

# Supplement: ATT Proofs and National JTPA Study Information

Adam N. Glynn<sup>\*</sup>    Konstantin Kashin<sup>†</sup>  
[aglynn@emory.edu](mailto:aglynn@emory.edu)    [kvkashin@gmail.com](mailto:kvkashin@gmail.com)

October 17, 2017

## A ATT Proofs

### A.1 Large-Sample Bias

The asymptotic bias in the front-door approach for  $E[Y(a_0)|a_1]$  is the following:

$$\begin{aligned} B_{0|a_1}^{fd} &= \mu_{0|a_1}^{fd} - \mu_{0|a_1} \\ &= \sum_x \sum_m P(m|a_0, x) \cdot E[Y|a_1, m, x] \cdot P(x|a_1) \\ &\quad - \sum_x \sum_u E[Y|a_0, x, u] \cdot P(u|x, a_1) \cdot P(x|a_1) \\ &= \sum_x \sum_m P(m|a_0, x) \sum_u E[Y|a_1, m, x, u] \cdot P(u|a_1, m, x) \cdot P(x|a_1) \\ &\quad - \sum_x \sum_u \sum_m E[Y|a_0, m, x, u] \cdot P(m|a_0, x, u) \cdot P(u|a_1, x) \cdot P(x|a_1) \\ &= \sum_x P(x|a_1) \sum_m \sum_u P(m|a_0, x) \cdot E[Y|a_1, m, x, u] \cdot P(u|a_1, m, x) \\ &\quad - \sum_x P(x|a_1) \sum_m \sum_u P(m|a_0, x, u) \cdot E[Y|a_0, m, x, u] \cdot P(u|a_1, x) \end{aligned} \tag{1}$$

---

<sup>\*</sup>Department of Political Science, Emory University, 327 Tarbuton Hall, 1555 Dickey Drive, Atlanta, GA 30322 (<http://scholar.harvard.edu/aglynn>).

<sup>†</sup>Data Scientist, Institute for Quantitative Social Science, Harvard University, 1737 Cambridge Street, Cambridge, MA 02138 (<http://konstantinkashin.com>).

The asymptotic bias of the back-door approach can be written as the following:

$$\begin{aligned}
B_{0|a_1}^{bd} &= \mu_{0|a_1}^{bd} - \mu_0 \\
&= \sum_x E[Y|a_0, x] \cdot P(x|a_1) - \sum_x \sum_u E[Y|a_0, x, u] \cdot P(u|x, a_1) \cdot P(x|a_1) \\
&= \sum_x \sum_u E[Y|a_0, x, u] \cdot P(u|a_0, x) \cdot P(x|a_1) \\
&\quad - \sum_x \sum_u E[Y|a_0, x, u] \cdot P(u|a_1, x) \cdot P(x|a_1) \\
&= \sum_x P(x|a_1) \sum_u E[Y|a_0, x, u] \cdot [P(u|a_0, x) - P(u|a_1, x)]
\end{aligned} \tag{2}$$

## A.2 Nonrandomized program evaluation with one-sided noncompliance

In the special case of nonrandomized program evaluations with one-sided noncompliance, the front-door and back-door bias can be written as the following, utilizing the fact that  $P(M = 0|a_0, u) = 1$  and  $P(M = 0|a_1, u) = 0$  for all  $u$ :

$$\begin{aligned}
B_{ATT}^{fd} &= \mu_1 - \mu_{0|a_1}^{fd} - (\mu_1 - \mu_{0|a_1}) \\
&= \mu_{0|a_1} - \mu_{0|a_1}^{fd} \\
&= -B_{0|a_1}^{fd} \\
&= \sum_x P(x|a_1) \sum_u E[Y|a_0, M = 0, x, u] P(u|a_1, x) \\
&\quad - \sum_x P(x|a_1) \sum_u E[Y|a_1, M = 0, x, u] P(u|a_1, M = 0, x)
\end{aligned}$$

Adding and subtracting  $\sum_x P(x) \sum_u E[Y|a_0, M = 0, u] \cdot P(u|a_1, M = 0)$ :

$$\begin{aligned}
&= \sum_x P(x|a_1) \sum_u E[Y|a_0, M = 0, x, u] \cdot [P(u|a_1, x) - P(u|a_1, M = 0, x)] \\
&\quad - \sum_x P(x|a_1) \sum_u \{E[Y|a_1, M = 0, x, u] - E[Y|a_0, M = 0, x, u]\} \cdot P(u|a_1, M = 0, x)
\end{aligned} \tag{3}$$

The bias can be re-written further if we note that the imbalance in  $U$  can be written in terms of the mediator,

$$\begin{aligned}
P(u|a_1, x) - P(u|a_1, m_0, x) &= P(u|a_1, m_1, x) \cdot P(m_1|a_1, x) + P(u|a_1, m_0, x) \cdot P(m_0|a_1, x) - P(u|a_1, m_0, x) \\
&= P(u|a_1, m_1, x) \cdot P(m_1|a_1, x) + P(u|a_1, m_0, x) \cdot [P(m_0|a_1, x) - 1] \\
&= P(u|a_1, m_1, x) \cdot P(m_1|a_1, x) - P(u|a_1, m_0, x) \cdot P(m_1|a_1, x) \\
&= P(m_1|a_1, x) \cdot [P(u|a_1, m_1, x) - P(u|a_1, m_0, x)].
\end{aligned}$$

and the difference in expectations over  $Y$  can be written in terms of the potential outcomes under control,

$$\begin{aligned}
& E[Y|a_0, m_0, x, u] - E[Y|a_1, m_0, x, u] \\
&= E[Y|a_0, x, u] - E[Y|a_1, m_0, x, u] \\
&= E[Y(a_0)|a_0, x, u] - E[Y|a_1, m_0, x, u] \\
&= E[Y(a_0)|a_1, x, u] - E[Y|a_1, m_0, x, u] \\
&= E[Y(a_0)|a_1, m_1, x, u] \cdot P(m_1|a_1, x, u) + E[Y(a_0)|a_1, m_0, x, u] \cdot P(m_0|a_1, x, u) - E[Y|a_1, m_0, x, u] \\
&= E[Y(a_0)|a_1, m_1, x, u] \cdot P(m_1|a_1, x, u) - E[Y(a_0)|a_1, m_0, x, u] \cdot \left[ \frac{E[Y|a_1, m_0, x, u]}{E[Y(a_0)|a_1, m_0, x, u]} - P(m_0|a_1, x, u) \right].
\end{aligned}$$

These equivalencies allow us to write the front-door bias under one-sided noncompliance as the following:

$$\begin{aligned}
B_{att}^{fd} &= \sum_x P(x|a_1) P(m_1|a_1, x) \sum_u E[Y|a_0, m_0, x, u] \cdot [P(u|a_1, m_1, x) - P(u|a_1, m_0, x)] \\
&+ \sum_x P(x|a_1) \sum_u \left\{ E[Y(a_0)|a_1, m_1, x, u] \cdot P(m_1|a_1, x, u) - E[Y(a_0)|a_1, m_0, x, u] \cdot \left[ \frac{E[Y|a_1, m_0, x, u]}{E[Y(a_0)|a_1, m_0, x, u]} - P(m_0|a_1, x, u) \right] \right\} P(u|a_1, m_0, x).
\end{aligned}$$

$$\begin{aligned}
B_{ATT}^{bd} &= \mu_1 - \mu_{0|a_1}^{bd} - (\mu_1 - \mu_{0|a_1}) \\
&= \mu_{0|a_1} - \mu_{0|a_1}^{bd} \\
&= -B_{0|a_1}^{bd} \\
&= \sum_x P(x|a_1) \sum_u E[Y|a_0, M=0, x, u] \cdot [P(u|a_1, x) - P(u|a_0, x)]
\end{aligned}$$

### A.3 Bias Simplification

In order to improve interpretability of the bias formulas and establish comparability with the results for back-door bias in [VanderWeele and Arah \(2011\)](#), we offer a simplification of the front-door bias formula under one-sided noncompliance and an exclusion restriction. Under Assumptions 3 and 4 from the main article we write front-door bias as:

$$\begin{aligned}
B_{ATT}^{fd} &= \sum_x P(x|a_1) P(m_1|a_1, x) \sum_u E[Y|a_0, m_0, x, u] \cdot \underbrace{[P(u|a_1, m_1, x) - P(u|a_1, m_0, x)]}_{\varepsilon} \\
&+ \sum_x P(x|a_1) \sum_u P(m_1|a_1, x, u) \underbrace{[E[Y(a_0)|a_1, m_1, x, u] - E[Y(a_0)|a_1, m_0, x, u]]}_{\eta} P(u|a_1, m_0, x)
\end{aligned}$$

Under Assumptions 6, 7, and 8 from the main article, we can simplify the above expression as:

$$\begin{aligned}
B_{ATT}^{fd} &= P(m_1|a_1) \cdot \varepsilon \cdot \sum_x P(x|a_1) \sum_u E[Y|a_0, m_0, x, u] \\
&+ P(m_1|a_1) \cdot \eta \cdot \sum_x P(x|a_1) \sum_u P(u|a_1, m_0, x)
\end{aligned}$$

Assuming  $U$  is binary as in [VanderWeele and Arah \(2011\)](#):

$$B_{ATT}^{fd} = P(m_1|a_1) \left[ \varepsilon \cdot \sum_x P(x|a_1) \underbrace{(E[Y|a_0, m_0, x, U = 1] - E[Y|a_0, m_0, x, U = 1])}_{\gamma} + \eta \right]$$

Assuming that  $\gamma$  does not depend on  $x$  as in [VanderWeele and Arah \(2011\)](#):

$$B_{ATT}^{fd} = P(m_1|a_1) [\gamma \cdot \varepsilon + \eta]$$

## B National JTPA Study

Our analysis makes use of the following samples in the National JTPA Study: experimental treatment group, experimental control group, and the nonexperimental / eligible nonparticipant (ENP) group. We restrict our attention to the 4 *service delivery areas* at which the ENP sample was collected: Fort Wayne, IN; Corpus Christi, TX; Jersey City, NJ; and Providence, RI. We also only examine 2 target groups: adult males and adult females. Note that the treatment group for our purposes means receiving any JTPA service, even though the services actually received from the JTPA varied across individuals.<sup>1</sup>

The raw data and edited analysis files are available as part of the National JTPA Study Public Use Data from the Upjohn Institute. The covariates for the experimental sample are available through the background information form (BIF) and the covariates for ENPs are available through the long baseline survey (LBS). The experimental samples completed the BIF, which contains demographic information, social program participation, and training and education histories, at the time of random assignment. The ENPs completed the LBS anywhere from 0 to 24 months following eligibility screening. Unlike the BIF which mostly covers the previous year in terms of labor market experiences, the LBS covers the past 5 years prior to the survey date and thus provides a much richer portrait of labor market participation. Moreover, experimental control units at the 4 ENP sites also received the long baseline survey, completed 1-2 months after random assignment. [Heckman et al. \(1998\)](#), [Heckman and Smith \(1999\)](#), and related works rely on the detailed labor force participation data and earnings histories in LBS to identify selection bias by comparing the experimental control units to the nonexperimental comparison units. Unfortunately, treated units were never administered the LBS and we have no detailed labor force participation data for multiple years prior to random assignment. Moreover, no one survey instrument was administered to all three of the samples we are using in this analysis, yielding issues of noncomparability. The limited set of covariates

---

<sup>1</sup>The National JTPA Study classified services received into the following 6 categories: classroom training in occupational skills, on-the-job training, job search assistance, basic education, work experience, and miscellaneous.

we use in the conditioning sets in our analysis have all been established to be comparable by verifying their values across the BIF and LBS for the experimental control group, which completed both surveys at the 4 ENP sites.

The dataset we end up using in our analysis was obtained in communication with Jeffrey Smith and Petra Todd. It is the dataset used in the estimates presented in Section 11 of Heckman et al. (1998) and contains all three samples we use in our analysis. It also contains compliance information for the experimental treated group sample. The covariates we utilize in our analysis have been cross-checked against the raw data from the Upjohn Institute. There are also additional covariates in the Heckman et al. (1998) data that have been imputed using a linear regression as described in Appendix B3 of their paper.

The outcome variable we use in the analysis is total 18-month earnings in the period following random assignment (for experimental units) or eligibility screening (for ENPs), calculated from the monthly `totearn` variable available from the public use data files for months 1-30 after random assignment (denoted as  $t + 1$  to  $t + 30$ , where  $t$  is the time of random assignment). The data also includes data for  $t$ , the month of random assignment. Note that this variable is not raw earnings data, but was constructed by Abt Associates from the First and Second Follow-up Surveys, as well as based on data from state unemployment agencies.<sup>2</sup> Please consult Appendix A of Orr et al. (1994) for description of the First Follow-up Survey, Second Follow-up Survey, and earnings data from state unemployment insurance agencies and Appendix B of the same report for construction and imputation of the 30-month earnings variables. The Narrative Description of the National JTPA Study Public Use Files also contains description of the imputation process (see <http://www.upjohninst.org/erdc/njtpa.html>).

In our analysis, we rely upon the monthly total earnings variable in the dataset we obtained from Jeffrey Smith and Petra Todd. We have verified the earnings data used in the calculation of the program impact from this dataset against the earnings variables in the public use data and they match exactly except for a few individuals where Heckman et al. (1998) have imputed missing monthly data. This applies to around 1% of observations and thus is unlikely to substantively change any results. A unit-by-unit comparison of earnings across the raw data and the data we are using can be obtained from us upon request.

The full dataset we obtained contains 1478 treated units, 649 experimental control units, and 667 ENPs for adult males. For adult females, there are 1734 treated units, 830 experimental control units, and 1340 ENPs. These numbers already exclude individuals without any earnings records. We follow the sample restrictions in Heckman et al. (1998) to reduce the full dataset to the final sample (see Appendix B1). We impose an age restriction of 22 to 54 years old on the experimental samples to match the ages of the ENP sample. We then omit individuals who are missing data on race and date of eligibility. Finally, we impose a *rectangular restriction* based on quarterly earnings.

---

<sup>2</sup>One of the major imputations was a decision to divide raw earnings by a `shares` variable which adjust earnings reported for incomplete months (due to the timing of the interviews) to full monthly earnings.

For experimental control and the ENP samples, we require (i) at least one month of valid earnings prior to random assignment (for experimental controls) or prior to eligibility screening (for ENPs), denoted as  $t = 0$ , (ii) valid earnings data at  $t = 0$ , and (iii) at least one month of valid earnings data in months  $t + 13$  to  $t + 18$ . For the treatment group, we impose only restriction iii. The final sample sizes are presented in Table 1.

Table 1: Sample Sizes Before and After Imposing Sample Restrictions. The treated units are broken up into compliers (C) and noncompliers (NC). Control denotes experimental control and ENP denotes the eligible nonparticipants.

	Adult Males				Adult Females			
	Treated		Control	ENP	Treated		Control	ENP
	C	NC			C	NC		
Pre-restriction	843	635	649	667	953	781	830	1340
Post-restriction	834	622	523	384	934	765	706	852

Even after imposing the rectangular restriction on earnings, some individuals had missing earnings data for some months. In the construction of the 18-month total earnings variable, we mean impute the missing months using the average of the individual's available monthly earnings. Details on the extent of missingness are available from authors upon request.

## References

- Heckman, J., Ichimura, H., Smith, J., and Todd, P. (1998). Characterizing selection bias using experimental data. *Econometrica*, 66:1017–1098. [4](#), [5](#)
- Heckman, J. J. and Smith, J. A. (1999). The pre-programme earnings dip and the determinants of participation in a social programme: implications for simple programme evaluation strategies. *Economic Journal*. [4](#)
- Orr, L. L., Bloom, H. S., Bell, S. H., Lin, W., Cave, G., and Doolittle, F. (1994). *The National JTPA Study: Impacts, Benefits, And Costs of Title IIA*. Abt Associates, Bethesda, MD. [5](#)
- VanderWeele, T. J. and Arah, O. A. (2011). Bias formulas for sensitivity analysis of unmeasured confounding for general outcomes, treatments, and confounders. *Epidemiology*, 22(1):42–52. [3](#), [4](#)