

History and Context
Within-Study Comparisons in Economics

Jeffrey Smith
Professor of Economics and Public Policy
University of Michigan
econjeff@umich.edu

Annual Meeting of the Impact Evaluation Network (IEN)
Harvard
September 19-20, 2012

Introduction

Thanks!

Who is that man?

Key points

What can and should be learned from within-study comparisons in economics?

Are we learning the correct lessons from existing within-study comparisons?

Outline / roadmap

Defining terms

LaLonde (1986)

Dehejia and Wahba (1999, 2002)

Defining concepts

Smith and Todd (2005)

Discussion

Defining terms

Experimental: involves random assignment to treated and untreated states

Non-experimental: everything else

Quasi-experimental: non-experimental

LaLonde (1986)

Compare the experimental estimates from the National Supported Work Demonstration evaluation with non-experimental estimates using comparison groups drawn from other data sets

Cook (and others) call this a “within study comparison”

LaLonde (1986) motivation

To quote LaLonde:

“The goal is to assess the likely ability of several *econometric methods* to accurately assess the economic benefits of employment and training programs”
(604 – italics mine)

Put differently, the goal is to find an estimator that solves the selection problem

LaLonde (1986) basic setup

Comparison group source one: Panel Study of Income Dynamics female heads (continuously) from 1975-1979

Comparison group source two: Current Population Survey persons in the March 1976 CPS in the labor force in 1976 with individual income < 20K and household income < 30K

Using representative samples for comparison groups was standard practice at the time – it is less so now

LaLonde (1986) NSW experiment

The National Supported Work Demonstration examined the impacts of an expensive treatment on four groups with labor market difficulties: long-term AFDC recipients, high school dropouts, ex-convicts and ex-addicts.

LaLonde looks at two groups: AFDC women and the men from the other three groups

Aside: why would you combine these groups?

Treatment group observations were randomly assigned from January 1976 to April 1977

Random assignment took place in 10 sites around the country, all of them in cities.

LaLonde (1986) estimators

Linear regression

First differences regression

Regression with lagged dependent variable

Bivariate normal selection model with exclusion restrictions, estimated using the Heckman (1979) two-step approach

LaLonde (1986) variables

Covariates: age, Black and Hispanic indicators, years of schooling, an indicator for married, and a high school completion indicator (and that is all!)

Outcome variable: Real earnings from NSW survey (treatment group), SSA earnings records (CPS comparison group), PSID survey (PSID comparison group)

Dependent variable: real earnings in 1978

Lagged dependent variable: real earnings in 1975

Exclusion restriction variables: urban residence (!), employment status in 1976 (!), AFDC status in 1975 (!), number of children (!)

TABLE 5—EARNINGS COMPARISONS AND ESTIMATED TRAINING EFFECTS FOR THE NSW MALE PARTICIPANTS USING COMPARISON GROUPS FROM THE PSID AND THE CPS-SSA^{a,b}

Name of Comparison Group ^d	Comparison Group Earnings Growth 1975-78 (1)	NSW Treatment Earnings Less Comparison Group Earnings				Difference in Differences: Difference in Earnings Growth 1975-78 Treatments Less Comparisons		Unrestricted Difference in Differences: Quasi Difference in Earnings Growth 1975-78		Controlling for All Observed Variables and Pre-Training Earnings (10)
		Pre-Training Year, 1975		Post-Training Year, 1978		Without Age (6)	With Age (7)	Unadjusted (8)	Adjusted ^c (9)	
		Unadjusted (2)	Adjusted ^c (3)	Unadjusted (4)	Adjusted ^c (5)					
Controls	\$2,063 (325)	\$39 (383)	\$-21 (378)	\$886 (476)	\$798 (472)	\$847 (560)	\$856 (558)	\$897 (467)	\$802 (467)	\$662 (506)
PSID-1	\$2,043 (237)	-\$15,997 (795)	-\$7,624 (851)	-\$15,578 (913)	-\$8,067 (990)	\$425 (650)	-\$749 (692)	-\$2,380 (680)	-\$2,119 (746)	-\$1,228 (896)
PSID-2	\$6,071 (637)	-\$4,503 (608)	-\$3,669 (757)	-\$4,020 (781)	-\$3,482 (935)	\$484 (738)	-\$650 (850)	-\$1,364 (729)	-\$1,694 (878)	-\$792 (1024)
PSID-3	(\$3,322) (780)	(\$455) (539)	(\$455) (704)	\$697 (760)	-\$509 (967)	\$242 (884)	-\$1,325 (1078)	\$629 (757)	-\$552 (967)	\$397 (1103)
CPS-SSA-1	\$1,196 (61)	-\$10,585 (539)	-\$4,654 (509)	-\$8,870 (562)	-\$4,416 (557)	\$1,714 (452)	\$195 (441)	-\$1,543 (426)	-\$1,102 (450)	-\$805 (484)
CPS-SSA-2	\$2,684 (229)	-\$4,321 (450)	-\$1,824 (535)	-\$4,095 (537)	-\$1,675 (672)	\$226 (539)	-\$488 (530)	-\$1,850 (497)	-\$782 (621)	-\$319 (761)
CPS-SSA-3	\$4,548 (409)	\$337 (343)	\$878 (447)	-\$1,300 (590)	\$224 (766)	-\$1,637 (631)	-\$1,388 (655)	-\$1,396 (582)	\$17 (761)	\$1,466 (984)

^a The columns above present the estimated training effect for each econometric model and comparison group. The dependent variable is earnings in 1978. Based on the experimental data an unbiased estimate of the impact of training presented in col. 4 is \$886. The first three columns present the difference between each comparison group's 1975 and 1978 earnings and the difference between the pre-training earnings of each comparison group and the NSW treatments.

^b Estimates are in 1982 dollars. The numbers in parentheses are the standard errors.

^c The exogenous variables used in the regression adjusted equations are age, age squared, years of schooling, high school dropout status, and race.

^d See Table 3 for definitions of the comparison groups.

The researchers who evaluated these federally sponsored programs devised both experimental and nonexperimental procedures to estimate the training effect, because they recognized that the difference between the trainees' pre- and post-training earnings was a poor estimate of the training effect. In a dynamic economy, the trainees' earnings may grow even without an effective program. The goal of these program evaluations is to estimate the earnings of the trainees had they not participated in the program. Researchers using experimental data take the earnings of the control group members to be an estimate of the trainees' earnings without the program. Without experimental data, researchers estimate the earnings of the trainees by using the regression-adjusted earnings of

a comparison group drawn from the population. This adjustment takes into account that the observable characteristics of the trainees and the comparison group members differ, and their unobservable characteristics may differ as well.

Any nonexperimental evaluation of a training program must explicitly account for these differences in a model describing the observable determinants of earnings and the process by which the trainees are selected into the program. However, unlike in an experimental evaluation, the nonexperimental estimates of the training effect depend crucially on the way that the earnings and participation equations are specified. If the econometric model is specified correctly, the nonexperimental estimates should be the same (within sampling error) as the training effect generated from the experimental data, but if there is a significant difference between the nonexperimental and the experi-

LaLonde (1986) results

The non-experimental impact estimates vary widely across estimators

The non-experimental impact estimates vary widely across comparison groups

Limited specification tests combined with a priori reasoning do not rule out all of the poorly-performing estimators

Bivariate normal model results are wrong for the reasons already described plus problems with choice-based sampling.

LaLonde (1986) conclusions

“... policymakers should be aware that the available non-experimental evaluations of employment and training programs may contain large and unknown biases resulting from specification errors.” (617)

This paper was widely interpreted to mean that only experiments could provide credible estimates of the impact of active labor market policies.

It directly resulted in the choice of an experimental design for the National Job Training Partnership Act Study

LaLonde (1986) alternative reading

The data are much (all?) of the problem, not the methods.

Why would you expect the handful of covariates here to solve the selection problem?

No measures of past AFDC participation or number of children for the women

No measures related to crime for the ex-convicts or to substance use for the ex-addicts

Lagged annual earnings not well aligned with the time of treatment and measured and aligned differently for the treated and untreated units.

And we blame the methods?

Heckman and Hotz (1989)

Use administrative earnings data for the NSW and CPS samples

The data are grouped, which rules out non-linear estimators

The NSW sample gets a lot larger as attrition is no longer an issue

Focus on specification tests (pre-program and extra lags)

Reject all models that differ in both sign and statistical significance from the experimental results (but point estimates vary among models not rejected)

Do more specification tests and research on specification tests

Keep in mind the “fallacy of alignment”

Dehejia and Wahba (1999, 2002)

Apply matching and weighting methods to LaLonde's data on men (the data on women having been lost due to a misbehaving magnetic field)

Use a sub-sample of the experimental data to allow greater conditioning on "pre" period earnings

Find that matching and weighting yield estimates quite close to those of the experiment in their sample (albeit with quite large standard errors)

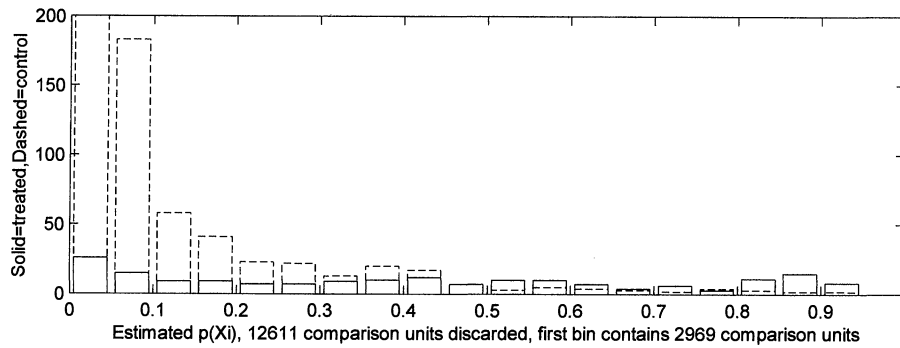


Figure 2. Histogram of the Estimated Propensity Score for NSW Treated Units and CPS Comparison Units. The 12,611 CPS units whose estimated propensity score is less than the minimum estimated propensity score for the treatment group are discarded. The first bin contains 2,969 CPS units. There is minimal overlap between the two groups, but the overlap is greater than in Figure 1; only one bin (.45-.5) contains no comparison units, and there are 35 treated and 7 comparison units with an estimated propensity score greater than .8.

treatment group (although the treatment impact still could be estimated in the range of overlap). With limited overlap, we can proceed cautiously with estimation. Because in our application we have the benchmark experimental estimate, we are able to evaluate the accuracy of the estimates. Even in the absence of an experimental estimate, we show in Section 5 that the use of multiple comparison groups provides another means of evaluating the estimates.

We use stratification and matching on the propensity score to group the treatment units with the small number of comparison units whose estimated propensity scores are greater than the minimum—or less than the maximum—propensity score for treatment units. We estimate the treatment effect by summing the within-stratum difference in means between the treatment and comparison observations (of earnings in 1978), where the sum is weighted by the

number of treated observations within each stratum [Table 3, column (4)]. An alternative is a within-block regression, again taking a weighted sum over the strata [Table 3, column (5)]. When the covariates are well balanced, such a regression should have little effect, but it can help eliminate the remaining within-block differences. Likewise for matching, we can estimate a difference in means between the treatment and matched comparison groups for earnings in 1978 [column (7)], and also perform a regression of 1978 earnings on covariates [column (8)].

Table 3 presents the results. For the PSID sample, the stratification estimate is \$1,608 and the matching estimate is \$1,691, compared to the benchmark randomized-experiment estimate of \$1,794. The estimates from a difference in means and regression on the full sample are -\$15,205 and \$731. In columns (5) and (8), controlling for covariates has little impact on the stratification and

Table 3. Estimated Training Effects for the NSW Male Participants Using Comparison Groups From PSID and CPS

	NSW earnings less comparison group earnings		NSW treatment earnings less comparison group earnings, conditional on the estimated propensity score					
	(1) Unadjusted	(2) Adjusted ^a	Quadratic in score ^b (3)	Stratifying on the score			Matching on the score	
				(4) Unadjusted	(5) Adjusted	(6) Observations ^c	(7) Unadjusted	(8) Adjusted ^d
NSW	1,794 (633)	1,672 (638)						
PSID-1 ^e	-15,205 (1,154)	731 (886)	294 (1,389)	1,608 (1,571)	1,494 (1,581)	1,255	1,691 (2,209)	1,473 (809)
PSID-2 ^f	-3,647 (959)	683 (1,028)	496 (1,193)	2,220 (1,768)	2,235 (1,793)	389	1,455 (2,303)	1,480 (808)
PSID-3 ^f	1,069 (899)	825 (1,104)	647 (1,383)	2,321 (1,994)	1,870 (2,002)	247	2,120 (2,335)	1,549 (826)
CPS-1 ^g	-8,498 (712)	972 (550)	1,117 (747)	1,713 (1,115)	1,774 (1,152)	4,117	1,582 (1,069)	1,616 (751)
CPS-2 ^g	-3,822 (670)	790 (658)	505 (847)	1,543 (1,461)	1,622 (1,346)	1,493	1,788 (1,205)	1,563 (753)
CPS-3 ^g	-635 (657)	1,326 (798)	556 (951)	1,252 (1,617)	2,219 (2,082)	514	587 (1,496)	662 (776)

^a Least squares regression: RE78 on a constant, a treatment indicator, age, age², education, no degree, black, Hispanic, RE74, RE75.
^b Least squares regression of RE78 on a quadratic on the estimated propensity score and a treatment indicator, for observations used under stratification; see note (g).
^c Number of observations refers to the actual number of comparison and treatment units used for (3)–(5); namely, all treatment units and those comparison units whose estimated propensity score is greater than the minimum, and less than the maximum, estimated propensity score for the treatment group.
^d Weighted least squares: treatment observations weighted as 1, and control observations weighted by the number of times they are matched to a treatment observation [same covariates as (a)]. Propensity scores are estimated using the logistic model, with specifications as follows:
^e PSID-1: Prob (T_i = 1) = F(age, age², education, education², married, no degree, black, Hispanic, RE74, RE75, RE74², RE75², u74*black).
^f PSID-2 and PSID-3: Prob (T_i = 1) = F(age, age², education, education², no degree, married, black, Hispanic, RE74, RE74², RE75, RE75², u74, u75).
^g CPS-1, CPS-2, and CPS-3: Prob (T_i = 1) = F(age, age², education, education², no degree, married, black, Hispanic, RE74, RE75, u74, u75, education*RE74, age³).

Dehejia and Wahba (1999, 2002) conclusions

“The methods we suggest are not relevant in all situations. There may be important unobservable covariates, for which the propensity score method cannot account. However, rather than giving up, or relying on assumptions about the unobserved variables, there is substantial reward in exploring first the information contained in the variables that *are* observed.”

This is all quite correct.

But: unobservable or unobserved?

The literature read this paper to mean that “matching works” even in weak data contexts.

Concerns about DW (1999, 2002) from HIST (1998)

Heckman, Ichimura, Smith and Todd (1998) draw some conclusions based on comparing the experimental estimates from the U.S. National Job Training Partnership Act Study to matching estimates obtained using “ideal comparison group data” from four of the experimental sites

Their key conclusions:

Conditioning variables matter a lot, particularly “pre” period outcomes

Putting comparison group members in the same local labor markets matters a lot

Measuring the dependent variable in the same way for the treatment group and comparison group members matters a lot

The LaLonde (1986) data fails to meet all these criteria!

Why then did Dehejia and Wahba (1999, 2002) get such good results?

Defining concepts: what are Dehejia and Wahba (1999, 2002) doing?

Study type 1: examining identification strategies

Example: examining whether unconfoundedness holds for a particular set of conditioning variables in a particular context using experimental estimates as benchmarks

Study type 2: examining the performance of alternative estimators that rely on the same identification strategy

Example: comparing three sets of estimates that all assume unconfoundedness generated by a parametric linear regression, nearest neighbor matching and inverse propensity weighting

Dehejia and Wahba (1999, 2002) combine these without being clear about it

Smith and Todd (2005)

Replicate Dehejia and Wahba (1999, 2002)

If you do what they did, you get what they got – not a minor feat!

If you use an equally (or more) plausible subset of the experimental data, the low bias results disappear, as they do with the full LaLonde (1986) sample

The estimates are quite sensitive to details of the specification

The estimates are quite sensitive to details of the estimation (even how ties are handled!)

Other (selected) recent literature

Fraker and Maynard (1987) *JHR*

Friedlander and Robins (1995) *AER*

Diaz and Handa (2006) *JHR*

Huber, Lechner and Wunsch (2010) *IZA*

Lechner and Wunsch (2011) *IZA*

Shadish, Clark and Steiner (2008) *JASA*

Jacob, Ludwig and Smith (2012) Unpublished (indeed, unwritten)

Conclusions (narrow)

LaLonde (1986) does not show that non-experimental estimators do not work.

Dehejia and Wahba (1999, 2002) do not show that propensity score matching “works”

Replication is useful and can be carried out peacefully

Conclusions (broad)

There is no magic bullet – no estimator that always solves the selection problem

The question is not “which estimator works”

Instead, we want to know the mapping from data, parameter of interest and institutional context to estimator choice

Sometimes there is no non-experimental estimator that solves the selection problem for a given data set and institutional context – a bummer indeed!

Clever econometrics will usually lose out to bad data

References mentioned in the talk

Dehejia, Rajeev and Sadek Wahba. 1999. “Causal Effects in Nonexperimental Studies: Reevaluating the Evaluation of Training Programs.” *Journal of the American Statistical Association*. 94(448): 1053-1062.

Dehejia, Rajeev and Sadek Wahba. 2002. “Propensity Score Matching Methods for Non-Experimental Causal Studies.” *Review of Economics and Statistics*. 84(1): 151-161.

Heckman, James, Hidehiko Ichimura, Jeffrey Smith and Petra Todd. 1998. “Characterizing Selection Bias Using Experimental Data.” *Econometrica* 66(5): 1017-1098.

LaLonde, Robert. 1986. “Evaluating the Econometric Evaluations of Training Programs with Experimental Data.” *American Economic Review* 76(4): 604-620.

Smith, Jeffrey and Petra Todd. 2005a. “Does Matching Overcome LaLonde’s Critique of Nonexperimental Methods?” *Journal of Econometrics* 125(1-2): 305-353.

Smith, Jeffrey and Petra Todd. 2005b. “Rejoinder.” *Journal of Econometrics* 125(1-2): 365-375.

Further reading of potential interest

Agodini, Roberto and Mark Dynarski. 2004. “Are Experiments the Only Option? A Look at Dropout Prevention Programs.” *Review of Economics and Statistics* 86(1): 180-194.

Bell, Stephen, Larry Orr, John Blomquist, and Glen Cain. 1995. *Program Applicants as a Comparison Group in Evaluating Training Programs*. Kalamazoo, MI: W.E. Upjohn Institute for Employment Research.

Cook, Thomas, William Shadish, and Vivian Wong. 2008. “Three Conditions under Which Experiments and Observational Studies Produce Comparable Causal Estimates: New Findings from Within-Study Comparisons.” *Journal of Policy Analysis and Management* 27(4): 724-750.

Diaz, Juan Jose and Sudhanshu Handa. 2006. “An Assessment of Propensity Score Matching as a Nonexperimental Impact Estimator: Evidence from Mexico’s PROGRESA Program.” *Journal of Human Resources* 41(2): 319-345.

Djebbari, Habiba and Jeffrey Smith. 2008. “Heterogeneous Program Impacts: Experimental Evidence from the PROGRESA Program.” *Journal of Econometrics* 145(1-2): 64-80.

Fraker, Thomas and Rebecca Maynard. 1987. "The Adequacy of Comparison Group Designs for Evaluation of Employment-Related Programs." *Journal of Human Resources*. 22(2): 194-227.

Friedlander, Daniel and Philip Robins. 1995. "Evaluating Program Evaluations: New Evidence on Commonly Used Nonexperimental Methods." *American Economic Review*. 85(4): 923-937.

Heckman, James and V. Joseph Hotz. 1989. "Choosing Among Alternative Methods of Estimating the Impact of Social Programs: The Case of Manpower Training." *Journal of the American Statistical Association*. 84(408): 862-874.

Lise, Jeremy, Shannon Seitz and Jeffrey Smith. 2004. "Equilibrium Policy Experiments and the Evaluation of Social Programs." NBER Working Paper No. 10283.

Michalopoulos, Charles, Howard Bloom, and Carolyn Hill. 2004. "Can Propensity Score Methods Match the Findings from a Random Assignment Evaluation of Mandatory Welfare-to-Work Programs?" *Review of Economics and Statistics* 86(1): 156-179.

Todd, Petra and Kenneth Wolpin. 2006. "Assessing the Impact of a School Subsidy Program in Mexico: Using a Social Experiment to Validate a Dynamic Behavioral Model of Child Schooling and Fertility." *American Economic Review* 96(5): 1384-1417.

Shadish, William, M. H. Clark, and Peter Steiner. 2008. "Can Nonrandomized Experiments Yield Accurate Answers? A Randomized Experiment Comparing Random and Nonrandom Assignments." *Journal of the American Statistical Association* 103(484): 1334-1356.

Wilde, Elizabeth and Robinson Hollister. 2007. "How Close is Close Enough? Evaluating Propensity Score Matching Using Data from a Class Size Reduction Experiment." *Journal of Policy Analysis and Management* 26: 455-477.