusted | | Adjusted | | | |
| | Low | High | Low | High | Low | High | Low | High | Low | High | Low | High | Low | High | Low | High | | |
| N = 200 | | | | | | | | | | | | | | | | | | |
| $\hat{\tau}_k$ | 2.94 | -2.97 | 2.97 | -2.99 | 4.57 | -4.60 | 4.63 | -4.69 | 6.86 | -6.87 | 6.86 | -6.87 | 7.20 | -7.24 | 7.20 | -7.24 | | |
| $\hat{\tau}_k^{LOO}$ | -0.41 | 0.40 | 0.01 | -0.02 | -0.26 | 0.23 | 0.76 | -0.79 | -0.14 | 0.17 | -0.14 | 0.17 | 2.48 | -2.48 | -0.00 | -0.00 | | |
| $\hat{\tau}_k^{RSS}$ | -0.03 | 0.01 | -0.02 | 0.01 | -0.01 | -0.01 | -0.02 | -0.02 | 0.00 | 0.01 | 0.00 | 0.01 | -0.02 | 0.01 | 0.01 | 0.02 | | |
| $\hat{\tau}_k^{UNF}$ | -0.01 | -0.01 | -0.02 | 0.00 | -0.01 | 0.02 | 0.00 | -0.01 | -0.03 | 0.04 | -0.03 | 0.04 | -0.03 | 0.02 | -0.03 | 0.01 | | |
| N = 1,000 | | | | | | | | | | | | | | | | | | |
| $\hat{\tau}_k$ | 0.95 | -0.94 | 0.95 | -0.94 | 1.64 | -1.64 | 1.64 | -1.64 | 2.64 | -2.64 | 2.64 | -2.64 | 2.66 | -2.65 | 2.66 | -2.65 | | |
| $\hat{\tau}_k^{LOO}$ | -0.11 | 0.12 | -0.10 | 0.11 | -0.10 | 0.11 | -0.03 | 0.04 | -0.07 | 0.08 | -0.07 | 0.08 | 0.14 | -0.13 | 0.00 | -0.01 | | |
| $\hat{\tau}_k^{RSS}$ | 0.01 | 0.01 | 0.01 | 0.01 | -0.00 | 0.00 | -0.00 | 0.01 | 0.00 | 0.01 | 0.00 | 0.01 | -0.00 | 0.01 | -0.00 | 0.01 | | |
| $\hat{\tau}_k^{UNF}$ | 0.01 | 0.01 | 0.01 | 0.01 | 0.00 | -0.00 | -0.00 | -0.01 | -0.01 | 0.00 | -0.01 | 0.00 | -0.01 | 0.03 | -0.01 | 0.03 | | |
| N = 5,000 | | | | | | | | | | | | | | | | | | |
| $\hat{\tau}_k$ | 0.23 | -0.22 | 0.23 | -0.22 | 0.43 | -0.43 | 0.43 | -0.43 | 0.81 | -0.81 | 0.81 | -0.81 | 0.81 | -0.81 | 0.81 | -0.81 | | |
| $\hat{\tau}_k^{LOO}$ | -0.02 | 0.04 | -0.02 | 0.03 | -0.03 | 0.03 | -0.03 | 0.03 | -0.03 | 0.03 | -0.03 | 0.03 | -0.02 | 0.02 | -0.01 | 0.02 | | |
| $\hat{\tau}_k^{RSS}$ | 0.00 | 0.00 | 0.00 | 0.00 | -0.00 | 0.00 | -0.00 | 0.00 | -0.00 | 0.00 | -0.00 | 0.00 | -0.00 | 0.00 | -0.00 | -0.00 | | |
| $\hat{\tau}_k^{UNF}$ | 0.00 | 0.01 | 0.00 | 0.01 | 0.00 | 0.01 | 0.00 | 0.00 | -0.01 | 0.00 | -0.01 | 0.00 | -0.00 | 0.00 | -0.00 | 0.00 | | |

Averages over 10,000 simulations. K is the number of included regressors. See section IV and the online appendix for details.

Overall, among the estimators that address the overfitting problem of full sample endogenous stratification, the repeated split sample estimator outperforms leave-one-out in the simulations. Moreover, in additional simulations (not reported for brevity), we found that the leave-one-out estimator can behave erratically in settings where the regressors take on only a few values and the variance of the regression error, $e_i = y_i - \mathbf{x}_i'\beta$ is large. This is the case, for example, in the STAR-based simulation if school indicators are excluded from the vector \mathbf{x}_i . In that case, \mathbf{x}_i includes only three indicator variables, for race, gender, and eligibility for a free lunch program. As a result, $\mathbf{x}_i'\beta$ takes on only eight different values. In this setting, we found that overfitting is not an issue and that the full sample and repeated split sample estimators produce small biases in simulations. However, the leave-one-out estimator generates extremely large biases, in fact, larger than the biases of the full-sample endogenous stratification estimator. The reason is that, in this setting, choosing c_1 and c_2 to be the quantiles to be one-third and two-thirds of the distribution of the predicted outcomes result in a large number of observations being located exactly at the boundaries of the values of $\mathbf{x}_i'\beta$ that define the predicted outcome groups. To be concrete, consider the untreated observations with $\mathbf{x}_i'\beta = c_1$. These observations are classified by the full-sample endogenous stratification estimator as members of the low predicted-outcome group. However, it is easy to see that if $\mathbf{x}_i'\beta = c_1$, then $\mathbf{x}_i'\beta_{(-i)} > c_1$ if $y_i < c_1$ and $\mathbf{x}_i'\beta_{(-i)} \leq c_1$ if $y_i \geq c_1$, which induces biases in the leave-one-out estimator because observations are reclassified into brackets of predicted outcomes without treatment on the basis of the value of their outcome variable. The reclassification pattern induced by the leave-one-out estimator creates a bias that is of the opposite sign to the bias of the full sample endogenous stratification estimator.¹⁷ The single split sample estimator has low bias and produces close-to-nominal coverage rates, but also large dispersion induced by the reduction in sample size. The increased variance of the single split sample estimator can make root-MSE of this estimator larger than the root-MSE of the full-sample endogenous stratification estimator (see table 3). All in all, the repeated split sample estimator displays the best performance in our simulations. It has low bias and accurate coverage rates, and it outperforms alternative estimators in terms of root-MSE.

V. Conclusion

In this paper, we have argued that the increasingly popular practice of stratifying experimental units on the basis of a prediction of the outcome without treatment estimated using full sample data from the control group leads to substantially biased estimates of treatment effects. We illustrate the magnitude of this bias using data from two well-known

¹⁷A “faint signal” of this bias is discernible in tables 3 to 6, in settings where the high dimensionality of the regressors makes the magnitude of this bias small.

social experiments: the National JTPA Study and the Project STAR class-size experiment. The full sample endogenous stratification approach is most problematic in studies with small sample sizes and many regressors, where the predictor of the outcome without treatment may be severely overfitted in the control sample. We demonstrate that in realistic scenarios, alternative endogenous stratification estimators based on leave-one-out and, especially, repeated split sample techniques display substantially improved small-sample behavior relative to the full-sample endogenous stratification estimator. In high-dimensional settings, the repeated split sample estimator dominates the leave-one-out estimator. Moreover, the repeated split sample estimator displays robust behavior in low-dimensional settings where overfitting is not a problem but the behavior of the leave-one-out estimator may be catastrophically bad.

Some questions remain open to future research. First, the methods described in this paper do not exhaust the possible approaches to the bias of the fullsample endogenous stratification estimator. Bootstrap/jackknife bias corrections on $\hat{\tau}_k$ and shrinkage estimation of β are potentially fruitful approaches that we are starting to explore. Another question of interest is whether the good small sample behavior of the repeated split sample estimator generalizes to other settings, like the two-step generalized method of moments setting analyzed by Altonji and Segal (1996). Finally, it would be useful to study the inferential properties of endogenous stratification estimators in settings where the first-step prediction problem is carried out using nonparametric/machine learning techniques.

REFERENCES

- Abadie, Alberto, Joshua D. Angrist, and Guido W. Imbens, "Instrumental Variables Estimates of the Effect of Subsidized Training on the Quantiles of Trainee Earnings," *Econometrica* 70 (2002), 91–117.
- Altonji, Joseph G., and Lewis M. Segal, "Small-Sample Bias in GMM Estimation of Covariance Structures," *Journal of Business and Economic Statistics* 14 (1996), 353–366.
- Angrist, Joshua D., and Victor Lavy, "The Effects of High Stakes High School Achievement Awards: Evidence from a Randomized Trial," *American Economic Review* 99 (2009), 1384–1414.
- Angrist, Joshua D., and Jorn-Steffen Pischke, "The Credibility Revolution in Empirical Economics: How Better Research Design Is Taking the Con Out of Econometrics," *Journal of Economic Perspectives* 24 (2010), 3–30.
- Billingsley, Patrick, *Probability and measure*, 3rd ed. (New York: Wiley, 1995).
- Bitler, Marianne P., Jonah B. Gelbach, and Hilary W. Hoynes, "Can Variation in Subgroups' Average Treatment Effects Explain Treatment Effect Heterogeneity? Evidence from a Social Experiment," NBER working paper 20142 (2013).
- Bloom, Howard S., Larry L. Orr, Stephen H. Bell, George Cave, Fred Doolittle, Winston Lin, and Johannes M. Bos, "The Benefits and Costs of JTPA Title II-A Programs," *Journal of Human Resources* 32 (1997), 549–576.
- Crump, Richard K., V. Joseph Hotz, Guido W. Imbens, and Oscar A. Mitnik, "Nonparametric Tests for Treatment Effect Heterogeneity," *THE REVIEW* 90 (2008), 389–405.
- , "Dealing with Limited Overlap in Estimation of Average Treatment Effects," *Biometrika* 96 (2009), 187–199.
- Djebbari, Habiba, and Jeffrey Smith, "Heterogeneous Impacts in PROGRESA," *Journal of Econometrics* 145 (2008), 64–80.
- Dynarski, Susan, Joshua M. Hyman, and Diane Schanzenbach, "Experimental Evidence on the Effect of Childhood Investments on Post-Secondary Attainment and Degree Completion," NBER working paper 17533 (2011).
- , "Experimental Evidence on the Effect of Childhood Investments on Postsecondary Attainment and Degree Completion," *Journal of Policy Analysis and Management* 32 (2013), 692–717.
- Giné, Xavier, Jessica Goldberg, and Dean Yang, "Credit Market Consequences of Improved Personal Identification," *American Economic Review* 102 (2012), 2923–2954.
- Goldrick-Rab, Sara, Douglas N. Harris, James Benson, and Robert Kelchen, "Conditional Cash Transfers and College Persistence: Evidence from a Randomized Need-Based Grant Program," Institute for Research on Poverty discussion paper 1393-11 (2011).
- , "Need-Based Financial Aid and College Persistence: Experimental Evidence from Wisconsin" (2012). SSRN 1887826.
- Hansen, Ben B., "The Prognostic Analogue of the Propensity Score," *Biometrika* 95 (2008), 481–488.
- Harris, Douglas N., and Sara Goldrick-Rab, "Improving the Productivity of Education Experiments: Lessons from a Randomized Study of Need-Based Financial Aid," *Education Finance and Policy* 7 (2012), 143–169.
- Harvill, Eleanor L., Laura R. Peck, and Stephen H. Bell, "On Overfitting in Analysis of Symmetrically Predicted Endogenous Subgroups from Randomized Experimental Samples: Part Three of a Method Note in Three Parts," *American Journal of Evaluation* 34 (2013), 545–566.
- Hausman, Jerry A., and David A. Wise, "Social Experimentation, Truncated Distributions, and Efficient Estimation," *Econometrica* 45 (1977), 919–938.
- , "Stratification on Endogenous Variables and Estimation: The Gary Income Maintenance Experiment," in C. F. Manski and D. L. McFadden, eds., *Structural Analysis of Discrete Data with Econometric Applications* (Cambridge, MA: MIT Press, 1981).
- Heckman, James J., Jeffrey Smith, and Nancy Clements, "Making the Most Out of Programme Evaluations and Social Experiments: Accounting for Heterogeneity in Programme Impacts," *Review of Economic Studies* 64 (1997), 487–535.
- Hemelt, Steven W., Kimberly B. Roth, and William W. Eaton, "Childhood Educational Interventions: Experimental Evidence on Post-Secondary Outcomes," paper presented at the Association for Education Finance and Policy Conference, Boston (2012).
- , "Elementary School Interventions: Experimental Evidence on Post-Secondary Outcomes," *Educational Evaluation and Policy Analysis* 35 (2013), 413–436.
- Kent, David M., and Rodney A. Hayward, "Limitations of Applying Summary Results of Clinical Trials to Individual Patients: The Need for Risk Stratification," *Journal of the American Medical Association* 298 (2007), 1209–1212.
- Krueger, Alan B., "Experimental Estimates of Education Production Functions," *Quarterly Journal of Economics* 114 (1999), 497–532.
- Murnane, Richard J., and John B. Willett, *Methods Matter: Improving Causal Inference in Educational and Social Science Research* (Oxford: Oxford University Press, 2011).
- Orr, Larry L., Howard S. Bloom, Stephen H. Bell, Fred Doolittle, Winston Lin, and George Cave, *Does Training for the Disadvantaged Work?* (Washington, DC: Urban Institute, 1996).
- Pane, John F., Beth Ann Griffin, Daniel F. McCaffrey, and Rita Karam, "Effectiveness of Cognitive Tutor Algebra I at Scale," *Educational Evaluation and Policy Analysis* 36 (2014), 127–144.
- Peck, Laura R., "Subgroup Analysis in Social Experiments: Measuring Program Impacts Based on Post-Treatment Choice," *American Journal of Evaluation* 24 (2003), 157–187.
- , "On Analysis of Symmetrically Predicted Endogenous Subgroups: Part One of a Method Note in Three Parts," *American Journal of Evaluation* 34 (2013), 225–236.
- Rodriguez-Planas, Nuria, "Schools and Drugs: Closing the Gap: Evidence from a Randomized Trial in the US," IZA discussion paper 6770 (2012).
- , "Schools and Drugs: Closing the Gap: Evidence from a Randomized Trial in the US," City University of New York working paper (2014).
- Sanbonmatsu, Lisa, Jeffrey R. Kling, Greg J. Duncan, and Jeanne Brooks-Gunn, "Neighborhoods and Academic Achievement: Results from the Moving to Opportunity Experiment," *Journal of Human Resources* 41 (2006), 649–691.

APPENDIX

Over-Fitting in a Regression Model

Consider the classical regression model,

$$y_i = \mathbf{x}'_i \boldsymbol{\beta} + u_i,$$

for $i = 1, \dots, N$, where $\{\mathbf{x}_1, \dots, \mathbf{x}_N\}$ are fixed regressors such that

$$\sum_{j=1}^N \mathbf{x}_j \mathbf{x}'_j$$

is nonsingular and $\{u_1, \dots, u_N\}$ are mutually independent mean-zero error terms independent of $\{\mathbf{x}_1, \dots, \mathbf{x}_N\}$. In particular, for $i, j = 1, \dots, N$ and $i \neq j$, we have $E[u_j | u_i] = E[u_j] = 0$. Notice that

$$\begin{aligned} \mathbf{x}'_i (\widehat{\boldsymbol{\beta}} - \boldsymbol{\beta}) &= \mathbf{x}'_i \left(\sum_{j=1}^N \mathbf{x}_j \mathbf{x}'_j \right)^{-1} \sum_{j=1}^N \mathbf{x}_j u_j \\ &= h_{ii} u_i + \sum_{j \neq i} h_{ij} u_j, \end{aligned}$$

where

$$h_{ij} = \mathbf{x}'_i \left(\sum_{j=1}^N \mathbf{x}_j \mathbf{x}'_j \right)^{-1} \mathbf{x}_j.$$

Therefore, we obtain

$$E[\mathbf{x}'_i (\widehat{\boldsymbol{\beta}} - \boldsymbol{\beta}) | u_i = u] = h_{ii} u.$$

Notice that $h_{ii} > 0$ as long as $\|\mathbf{x}_i\| > 0$. So $E[\mathbf{x}'_i (\widehat{\boldsymbol{\beta}} - \boldsymbol{\beta}) | u_i = u]$ is equal to a positive constant times u . That is, conditional on $u_i < 0$, the expectation of $\mathbf{x}'_i \widehat{\boldsymbol{\beta}} - \mathbf{x}'_i \boldsymbol{\beta}$ is negative, and conditional on $u_i > 0$, the expectation of $\mathbf{x}'_i \widehat{\boldsymbol{\beta}} - \mathbf{x}'_i \boldsymbol{\beta}$ is positive.

To simplify notation, in this appendix, we did not distinguish between treated and control units. Notice, however, that in the context of the endogenous stratification estimators studied in this paper, the regression of the outcome, y_i , on the covariates, \mathbf{x}_i , is estimated using the sample of experimental controls only.