

Training, Wages, and Sample Selection: Estimating Sharp Bounds on Treatment Effects

DAVID S. LEE

Princeton University and NBER

First version received October 2005; final version accepted August 2008 (Eds.)

This paper empirically assesses the wage effects of the Job Corps program, one of the largest federally funded job training programs in the U.S. Even with the aid of a randomized experiment, the impact of a training program on wages is difficult to study because of sample selection, a pervasive problem in applied microeconomic research. Wage rates are only observed for those who are employed, and employment status itself may be affected by the training program. This paper develops an intuitive trimming procedure for bounding average treatment effects in the presence of sample selection. In contrast to existing methods, the procedure requires neither exclusion restrictions nor a bounded support for the outcome of interest. Identification results, estimators, and their asymptotic distribution are presented. The bounds suggest that the program raised wages, consistent with the notion that the Job Corps raises earnings by increasing human capital, rather than solely through encouraging work. The estimator is generally applicable to typical treatment evaluation problems in which there is nonrandom sample selection/attrition.

1. INTRODUCTION

For decades, many countries around the world have administered government-sponsored employment and training programs, designed to help improve the labour market outcomes of the unemployed or economically disadvantaged.¹ To do so, these programs offer a number of different services, ranging from basic classroom education and vocational training to various forms of job search assistance. The key question of interest to policymakers is whether or not these programs are actually effective, sufficiently so to justify the cost to the public. The evaluation of these programs has been the focus of a large substantive and methodological literature in economics. Indeed, Heckman, LaLonde and Smith (1999) observe that “[f]ew U.S. government programs have received such intensive scrutiny, and been subject to so many different types of evaluation methodologies, as governmentally-supplied job training”.

Econometric evaluations of these programs typically focus on their reduced-form impacts on total earnings, a first-order issue for cost–benefit analysis. Unfortunately, exclusively studying the effect on total earnings leaves open the question of whether any earnings gains are achieved through raising individuals’ *wage rates* (price effects) or hours of work (quantity effects). That is, a training program may lead to a meaningful increase in human capital, thus raising participants’ wages. Alternatively, the program may have a pure labour supply effect: through career counselling and encouragement of individuals to enter the labour force, a training program may simply be raising incomes by increasing the likelihood of employment, without any increase in wage rates.

1. See Heckman *et al.* (1999) for figures on expenditures on active labour market programs in OECD countries. See also Martin (2000).

But assessing the impact of training programs on wage rates is not straightforward, due to the well-known problem of sample selection, which is pervasive in applied microeconomic research. That is, wages are only observed for individuals who are employed. Thus, even if there is random assignment of the “treatment” of a training program, there may not only be an effect on wages but also on the probability that a person’s wage will even be observed. Even a randomized experiment cannot guarantee that treatment and control individuals will be comparable *conditional on being employed*. Indeed, standard labour supply theory predicts that wages will be correlated with the likelihood of employment, resulting in sample selection bias (Heckman, 1974). This missing data problem is especially relevant for analysing public job training programs, which typically target individuals who have low employment probabilities.

This paper empirically assesses the *wage* effects of the Job Corps program, one of the largest federally funded job training programs in the U.S.² The Job Corps is a comprehensive program for economically disadvantaged youth aged 16–24 years and is quite intensive: the typical participant will live at a local Job Corps centre, receiving room, board, and health services while enrolled, for an average of about 8 months. During the stay, the individual can expect to receive about 1100 hours of vocational and academic instruction, equivalent to about 1 year in high school. The Job Corps is also expensive, with the average cost at about \$14,000 per participant.³ This paper uses data from the National Job Corps Study, a randomized evaluation funded by the U.S. Department of Labor.

Standard parametric or semiparametric methods for correcting for sample selection require exclusion restrictions that have little justification in this case. As shown below, the data include numerous baseline variables, but all those that are found to be related to employment probabilities (*i.e.*, sample selection) could also potentially have a direct impact on wage rates.

Thus, this paper develops an alternative method, a general procedure for bounding the treatment effects. The method amounts to first identifying the excess number of individuals who are induced to be selected (employed) because of the treatment and then “trimming” the upper and lower tails of the outcome (*e.g.*, wage) distribution by this number, yielding a worst-case scenario bound. The assumptions for identifying the bounds are already assumed in conventional models for sample selection: (1) the regressor of interest is independent of the errors in the outcome and selection equations and (2) the selection equation can be written as a standard latent variable binary response model. In the case of an experiment, random assignment ensures that the first assumption holds. It is proven that the trimming procedure yields the tightest bounds for the average treatment effect that are consistent with the observed data. No exclusion restrictions are required, nor is a bounded support for the outcome variable.

An estimator for the bounds is introduced and shown to be \sqrt{n} consistent and asymptotically normal with an intuitive expression for its asymptotic variance. It not only depends on the variance of the trimmed outcome variable but also on the trimming threshold, which is an estimated quantile. There is also an added term that accounts for the estimation of *which* quantile (the 10th, 11th, 12th, etc. percentile) of the distribution to use as the trimming threshold.

For the analysis of Job Corps, the trimming procedure is instrumental to measuring the wage effects, producing bounds that are somewhat narrow. For example, at week 90 after random assignment, the estimated interval for the treatment effect is 4.2–4.3%, even when wages

2. In the 2004 fiscal year, the U.S. Department of Labor’s Employment and Training Administration spent \$1.54 billion for the operation of the Job Corps. By comparison, it spent about \$893 million on “Adult Employment and Training Activities” (job search assistance for anyone and job training available to anyone if such training is needed for obtaining or retaining employment) and about \$1.44 billion on “Dislocated Workers Employment and Training Activities” (employment and training services for unemployment and underemployed workers) (U.S. Department of Labor, 2005a).

3. A summary of services provided and costs can be found in Burghardt, Schochet, McConnell, Johnson, Gritz, Glazerman, Homrighausen and Jackson (2001).

are missing for about 54% of individuals. By the end of the 4-year follow-up period, the interval is still somewhat informative, statistically rejecting effects more negative than -3.7% and more positive than 11.2% . By comparison, the assumption-free, “worst-case scenario” bounds proposed by Horowitz and Manski (2000a) produce a lower bound of -75% effect and an upper bound of 80% .

Adjusting for the reduction in potential work experience likely caused by the program, the evidence presented here points to a positive causal effect of the program on wage rates. This is consistent with the view that the Job Corps program represents a human capital investment rather than a means to improve earnings through raising work effort alone.

The proposed trimming procedure is neither specific to this application nor to randomized experiments. It will generally be applicable to treatment evaluation problems when outcomes are missing, a problem that often arises in applied research. Reasons for missing outcomes range from survey nonresponse (e.g., students not taking tests), to sample attrition (e.g., inability to follow individuals over time), to other structural reasons (e.g., mortality). Generally, this estimator is well suited for cases where the researcher is uncomfortable imposing exclusion restrictions in the standard two-equation sample selection model and when the support of the outcome variable is too wide to yield informative bounds on treatment effects.

This paper is organized as follows. It begins, in Section 2, with a description of the Job Corps program, the randomized experiment, and the nature of the sample selection problem. After this initial analysis, the proposed bounding procedure is described in Sections 3 and 4. Section 3 presents the identification results, while Section 4 introduces a consistent and asymptotically normal estimator of the bounds and discusses inference. Section 5 reports the results from the empirical analysis of the Job Corps. Section 6 concludes.

2. THE NATIONAL JOB CORPS STUDY AND SAMPLE SELECTION

This section describes both the Job Corps program and the data used for the analysis, replicates the main earnings results of the recently completed randomized evaluation, and illustrates the nature of the sample selection problem. It is argued below that standard sample selection correction procedures are not appropriate for this context. Also, to provide an initial benchmark, the approach of Horowitz and Manski (2000a) is used to provide bounds on the Job Corps’ effect on wages. They are to be compared with the “trimming” bounds presented in Section 5, which implements the estimator developed in Sections 3 and 4.

2.1. *The Job Corps program and the randomized experiment*

The U.S. Department of Labor describes the Job Corps program today as “a no-cost education and vocational training program ... that helps young people ages 16 through 24 get a better job, make more money and take control of their lives” (U.S. Department of Labor, 2005b). To be eligible, an individual must be a legal resident of the U.S., be between the ages of 16 and 24, and come from a low-income household (Schochet, Burghardt and Glazerman, 2001). The administration of the Job Corps is considered to be somewhat uniform across the 110 local Job Corps centres in the U.S.

Perhaps the most distinctive feature of the program is that most participants live at the local Job Corps centre while enrolled. This residential component of the program includes formal social skills training, meals, and a dormitory-style life. During the stay, with the help of counsellors, the participants develop individualized, self-paced programs, which will consist of a combination of remedial high school education, including consumer and driver education, as well as vocational training in a number of areas, including clerical work, carpentry, automotive repair,

building and apartment maintenance, and health-related work. On average, enrollees can expect to receive about 440 hours of academic instruction and about 700 hours of vocational training, over an average of 30 weeks. Centres also provide health services as well as job search assistance upon the students' exit from the Job Corps.

In the mid-1990's, three decades after the creation of Job Corps, the U.S. Department of Labor funded a randomized evaluation of the program, which was carried out by Mathematica Policy Research, Inc. Persons who applied for the program for the first time between November 1994 and December 1995 and were found to be eligible (80,883 persons) were randomized into a "program" group and a "control" group. The control group of 5977 individuals was essentially embargoed from the program for 3 years, while the remaining applicants could enrol in the Job Corps as usual. Since those who were still eligible after randomization were not compelled to participate, the differences in outcomes between program and control group members represent the reduced-form effect of eligibility or the "intent-to-treat" effect. This treatment effect is the focus of the empirical analysis presented below. Throughout the paper, when I use the phrase "effect of the program", I am referring to this reduced-form treatment effect.

Of the program group, 9409 applicants were randomly selected to be followed for data collection. The research sample of 15,386 individuals was interviewed at random assignment, and at three subsequent points in time: 12, 30, and 48 months after random assignment. Due to programmatic reasons, some subpopulations were randomized into the program group with differing, but known, probabilities. Thus, analysing the data requires the use of the design weights (the variable DSGN_WGT as described in Schochet, Cao, Glazerman, Grady, Gritz, McConnell, Johnson and Burghardt, 2003).

This paper uses the public-release data of the National Job Corps Study. Table 1 provides descriptive statistics for the data used in the analysis below. For baseline as well as postassignment variables, it reports the treatment and control group means, S.D., proportion of the observations with nonmissing values for the specified variable, as well as the difference in the means and associated S.E. The table shows that the proportion nonmissing and the means for the demographic variables (the first 12 rows), education and background variables (the next 4 rows), income at baseline (the next 9 rows), and employment information (the next 6 rows) are quite similar. For only one of the variables—usual weekly hours on the most recent job at the baseline—is the difference (0.91 hours) statistically significant. A logit of the treatment indicator on all baseline characteristics in Table 1 was estimated; the chi-square test of all coefficients equalling zero yielded a p value of 0.577.⁴ The overall comparability between the treatment and the control groups is consistent with successful randomization of the treatment.

It is important to note that the analysis in this paper, abstracts from missing values due to interview nonresponse and sample attrition over time. Thus, only individuals who had nonmissing values for weekly earnings and weekly hours *for every week* after the random assignment are used; the estimation sample is thus somewhat smaller (9145 vs. 15,386). It will become clear below that the trimming procedure could be applied exclusively to the attrition/nonresponse problem, which is a mechanism for sample selection that is quite distinct from the selection into employment status. More intensive data collection can solve the attrition/nonresponse problem but not the problem of sample selection on wages caused by employment. For this reason, the analysis below focuses exclusively on the latter problem and analyses the data conditional on individuals having continuously valid earnings and hours data.⁵

4. Missing values for each of the baseline variables were imputed with the mean of the variable. The analysis below uses this imputed data.

5. Although the analysis here abstracts from the nonresponse problem, there is some evidence that it is a second-order issue, as mentioned in Remark 2 of Subsection 3.1.

TABLE 1
Summary statistics, by treatment status, National Job Corps Study

Variable	Control			Program			Difference	
	Proportion of nonmissing	Mean	S.D.	Proportion of nonmissing	Mean	S.D.	Difference	S.E.
Female	1.00	0.458	0.498	1.00	0.452	0.498	-0.006	0.011
Age at baseline	1.00	18.351	2.101	1.00	18.436	2.159	0.085	0.045
White, non-Hispanic	1.00	0.263	0.440	1.00	0.266	0.442	0.002	0.009
Black, non-Hispanic	1.00	0.000	0.500	1.00	0.493	0.500	0.003	0.011
Hispanic	1.00	0.172	0.377	1.00	0.169	0.375	-0.003	0.008
Other race/ethnicity	1.00	0.074	0.262	1.00	0.072	0.258	-0.002	0.006
Never married	0.98	0.916	0.278	0.98	0.917	0.275	0.002	0.006
Married	0.98	0.023	0.150	0.98	0.020	0.139	-0.003	0.003
Living together	0.98	0.040	0.197	0.98	0.039	0.193	-0.002	0.004
Separated	0.98	0.021	0.144	0.98	0.024	0.154	0.003	0.003
Has child	0.99	0.193	0.395	0.99	0.189	0.392	-0.004	0.008
Number of children	0.99	0.268	0.640	0.99	0.270	0.650	0.002	0.014
Education	0.98	10.105	1.540	0.98	10.114	1.562	0.009	0.033
Mother's education	0.81	11.461	2.589	0.82	11.483	2.562	0.022	0.061
Father's education	0.61	11.540	2.789	0.62	11.394	2.853	-0.146	0.077
Ever arrested	0.98	0.249	0.432	0.98	0.249	0.432	-0.001	0.009
Household income								
<3000	0.65	0.251	0.434	0.63	0.253	0.435	0.002	0.012
3000-6000	0.65	0.208	0.406	0.63	0.206	0.405	-0.002	0.011
6000-9000	0.65	0.114	0.317	0.63	0.117	0.321	0.003	0.008
9000-18,000	0.65	0.245	0.430	0.63	0.245	0.430	0.000	0.011
>18,000	0.65	0.182	0.386	0.63	0.179	0.383	-0.003	0.010
Personal income								
<3000	0.92	0.789	0.408	0.92	0.789	0.408	-0.001	0.009
3000-6000	0.92	0.131	0.337	0.92	0.127	0.334	-0.003	0.007
6000-9000	0.92	0.046	0.209	0.92	0.053	0.223	0.007	0.005
>9000	0.92	0.034	0.181	0.92	0.031	0.174	-0.003	0.004
At baseline								
Have job	0.98	0.192	0.394	0.98	0.198	0.398	0.006	0.009
Months employed, previous year	1.00	3.530	4.238	1.00	3.596	4.249	0.066	0.091
Had job, previous year	0.98	0.627	0.484	0.98	0.635	0.482	0.007	0.010
Earnings, previous year	0.93	2810.482	4435.616	0.94	2906.453	6401.328	95.971	117.097
Usual hours/week	1.00	20.908	20.704	1.00	21.816	21.046	0.908*	0.446
Usual weekly earnings	1.00	102.894	116.465	1.00	110.993	350.613	8.099	5.093
After random assignment								
Week 52 weekly hours	1.00	17.784	23.392	1.00	15.297	22.680	-2.487*	0.495
Week 104 weekly hours	1.00	21.977	26.080	1.00	22.645	26.252	0.668	0.560
Week 156 weekly hours	1.00	23.881	26.151	1.00	25.879	26.574	1.997*	0.563
Week 208 weekly hours	1.00	25.833	26.250	1.00	27.786	25.745	1.953*	0.558
Week 52 weekly earnings	1.00	103.801	159.893	1.00	91.552	149.282	-12.249*	3.335
Week 104 weekly earnings	1.00	150.407	210.241	1.00	157.423	200.266	7.015	4.417
Week 156 weekly earnings	1.00	180.875	224.426	1.00	203.714	239.802	22.839*	4.936
Week 208 weekly earnings	1.00	200.500	230.661	1.00	227.912	250.222	27.412*	5.106
Total earnings (4 years)	1.00	30,007	26,894	1.00	30,800	26,437	794	572
Number of observations	3599			5546				

Notes: $N = 9145$. Computations use design weights. Chi-square test of all coefficients equalling zero, from a logit of the treatment indicator on all baseline characteristics (where mean values were imputed for missing values) yields 24.95; associated p value from a chi-squared (27df) distribution is 0.577.

*Indicates difference is statistically significant from 0 at the 5% (or less) level.

The bottom of Table 1 shows that the only set of variables that show important (and statistically significant) differences between treatment and control are the postassignment labour market outcomes. The treatment group has lower weekly hours and earnings at week 52 but higher hours and earnings at the 3-year and 4-year marks. At week 208, the earnings gain is about \$27, with the control mean of about \$200. This is consistent with Mathematica's final report, which showed that the program had about a 12% positive effect on earnings by the fourth year after enrolment and suggested that lifetime gains in earnings could very well exceed the program's costs (Burghardt *et al.*, 2001). The effect on weekly hours at that time is a statistically significant 1.95 hours.

Figure 1 illustrates the treatment effects on earnings for each week subsequent to random assignment. It shows an initial negative impact on earnings for the first 80 weeks, after which

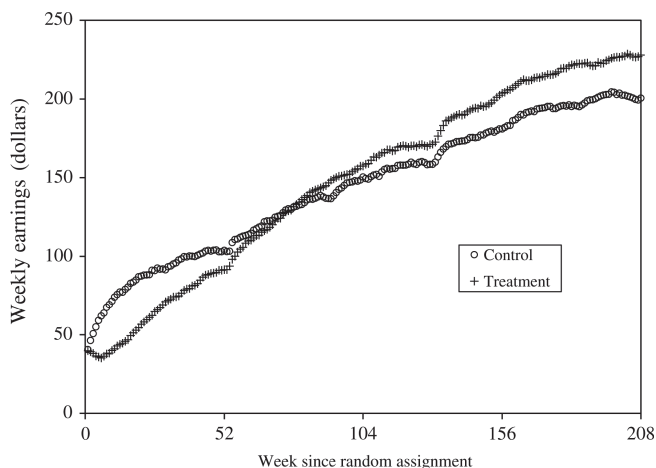


FIGURE 1
Impact of Job Corps on weekly earnings

point a positive treatment effect appears and grows. The estimates in the bottom of Table 1 and plotted in Figure 1 are similar qualitatively and quantitatively to the impact estimates reported in Schochet *et al.* (2001).⁶

2.2. The effect on wages and the sample selection problem

It seems useful to assess the impact of the program on *wage rates*, as distinct from total earnings, which is a product of both the price of labour (the wage) and labour supply (whether the person works, and if so, how many hours). Distinguishing between price and quantity effects is important for better understanding the mechanism through which the Job Corps leads to more favourable labour market outcomes.

On the one hand, one of the goals of the Job Corps is to encourage work and self-sufficiency; thus, participants' total earnings might rise simply because the program succeeds in raising the likelihood that they will be employed, while at the same time leaving the market wage for their labour unaffected. On the other hand, the main component of the Job Corps is significant academic and vocational training, which could be expected to raise wages. There is a great deal of empirical evidence to suggest a positive causal effect of education on wages (see Card, 1999).

Unfortunately, even though the National Job Corps study was a randomized experiment, one cannot use simple treatment–control differences to estimate the effect of the program on wage rates. This is because the effective “prices” of labour for these individuals are only observed to the econometrician when the individuals are employed. This gives rise to the classic sample selection problem (*e.g.*, see Heckman, 1979).

Figure 2 suggests that sample selection may well be a problem for the analysis of wage effects of the Job Corps. It reports employment rates (the proportion of the sample that has positive work hours in the week) for both treated and control individuals, for each week following

6. In Schochet *et al.* (2001), the reported estimates used a less stringent sample criterion. Instead of requiring nonmissing values for 208 consecutive weeks, individuals only needed to complete the 48-month interview (11,313 individuals). Therefore, for that sample, some weeks' data will be missing. Despite the difference in the samples, the levels, impact estimates, and time profile reported in Schochet *et al.* (2001) are also quite similar to those found in Figures 2 and 3 (below).

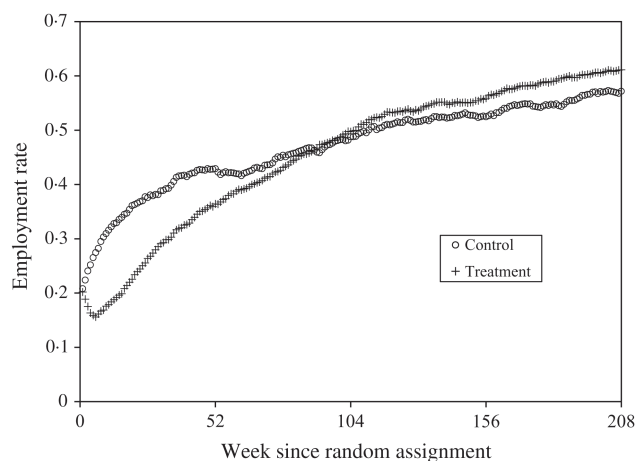


FIGURE 2
Impact of Job Corps on employment rates

random assignment. The results show that the program had a negative impact on employment propensities in the first half of the follow-up and a positive effect in the latter half. This shows that the Job Corps itself affected whether individuals would have a nonmissing wage rate.

Put another way, Figure 2 illustrates that even though proper random assignment will imply that the treatment and control groups are comparable at the baseline, they may well be systematically different *conditional on being employed* in a given period subsequent to the random assignment. As a result, the treatment–control difference in mean log hourly wages, as plotted in Figure 3 (with pointwise 95% confidence intervals), may not represent the true causal effect of the program.⁷

There are two other reasons why sample selection can potentially be important in this case. As shown in Figure 2, a large fraction of individuals are not employed: employment rates start at about 20% and grow to at most 60% at the 4-year mark. Second, nonemployed and employed individuals appear to be systematically different on a number of important observable dimensions. Table 2 reports log-odds coefficients from a logit of employment in week 208 on the treatment dummy and the baseline characteristics listed in Table 1. As might be expected, gender, race, education, criminal history, and employment status at the baseline are all very strong predictors of employment in week 208.

The problem of nonrandom sample selection is well understood in the training literature; it may be one of the reasons why most evaluations of job training programs focus on total earnings, including zeros for those without a job, rather than on wages conditional on employment. Of the 24 studies referenced in a survey of experimental and nonexperimental studies of U.S. employment and training programs (Heckman *et al.*, 1999), most examine annual, quarterly, or monthly earnings without discussing the sample selection problem of examining wage rates.⁸ As for the

7. Hourly wage is computed by dividing weekly earnings by weekly hours worked, for the treatment and control groups. Note the pattern of “kinks” that occur at the 12- and 30-month marks, which is also apparent in Figure 1. This could be caused by the retrospective nature of the interviews that occur at 12-, 30-, and 48-months postrandom assignment. This pattern would be found if there were systematic overestimation of earnings on employment that was further away from the interview date. The lines would “connect” if respondents were reminded of their answer from the previous interview. Note that these potential errors do not seem to be too different between the treatment and the control groups, as there are no obvious kinks in the difference (solid squares).

8. The exceptions include Kiefer (1979), Hollister, Kemper and Maynard (1984), and Barnow (1987). The sources from tables 22 and 24 in Heckman *et al.* (1999) were surveyed.

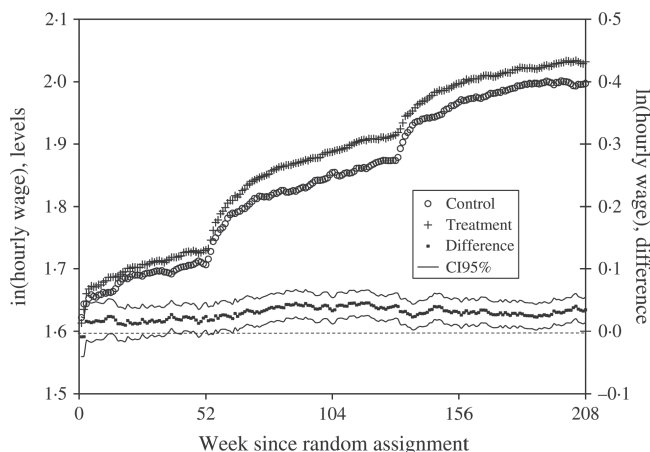


FIGURE 3

Differences in log(hourly wage), conditional on employment

Job Corps, when reporting results on hourly wages for the working, Schochet *et al.* (2001) is careful to note that because of the selection into employment, the treatment–control differences cannot be interpreted as impact estimates.

2.3. Existing approaches

Currently, there are two general approaches to addressing the sample selection problem. The first is to explicitly model the process determining selection. The conventional setup, following Heckman (1979), models the wage determining process as:

$$Y^* = D\beta + X\pi_1 + U \quad (1)$$

$$Z^* = D\gamma + X\pi_2 + V$$

$$Y = 1[Z^* \geq 0] \cdot Y^*,$$

where Y^* is the offered market wage as of a particular point in time (*e.g.*, week 208 after randomization), D is the indicator variable of receiving the treatment of being given access to the Job Corps program, and X is a vector of baseline characteristics. Z^* is a latent variable representing the propensity to be employed. γ represents the causal effect of the treatment on employment propensities, while β is the (constant) treatment effect of interest.⁹ Both Y^* and Z^* are unobserved, but the wage conditional on employment Y is observed, where $1[\cdot]$ is the indicator variable. (U, V) are assumed to be jointly independent of the regressors (D, X) .¹⁰ Within a standard labour supply framework, it is easy to imagine the possibility that job training could raise the market wage for individuals, leading to a positive β , and at the same time raise the probability of participating in the labour force ($\gamma > 0$) since a higher wage will more likely exceed the reservation wage for participating.¹¹

9. In this specification, the treatment effect is constant.

10. This assumption, which is stronger than necessary, is invoked now for expositional purposes. It will be shown below that what is required is instead independence of (U, V) and D , conditional on X .

11. Of course, it should be noted that since the goal here is to estimate a reduced-form treatment effect, we do not adopt a particular labour supply model or prohibit ways in which the treatment could affect participation. For example, γ could be positive if the program's job search assistance component was important.

TABLE 2
Logit of employment in week 208 on baseline characteristics

Variable	Estimate
Treatment status	0.172* (0.046)
Female	-0.253* (0.051)
Age at baseline	0.027 (0.014)
Black, non-Hispanic	-0.471* (0.060)
Hispanic	-0.225* (0.077)
Other race/ethnicity	-0.412* (0.099)
Married	-0.193 (0.175)
Living together	0.106 (0.130)
Separated	-0.261 (0.165)
Has child	0.121 (0.114)
Number of children	-0.031 (0.070)
Education	0.104* (0.019)
Mother's education	0.007 (0.012)
Father's education	-0.006 (0.012)
Ever arrested	-0.223* (0.055)
Household income	
3000–6000	0.033 (0.085)
6000–9000	0.213* (0.104)
9000–18,000	0.149 (0.086)
>18,000	0.103 (0.095)
Personal income	
3000–6000	0.105 (0.080)
6000–9000	0.180 (0.127)
>9000	0.197 (0.162)
At baseline	
Have job	0.218* (0.071)
Months employed, previous year	0.049* (0.011)
Had job, previous year	0.306* (0.091)
Earnings, previous year (*10,000)	0.012 (0.120)
Usual hours/week (*10,000)	-26.580 (19.508)
Usual weekly earnings (*10,000)	0.845 (1.990)
Constant	-1.288* (0.285)

Notes: $N = 9145$. Robust S.E. are given in parentheses. Table reports are (log-odds) coefficients from a logit of employment (positive hours) in week 208 on treatment status and baseline characteristics.

*Indicates statistically significant at the 0.05 (or less) level.

As in Heckman (1979), sample selection bias can be seen as specification error in the conditional expectation

$$E[Y|D, X, Z^* \geq 0] = D\beta + X\pi_1 + E[U|D, X, V \geq -D\gamma - X\pi_2].$$

One modelling approach is to assume that data are missing at random, perhaps conditional on a set of covariates (Rubin, 1976). This amounts to assuming that U and V are independent of one another or that employment status is unrelated to the determination of wages. This assumption is strictly inconsistent with standard models of labour supply that account for the participation decision (e.g., see Heckman, 1974).

A more common modelling assumption is that some of the exogenous variables determine sample selection but do not have their own direct impact on the outcome of interest; i.e., some of the elements of π_1 are zero, while corresponding elements of π_2 are nonzero. Such exclusion restrictions are used in parametric and semiparametric models of the censored selection process (e.g., Heckman, 1979, 1990; Ahn and Powell, 1993; Andrews and Schafgans, 1998; Das, Newey and Vella, 2003).

TABLE 3
*Bounds on treatment effects for week 208 $\ln(\text{wage})$ using bounds of support
 (Horowitz and Manski)*

Control group		
(i)	Observations	3599
(ii)	Employment rate	0.566
(iii)	Mean $\log(\text{wage})$	1.997
(iv)	Upper bound	2.332
(v)	Lower bound	1.520
Treatment group		
(vi)	Observations	5546
(vii)	Employment rate	0.607
(viii)	Mean $\log(\text{wage})$	2.031
(ix)	Upper bound	2.321
(x)	Lower bound	1.586
Difference		
(xi)	Upper bound: (ix) – (v)	0.802
(xii)	Lower bound: (x) – (iv)	–0.746

Notes: 0.90 and 2.77 are the lower and upper bounds of the support of $\ln(\text{hourly wage})$ in week 208 after random assignment; (iv) = (ii) \times (iii) + $[1 - (\text{ii})] \times 2.77$; (v) = (ii) \times (iii) + $[1 - (\text{ii})] \times (0.90)$. Rows (ix) and (x) are defined analogously.

The practical limitation to relying on exclusion restrictions for the sample selection problem is that there may not exist credible “instruments” that can be excluded from the outcome equation. This seems to be true for an analysis of the Job Corps experiment. There are many variables available to the researcher from the Job Corps evaluation and many of the key variables are listed in Tables 1 and 2. But for each of the variables in Table 2 that have significant associations with employment, there is a well-developed literature suggesting that those variables may also influence wage offers. For example, race, gender, education, and criminal histories all could potentially impact wages. Household income and past employment experiences are also likely to be correlated with unobserved determinants of wages.

Researchers’ reluctance to rely upon specific exclusion restrictions motivates a second, general approach to addressing the sample selection problem: the construction of worst-case scenario bounds of the treatment effect. When the support of the outcome is bounded, the idea is to impute the missing data with either the largest or the smallest possible values to compute the largest and smallest possible treatment effects consistent with the data that are observed. Horowitz and Manski (2000a) use this notion to provide a general framework for constructing bounds for treatment effect parameters when outcome and covariate data are nonrandomly missing in an experimental setting.¹² This strategy is discussed in detail in Horowitz and Manski (2000a), which shows that the approach can be useful when Y is a binary outcome.

This imputation procedure cannot be used when the support is unbounded. Even when the support is bounded, if it is very wide, so too will be the width of the treatment effect bounds. In the context of the Job Corps program, the bounds are somewhat uninformative. Table 3 computes the Horowitz and Manski (2000a) bounds for the treatment effect of the Job Corps program on log wages in week 208. Specifically, it calculates the upper bound of the treatment effect as:

12. An early example of sensitivity analysis that imputed missing values is found in the work of Smith and Welch (1986). Others (Balke and Pearl, 1997; Heckman and Vytlacil, 1999, 2000a,b) have constructed such bounds to address a very different problem—that of imperfect compliance of the treatment, even when “intention” to treat is effectively randomized (Bloom, 1984; Robins, 1989; Angrist *et al.*, 1996).

$$\Pr[Z^* \geq 0|D=1] E[Y|D=1] + \Pr[Z^* < 0|D=1] Y^{UB} \\ - (\Pr[Z^* \geq 0|D=0] E[Y|D=0] + \Pr[Z^* < 0|D=0] Y^{LB}),$$

where all population quantities can be estimated, and Y^{UB} and Y^{LB} are the upper and lower bounds of the support of log wages. As reported in the table, Y^{UB} and Y^{LB} are taken to be 2.77 and 0.90 (\$15.96 and \$2.46 an hour), respectively.¹³

Table 3 shows that the lower bound for the treatment effect on week 208 log wages is -0.75 and the upper bound is 0.80 . Thus, the interval is almost as consistent with extremely large negative effects as it is with extremely large positive effects. The reason for this wide interval is that more than 40% of the individuals are not employed in week 208. In this context, imputing the missing values with the maximal and minimal values of Y is so extreme as to yield an interval that includes effect sizes that are arguably implausible. Nevertheless, the Horowitz and Manski (2000a) bounds provide a useful benchmark and highlight that some restrictions on the sample selection process are needed to produce tighter bounds (Horowitz and Manski, 2000b).

The procedure proposed below is a kind of “hybrid” of the two general approaches to the sample selection problem. It yields bounds on the treatment effect, even when the outcome is unbounded. It does so by imposing some structure on the sample selection process but without requiring exclusion restrictions.

3. IDENTIFICATION OF BOUNDS ON TREATMENT EFFECTS

This section first uses a simple case to illustrate the intuition behind the main identification result and then generalizes it for a very unrestrictive sample selection model.

Consider the case where there is only the treatment indicator, with no other covariates. That is, X is a constant, so that π_1 and π_2 will be intercept terms. It will become clear that the result below is also valid conditional on any value of X . Describing the identification result in this simple case makes clear that the proposed procedure does not rely on exclusion restrictions.¹⁴ In addition, this section and the next assume that U (and hence Y) has a continuous distribution. Doing so will simplify the exposition; it can be shown that the proposed procedure can be applied to discrete outcome variables as well (see Lee, 2002). Without loss of generality, assume $\gamma > 0$, so that the treatment causes an increase in the likelihood of the outcome being observed.

From Equation (1), the observed population means for the control and treatment groups can be written as

$$E[Y|D=0, Z^* \geq 0] = \pi_1 + E[U|D=0, V \geq -\pi_2] \quad (2)$$

and

$$E[Y|D=1, Z^* \geq 0] = \pi_1 + \beta + E[U|D=1, V \geq -\pi_2 - \gamma], \quad (3)$$

respectively. This shows that when U and V are correlated, the difference in the means will generally be different from β .

13. The wage variable was transformed before being analysed to minimize the effect of outliers so that the Horowitz and Manski (2000a) bounds would not have to rely on these outliers. Specifically, the entire observed wage distribution was split into 20 categories, according to the 5th, 10th, 15th, ... 95th percentile wages, and the individual was assigned the mean wage within each of the 20 groups. Thus, the upper “bound” of the support, *e.g.*, is really the mean log wage for those earning more than the 95th percentile. The same data are used for the trimming procedure described below. Strictly speaking, the Horowitz and Manski (2000a) bounds would use the theoretical bounds of the support of the population log-wage distribution. Since these population maximums and minimums are not observed, one could instead use the log of the minimum and maximum log-wage observed in the sample. It is clear that doing so would produce wider bounds than that given by the implementation here.

14. Note that while existing procedures for point identification require an instrument that satisfies an exclusion restriction, the existence of such an instrument is not sufficient for identification. For example, a single binary instrument will not allow identification without imposing further assumptions.

Identification of β would be possible if we could estimate

$$E[Y|D = 1, V \geq -\pi_2] = \pi_1 + \beta + E[U|D = 1, V \geq -\pi_2] \quad (4)$$

because (2) could be subtracted to yield the effect β (since D is independent of (U, V) by assumption). But the mean in (4) is not observed.

But this mean can be bounded. This is because all observations on Y needed to compute this mean are a subset of the selected population ($V \geq -\pi_2 - \gamma$). For example, we know that

$$E[Y|D = 1, Z^* \geq 0] = (1 - p) E[Y|D = 1, V \geq -\pi_2] + p E[Y|D = 1, -\pi_2 - \gamma \leq V < -\pi_2],$$

where $p = \frac{\Pr[-\pi_2 - \gamma \leq V < -\pi_2]}{\Pr[-\pi_2 - \gamma \leq V]}$. The observed treatment mean is a weighted average of (4) and the mean for a subpopulation of “marginal” individuals ($-\pi_2 - \gamma \leq V < -\pi_2$) that are induced to be selected into the sample because of the treatment.

$E[Y|D = 1, V \geq -\pi_2]$ is therefore bounded above by $E[Y|D = 1, Z^* \geq 0, Y \geq y_p]$, where y_p is the p th quantile of the treatment group’s observed Y distribution. This is true because among the selected population with $V \geq -\pi_2 - \gamma$, $D = 1$, no subpopulation with proportion $(1 - p)$ can have a mean that is larger than the average of the largest $(1 - p)$ values of Y .

Put another way, we cannot identify which observations are inframarginal ($V \geq -\pi_2$) and which are marginal ($-\pi_2 - \gamma \leq V < -\pi_2$). But the worst-case scenario is that the smallest p values of Y belong to the marginal group and the largest $1 - p$ values belong to the inframarginal group. Thus, by trimming the lower tail of the Y distribution by the proportion p , we obtain an upper bound for the inframarginal group’s mean in (4). Consequently, $E[Y|D = 1, Z^* \geq 0, Y \geq y_p] - E[Y|D = 0, Z^* \geq 0]$ is an upper bound for β . Note that the trimming proportion p is equal to

$$\frac{\Pr[Z^* \geq 0|D = 1] - \Pr[Z^* \geq 0|D = 0]}{\Pr[Z^* \geq 0|D = 1]},$$

where each of these probabilities is identified by the data.

To summarize, a standard latent variable sample selection model implies that the observed outcome distribution for the treatment group is a mixture of two distributions: (1) the distribution for those who would have been selected irrespective of the treatment (the inframarginal group) and (2) the distribution for those induced into being selected because of the treatment (the marginal group). It is possible to quantify the proportion of the treatment group that belongs to this second group, using a simple comparison of the selection probabilities of the treatment and control groups. Although it is impossible to identify specifically *which* treated individuals belong to the second group, worst-case scenarios can be constructed by assuming that they are either at the very top or at the very bottom of the distribution. Thus, trimming the data by the known proportion of excess individuals should yield bounds on the mean for the inframarginal group.

3.1. Identification under a generalized sample selection model

This identification result applies to a much wider class of sample selection models. It depends neither on a constant treatment effect nor on homoskedasticity, which are both implicitly assumed in Equation (1).

To see this, consider a general sample selection model that allows for heterogeneity in treatment effects:

$$(Y_1^*, Y_0^*, S_1, S_0, D) \text{ is i.i.d. across individuals} \quad (5)$$

$$S = S_1 D + S_0 (1 - D)$$

$$Y = S \cdot \{Y_1^* D + Y_0^* (1 - D)\}$$

(Y, S, D) is observed,

where D , S , S_0 , and S_1 are all binary indicator variables. D denotes treatment status; S_1 and S_0 are “potential” sample selection indicators for the treated and control states, respectively. For example, when an individual has $S_1 = 1$ and $S_0 = 0$, this means that there will be nonmissing data on the outcome ($S = 1$) if treatment is given and there will be missing data on the outcome ($S = 0$) if treatment is denied. The second line highlights the fact that for each individual, we only observe S_1 or S_0 . Y_1^* and Y_0^* are latent potential outcomes for the treated and control states, and the third line points out that we observe only one of the latent outcomes Y_1^* or Y_0^* and only if the individual is selected into the sample $S = 1$. It is assumed throughout the paper that $E[S|D = 1], E[S|D = 0] > 0$.

Assumption 1. (*Independence*): (Y_1^*, Y_0^*, S_1, S_0) is independent of D .

This assumption corresponds to the independence of (U, V) and (D, X) in the previous section. In the context of experiments, random assignment will ensure this assumption will hold.

Assumption 2a. (*Monotonicity*): $S_1 \geq S_0$ with probability 1.

This assumption implies that treatment assignment can only affect sample selection in “one direction”. Some individuals will never be observed, regardless of treatment assignment ($S_0 = S_1 = 0$), others will always be observed ($S_0 = 1, S_1 = 1$), and others will be selected into the sample *because* of the treatment ($S_0 = 0, S_1 = 1$). This assumption is commonly invoked in studies of imperfect compliance of treatment (Imbens and Angrist, 1994; Angrist, Imbens and Rubin, 1996); the difference is that in those studies, monotonicity is for how an instrument affects *treatment status*. Here, the monotonicity is for how treatment affects *sample selection*.

In the context of the Job Corps program, the monotonicity assumption essentially limits the degree of heterogeneity in the effect of the program on labour force participation. It does not allow, *e.g.*, the job search assistance services provided by Job Corps to induce some to become employed while simultaneously causing others to drop out of the labour force. A negative impact could occur, *e.g.*, if the job search counselling induced some to pursue further education (and hence drop out of the labour force). Similar to the case of LATE, with only information on the outcome, treatment status, and selection status, the monotonicity assumption is fundamentally untestable. It should be noted that monotonicity has been shown to be equivalent to assuming a latent variable threshold-crossing model (Vytlacil, 2002), which is the basis for virtually all sample selection models in econometrics.

Proposition 1a. Let Y_0^* and Y_1^* be continuous random variables. If Assumptions 1 and 2a hold then Δ_0^{LB} and Δ_0^{UB} are sharp lower and upper bounds for the average treatment effect $E[Y_1^* - Y_0^* | S_0 = 1, S_1 = 1]$, where

$$\Delta_0^{LB} \equiv E[Y|D = 1, S = 1, Y \leq y_{1-p_0}] - E[Y|D = 0, S = 1]$$

$$\Delta_0^{UB} \equiv E[Y|D = 1, S = 1, Y \geq y_{p_0}] - E[Y|D = 0, S = 1]$$

$$y_q \equiv G^{-1}(q), \text{ with } G \text{ the c.d.f. of } Y, \text{ conditional on } D = 1, S = 1$$

$$p_0 \equiv \frac{\Pr[S = 1|D = 1] - \Pr[S = 1|D = 0]}{\Pr[S = 1|D = 1]}.$$

The bounds are sharp in the sense that Δ_0^{LB} (Δ_0^{UB}) is the largest (smallest) lower (upper) bound that is consistent with the observed data. Furthermore, the interval $[\Delta_0^{LB}, \Delta_0^{UB}]$ is contained in any other valid bounds that impose the same assumptions. (If $S_0 \geq S_1$ with probability 1 then the control group's, rather than the treatment group's, outcome distribution must be trimmed.)

Obviously, this result is equally valid if one were to assume monotonicity in the opposite direction ($S_0 \geq S_1$ with probability 1).

Remark 1. The sharpness of the bound Δ_0^{UB} means that it is the “best” upper bound that is consistent with the data. A specific example of where this proposition can be applied is in Krueger and Whitmore (2001), who study the impact of the Tennessee STAR (Student Teacher Achievement Ratio) class-size experiment. In that study, students are randomly assigned to a regular or small class and the outcome of interest is the SAT (or ACT) scores, but not all students take the exam. Krueger and Whitmore (2001, p. 25) use Assumptions 1 and 2a to derive a different upper bound, given by $B \equiv E[Y | D = 1, S = 1] \cdot \frac{\Pr[S=1|D=1]}{\Pr[S=1|D=0]} - E[Y | D = 0, S = 1]$. Proposition 1a implies that this bound B , like any other proposed bound using these assumptions, cannot be smaller than Δ_0^{UB} .¹⁵

Remark 2. An important practical implication of Assumptions 1 and 2a is that as p_0 vanishes, so does the sample selection bias.¹⁶ The intuition is that if $p_0 = 0$ then under the monotonicity assumption, both treatment and control groups are comprised of individuals whose sample selection was unaffected by the assignment to treatment, and therefore, the two groups are comparable. These individuals can be thought of as the “always-takers” subpopulation (Angrist *et al.*, 1996), except that “taking” is not the taking of the treatment, but rather selection into the sample. It follows that when analysing randomized experiments, if the sample selection rates in the treatment and control groups are similar and if the monotonicity condition is believed to hold then a comparison of the treatment and control means is a valid estimate of an average treatment effect.¹⁷ As an example, the proportion of control group individuals, at week 90, that have continuously nonmissing earnings and hours data is 0.822 and the proportion is 0.003 smaller (S.E. of 0.006) for the treatment group. Thus, if the above assumptions are invoked to examine the nonresponse/attrition problem (as opposed to the focus of this study, missing wages due to nonemployment) then the data suggest little bias due to nonresponse/attrition.

Remark 3. Assumptions 1 and 2a are minimally sufficient for computing the bounds. First, the independence assumption is also important since it is what justifies the contrast between the trimmed treatment group and the control group.

Second, monotonicity ensures that the sample-selected control group consists only of those individuals with $S_0 = 1, S_1 = 1$. Without monotonicity, the control group could consist solely of observations with $S_0 = 1, S_1 = 0$ and the treatment group solely of observations with $S_0 = 0,$

15. Thus, in the context of Krueger and Whitmore (2001), Proposition 1a implies that computing the bound B is unnecessary after already computing a very different estimate T , their “linear truncation” estimate. They justify T under a different set of assumptions that (1) the additional small-class students induced to take the ACT exam are from the left tail of the distribution and (2) if attending a small class did not change the ranking of students in small classes. Their estimate T is mechanically equivalent to the bound Δ_0^{UB} . Therefore, Proposition 1a implies that their estimate T is actually the sharp upper bound given the mild assumptions that were used to justify their bound B .

16. A vanishing p corresponds to individuals with the same value of the sample selection correction term, and it is well known that there is no selection bias, conditional on the correction term (see, *e.g.*, Heckman and Robb, 1986; Heckman, 1990; Ahn and Powell, 1993; Angrist, 1997).

17. Note that p_0 here is proportional to the difference in the fraction that are sample selected between the treatment and the control groups. Thus, the notion of a vanishing p should not be confused with “identification at infinity” in Heckman (1990), in which the bias term vanishes as the fraction that is selected into the sample tends to 1.

$S_1 = 1$. Since the two subpopulations do not “overlap”, the difference in the means could not be interpreted as a causal effect.

An interesting exception to this arises in the special case that $E[S|D=0] + E[S|D=1] > 1$, in which case informative bounds can be constructed without invoking monotonicity, as demonstrated in Zhang and Rubin (2003). There, the insight is that the proportion of those who are $S_0 = 1, S_1 = 0$ can be no larger than the proportion in the treatment group who have missing values, $1 - E[S|D=1]$. It follows that within the control group, the fraction of $S_0 = 1, S_1 = 1$ individuals cannot be less than $E[S|D=0] - (1 - E[S|D=1])$, which is positive, as assumed. It thus follows that, e.g., the upper bound for the mean of Y_0^* for $S_0 = 1, S_1 = 1$ is the mean after trimming the bottom $\frac{1-E[S|D=1]}{E[S|D=0]}$ fraction of the observed control group distribution. A symmetric argument can be made for bounding the mean of Y_1^* for $S_0 = 1, S_1 = 1$. This idea is formalized in Zhang and Rubin (2003) and also discussed in Zhang, Rubin and Mealli (2008). It should be noted, however, that the procedure of Zhang and Rubin (2003) will not produce informative bounds for a general sample selection model, as the assumption $E[S|D=0] + E[S|D=1] > 1$ is crucial.¹⁸ Specifically, if $E[S|D=0] + E[S|D=1] \leq 1$ then the worst-case scenario would involve trimming *all* the observed treatment and control observations, resulting in noninformative (or “vacuous”) bounds.¹⁹

Remark 4. When $p_0 = 0$ in a randomized experimental setting, there is a limited test of whether monotonicity holds (and therefore whether the simple difference in means in the outcome suffers from sample selection bias). If $p_0 = 0$ and monotonicity holds then the selected subsets of both the treatment and the control groups will consist solely of individuals with $(S_0 = 1, S_1 = 1)$. Under randomization, the treatment control difference in the outcome should represent a causal effect. In addition, the distribution of the X s should be the same in the treatment and control groups, *conditional on being selected*. This can be tested empirically.

In order for this test to have power, the two subpopulations $(S_0 = 0, S_1 = 1)$ and $(S_0 = 1, S_1 = 0)$, need to have different distributions of baseline characteristics X . Recall that without monotonicity, the selected treated group will be comprised of two subpopulations, $(S_0 = 1, S_1 = 1)$ and $(S_0 = 0, S_1 = 1)$, while the selected control group will be comprised of the groups $(S_0 = 1, S_1 = 1)$ and $(S_0 = 1, S_1 = 0)$. So if the distribution of the X is the same for $(S_0 = 0, S_1 = 1)$ and $(S_0 = 1, S_1 = 0)$ then the selected treatment and control groups will have the same distribution of X , whether or not monotonicity holds.

Finally, the trimming procedure described above places sharp bounds on the average treatment effect for a particular subpopulation—those individuals who will be selected irrespective of the treatment assignment $(S_0 = 1, S_1 = 1)$. It should be noted, however, that this subpopulation is the only one for which it is possible to learn about treatment effects, given Assumptions 1 and 2a (at least, in this missing data problem). For the marginal $(S_0 = 0, S_1 = 1)$ observations, the outcomes are missing in the control regime. For the remaining $(S_0 = 0, S_1 = 0)$ observations, outcomes are missing in both the treatment and the control regimes. It would still be possible to appeal to the bounds of Horowitz and Manski (2000a) to construct bounds on this remaining population of the “never observed”, but this interval (whose width would be two times the width of the outcome variable’s support) would not require any data. Whether or not the subpopulation

18. For example, the procedure will not work if 49% of the treatment group is missing and 52% of the control group is missing.

19. Although $E[S|D=0] + E[S|D=1] > 1$ is not formally stated as an assumption in Zhang and Rubin (2003) or in Zhang *et al.* (2008), it is clear that it is a necessary one to produce informative bounds. Using the notation of Zhang and Rubin (2003), P_{CG} and P_{TG} are equivalent to $E[S|D=0]$ and $E[S|D=1]$, respectively. If $P_{CG} + P_{TG} < 1$, this means that π_{DG} is bounded above by P_{CG} (the line below their equation (12)), which means that their equations (11) and (12) yield $(-\infty, \infty)$ as bounds (if the dependent variable has unbounded support).

of the “always observed” is of interest will depend on the context. In the case of the Job Corps program, *e.g.*, it is useful to assess the impact of the program on wage rates for those whose employment status was not affected by the program.

3.2. Narrowing bounds using covariates

A straightforward extension to the above analysis is to produce bounds of the treatment effect, stratified by observed “baseline” characteristics X (those determined prior to the assignment of treatment). Examples of such covariates in the case of Job Corps might include gender or race. It is clear that the above analysis can all be conditioned on covariates X . It is possible to estimate bounds for the average treatment effect for each value of X .

Alternatively, one can use these covariates to reduce the width of the bounds for the same estimand that has been discussed so far (the average treatment effect for those who would always be observed). To gain intuition for this, suppose half of the workers in the treatment group earns the wage w^H , while the other half earns the lower wage of w^L . The trimming procedure described in the previous sections suggests removing only low-wage individuals, by a proportion p_0 to obtain an upper bound of the mean for the “inframarginally” selected. The trimmed mean will necessarily be larger.

Suppose now there is a baseline covariate X that perfectly predicts whether an individual will earn w^H or w^L . Then, due to the random assignment of treatment, Assumptions 1 and 2a also hold conditional on X . Therefore, the results in the previous section can be applied separately for the two types of workers. If, for both groups, the same proportion of observations is trimmed, the overall mean will not be altered by this trimming procedure.

More formally, consider the following alternative to Assumption 1,

Assumption 3. (*Independence*): Let X be a vector of covariates, and let $(Y_1^*, Y_0^*, S_1, S_0, X)$ be independent of D .

As an example, this would hold in the case of the Job Corps experiment due to random assignment.

Proposition 1b. Let Y_0^* and Y_1^* be continuous random variables. If Assumptions 3 and 2a hold then Δ_0^{LB} and Δ_0^{UB} are sharp lower and upper bounds for the average treatment effect $E[Y_1^* - Y_0^* | S_0 = 1, S_1 = 1]$, where

$$\begin{aligned}\Delta_0^{\text{LB}} &\equiv \int \Delta_x^{\text{LB}} dH(x) \\ \Delta_0^{\text{UB}} &\equiv \int \Delta_x^{\text{UB}} dH(x), \text{ where } H \text{ is the c.d.f. of } X \text{ conditional on } D = 0, S = 1 \\ \Delta_x^{\text{LB}} &\equiv E[Y | D = 1, S = 1, Y \leq y_{1-p_x}, X = x] - E[Y | D = 0, S = 1, X = x] \\ \Delta_x^{\text{UB}} &\equiv E[Y | D = 1, S = 1, Y \geq y_{p_x}, X = x] - E[Y | D = 0, S = 1, X = x] \\ y_q &\equiv G_x^{-1}(q), \text{ with } G_x \text{ the c.d.f. of } Y, \text{ conditional on } D = 1, S = 1, X = x \\ p_x &\equiv \frac{\Pr[S = 1 | D = 1, X = x] - \Pr[S = 1 | D = 0, X = x]}{\Pr[S = 1 | D = 1, X = x]}.\end{aligned}$$

The bounds are sharp in the sense that Δ_0^{LB} (Δ_0^{UB}) is the largest (smallest) lower (upper) bound that is consistent with the observed data. Furthermore, $\Delta_0^{\text{LB}} \geq \Delta_0^{\text{LB}}$ and $\Delta_0^{\text{UB}} \leq \Delta_0^{\text{UB}}$.

The first part of the proposition follows from applying Proposition 1a conditionally on $X = x$. The second claim, that the width of the bounds must be narrower after using the covariates, is seen by noting that any treatment effect that is consistent with an observed population distribution of (Y, S, D, X) , must also be consistent with the data after throwing away information on X , and observing only the distribution of (Y, S, D) . This necessity is strictly inconsistent with $\Delta_0^{\text{UB}} > \Delta_0^{\text{UB}}$.

4. ESTIMATION AND INFERENCE

This section proposes and discusses an estimator for the bounds. The estimator can be shown to be \sqrt{n} consistent and asymptotically normal. The asymptotic variance comprises three components, reflecting (1) the variance of the trimmed distribution, (2) the variance of the estimated trimming threshold, and (3) the variance in the estimate of how much of the distribution to trim. To minimize redundancies, the discussion below continues to consider the case that $S_1 \geq S_0$ with probability 1 (from Assumption 2a); the results are also analogously valid for the reverse case of $S_0 \geq S_1$.

4.1. Estimation

The estimates of the bounds are sample analogs to the parameters defined in Proposition 1a. First, the trimming proportion \hat{p} is estimated by taking the treatment control difference in the proportion with nonmissing outcomes and dividing by the proportion that is selected in the treatment group. Next, the \hat{p} th (or the $(1 - \hat{p})$ th) quantile of the treatment group's outcome distribution is calculated. Finally, these quantiles are used to trim the data for the treatment group's outcomes and compute the bounds Δ^{LB} and Δ^{UB} .

Formally, we have

Definition of estimator:

$$\begin{aligned}\widehat{\Delta}^{\text{LB}} &\equiv \frac{\sum Y \cdot S \cdot D \cdot 1[Y \leq \widehat{y}_{1-\hat{p}}]}{\sum S \cdot D \cdot 1[Y \leq \widehat{y}_{1-\hat{p}}]} - \frac{\sum Y \cdot S \cdot (1 - D)}{\sum S \cdot (1 - D)} \\ \widehat{\Delta}^{\text{UB}} &\equiv \frac{\sum Y \cdot S \cdot D \cdot 1[Y \geq \widehat{y}_{\hat{p}}]}{\sum S \cdot D \cdot 1[Y \geq \widehat{y}_{\hat{p}}]} - \frac{\sum Y \cdot S \cdot (1 - D)}{\sum S \cdot (1 - D)} \\ \widehat{y}_q &\equiv \min \left\{ y : \frac{\sum S \cdot D \cdot 1[Y \leq y]}{\sum S \cdot D} \geq q \right\} \\ \widehat{p} &\equiv \left(\frac{\sum S \cdot D}{\sum D} - \frac{\sum S \cdot (1 - D)}{\sum (1 - D)} \right) / \left(\frac{\sum S \cdot D}{\sum D} \right),\end{aligned}\tag{6}$$

where the summation is over the entire sample of size n .

4.2. Consistency, asymptotic normality, variance estimation, and inference

The estimators $\widehat{\Delta}^{\text{LB}}$ and $\widehat{\Delta}^{\text{UB}}$ are consistent for Δ_0^{LB} and Δ_0^{UB} under fairly standard conditions:

Proposition 2. (Consistency): Let Y have bounded support (i.e. \exists finite L, U such that $\Pr[Y \leq L]$ and $\Pr[Y \geq U] = 0$), and suppose $E[S|D = 0] > 0$ and $p_0 \geq 0$, then $\widehat{\Delta}^{\text{LB}} \xrightarrow{p} \Delta_0^{\text{LB}}$ and $\widehat{\Delta}^{\text{UB}} \xrightarrow{p} \Delta_0^{\text{UB}}$.

As shown in the Appendix, the proof involves showing that the estimator is a solution to a GMM (Generalized Method of Moments) problem, showing that the moment function vector is, with probability 1, continuous at each possible value of Δ_0^{LB} and Δ_0^{UB} and applying theorem 2.6 of Newey and McFadden (1994).²⁰

The estimators Δ^{LB} and Δ^{UB} are also asymptotically normal, with an intuitive expression for the variance.

Proposition 3. (Asymptotic normality): Define $\mu^{\text{LB}} \equiv E[Y | D = 1, S = 1, Y \leq y_{1-p_0}]$ and $\mu^{\text{UB}} \equiv E[Y | D = 1, S = 1, Y \geq y_{p_0}]$. In addition to the conditions in Proposition 2, assume $E[S|D=0] < E[S|D=1] < 1$. Then, $\sqrt{n}(\Delta^{\text{LB}} - \Delta_0^{\text{LB}}) \xrightarrow{d} N(0, V^{\text{LB}} + V_C)$ and $\sqrt{n}(\Delta^{\text{UB}} - \Delta_0^{\text{UB}}) \xrightarrow{d} N(0, V^{\text{UB}} + V_C)$, where

$$\begin{aligned} V^{\text{LB}} &= \frac{\text{Var}[Y|D=1, S=1, Y \leq y_{1-p_0}]}{E[SD](1-p_0)} + \frac{(y_{1-p_0} - \mu^{\text{LB}})^2 p_0}{E[SD](1-p_0)} \\ &\quad + \left(\frac{y_{1-p_0} - \mu^{\text{LB}}}{1-p_0} \right)^2 \cdot V^p \\ V^{\text{UB}} &= \frac{\text{Var}[Y|D=1, S=1, Y \geq y_{p_0}]}{E[SD](1-p_0)} + \frac{(y_{p_0} - \mu^{\text{UB}})^2 p_0}{E[SD](1-p_0)} \\ &\quad + \left(\frac{y_{p_0} - \mu^{\text{UB}}}{1-p_0} \right)^2 \cdot V^p \\ V^p &= (1-p_0)^2 \left(\frac{\left(1 - \frac{\alpha_0}{1-p_0}\right)}{E[D]\left(\frac{\alpha_0}{1-p_0}\right)} + \frac{(1-\alpha_0)}{(1-E[D])\alpha_0} \right) \end{aligned} \quad (7)$$

and $V_C = \text{Var}[Y|D=0, S=1]/E[S(1-D)]$.²¹

Consider the three terms in V^{LB} . The first term would be the variance of the estimate if the trimming threshold y_{1-p_0} were known. The term $\frac{1}{E[SD](1-p_0)}$ exists because n is the size of the entire sample (both treatment and control, and all observations including those with missing outcomes). The second term reflects the fact that the threshold is a quantile that needs to be estimated. Taken together, the first two terms are exactly equivalent to the expression given in Stigler (1973), which derives the asymptotic distribution of a one-sided “ p_0 -trimmed” mean, when p_0 is known. But p_0 is not known, and must be estimated, which is reflected in the third term. The third term itself includes the asymptotic variance of \hat{p} multiplied by the square of the gradient of the population trimmed mean with respect to p_0 . Note that $\frac{1-\alpha_0}{\alpha_0}$ and $\frac{\left(1 - \frac{\alpha_0}{1-p_0}\right)}{\left(\frac{\alpha_0}{1-p_0}\right)}$ are the odds of an observation being missing conditional on being in the control group and the treatment group, respectively. The Appendix contains the proposition’s proof, which involves applying theorem 7.2 of Newey and McFadden (1994), an asymptotic normality result for GMM estimators when the moment function is not smooth.

20. Recall that boundedness of the support of Y is unnecessary for identification. Furthermore, consistency can be proven without boundedness (see Lee, 2005).

21. Note that the usual asymptotic variance of the estimated mean for the control group is divided by $E[S(1-D)]$ because n here is the total number of observations (selected and nonselected and treated and control).

Estimation of the variances is easily carried out by replacing all of the above quantities (e.g., $E[SD]$, y_{p_0}) with either of their sample analogs (e.g., $\frac{1}{n} \sum SD$, $\hat{y}_{\hat{p}}$). After assuming a finite second moment for Y , consistency follows because the resulting estimator is a continuous function of consistent estimators for each part.

There are two simple ways to compute confidence intervals. First, one can compute the interval $[\widehat{\Delta}^{\text{LB}} - 1.96 \frac{\widehat{\sigma}_{\text{LB}}}{\sqrt{n}}, \widehat{\Delta}^{\text{UB}} + 1.96 \frac{\widehat{\sigma}_{\text{UB}}}{\sqrt{n}}]$, $\widehat{\sigma}_{\text{LB}} \equiv \sqrt{V(\widehat{\Delta}^{\text{LB}})}$, $\widehat{\sigma}_{\text{UB}} \equiv \sqrt{V(\widehat{\Delta}^{\text{UB}})}$. This interval will asymptotically contain the region $[\Delta_0^{\text{LB}}, \Delta_0^{\text{UB}}]$ with at least 0.95 probability.²² Imbens and Manski (2004) point out that this same interval will contain the parameter $E[Y_1^* - Y_1^* | S_0 = 1, S_1 = 1]$ with an even greater probability, suggesting that the confidence interval for the parameter will be narrower for the same coverage rate. The results of Imbens and Manski (2004) imply that a (smaller) interval of $[\widehat{\Delta}^{\text{LB}} - \bar{C}_n \frac{\widehat{\sigma}_{\text{LB}}}{\sqrt{n}}, \widehat{\Delta}^{\text{UB}} + \bar{C}_n \frac{\widehat{\sigma}_{\text{UB}}}{\sqrt{n}}]$, where \bar{C}_n satisfies

$$\Phi\left(\bar{C}_n + \sqrt{n} \frac{\widehat{\Delta}^{\text{UB}} - \widehat{\Delta}^{\text{LB}}}{\max(\widehat{\sigma}_{\text{LB}}, \widehat{\sigma}_{\text{UB}})}\right) - \Phi(-\bar{C}_n) = 0.95,$$

can be computed, and will contain the parameter $E[Y_1^* - Y_1^* | S_0 = 1, S_1 = 1]$ with a probability of at least 0.95.

The interval of Imbens and Manski (2004) is more appropriate here since the object of interest is the treatment effect and not the *region* of all rationalizable treatment effects. Nevertheless, for completeness, both intervals are reported in the presentation of the results.

4.3. Inference with unknown $\text{sgn}(p_0)$

The discussion to this point has presumed that $p_0 > 0$ and therefore the procedure described so far is appropriate when the researcher has reason to impose the assumption that the treatment status has a (strictly) positive impact on the outcome being observed. But a researcher may want to remain agnostic about the sign of p_0 . Specifically, we have so far assumed that $S_1 \geq S_0$ with probability 1. But the researcher—still concerned about sample selection—may instead want to adopt the following assumption.

Assumption 2b. (*Monotonicity*): Either $S_1 \geq S_0$ with probability 1 or $S_0 \geq S_1$ with probability 1.

This means that monotonicity is maintained, but the *direction* in which treatment affects selection is unknown.

The above identification, estimation, and inference procedure readily generalizes to this case. First, from an identification standpoint, it is clear that the sharp lower bound is given by

$$\begin{aligned} \Delta_0^{\text{LB}} \equiv & 1[p_0 \geq 0] \{E[Y|D = 1, S = 1, Y \leq y_{1-p_0}] - E[Y|D = 0, S = 1]\} \\ & + 1[p_0 < 0] \{E[Y|D = 1, S = 1] - E[Y|D = 0, S = 1|Y \geq y_{p_0^*}]\}, \end{aligned}$$

22. To see this, note that $\Pr[\widehat{\Delta}^{\text{LB}} - 1.96 \frac{\widehat{\sigma}_{\text{LB}}}{\sqrt{n}} < \Delta_0^{\text{LB}}, \widehat{\Delta}^{\text{UB}} + 1.96 \frac{\widehat{\sigma}_{\text{UB}}}{\sqrt{n}} > \Delta_0^{\text{UB}}]$ is equivalent to $\Pr\left[\frac{\widehat{\Delta}^{\text{LB}} - \Delta_0^{\text{LB}}}{\widehat{\sigma}_{\text{LB}}} < -1.96, \frac{\widehat{\Delta}^{\text{UB}} - \Delta_0^{\text{UB}}}{\widehat{\sigma}_{\text{UB}}} > 1.96\right] = 1 - \Pr\left[\frac{\widehat{\Delta}^{\text{LB}} - \Delta_0^{\text{LB}}}{\widehat{\sigma}_{\text{LB}}} > 1.96\right] - \Pr\left[\frac{\widehat{\Delta}^{\text{UB}} - \Delta_0^{\text{UB}}}{\widehat{\sigma}_{\text{UB}}} < -1.96\right] + \Pr\left[\frac{\widehat{\Delta}^{\text{LB}} - \Delta_0^{\text{LB}}}{\widehat{\sigma}_{\text{LB}}} > 1.96, \frac{\widehat{\Delta}^{\text{UB}} - \Delta_0^{\text{UB}}}{\widehat{\sigma}_{\text{UB}}} < -1.96\right]$, which is equal to $1 - 0.025 - 0.025 + \Pr\left[\frac{\widehat{\Delta}^{\text{LB}} - \Delta_0^{\text{LB}}}{\widehat{\sigma}_{\text{LB}}} > 1.96, \frac{\widehat{\Delta}^{\text{UB}} - \Delta_0^{\text{UB}}}{\widehat{\sigma}_{\text{UB}}} < -1.96\right]$, when $\frac{\widehat{\Delta}^{\text{LB}} - \Delta_0^{\text{LB}}}{\widehat{\sigma}_{\text{LB}}}, \frac{\widehat{\Delta}^{\text{UB}} - \Delta_0^{\text{UB}}}{\widehat{\sigma}_{\text{UB}}}$ is standard bivariate normal.

where $y_{p_0}^*$ is the p_0^{th} quantile of the control group's observed distribution of Y . In other words, when $p_0 > 0$, the upper tail of the treatment group's Y distribution is trimmed, as described above; but when $p_0 < 0$, the *lower* tail of the *control* group is trimmed for exactly same reasoning as described in the previous section. There is an analogous expression for Δ_0^{UB} .

Replacing the above population quantities with their sample analogues, an estimator for the bounds in this less restrictive model becomes

$$\begin{aligned}\widehat{\Delta}^{\text{LB}} &= 1[\hat{p} \geq 0] \cdot \widehat{\Delta}^{\text{LB}} + 1[\hat{p} < 0] \cdot \widehat{\Delta}^{\text{LB}*} \\ \widehat{\Delta}^{\text{UB}} &= 1[\hat{p} \geq 0] \cdot \widehat{\Delta}^{\text{UB}} + 1[\hat{p} < 0] \cdot \widehat{\Delta}^{\text{UB}*},\end{aligned}$$

where $\widehat{\Delta}^{\text{LB}*}$ and $\widehat{\Delta}^{\text{UB}*}$ are the analogous bounds when the control groups are trimmed.²³ As long as $p_0 \neq 0$, $\widehat{\Delta}^{\text{LB}}$ is consistent because it is a function of consistent estimators \hat{p} , $\widehat{\Delta}^{\text{LB}}$, and $\widehat{\Delta}^{\text{LB}*}$, and the function is continuous at the true parameter values of those estimators.

It follows from the delta method that the above estimator is also asymptotically normal with

$$\begin{aligned}\sqrt{n}(\widehat{\Delta}^{\text{LB}} - \Delta_0^{\text{LB}}) &\xrightarrow{d} N(0, 1[p_0 \geq 0]\{V^{\text{LB}} + V_C\} + 1[p_0 < 0]\{V_T + V_C^{\text{UB}}\}) \\ \sqrt{n}(\widehat{\Delta}^{\text{UB}} - \Delta_0^{\text{UB}}) &\xrightarrow{d} N(0, 1[p_0 \geq 0]\{V^{\text{UB}} + V_C\} + 1[p_0 < 0]\{V_T + V_C^{\text{LB}}\}),\end{aligned}$$

where the variance for the untrimmed treatment mean V_T is analogous to V_C defined previously and V_C^{UB} and V_C^{LB} use the analogous expressions in Proposition 3 but for the control group.

To summarize, suppose the researcher is unsure about the sign of p_0 , but knows that p_0 is nonzero. As an overall procedure, it is asymptotically valid to estimate \hat{p} , and if positive, trim the treatment group and conduct inference as discussed in Subsections 4.1 and 4.2. And if negative, trim the control group instead and conduct inference using the same formulas (*i.e.*, let $D^* = 1 - D$ and replace D everywhere with D^*). The intuition behind this is that as sample size increases, and the sampling variability of \hat{p} shrinks, the probability that the “wrong” group (treatment or control) is trimmed, leading to the wrong asymptotic variance being used, vanishes.

It is useful to consider the asymptotic behaviour of this estimator when $p_0 = 0$. In the Appendix, the estimator is shown to remain consistent, even without bounded support. Intuitively, the amount of trimming vanishes with sample size and so the trimmed mean converges to the (unbiased) untrimmed mean. On the other hand, it is clear that conventional first-order asymptotics will not apply. Close inspection of the above expressions reveals that keeping all other parameters constant, the asymptotic variance of either of the bounds is in general discontinuous at $p_0 = 0$. Specifically, when p_0 approaches zero from the right, the third component of the variance of the trimmed treatment mean will in general converge to a quantity that differs from the third component that must appear for the variance of the trimmed control mean when p_0 becomes negative.

This leads to two practical implications. First, when the researcher knows p_0 to be exactly zero, the above asymptotic expressions do not apply. Second, in the case when $p_0 \neq 0$, even

23. That is, more formally, $\widehat{\Delta}^{\text{LB}*} \equiv \frac{\sum Y \cdot S \cdot D}{\sum S \cdot D} - \frac{\sum Y \cdot S \cdot (1-D) \cdot 1[Y \geq \widehat{y}_{p^*}^*]}{\sum S \cdot (1-D) \cdot 1[Y \geq \widehat{y}_{p^*}^*]}$ and $\widehat{\Delta}^{\text{UB}*} \equiv \frac{\sum Y \cdot S \cdot D}{\sum S \cdot D} - \frac{\sum Y \cdot S \cdot (1-D) \cdot 1[Y \leq \widehat{y}_{1-p^*}^*]}{\sum S \cdot (1-D) \cdot 1[Y \leq \widehat{y}_{1-p^*}^*]}$, with $\widehat{y}_q^* \equiv \min\{y : \frac{\sum S \cdot (1-D) \cdot 1[Y \leq y]}{\sum S \cdot (1-D)} \geq q\}$ and $\widehat{p}^* \equiv (\frac{\sum S \cdot (1-D)}{\sum (1-D)} - \frac{\sum S \cdot D}{\sum D}) / (\frac{\sum S \cdot (1-D)}{\sum (1-D)}).$

though coverage rates for confidence intervals are asymptotically correct, a large discontinuity in the asymptotic variance suggests coverage rates may be inaccurate when sample sizes are small and p_0 is “close” to zero, which would imply that the “wrong” group is being trimmed with nontrivial probability in repeated samples.²⁴

It is useful to note, however, that for any finite sample size, as p_0 approaches zero, the confidence interval constructed from the *untrimmed* estimator will have coverage for the parameter of interest that approaches the correct rate since the bias (the difference between the untrimmed population mean and the population trimmed mean) is continuous in p_0 and equal to zero at $p_0 = 0$. Therefore, the untrimmed estimator for the treatment effect may have better coverage rates in a finite sample, even though its coverage will be zero asymptotically. Thus, at a minimum, it seems worthwhile for the researcher to additionally report the untrimmed estimator and S.E. A simple, conservative approach to combining the trimmed and untrimmed intervals is to compute their union. In repeated finite samples, at p_0 arbitrarily close to zero, this guarantees at least nominal coverage.²⁵

The issue of the estimator’s finite sample behaviour when p_0 is close to zero has some similarities to that regarding inference in instrumental variables when the first-stage coefficient is close to zero. Just as instrumental variables presumes the existence of a first stage, here, we presume that there is a nontrivial selection problem (p_0 nonzero). In both cases, first-order asymptotic approximations may be inadequate in finite samples when the nuisance parameter (here, p_0) is close to zero. The problem for instrumental variables is indeed nontrivial and has motivated a number of theoretical papers focusing on inference with weak instruments.²⁶

5. EMPIRICAL RESULTS

This section uses the trimming estimator to compute bounds on the treatment effect of the Job Corps on wage rates. The procedure is first employed for wages at week 208, 4 years after the date of random assignment. The widths of the bounds are reasonably narrow and are suggestive of positive wage effects of the program. The bounds for the effect at week 208 do contain zero, but the bounds at week 90 do not. Overall, the evidence presented below points towards a positive treatment effect, but not significantly more than a 10% effect.

5.1. Main results at week 208

Table 4 reports the estimates of the bounds of the treatment effect on wages at week 208. The construction of the bounds and their S.E. are illustrated in the table. Rows (iii) and (vi) report the means of log wages for the treated and control groups. Rows (ii) and (v) report that about 61% of the treated group have nonmissing wages, while about 57% of the control group have nonmissing wages. This implies a trimming proportion of about 6-8% of the treated group sample. The p th quantile is about 1.64, and therefore, the upper bound for the treated group is the mean after

24. As can be seen from the asymptotic expressions above, the discontinuity in the asymptotic variance disappears when the treatment and control groups have similar scale, in the sense that $\bar{y} - \mu_T$ for the treatment group is equal to $\mu_C - \underline{y}$ for the control group, where μ_T and μ_C are the untrimmed treatment and control means, and \bar{y} and \underline{y} are the population maximum and minimum for the treatment and control groups, respectively.

25. It should also be recalled that the untrimmed estimator lies between the point estimators of the two bounds with probability 1, and therefore, it may well be with many applications and sample sizes the untrimmed confidence interval may be contained in the trimmed confidence interval with high probability, meaning that inferences based on the trimming bounds would be too conservative.

26. See, e.g., Staiger and Stock (1997) and Andrews, Moreira and Stock (2007) and the references therein. Although there are some similarities, the trimming problem presented here is quite distinct from the IV case. For one, the bounds are still identified and the proposed estimator is still consistent (with bounded support) even when $p_0 = 0$.

TABLE 4
Bounds on treatment effects for ln(wage) in week 208 using trimming procedure

Control	(i) Number of observations	3599	Control S.E.	
	(ii) Proportion of nonmissing	0.566	S.E.	0.0082
	(iii) Mean ln(wage) for employed	1.997	Treatment upper bound S.E.	
Treatment	(iv) Number of observations	5546	Component 1	0.0053
	(v) Proportion nonmissing	0.607	Component 2	0.0021
	(vi) Mean ln(wage) for employed	2.031	Component 3	0.0083
			Total	0.0100
	$p = [(v) - (ii)] / (v)$	0.068	Treatment lower bound S.E.	
	(vii) p th quantile	1.636	Component 1	0.0058
	(viii) Trimmed mean: $E[Y Y > y_p]$	2.090	Component 2	0.0037
	(ix) $(1-p)$ th quantile	2.768	Component 3	0.0144
	(x) Trimmed mean: $E[Y Y < y_{1-p}]$	1.978	Total	0.0159
			Effect	
Effect	(xi) Upper bound estimate = (viii) - (iii)	0.093	(xiii) Upper bound S.E.	0.0130
	(xii) Lower bound estimate = (x) - (iii)	-0.019	(xiv) Lower bound S.E.	0.0179
Confidence interval 1 = [(xii) - 1.96 × (xiv), (xi) + 1.96 × (xiii)]				[-0.055 to 0.119]
Confidence interval 2 (Imbens and Manski) = [(xii) - 1.645 × (xiv), (xi) + 1.645 × (xiii)]				[-0.049 to 0.114]
Heckman two-step estimator				0.0148 (0.0117)
Das <i>et al.</i> (2003)				0.0140 (0.0122)

Notes: Before trimming, there are 3371 nonmissing observations in the treatment group. After trimming, there are 3148 (3142) observations remaining in the treatment group after trimming the lower p (upper $1-p$) of the distribution (these numbers are not equal due to using the design weights). For the upper bound S.E., component 1 is the usual S.E. of the mean, using the trimmed sample. Component 2 is the square root of $(1/3371) \times (p/(1-p)) \times \{(viii) - (vii)\}^2$. Component 3 is the square root of $\{((viii) - (vii))/(1-p)\}^2 \times \text{Var}(p)$, where $\text{Var}(p) = (1-p)^2 \times \{(1/5546) \times ((1-(v))/(v)) + (1/3599) \times ((1-(ii))/(ii))\}$. "Total" refers to the square root of the sum the squared components. The entries for the treatment lower bound S.E. are defined analogously. (xiii) and (xiv) are the square root of the sum of the squared S.E. for the treatment upper bound (or lower bound) and control group. For the Imbens and Manski confidence interval 1.645 satisfies $F(1.645 + ((xi) - (xii))/(\max((xiii), (xiv)))) - F(-1.645) = 0.95$, where F is the standard normal c.d.f. See Imbens and Manski (2004) for details. The Heckman two-step estimator uses months employed in the previous year and treatment status in the first-stage probit. The Das *et al.* (2003) estimator is described in the text.

trimming the tail of the distribution below 1.64.²⁷ After trimming, the resulting mean is about 2.09, and so the upper bound of the treatment effect Δ^{UB} is 0.093 (row (xi)). A symmetric procedure yields Δ^{LB} of -0.019 (row (xii)).

The width of these bounds is about 0.11. Note that this is 1/14th the width of the bounds yielded by existing "imputation" procedures as reported in Table 3 (calculate 1.55 from rows (xi) and (xii)). The much larger interval in Table 3 is clearly driven by the relatively wide support of the outcome variable.²⁸ The difference between the two sets of bounds makes an important difference in gauging the magnitude of the effects of the program. From Table 3, the negative region covered by the bounds is almost as large as the positive region contained by the bounds. In this sense, the bounds from Table 3 are almost as consistent with large negative effects as they are with large positive effects.

27. The procedure can be easily adapted to the case of a dependent variable with discrete support, which can generate "ties" in the data. After sorting the data by the dependent variable, unique ranks can be imposed (*i.e.*, so that individuals with the exact same wage level all have different ranks). The correct proportion of data can be trimmed based on those ranks, before calculating the trimmed mean, which is based on the remaining data. This procedure was used here, with the slight modification that the design weights were used, so the observations were dropped until the accumulated sum of the weights equalled the trimming proportion times the total sum of the weights in the treatment group.

28. For a detailed theoretical discussion of how the imputation bounds (*e.g.*, Table 3) compare to the trimming bounds (*e.g.*, Table 4) when the outcome is binary, see Lee (2002).

The width of the trimming bounds in Table 4 is also narrow enough to rule out plausible effect sizes. For example, suppose the training component of the Job Corps program was ineffective at raising the marketable skills of the participants. We would then expect Job Corps to have a negative impact on wages, insofar as the time spent in the program caused a delay in accumulating labour market experience.

Suppose annual wage growth is about 8% a year, and the program group spent more time in education and training programs than the control group by an amount equivalent to 0.72 of a school year.²⁹ If a full school year in training causes a year delay in earnings growth, this would imply Job Corps impact of about -0.058 . The lower bound in Table 4 is -0.019 . Thus, the scenario described above is ruled out by the trimming bounds computed in Table 4. By contrast, an impact of -0.058 is easily contained by the support-dependent interval $[-0.746, 0.802]$ of Table 3.

An impact of -0.058 is also outside the interval after accounting for sampling errors of the estimated bounds. The right side of Table 4 illustrates the construction of these S.E. For the estimate of the upper bound for the treatment group, component 1 is the S.E. associated with the first term in Equation (7).³⁰ Component 2 reflects sampling error in estimating the trimming threshold.³¹ Component 3 reflects sampling error in estimating the trimming proportion.³² In this case, the largest source of the variance in the upper bound comes from the estimation of the trimming proportion. The total of 0.010 is the square root of the sum of the squared components.

Doing a similar calculation for the lower bound, and then using the S.E. on the mean for the control group, yields S.E. for $\widehat{\Delta}^{UB}$ and $\widehat{\Delta}^{LB}$ of 0.0130 and 0.0179, as shown in the bottom of Table 4. These S.E. can then be used to compute two types of 95% confidence intervals. The first covers the entire set of possible treatment effects with at least 0.95 probability, while the second interval, using the result from Imbens and Manski (2004), covers the true treatment effect at least 95% of the time. A plausible negative impact of -0.058 is outside both of these intervals.

As argued previously, the Job Corps data do not seem to include a plausible instrument for selection. Nevertheless, it is useful to compare the bounding inference to conventional parametric and nonparametric sample selection estimators that do rely on exclusion restrictions. The bottom of Table 4 presents both a Heckman two-step estimator and the nonparametric estimator of Das *et al.* (2003). Both use the “Months Employed in Previous Year” variable to predict sample selection.³³

29. From Figure 2, there appears to be about 40% nominal wage growth more than 4 years. Inflation over that length of time in the late 1990's was about 9% (CPI-U (Consumer Price Index for all Urban Consumers) for 1995: 152.4 and for 1999: 166.6). Schochet *et al.* (2001) find that the Job Corps impact on time spent in any education and training programs amounted to about 1 school year per participant. The estimated impact per eligible applicant was 28% lower.

30. Specifically, it is the square root of the sample analog of $\frac{1}{n^{TRIM}} \text{Var}[Y|D=1, S=1, Y \geq y_{p_0}]$, where n^{TRIM} is the number of observations after trimming.

31. It is the square root of the sample analog of $\frac{1}{n^{UNTRIM}} \frac{(y_{p_0} - \mu^{UB})^2 p_0}{(1-p_0)}$, where n^{UNTRIM} is the number of nonmissing observations before trimming.

32. It is the square root of the sample analog of $(y_{p_0} - \mu^{UB})^2 \left(\frac{1}{n^T} \frac{1 - \frac{\alpha_0}{1-p_0}}{\frac{\alpha_0}{1-p_0}} + \frac{1}{n^C} \frac{1 - \alpha_0}{\alpha_0} \right)$, where n^T and n^C are the number of treatment and control observations (missing and nonmissing) in the sample.

33. Specifically, for the Heckman two-step estimator, selection status was the dependent variable in a first-step probit including the treatment status and months employed. The predicted inverse Mill's ratio was used as an additional regressor in a regression of wages at week 208 on treatment status. For the estimator of Das *et al.* (2003), the probability of selection was predicted from a regression of selection status on treatment months employed, their interaction and the square of months employed. The second-stage regressed wages at week 208 on treatment status and the predicted probability. As in Das *et al.* (2003), the orders of the polynomials and interactions for both first and second stages were determined by cross-validation.

TABLE 5

Bounds on treatment effects for $\ln(\text{wage})$ in week 208 trimming procedure using baseline covariates

Group	Lower bound for treatment mean			Upper bound for treatment			Weight
	Estimate	S.E.	Observation	Estimate	S.E.	Observation	
1	1.795	0.030	343	1.979	0.025	348	0.107
2	1.938	0.052	248	1.963	0.065	250	0.131
3	1.934	0.020	931	2.051	0.017	935	0.291
4	2.025	0.028	745	2.127	0.020	748	0.238
5	2.121	0.025	712	2.204	0.022	715	0.234
Total	1.985	0.013	2979	2.086	0.012	2996	1.000
Effect	Lower bound for effect			Upper bound for effect			
	-0.0118	0.0151		0.0889	0.0142		

Notes: Trimming procedure from Table 3 applied separately to each group (defined in text). “Total” estimates are means of the five groups using the “Weight” as weights. Asymptotic variance for “Total” is computed according to Chamberlain (1993): it is the (weighted, using “Weight”) average of the asymptotic variance for each group (each group’s sampling variance times the number of observations for the group) plus the (weighted by “Weight”) average squared deviation of each group’s estimate from the “Total” mean. Control mean, (iii) in Table 4, is then subtracted to obtain bounds on the treatment effect.

5.2. Using covariates to narrow bounds

The construction of bounds that use the baseline covariates, as presented in Proposition 1b, is illustrated using a variable that splits the sample into five mutually exclusive groups, based on their observed baseline characteristics. Any baseline covariate will do, as will any function of all the baseline covariates. In the analysis here, a single baseline covariate—which is meant to be a proxy for the predicted wage potential for each individual—is constructed from a linear combination of all observed baseline characteristics. This single covariate is then discretized, so that effectively five groups are formed according to whether the predicted wage is within intervals defined by \$6.75, \$7, \$7.50, and \$8.50.³⁴

Then, a trimming analysis is conducted for each of the five groups separately. Note that for each of the five groups, there is a different trimming proportion. The lower and upper bounds of the treatment group means, by each of the five groups, are given in the left and right columns of Table 5, respectively. The lower bounds range from 1.80 to 2.12, while the upper bounds range from 1.96 to 2.20. The S.E. are computed for each group separately in the same manner as in Table 4.

To compute the bounds for the overall average $E[Y_1^* | S_0 = 1, S_1 = 1]$, the group-specific bounds must be averaged, weighted by the proportions $\Pr[\text{Group } J | S_0 = 1, S_1 = 1]$. This is provided in the row labelled “Total”.³⁵ This leads to an interval of $[-0.0118, 0.0889]$. This interval is about 11% narrower than that reported in Table 4. The estimated asymptotic variance for these overall averages is the sum of (1) a weighted average of the group-specific variances and (2) the (weighted-) mean squared deviation of the group-specific estimates from the overall mean. This second term takes into account the sampling variability of the weights, as described in

34. Specifically, the coefficients from the linear combination of the X s are the coefficients from a regression of week 208 wages on all baseline characteristics in Table 1. The coefficients were then applied to *all* individuals to impute a predicted wage.

35. There are slight differences in the number of observations in each group after trimming, for the upper and lower bounds. This is due to the use of the design weights.

TABLE 6
Treatment effect estimates and bounds, by week

	Fraction nonmissing		Trimming proportion	Effect		
	Control	Treatment		Untrimmed	Lower bound	Upper bound
Week 45	0.4223	0.3424	0.1892 (0.0219)	0.022 (0.011)	-0.074 (0.014)	0.127 (0.015)
Week 90	0.4600	0.4601	0.0003 (0.0232)	0.043 (0.011)	0.042 (0.024)	0.043 (0.025)
Week 135	0.5173	0.5451	0.0509 (0.0192)	0.028 (0.011)	-0.016 (0.021)	0.076 (0.014)
Week 180	0.5403	0.5825	0.0724 (0.0177)	0.026 (0.011)	-0.033 (0.019)	0.087 (0.013)

Notes: ($N = 9145$ for each row). S.E. are given in parentheses. S.E. for trimming proportion given by formula in note to Table 4. Bounds computed according to Table 4. See text for details.

Chamberlain (1994).³⁶ These sampling errors lead to a 95% Imbens–Manski interval of $[-0.037, 0.112]$.

By statistically ruling out any effect more negative than -0.037 , this suggests that after 4 years, the Job Corps enabled program group members to offset at least 35% (and perhaps more) of the potential 0.058 loss in wages due to lost labour market experience that could have been caused by the program.

5.3. Effects by time horizon and testable implications

An analysis of the bounds at different time horizons provides further evidence that the Job Corps program had a positive impact on wage rates. The analysis of Table 4 was performed for impacts on wage rates at weeks 45, 90, 135, and 180, and these results are reported in Table 6.

As would be expected, the widths of the intervals are directly related to the treatment control difference in the proportion missing. When the proportion is the largest, as at week 45, the range is $[-0.074, 0.127]$. At week 180, when the proportion is 0.0724, the interval is $[-0.033, 0.087]$.

At week 90, the estimated trimming proportion is close to zero, and the resulting bounds are given by the interval $[0.042, 0.043]$. Maintaining the assumption that the true $p_0 \neq 0$, we note that the S.E. are larger for these bounds, even though they are quite similar to the untrimmed treatment control difference. This is partly due to the sampling error in the trimming proportion. Using these S.E. and the Imbens and Manski (2004) confidence interval for the treatment effect, *parameter* is computed to be $[-0.004, 0.092]$. As noted above, if the true trimming proportion p_0 is arbitrarily close to zero then the untrimmed confidence interval will have almost accurate coverage in a finite sample. This untrimmed treatment effect confidence interval is $[0.020, 0.065]$. Thus, both procedures can rule out effects more negative than -0.004 at conventional levels of significance.

If we were to alternatively assume that $p_0 = 0$ at week 90 then one can provide limited evidence on the plausibility of the monotonicity condition (Assumption 2b). If at week 90, $E[S|D = 1] - E[S|D = 0]$ is truly zero then the average causal effect on sample selection

36. The weighted mean of the five group-specific means can be seen as a minimum distance estimator where the weights are the estimated proportions in each group. Chamberlain (1994) gives the asymptotic variance for this estimator even when the moment vector is misspecified, as would be the case if the group-specific means are different. The asymptotic variance is the sum of two components: (1) the (weighted) average of the asymptotic variance for each group (Λ_1 in Chamberlain, 1994) and (2) the (weighted) average squared deviation of each group's estimate from the "Total" mean (Λ_2 in Chamberlain, 1994).

$E[S_1 - S_0]$ is zero. If monotonicity holds then this can only be true if $S_1 = S_0$ with probability 1.³⁷

If the only observed data are the triple (Y, S, D) then it is impossible to test this monotonicity assumption. On the other hand, if there exist baseline characteristics X , as in the case of the Job Corps experiment, then it is possible to test whether $S_0 = S_1$ with probability 1. That is, it is possible to test whether for each value of X , $\Pr[S = 1|D = 1, X = x] = \Pr[S_1 = 1|X = x]$ is equal to $\Pr[S = 1|D = 0, X = x] = \Pr[S_0 = 1|X = x]$, which should be the case for all x if $S_0 = S_1$ with probability 1. Intuitively, if it was found that for some values of X , the treatment caused wages to be observed, while for other values of X , the treatment was found to cause wages to be missing, then Assumption 2a must not hold.

By Bayes' Rule and independence (Assumption 1), $\Pr[S = 1|D = 1, X = x] = \Pr[S = 1|D = 0, X = x]$ for all x implies that the distribution of X conditional on $S = 1, D = 1$ should be the same as the distribution conditional on $S = 1, D = 0$. This is because the density of X , conditional on D , does not depend on the value of D , and the probability of $S = 1$ conditional on D also does not depend on D , by assumption.

A simple way to check this empirically is to examine the means of the variables in Table 1, but *conditional* on having nonmissing wages. This is done for week 90 and is reported in the Appendix, Table A1. The differences between the treatment and the control means for each variable are small and consistently statistically insignificant. A joint test of significance is given by a logistic regression of the treatment indicator on the baseline characteristics X , using a sample of all those with nonmissing wages at week 90.³⁸ The resulting test of all coefficients equalling zero yields a p value of 0.851. Thus, the data are consistent with the monotonicity condition holding at week 90.

6. CONCLUSIONS: IMPLICATIONS AND APPLICATIONS

This paper focuses on an important issue in evaluating the impact of a job training program on wage rates—the sample selection problem. It is a serious issue even when the treatment of a training program is believed to be independent of all other factors, as was the case in the randomized experimental evaluation of the U.S. Job Corps. Existing sample selection correction methods are infeasible due to the absence of plausible exclusion restrictions, and in this case, one cannot rely upon the boundedness of the outcome variable's support to yield informative bounds on the treatment effect of interest.

To estimate the impact of the Job Corps on wages, this paper develops a new method for bounding treatment effects in the presence of sample selection in the outcome. An appealing feature of the method is that the assumptions for identification, independence, and monotonicity are typically already assumed in standard models of the sample selection process, such as in Equation (1). In the case of randomized experiments, the independence assumption is satisfied, and as illustrated in the previous section, the existence of baseline characteristics suggests a limited test of monotonicity. More importantly, the bounding approach does not require any exclusion restrictions for the outcome equation. Nor do the trimming bounds rely on the bounds of the support of the outcome variable.

The analysis using the proposed “trimming” bounds points to two substantive conclusions about the Job Corps. First, the evidence casts doubt on the notion that the program only raised earnings through raising labour force participation. Effects more negative than -0.037 can be

37. If $S_1 = S_0$ with less than probability 1 then there would be a nonzero probability of $S_1 < S_0$ and it would be equal to the probability of $S_0 > S_1$ (in order for $E[S_1 - S_0] = 0$). This would contradict monotonicity.

38. This is a valid test since in this context, $\Pr[S = 1|D = 1, X = x] = \Pr[S = 1|D = 0, X = x]$ for all x and is equivalent to the test $\Pr[D = 1|S = 1, X = x] / \Pr[D = 0|S = 1, X = x] = \Pr[D = 0] / \Pr[D = 1]$.

statistically ruled out. If there were literally no wage effect, one might expect to see a more negative impact (perhaps around a -0.058 effect) due to lost labour market experience since these young applicants are on the steep part of their wage profile.

Another reason to interpret the evidence as pointing to positive wage effects is that the lower bound is based on an extreme and unintuitive assumption—that wage outcomes are perfectly *negatively* correlated with the propensity to be employed. From a purely theoretical standpoint, a simple labour supply model suggests that, all other things equal, those on the margin of being employed will have lowest wages not the highest wages (*i.e.*, the “reservation wage” will be the smallest wage that draws the individual into the labour force).³⁹ In addition, the empirical evidence in Table 2 suggests that there is positive selection into employment: those who are predicted to have higher wages are more likely to be employed (*i.e.*, U and V are positively correlated). If this is true, it seems relatively more plausible to trim the lower rather than the upper tail of the distribution to get an estimate of the treatment effect.

Second, the intervals provided here are comparable to rates of return found in the returns to education literature. At week 208, the point estimates an interval of $[-0.0118, 0.0889]$. Program participants may be lagging behind their control counterparts by as much as 8 months in labour market experience due to enrolment in the program. As argued above, this could translate to as much as a 5.8% wage disadvantage even 4 years after random assignment because many of the individuals in this sample are still on the steep part of their age-earnings profiles. Projecting to ages when the wage profile flattens leads to an interval of $[0.047, 0.145]$. A similar adjustment for week 90 wages yields an interval tightly centred around 0.10. As found in a survey of studies that exploit institutional features of school systems (Card, 1999), point estimates of the return to a single year of schooling range from 0.060 to 0.153.⁴⁰ Thus, the magnitudes found in this analysis of the Job Corps are roughly consistent with viewing the program as a human capital investment of 1 year of schooling.

It should be emphasized that the trimming bounds introduced here are specific neither to selection into employment nor to randomized experiments. For example, outcomes can be missing due to survey nonresponse (*e.g.*, students not taking tests), sample attrition (*e.g.*, inability to follow individuals over time), or other structural reasons (*e.g.*, mortality). As long as the researcher believes that the sample selection process can be written as a model like Equation (1) or (5), the same trimming method can be applied. Also, the basis for matching estimators for evaluations is the weaker assumption that (Y_1^*, Y_0^*) is independent of D , conditional on X , rather than Assumption 3. It is immediately clear that the trimming bounds proposed here can be applied even when (Y_1^*, Y_0^*, S_0, S_1) is independent of D , but only conditional on X , as long as Assumption 2b holds conditional on X . In this situation, the procedure described in Subsection 5.2 can be applied.⁴¹

MATHEMATICAL APPENDIX

Lemma 1. Let Y be a continuous random variable and a mixture of two random variables, with c.d.f.s $M^*(y)$ and $N^*(y)$, and a known mixing proportion $p^* \in [0, 1]$, so that we have $F^*(y) = p^*M^*(y) + (1 - p^*)N^*(y)$. Consider $G^*(y) = \max\left[0, \frac{F^*(y) - p^*}{1 - p^*}\right]$, which is the c.d.f. of Y after truncating the p^* lower tail of Y . Then, $\int_{-\infty}^{\infty} y dG^*(y) \geq \int_{-\infty}^{\infty} y dN^*(y)$. $\int_{-\infty}^{\infty} y dG^*(y)$ is a sharp (in the sense of Horowitz and Manski, 1995) upper bound for $\int_{-\infty}^{\infty} y dN^*(y)$.

39. A recent example of the use of more restrictive assumptions that rule out extreme or arguably unintuitive (negative) correlations in the labour supply context is found in the analysis of Blundell, Gosling, Ichimura and Meghir (2007).

40. See table 4 in Card (1999).

41. But it should be noted that since the baseline characteristics X would no longer be independent of the treatment, one could no longer use Remark 4 to test the monotonicity assumption.

TABLE A1

Summary statistics, by treatment status, National Job Corps Study conditional on positive earnings in week 90

Variable	Control		Program		Difference	
	Proportion of nonmissing	Mean	Proportion of nonmissing	Mean	Difference	S.E.
Female	1.00	0.429	1.00	0.419	-0.009	0.016
Age at baseline	1.00	18.691	1.00	18.729	0.038	0.068
White, non-Hispanic	1.00	0.310	1.00	0.328	0.018	0.015
Black, non-Hispanic	1.00	0.447	1.00	0.443	-0.004	0.016
Hispanic	1.00	0.171	1.00	0.167	-0.004	0.012
Other race/ethnicity	1.00	0.072	1.00	0.063	-0.009	0.008
Never married	0.99	0.909	0.99	0.909	0.000	0.009
Married	0.99	0.030	0.99	0.023	-0.007	0.005
Living together	0.99	0.039	0.99	0.045	0.006	0.006
Separated	0.99	0.022	0.99	0.022	0.001	0.005
Has child	0.99	0.188	1.00	0.178	-0.009	0.012
Number of children	0.99	0.247	0.99	0.241	-0.007	0.019
Education	0.99	10.381	0.98	10.371	-0.010	0.050
Mother's education	0.83	11.506	0.84	11.579	0.072	0.090
Father's education	0.66	11.644	0.67	11.458	-0.186	0.111
Ever arrested	0.99	0.238	0.99	0.232	-0.006	0.013
Household income						
<3000	0.68	0.188	0.66	0.202	0.014	0.015
3000-6000	0.68	0.188	0.66	0.182	-0.006	0.015
6000-9000	0.68	0.116	0.66	0.119	0.003	0.012
9000-18,000	0.68	0.289	0.66	0.270	-0.019	0.017
>18,000	0.68	0.219	0.66	0.227	0.008	0.016
Personal income						
<3000	0.95	0.726	0.93	0.732	0.005	0.014
3000-6000	0.95	0.164	0.93	0.154	-0.010	0.012
6000-9000	0.95	0.065	0.93	0.068	0.003	0.008
>9000	0.95	0.045	0.93	0.047	0.002	0.007
At baseline						
Have job	0.98	0.251	0.98	0.254	0.002	0.014
Months employed, previous year	1.00	4.572	1.00	4.558	-0.013	0.143
Had job, previous year	0.99	0.725	0.99	0.727	0.002	0.014
Earnings, previous year	0.94	3783.940	0.94	3699.524	-84.416	159.333
Usual hours/week	1.00	24.600	1.00	25.165	0.565	0.642
Usual weekly earnings	1.00	125.147	1.00	126.297	1.150	3.838
After random assignment						
Week 90 ln(wage)	1.00	1.827	1.00	1.870	0.043*	0.011
Number of observations	1660		2564			

Notes: $N = 4224$. Computations use design weights. Chi-square test of all coefficients equalling zero, from a logit of the treatment indicator on all baseline characteristics (where mean values were imputed for missing values) yields 19.50; associated p value from a chi-squared (27 df) distribution is 0.851.

*Indicates difference is statistically significant from 0 at the 5% level.

Proof of Lemma. See Horowitz and Manski (1995), corollary 4.1.

Proof of Proposition 1a. It suffices to show that $\mu^{UB} \equiv E[Y|D=1, S=1, Y \geq y_{p_0}]$ is a sharp upper bound for $E[Y_1^*|S_0=1, S_1=1]$. A similar argument for the sharp lower bound would follow. Assumptions 1 and 2a imply that $p_0 = \frac{\Pr[S=1|D=1] - \Pr[S=1|D=0]}{\Pr[S=1|D=1]} = \frac{\Pr[S_0=0, S_1=1|D=1]}{\Pr[S=1|D=1]}$. Let $F(y)$ be the c.d.f. of Y conditional on $D=1, S=1$. Assumption 2a implies that $F(y) = p_0 M(y) + (1-p_0)N(y)$, where $M(y)$ denotes the c.d.f. of Y_1^* , conditional on $D=1, S_0=0, S_1=1$, and $N(y)$ denotes the c.d.f. of Y_1^* , conditional on $D=1, S_0=1, S_1=1$. By Assumption 1, $N(y)$ is also the c.d.f. of Y_1^* , conditional on $S_0=1, S_1=1$. By the lemma, $\mu^{UB} \equiv \frac{1}{1-p_0} \int_{y_{p_0}}^{\infty} y dF(y) \geq \int_{-\infty}^{\infty} y dN(y) = E[Y_1^*|S_0=1, S_1=1]$.

To show that μ^{UB} equals the maximum possible value for $E[Y_1^*|S_0=1, S_1=1]$ that is consistent with the distribution of the observed data on (Y, S, D) , it must be shown that (1) conditional on p_0 , μ^{UB} is a sharp upper bound and (2) p_0 is uniquely determined by the data. (1) follows from the lemma. (2) is true because the data yield a unique probability function $\Pr[S=s, D=d], s, d=0, 1$, which uniquely determines p_0 .

To show that $[\Delta_0^{LB}, \Delta_0^{UB}]$ is contained in any other valid bounds that impose the same assumptions, it suffices to show that any Δ strictly within the interval $[\Delta_0^{LB}, \Delta_0^{UB}]$ cannot be ruled out by the observed data, note

that $\Delta_0^{\text{UB}} \geq E[Y|D=1, S=1] - E[Y|D=0, S=1] \geq E[Y|D=1, S=1, Y < y_{p_0}] - E[Y|D=0, S=1]$. Therefore, for any Δ between $E[Y|D=1, S=1] - E[Y|D=0, S=1]$ and Δ_0^{UB} , there exists $\lambda \in [0, 1]$ such that $\Delta = \lambda \Delta_0^{\text{UB}} + (1-\lambda)\{E[Y|D=1, S=1, Y < y_{p_0}] - E[Y|D=0, S=1]\}$. With this λ , we can construct (1) a density of Y_1^* conditional on $S_0=1, S_1=1$ as $\lambda g(y) + (1-\lambda)h(y)$ and (2) a density of Y_1^* conditional on $S_0=0, S_1=1$ being $\left(\frac{1-p_0}{p_0} - \frac{1-p_0}{p_0}\lambda\right)g(y) + \left(1 - \frac{1-p_0}{p_0}(1-\lambda)\right)h(y)$, where $g(y)$ is the density of Y conditional on $Y \geq y_{p_0}$ and $h(y)$ is the density of Y conditional on $Y < y_{p_0}$. The mixture of these two latent densities, by construction, replicates the observed density of Y conditional on $D=1, S=1$; furthermore, by construction the mean of the constructed density of Y_1^* conditional on $S_0=1, S_1=1$ minus the control mean yields the proposed Δ . A symmetric argument can be made about any Δ in between $E[Y|D=1, S=1] - E[Y|D=0, S=1]$ and Δ_0^{LB} . Therefore, each Δ within the interval $[\Delta_0^{\text{LB}}, \Delta_0^{\text{UB}}]$ cannot be ruled out by the observed data. Q.E.D.

Proof of Proposition 2. It is sufficient to prove consistency for the trimmed mean for the treatment group, and only for the lower bound, since a symmetric argument will follow for the upper bound. Denote $\mu_0 \equiv E[Y|D=1, S=1, Y \leq y_{1-p_0}]$ as the true lower bound of interest. Consistency follows from applying theorem 2.6 of Newey and McFadden (1994), which applies to GMM estimators. Define the moment function

$$g(z, \theta) \equiv \begin{pmatrix} (Y - \mu)SD \cdot 1[Y \leq y_{1-p}] \\ (1[Y > y_{1-p}] - p)SD \\ \left(S - \alpha \frac{1}{1-p}\right)D \\ (S - \alpha)(1 - D) \end{pmatrix},$$

where $\theta' = (\mu, y_{1-p}, p, \alpha)'$, $\theta'_0 = (\mu_0, y_{1-p_0}, p_0, \alpha_0)'$, $\alpha_0 \equiv \Pr[S=1|D=0]$, and $z' = (Y, S, D)'$. The estimator of μ_0 , the lower bound of $E[Y_1^*|S_0=1, S_1=1]$, as provided in Equation (6), is a solution to $\min_{\theta} \left(\sum g(z, \theta)\right)' \cdot \left(\sum g(z, \theta)\right)$. From theorem 2.6, (i) holds because as long as $E[S|D=0] > 0$, this just-identified system yields only one solution, (ii) holds if we take the parameter space to be the bounds of the support for the trimmed mean and quantiles, and $[0, 1]$ to be the parameter space for the two probabilities α and p , (iii) continuity holds, and bounded support implies (iv). Q.E.D.

Proof of Proposition 3. As in the proof above, it is sufficient to focus only on the asymptotic properties of the estimator of μ_0 . This estimator will be independent of that for the (untrimmed) control group mean. The proof follows by showing that the conditions of theorem 7.2 of Newey and McFadden (1994) are satisfied.

Define $g_0(\theta) \equiv E[g(z, \theta)]$ and $\hat{g}_n(\theta) \equiv n^{-1} \sum g(z, \theta)$. (i) of theorem 7.2 holds. (iii) holds because by assumption, each of the parameters is in the interior of the parameter space defined in Proposition 2. (iv) holds by the central limit theorem. Let G be the derivative of $g_0(\theta)$ at $\theta = \theta_0$. An explicit expression for G , a square matrix, is given below and will be shown to be nonsingular; hence, (ii) holds as well.

The stochastic equicontinuity condition in (v) can be shown to hold using theorem 1 of Andrews (1994). Assumption C of this theorem holds, and Assumption A holds with envelope $\overline{M} = |Y - D\mu_0| + |D| \sup_{\mu} \|\mu_0 - \mu\|$ for the first element and 1 for the remaining elements of $g(z, \theta)$. Boundedness of the support implies $E|Y|^{2+\delta} < \infty$ for some $\delta > 0$, which implies that $E|\overline{M}|^{2+\delta} < \infty$ for some $\delta > 0$, and therefore, Assumption B holds as well.

From theorem 7.2 of Newey and McFadden (1994), the asymptotic variance is $V^{\text{LB}} = G^{-1} \Sigma (G')^{-1}$, where Σ is the asymptotic variance of $\hat{g}_n(\theta_0)$. After letting $\gamma' \equiv (\mu, y_{1-p})'$ and $\delta \equiv (p, \alpha)'$, it can be shown that G can be written as the partitioned matrix $\begin{pmatrix} G_{\gamma} & G_{\delta} \\ 0 & M_{\delta} \end{pmatrix}$ and Σ can be partitioned as $\begin{pmatrix} \Sigma_1 & 0 \\ 0 & \Sigma_2 \end{pmatrix}$. The upper left 2×2 block of V^{LB} can then be shown to be equal to $G_{\gamma}^{-1} \Sigma_1 (G_{\gamma}^{-1})' + G_{\gamma}^{-1} G_{\delta} M_{\delta}^{-1} \Sigma_2 \cdot (M_{\delta}^{-1})' G_{\delta}' (G_{\gamma}^{-1})'$. The first term contains the variance of the trimmed mean if the trimming proportion p_0 is known. The second term captures the variance due to the estimation of the trimming proportion.

Consider the first term. After computing $g_0(\theta)$, G_{γ} can be shown to equal

$$E[SD] \begin{pmatrix} -(1-p_0) & (y_{1-p_0} - \mu_0) f(y_{1-p_0}) \\ 0 & -f(y_{1-p_0}) \end{pmatrix},$$

where $f(\cdot)$ is the density of Y conditional on $D=1, S=1$. Σ_1 is equal to

$$\begin{pmatrix} \int_{-\infty}^{y_{1-p_0}} (y - \mu_0)^2 f(y) dy \cdot E[SD] & 0 \\ 0 & p_0(1-p_0)E[SD] \end{pmatrix}.$$

It follows that the upper left element of $G_{\gamma}^{-1} \Sigma_1 (G_{\gamma}^{-1})'$ is

$$\frac{1}{E[SD](1-p_0)} \left\{ \text{Var}[Y|D=1, S=1, Y \leq y_{1-p_0}] + (y_{1-p_0} - \mu_0)^2 p_0 \right\},$$

as stated in Equation (7).

Consider the second term. Direct calculation of G_{δ} , M_{δ} , and Σ_2 yields

$$G_{\delta} = E[SD] \begin{pmatrix} 0 & 0 \\ -1 & 0 \end{pmatrix}, M_{\delta} = \begin{pmatrix} -E[D] \alpha_0 \frac{1}{(1-p_0)^2} & -E[D] \frac{1}{1-p_0} \\ 0 & -(1-E[D]) \end{pmatrix}$$

$$\Sigma_2 = \begin{pmatrix} \frac{\alpha_0}{1-p_0} \left(1 - \frac{\alpha_0}{1-p_0}\right) E[D] & 0 \\ 0 & \alpha_0(1-\alpha_0)(1-E[D]) \end{pmatrix}.$$

After simplifying terms, it follows that the upper left element of $G_{\gamma}^{-1} G_{\delta} M_{\delta}^{-1} \Sigma_2 (M_{\delta}^{-1})' G'_{\delta} (G_{\gamma}^{-1})'$ is equal to

$$(y_{1-p_0} - \mu_0)^2 \left(\frac{\left(1 - \frac{\alpha_0}{1-p_0}\right)}{E[D] \left(\frac{\alpha_0}{1-p_0}\right)} + \frac{(1-\alpha_0)}{(1-E[D]) \alpha_0} \right),$$

as stated in Equation (7), after substituting in V^P .

Finally, direct computation of the upper left element of $M_{\delta}^{-1} \Sigma_2 (M_{\delta}^{-1})'$ yields the expression for V^P . Q.E.D.

Proof of Consistency When $p_0 = 0$. Assume $p_0 = 0$. We know that as long as $E|Y| < \infty$, the untrimmed treatment effect estimator $\hat{\Delta}$ converges to the true treatment effect Δ_0 . It is thus sufficient to show that for any $\delta > 0$, we have $\lim_{n \rightarrow \infty} \Pr[|\hat{\Delta}^{LB} - \hat{\Delta}| < \delta] = 1$. First, note that $\Pr[|\hat{\Delta}^{LB} - \hat{\Delta}| < \delta] = \Pr[|\hat{\Delta}^{LB} - \hat{\Delta}| < \delta | 0 \leq \hat{p} \leq \bar{p}] \Pr[0 \leq \hat{p} \leq \bar{p}] + \Pr[|\hat{\Delta}^{LB} - \hat{\Delta}| < \delta | \hat{p} > \bar{p}] \Pr[\hat{p} > \bar{p}] + \Pr[|\hat{\Delta}^{LB} - \hat{\Delta}| < \delta | 0 > \hat{p} \geq \bar{p}^*] \Pr[0 > \hat{p} \geq \bar{p}^*] + \Pr[|\hat{\Delta}^{LB} - \hat{\Delta}| < \delta | \hat{p} < \bar{p}^*] \Pr[\hat{p} < \bar{p}^*]$. Since $p_0 = 0$, for any positive \bar{p} and negative \bar{p}^* , the second and fourth terms converge to zero. Now consider the first term. Let \bar{p} be any positive value such that $\Delta_0 - \frac{\Delta^{LB}}{\bar{p}} < \delta$, where $\frac{\Delta^{LB}}{\bar{p}}$ is the population trimmed mean after trimming the top tail by the proportion \bar{p} . Now note that for any sample indexed by N , we have $\Pr[|\hat{\Delta}^{LB} - \hat{\Delta}| < \delta | 0 \leq \hat{p} \leq \bar{p}] = \int_0^{\bar{p}} \Pr[|\hat{\Delta}^{LB} - \hat{\Delta}| < \delta | \hat{p} = p] dF_N(p)$, where F_N is the c.d.f. of \hat{p} conditional on $0 \leq \hat{p} \leq \bar{p}$. For any realization of the data, $|\hat{\Delta}^{LB} - \hat{\Delta}|$ is nondecreasing in \hat{p} . Therefore, $\Pr[|\hat{\Delta}^{LB} - \hat{\Delta}| < \delta | \hat{p} = p]$ is nonincreasing in \hat{p} . It follows that $\int_0^{\bar{p}} \Pr[|\hat{\Delta}^{LB} - \hat{\Delta}| < \delta | \hat{p} = p] dF_N(p) \geq \int_0^{\bar{p}} \Pr[|\hat{\Delta}^{LB} - \hat{\Delta}| < \delta | \hat{p} = \bar{p}] dF_N(p) = \Pr[|\hat{\Delta}^{LB} - \hat{\Delta}| < \delta | \hat{p} = \bar{p}]$, which converges to 1, by construction of \bar{p} . $\Pr[0 \leq \hat{p} \leq \bar{p}]$ converges to 0.5 and therefore so does the first term above. A parallel argument shows the third term converges to 0.5 as well. Q.E.D.

Acknowledgements. Earlier drafts of this paper were circulated as “Trimming for Bounds on Treatment Effects with Missing Outcomes”, Center for Labor Economics Working Paper No. 51, March 2002, and NBER Technical Working Paper No. 277, June 2002, as well as a revision with the above title, as NBER Working Paper No. 11721, October 2005. Emily Buchsbaum, Vivian Hwa, Xiaotong Niu, and Zhuan Pei provided excellent research assistance. I thank David Card, Guido Imbens, Justin McCrary, Marcelo Moreira, Enrico Moretti, Jim Powell, Jesse Rothstein, Mark Watson, and Edward Vytlacil for helpful discussions and David Autor, Josh Angrist, John DiNardo, Jonah Gelbach, Alan Krueger, Doug Miller, Aviv Nevo, Jack Porter, Diane Whitmore, and participants of the UC Berkeley Econometrics and Labor Lunches for useful comments and suggestions.

REFERENCES

- AHN, H. and POWELL, J. (1993), “Semiparametric Estimation of Censored Selection Models with a Nonparametric Selection Mechanism”, *Journal of Econometrics*, **58**, 3–29.
- ANDREWS, D. W. K. (1994), “Empirical Process Methods in Econometrics”, in R. F. Engle and D. L. McFadden (eds.) *Handbook of Econometrics*, Vol. 4 (Amsterdam: North Holland) 2248–2296.
- ANDREWS, D., MOREIRA, M. J. and STOCK, J. H. (2007), “Performance of Conditional Wald Tests in IV Regression with Weak Instruments”, *Journal of Econometrics*, **139**, 116–132.
- ANDREWS, D. and SCHAFFGANS, M. (1998), “Semiparametric Estimation of the Intercept of a Sample Selection Model”, *Review of Economic Studies*, **65**, 497–517.
- ANGRIST, J. (1997), “Conditional Independence in Sample Selection Models”, *Economics Letters*, **54**, 103–112.

- ANGRIST, J., IMBENS, G. and RUBIN, D. (1996), "Identification of Causal Effects Using Instrumental Variables", *Journal of the American Statistical Association*, **91**, 444–455.
- BALKE, A. and PEARL, J. (1997), "Bounds on Treatment Effects from Studies with Imperfect Compliance", *Journal of the American Statistical Association*, **92**, 1171–1177.
- BARNOW, B. (1987), "The Impact of CETA on the Post-Program Earnings of Participants", *Journal of Human Resources*, **22**, 157–193.
- BLOOM, H. (1984), "Accounting for No-Shows in Experimental Evaluation Designs", *Evaluation Review*, **8**, 225–246.
- BLUNDELL, R., GOSLING, A., ICHIMURA, H. and MEGHIR, C. (2007), "Changes in the Distribution of Male and Female Wages Accounting for Employment Composition Using Bounds", *Econometrica*, **75**, 323–363.
- BURGHARDT, J., SCHOCHET, P. Z., MCCONNELL, S., JOHNSON, T., GRITZ, R. M., GLAZERMAN, S., HOMRIGHAUSEN, J. and JACKSON, R. (2001), *Does Job Corps Work? Summary of the National Job Corps Study*, Report (Washington, DC: Mathematica Policy Research, Inc.).
- CARD, D. (1999), "The Causal Effect of Education on Earnings", in O. Ashenfelter and D. Card (eds.) *Handbook of Labor Economics*, Vol. 3A (Amsterdam: North Holland) 1801–1863.
- CHAMBERLAIN, G. (1994), "Quantile Regression, Censoring, and the Structure of Wages", in C. A. Sims (ed.) *Advances in Econometrics, Sixth World Congress*, Vol. 1 (Cambridge: Cambridge University Press) 171–220.
- DAS, M., NEWWEY, W. K. and VELLA, F. (2003), "Nonparametric Estimation of Sample Selection Models", *Review of Economic Studies*, **70**, 33–58.
- HECKMAN, J. J. (1974), "Shadow Prices, Market Wages, and Labor Supply", *Econometrica*, **42**, 679–694.
- HECKMAN, J. J. (1979), "Sample Selection Bias as a Specification Error", *Econometrica*, **47**, 153–161.
- HECKMAN, J. J. (1990), "Varieties of Selection Bias", *American Economic Review*, **80**, 313–318.
- HECKMAN, J. J., LALONDE, R. J. and SMITH, J. A. (1999), "The Economics and Econometrics of Active Labor Market Programs", in O. Ashenfelter and D. Card (eds.) *Handbook of Labor Economics*, Vol. 3A (Amsterdam: North Holland) 1865–2097.
- HECKMAN, J. J. and ROBB, R. (1986), "Alternative Methods for Solving the Problem of Selection Bias in Evaluating the Impact of Treatments on Outcomes", in H. Wainer (ed.) *Drawing Inferences from Self-selected Samples* (New York: Springer) 63–107.
- HECKMAN, J. J. and VYTLACIL, E. (1999), "Local Instrumental Variables and Semiparametric Estimation and Latent Variable Models for Identifying and Bounding Treatment Effects", *Proceedings of the National Academy of Sciences*, **96**, 4730–4734.
- HECKMAN, J. J. and VYTLACIL, E. (2000a), "Instrumental Variables, Selection Models, and Tight Bounds on the Average Treatment Effect" (Technical Working Paper No. 259, National Bureau of Economic Research).
- HECKMAN, J. J. and VYTLACIL, E. (2000b), "Local Instrumental Variables" (Technical Working Paper No. 252, National Bureau of Economic Research).
- HOLLISTER, R., KEMPER, P. and MAYNARD, R. (1984), *The National Supported Work Demonstration* (Madison, WI: University of Wisconsin Press).
- HOROWITZ, J. L. and MANSKI, C. F. (1995), "Identification and Robustness with Contaminated and Corrupted Data", *Econometrica*, **63**, 281–302.
- HOROWITZ, J. L. and MANSKI, C. F. (2000a), "Nonparametric Analysis of Randomized Experiments with Missing Covariate and Outcome Data", *Journal of the American Statistical Association*, **95**, 77–84.
- HOROWITZ, J. L. and MANSKI, C. F. (2000b), "Rejoinder: Nonparametric Analysis of Randomized Experiments with Missing Covariate and Outcome Data", *Journal of the American Statistical Association*, **95**, 87.
- IMBENS, G. W. and ANGRIST, J. (1994), "Identification and Estimation of Local Average Treatment Effects", *Econometrica*, **62**, 467–476.
- IMBENS, G. W. and Manski, C. F. (2004), "Confidence Intervals for Partially Identified Parameters", *Econometrica*, **72**, 1845–1857.
- KIEFER, N. (1979), *The Economic Benefits of Four Employment and Training Programs* (New York: Garland Publishing).
- KRUEGER, A. B. and WHITMORE, D. M. (2001), "The Effect of Attending a Small Class in the Early Grades on College-Test Taking and Middle School Test Results: Evidence from Project STAR", *Economic Journal*, **111**, 1–28.
- LEE, D. S. (2002), "Trimming for Bounds on Treatment Effects with Missing Outcomes" (Center for Labor Economics Working Paper No. 38, University of California, Berkeley).
- LEE, D. S. (2005), "Training, Wages, and Sample Selection: Estimating Sharp Bounds on Treatment Effects" (NBER Working Paper No. 11721, National Bureau of Economic Research).
- MARTIN, J. P. (2000), "What Works among Active Labour Market Policies: Evidence from OECD Countries' Experiences", *OECD Economic Studies*, **30**, 79–113.
- NEWWEY, W. K. and MCFADDEN, D. (1994), "Large Sample Estimation and Hypothesis Testing", in R. F. Engle and D. L. McFadden (eds.) *Handbook of Econometrics*, Vol. 4 (Amsterdam: North Holland) 2113–2245.
- ROBINS, J. (1989), "The Analysis of Randomized and Non-Randomized AIDS Treatment Trials Using a New Approach to Causal Inference in Longitudinal Studies", in L. Sechrest, H. Freeman and A. Mulley (eds.) *Health Service Research Methodology: A Focus on AIDS* (Washington, DC: U.S. Public Health Service) 113–159.
- RUBIN, D. (1976), "Inference and Missing Data", *Biometrika*, **63**, 581–592.

- SCHOCHET, P. Z., BURGHARDT, J. and GLAZERMAN, S. (2001), "National Job Corps Study: The Impacts of Job Corps on Participants' Employment and Related Outcomes", Report (Washington, DC: Mathematica Policy Research, Inc.).
- SCHOCHET, P. Z., CAO, J. B. R.-J., GLAZERMAN, S., GRADY, A., GRITZ, M., MCCONNELL, S., JOHNSON, T. and BURGHARDT, J. (2003), "National Job Corps Study: Data Documentation and Public Use Files, Volume I," Documentation (Washington, DC: Mathematica Policy Research, Inc.).
- SMITH, J. P. and WELCH, F. R. (1986), *Closing the Gap: Forty Years of Economic Progress for Blacks* (Santa Monica, CA: Rand Corporation).
- STAIGER, D. and STOCK, J. H. (1997), "Instrumental Variables regression with Weak Instruments", *Econometrica*, **65**, 557–586.
- STIGLER, S. M. (1973), "The Asymptotic Distribution of the Trimmed Mean", *Annals of Statistics*, **1**, 472–477.
- U.S. DEPARTMENT OF LABOR (2005a), "Summary of Budget Authority: Fiscal Years 2004–2005" (Table, Employment and Training Administration 2).
- U.S. DEPARTMENT OF LABOR (2005b), "What is Job Corps?" (Web Page, Employment and Training Administration, <http://jobcorps.doleta.gov/about.cfm>).
- VYTLACIL, E. (2002), "Independence, Monotonicity, and Latent Index Models: An Equivalence Result", *Econometrica*, **70**, 331–341.
- ZHANG, J. L. and RUBIN, D. B. (2003), "Estimation of Causal Effects via Principal Stratification When Some Outcomes Are Truncated by 'Death'", *Journal of Educational and Behavioral Statistics*, **28**, 353–368.
- ZHANG, J. L., RUBIN, D. B. and MEALLI, F. (2008), "Evaluating the Effects of Job Training Programs on Wages through Principal Stratification", in D. Millimet, J. Smith and E. Vytlacil (eds.) *Modelling and Evaluating Treatment Effects in Econometrics*, Vol. 21 (Amsterdam: Elsevier) 119–147.