

ARE PROGRAM PARTICIPANTS GOOD EVALUATORS?

Jeffrey Smith
University of Michigan
UWO, IFS, NBER, IZA and PSI
econjeff@umich.edu

Alexander Whalley
University of California - Merced
awhalley@ucmerced.edu

Nathaniel Wilcox
University of Houston

Version of August 20, 2006

We thank the W.E. Upjohn Institute for funding this research. We are very grateful to Dan Black, Hilary Hoynes and audience members and discussants at University College London (especially Richard Blundell, Hidehiko Ichimura and Costas Meghir), the 2004 APPAM meetings in Atlanta, the 2004 SEA meetings in New Orleans, the 2006 SOLE meeting in Boston for helpful comments. We thank participants in a 2002 OECD conference on evaluating local economic development programs for inspiration. Any and all scandalous mistakes and shocking opinions remain our own.

Abstract

Participants, like econometricians, may have difficulty in constructing the counterfactual outcome required to estimate the impact of a program. In this paper, we directly assess this question by examining the extent to which program participants are able to estimate their individual program impacts ex-post. Utilizing experimental data from the National Job Training Partnership Act (JTPA) Study (NJS) we compare experimentally estimated program impacts to individual self-reports of program effectiveness after the completion of the program. We estimate individual experimental impacts by two methods: (1) subgroup variation; and (2) the assumption of perfect rank correlation in impacts. Little evidence of a relationship between these experimentally estimated program impacts and self-reported program effectiveness is found. We do find evidence that cognitively inexpensive potential proxies for program impacts such as before-after differences in earnings, the type of training received, and labor market outcomes are correlated with self-reported program effectiveness.

1.0 Introduction

Evidence-based policy choice requires systematic and rigorous empirical program evaluation. Data and methods for this have improved rapidly. Social experiments, virtually unknown before 1970, now provide some random assignment of individuals to participation or nonparticipation in programs, and also measure individual outcomes of policy interest, for a wide variety of economic, social and criminal justice programs; see Greenberg and Shroder (2004). Estimation methods have also improve over time; see e.g., Heckman, LaLonde and Smith (1999) and Abbring, Heckman and Vytlacil (2005).

At the same time, participants' own program evaluations have gained attention as a complement to, or substitute for, experimental and/or econometric evaluation. Many surveys ask participants whether they believe a program helped them in some way, though the specifics, wording and number of these questions varies substantially across surveys. Many, if not most, surveys that contain actual outcome measures of policy interest (that is, the data useful for experimental and/or econometric evaluation—instead of, or in addition to, administrative data) also include some participant evaluations.

We compare econometric estimates of program impacts at the individual or subgroup level to participant evaluations, using the rich data from the U.S. National Job Training Partnership Act (JTPA) Study (NJS). JTPA was the major employment and training program for the disadvantaged in the U.S. during the 1980s and 1990s. Section 2 describes this experimental evaluation and its resulting data in detail, including the specific structure and wording of the participant evaluation survey question. Later sections develop two different econometric estimates of program impacts, and we find that these have little systematic relationship to participant evaluations.

We consider two broad interpretations of these results. The “subjective rationality” (Simon (1956)) interpretation is that both participants and econometricians make rational judgments about program success, given their own evaluational premises and definitions of success; but their definitions and premises simply differ. Consistent econometric impact estimates measure program effects on specific outcomes of policy interest (earnings and/or employment, over some specific time period). Participant evaluations instead measure program effects on outcomes, and over time periods, defined by participants in an unobserved way (these may depend both on question wording and participant-idiosyncratic interpretations of that wording). Under this interpretation, weak relationships between participant evaluations and econometric impact estimates can occur even if participants care a lot about the specific outcomes and time periods analyzed by econometricians, since these may still be but a small subset of the outcomes and/or time periods participants consider in their own evaluations.

Our second interpretation of the results is borrowed from cognitive psychology. Nisbett and Ross (1980) discuss the idea that agents act as “lay scientists” when they are asked to produce verbal judgments about the causal structure of their social environment, their own behavior or that of other agents. Like real scientists, lay scientists make these judgments using either empirical or theoretical reasoning, or some mixture of these, depending on how they interpret questions and, perhaps, which approach appears reasonable or easy to them. Yet lay scientists are not real scientists, in two critical senses. First, when acting as “lay empiricists,” they are not compelled to follow canons of formal inference on pain of professional embarrassment if they do not; nor, when acting as “lay theorists,” are they necessarily well-informed as to which theories are well-supported by

evidence based on those canons, but may instead subscribe to a stock of “folk theories” they share with other lay scientists. These two possibilities, of course, can interact: If lay scientists use a poorly supported folk theory as an identifying restriction for empirical inference, their inferences will very likely be flawed (Ross 1989).

While lay empiricists may depend on judgment heuristics (cognitively inexpensive shortcuts) that are generally adaptive, these may occasionally lead to predictable biases of judgment. Memory is fragmentary, and sometimes biased, so lay empiricists may frequently depend on inherently limited data, and may not always correct their judgments for such potentially knowable data limitations (Nisbett and Ross 1980). They may also fail to correct for potential confounds, particularly those that are not salient to them at the time of judgment (Nisbett and Wilson 1977). Put differently, participants may not appropriately construct the counterfactual outcome required to estimate the impact a program had on them.

For instance, participants may wholly depend on relatively crude proxies, such as simple before-after comparisons, in order to make judgments, without accounting for other things that may have changed during treatment periods. Interestingly, before-after comparisons and other crude impact proxies are commonly collected and used in administrative performance standards systems for employment and training programs, perhaps because they are quick and low cost bureaucratic alternatives to the more difficult construction of consistent impact estimates. Participants may rely on the very same proxies to construct their program evaluations. And in fact, we find evidence consistent with this: Participants’ evaluations are predicted both by such crude impact proxies, and by imperfect impact estimates based solely on such proxies. This analysis

has independent interest since it suggests the extent to which participant evaluations might substitute for these measures in administrative performance systems.

The pure “lay theorist” attempts no empirical evaluation at all, but instead consults one of her folk theories and provides a program evaluation (actually, a program outcome prediction, instead of the evaluation of her actual outcome) based on it. Evidence of theory-driven retrospective evaluations by subjects, from laboratory experiments, is long-standing and extensive; see e.g. Nisbett and Wilson (1977) and Ross (1989). In our conclusions, we will discuss the close relationship between our findings and a very direct laboratory demonstration of theory-driven evaluation by Conway and Ross (1984)—a controlled study of participant evaluation. For a simple and pertinent example, however, suppose participants have a folk theory that output is generally increasing in input expense or resource intensity. They may then be more likely to say a program service had a positive impact on them if it seemed relatively expensive or resource-intensive, *ceteris paribus*. We also find evidence consistent with this sort of judgment process in participants’ evaluations.

In addition to informing decisions about how best to evaluate policies, our research may have broader implications. First, whether or not individuals can accurately assess their program impacts, and how they go wrong if they cannot, may have implications for the interpretation of instrumental variables estimates in the context of the correlated random coefficient model, as in Angrist (2004), Heckman (1997a), Heckman and Vytlacil (2005) and Carniero, Heckman and Vytlacil (2005). In that model, complications arise when using instruments that are correlated with the individual-specific component of impacts. Those problems go away if individuals do not know their

impacts (that is, if they make decisions based on “noise”). Of course, if individuals actually use biased estimates of their impacts in making decisions, the problems may return in a different form, depending on how the bias relates to the instrument.

Still more broadly, individuals’ ability to accurately represent outcomes associated with available alternatives lies at the heart of rational models of human behavior. From this perspective, the lay science interpretation of our findings may appear to challenge the notion that people accurately represent alternatives and their future consequences when making decisions. Yet the heart of Nisbett and Wilson’s (1977) and Ross’ (1989) surveys is a dissociation between experimentally measured causes of subject behavior and subjects’ own verbal reports on those causes. Nisbett and Wilson remind us that though subjects sometimes tell more than they can know, they also clearly know more than they can tell (Polanyi 1964). There is no paradox here: skilled performance, which can depend crucially on neural processes hidden from consciousness, does not imply a capacity for accurate verbal description of processes underlying such performance. A neoclassicist gives no important ground by embracing a lay science interpretation of participants’ inability to accurately report program impacts. The capacity for verbal report, and the capacity for decision making, are simply two different things.

In contrast to the extensive literature from laboratory experiments in social psychology referred to above, few studies based on survey data from social experiments examine relationships between consistent impact estimates and participant evaluations. However, Heckman and Smith (1998) and Philipson and Hedges (1998) do, using treatment group dropout (rather than responses to actual evaluation questions) as an indicator of participants’ evaluations. More broadly, Jacob and Lefgren (2005) compare

principals' subjective evaluations of teachers to econometric estimates of teacher value-added, but do not consider the teachers' own evaluations of their value-added.

Prendergast (1999) reviews the literature on subjective performance evaluation, but that literature primarily views subjective evaluations as a way to deal with situations in which agents have multiple tasks (the outputs from some but not all of which allow objective measurement), not as a potentially cost-saving alternative to objective evaluation. That literature is also focused mainly on performance evaluation of workers within firms, not evaluation of the effects of programs on participant labor market outcomes.

We organize the remainder of the paper as follows. Section 2 describes the basic structure of the JTPA program, the NJS experiment and the resulting data. Section 3 presents the conceptual framework that guides our econometric analysis and our interpretation of results. Section 4 discusses the construction and interpretation of the alternative econometric estimates of program impact on employment and earnings that we construct using the experimental data. Section 5 presents results on the relationship between participants' self-reported impacts and impacts estimated using the experimental data. Section 6 examines the relationship between self-reported impacts and before-after employment and earnings changes, as well as proxies such as inputs and outcomes, while Section 7 examines the relationship between self-reported evaluations and performance measures. Finally, Section 8 lays out conclusions we draw from our analysis.

2.0 Data and institutions

2.1 The JTPA program

The U.S. Job Training Partnership Act program was the primary federal program providing employment and training services to the disadvantaged from 1982, when it replaced the Comprehensive Employment and Training Act (CETA) program, to 1998, when it was replaced by the Workforce Investment Act (WIA) program. All of these programs share more or less the same set of services (though the latter two omit the public sector jobs that led to scandal under CETA) and serve the same basic groups. They differ primarily in their organizational details (i.e. do cities or counties play the primary role) and in the emphasis on, and ordering of, the various services provided. Nonetheless, the commonalities dominate with the implication that our results for JTPA likely generalize to WIA (and CETA).¹

The JTPA eligibility rules included categorical eligibility for individuals receiving means tested transfers such as Aid to Families with Dependent Children (AFDC) or its successor Temporary Aid to Needy Families (TANF) as well as food stamps. In addition, individuals were eligible if their family income in the preceding six months fell below a specific cutoff value. There were also special eligibility rules for a number of small groups and a 10 percent “audit” window that basically allowed local sites to enroll individuals at their own discretion. See Devine and Heckman (1996) for more details on the JTPA eligibility rules and Kemple, Doolittle and Wallace (1993) for detailed descriptive statistics on the experimental sample in the NJS. Heckman and Smith (1999, 2004) provide thorough analyses of the determinants of participation in JTPA conditional on eligibility.

¹ One possible caveat is that potential participants may have better information to guide them in making participation decisions about a relatively old program, as JTPA was at the time of the experiment, than about a relatively new program. This reasoning suggests greater selection on impacts over time as a program matures.

The JTPA program provided five major services: classroom training in occupational skills (CT-OS), subsidized on-the-job training (OJT), job search assistance (JSA), adult basic education (ABE) and subsidized work experience (WE). Local sites had the flexibility to emphasize or de-emphasize particular services in response to the needs of the local population and the availability of local service providers. In general, CT-OS was the most expensive service, followed by OJT, ABE and WE. JSA costs a lot less. See Heinrich, Marschke and Zhang (1998) for a detailed study of costs in JTPA and Wood (1995) for information on costs at the NJS study sites.

Services get assigned to individuals by caseworkers, typically as the result of a decision process that incorporates the participant's abilities and desires. This process leads to clear patterns in terms of the observable characteristics of participants assigned to each service. The most job ready individuals typically get assigned to JSA or OJT, while less job ready individuals typically get assigned to CT-OS, BE or WE, where CT-OS often gets followed by JSA. See Kemple, Doolittle and Wallace (1993) for more about the service assignment process. This strongly non-random assignment process has implications for our analyses below in which we examine the relationship between the participant evaluations and types of services they receive.

2.2 The National JTPA Study data

The National JTPA Study (NJS) evaluated the JTPA program using a random assignment design. It was the first major social experiment to evaluate an ongoing program rather than a demonstration program brought into existence solely for the purposes of the experiment. Random assignment in the NJS took place at a non-random sample of 16 of

the more than 600 JTPS Service Delivery Areas (SDAs). Each SDA had a local geographic monopoly on the provision of employment and training services funded under the JTPA. The exact period of random assignment varied among the sites, but in most cases random assignment ran from late 1987 or early 1988 until sometime in spring or summer of 1989. A total of 20,601 individuals were random assigned, usually but not always with the probability of assignment to the treatment group set at 0.67.

The NJS data come from multiple sources. First, respondents completed a Background Information Form (BIF) at the time of random assignment. The BIF collected basic demographic information along with information on past schooling and training and on labor market outcomes at the time of random assignment and earlier. Second, all experimental sample members were asked to complete the first follow-up survey around 18 months after random assignment. This survey collected information on employment and training services (and any formal schooling) received in the period since random assignment, as well as monthly information on employment, hours and wages, from which a monthly earnings measure was constructed. Third, a random subset (for budgetary reasons) of the experimental sample members was asked to complete a second follow-up survey around 32 months after random assignment. This survey collected similar information for the period since the completion of the first follow-up survey or, in the case of individuals who did not complete the first follow-up survey, over the period since random assignment. Response rates to both follow-up surveys were around 80 percent. Finally, administrative data on quarterly earnings and unemployment from state

UI records in the states corresponding to the 16 NJS states were collected.² See Doolittle and Traeger (1990) on the design of the NJS, Orr et al. (1996) and Bloom et al. (1997) for the official impact reports and Heckman and Smith (2000) and Heckman, Hohmann, Smith and Khoo (2000) for further interpretation. Appendix 1 describes the data used in this study in greater detail.

2.3 The self-evaluation questions

Exhibit 1 presents the two survey questions that, taken together, define the participant evaluation measure we use in this paper. The question appears on both the first follow-up survey and the second follow-up survey. In both cases, the skip pattern in the survey excluded control group members from both questions. Respondents in the treatment group were asked these questions in the second follow-up survey only if they did not complete the first follow-up survey.

The first question asks treatment group members whether or not they participated in JTPA. The question assumes application because it is implied by the respondent having been randomly assigned. The JTPA program had different names in the various sites participating in the evaluation; the interviewer included the appropriate local name in each site as indicated in the question.

In the second question, respondents who self-report having participated in the program get asked whether the program helped them get a job or perform better on the job. This is not the ideal question from our point of view, as it focuses more on a specific outcome than on an overall impact, but it is what we have to work with in the JTPA

² These data were collected twice, once for 12 of the 16 sites by Abt Associates, one of the prime contractors on the original experiment, and then for all 16 sites later on by Westat under a separate contract. We use the latter dataset in our analysis.

evaluation. However, to the extent that it focuses respondents' attention specifically on the effect of program participation on labor market outcomes, it should increase the strength of the relationship between the participant evaluations and the econometric estimates of labor market impacts, relative to a broader question that asked about generic program benefits.

We code the responses to both questions as dummy variables. The participant evaluation measure employed in our empirical work consists of the product of the two dummy variables. Put differently, our self-reported evaluation measure equals one if the respondent replies "YES" to question (D7), and "YES" to question (D9). Otherwise, it equals zero.

3.0 Conceptual framework

3.1 A simple model of participants' self-reported evaluations

In this section, we lay out a model of how individual participants might respond to a question regarding whether or not they benefited from a program. The discussion here is inspired by those in Manski (1990) and Dominitz and Manski (1994), who provide careful economic (and econometric) analyses of responses to questions about fertility intentions and returns to schooling, respectively. Our (very) simple model helps to structure the design and interpretation of our empirical work.

To begin, we suppose that respondents compare their observed utility given participation with the utility they would have experienced had they not participated. Let U_1 denote utility given participation, U_0 denote participation given non-participation and let $\Delta_{SR} \in \{0,1\}$ denote the response to the self-evaluation question. Then if respondents generate their answer by comparing the two utilities, we have

$$\Delta_{SR} = 1(U_1 > U_0).$$

This formulation ignores the timing of any affects of participation on utility relative to the survey response. Depending on the wording of the survey question and the respondents' interpretation thereof, respondents may focus on impacts during the period up to the survey response, after the survey response, or some combination of the two. In the JTPA context, we expect them to focus primarily on the effects of the program in the period leading up to the survey response. Expanding our notation, let the subscript "b" denote the period before the survey response and the subscript "a" denote the period following the survey response. We can then write

$$\Delta_{SR} = 1(U_1 > U_0) = 1(f(U_{1b}, E(U_{1a})) > f(U_{0b}, E(U_{0a}))),$$

where $f(\square)$ is an increasing function of both its arguments that maps the utility associated with participation or non-participation, both before and after the self-reported evaluation, into an overall valuation.

Next we consider what aspects of participation affect the utility levels of individuals. In particular, we can decompose the impacts that individuals experience into components related to earnings or employment and a residual component that includes other direct costs and benefits as well as psychic costs and benefits. Denote labor market impacts in the standard notation in the evaluation literature as

$$\Delta_Y = Y_1 - Y_0,$$

where Y_1 denotes the labor market outcome in the treated state and Y_0 denotes the labor market outcome in the untreated state. Similarly, denote the impact on all other determinants of participant utility by

$$\Delta_B = B_1 - B_0,$$

where B_1 and B_0 parallel Y_1 and Y_0 in their interpretation. In what follows, we will further distinguish between impacts realized before and after the survey response.

This decomposition into impacts on labor market outcomes and on all other outcomes that individuals care about corresponds to the components of the impacts that we can and cannot estimate econometrically using our data. The outcomes we (and hopefully the respondents) have in mind other than labor market outcomes include direct costs of participating in training, such as transportation, childcare expenses, alternative time use as in Greenberg (1997), as well as any psychic costs and benefits from participating. Rewriting the survey response function in terms of this additional notation yields

$$\Delta_{SR} = 1(U_1 > U_0) = 1(U(Y_{1b}, B_{1b}, E(Y_{1a}), E(B_{1a})) > U(Y_{0b}, B_{0b}, E(Y_{0a}), E(B_{0a}))),$$

or, alternatively

$$(1) \quad \Delta_{SR} = 1(U_1 > U_0) = 1(g(Y_{1b} - Y_{0b}, B_{1b} - B_{0b}, E(Y_{1a} - Y_{0a}), E(B_{1a} - B_{0a}))).$$

We estimate two variants of equation (1), one suitable for cases where we have econometric estimates of $(Y_{1b} - Y_{0b})$ and another in cases where we examine simple proxies for $(Y_{1b} - Y_{0b})$. The next two subsections define these variants.

3.2 Econometric specification: relationship between econometric impact estimates and participant evaluations.

To examine relationships between our preferred, formal econometric impact estimates (discussed shortly in section 4) and participant evaluations, we begin by assuming additive separability of the $g(\square)$ function into components related to the labor market

impact in the period prior to the survey and the remainder of the function.³ Assuming that the utility function is monotonic in its arguments, we can put an actual labor market impact on the left hand side, yielding

$$(2) \quad Y_{1b} - Y_{0b} = h(\Delta_{SR}, B_{1b} - B_{0b}, E(Y_{1a} - Y_{0a}), E(B_{1a} - B_{0a})).$$

We actually estimate linear versions of (2), given by

$$(3) \quad \hat{Y}_{1b} - Y_{0b} = \beta_0 + \beta_1 \Delta_{SR} + \varepsilon,$$

where the hat on the impact denotes an estimate and where ε includes the idiosyncratic pieces of $B_{1b} - B_{0b}$, $E(Y_{1a} - Y_{0a})$ and $E(B_{1a} - B_{0a})$ (the means are captured in the intercept) as well as the estimation error in the impact estimate and any approximation error due to inappropriate linearization.

We adopt this formulation in the case of the econometric impact estimates for two reasons. First, because econometric impact estimates include estimation error, we want them on the left hand side for the usual reasons associated with measurement error. In contrast, the participant evaluation has no measurement error; the variable is defined as the response to the survey question.⁴ We include no additional covariates on the right hand side because one of our two econometric impact estimates (described in detail in Section 4.1) consists of predicted subgroup experimental impact estimates. To include both these predicted impacts and a set of observables would require excluding at least one observable from this equation, but including it among the observables used to construct

³ Additive separability is not innocuous here; it implies no complementarities between the component of the impacts we estimate and the other components of the impacts.

⁴ The counter-argument in favor of making the participant evaluation the dependent variable despite the estimation error in our impact estimates relies on the econometric impacts having the larger variance of the two variables.

the subgroup impacts. The observables available to us lack an obvious candidate for exclusion.

Under the first interpretation of our analysis (subjectively rational participant evaluations), a weak estimated relationship in equation (3) indicates that earnings or employment impacts, strictly in the period before the evaluation question, are a relatively small part of what participants evaluate. This conclusion requires the qualification that we should not forget what lies in the error term. Among the items in the error term, we would expect long term impacts to correlate positively with short term impacts; in contrast, impacts on leisure likely correlate negatively with impacts on labor market outcomes prior to the survey. A weak relationship in (3) could thus also result from a combination of a positive direct effect of impacts on employment or earnings and a negative indirect effect on leisure, working through the correlation between the omitted impact on leisure and the included impact on employment or earnings. Finally, in a common effect world in which the program has the same impact on everyone, or in which the impact varies but the idiosyncratic component is unknown even to participants, the true coefficient on the econometric estimate in (3) equals zero.

Under the second interpretation of our analysis, the absence of a relationship between the participant impacts and our econometric estimates has an additional possible meaning, namely that participants have used the less-than-formal inferential methods of lay science to construct flawed impact estimates. However, the lay science explanation has additional implications: Variables used by participants to create their estimates, such as outcomes or before/after outcome differences (in the case of lay empiricists) or measures of program inputs (in the case of lay theorists depending on a theory of positive

marginal products of training and education inputs), should simultaneously display a relatively strong relationship with the participants' self-reported evaluations.

Finally, regardless of our interpretation, large estimated standard errors suggest that our econometric impact estimates embody substantial estimation error.

3.3. *Econometric specification: lay science, impact proxies and performance measures*

To further explore the lay scientist idea—that participants' evaluations depend on relatively crude proxy measures of program inputs and outputs—we adopt a more direct analog to equation (1) as our econometric specification. In particular, we assume that

$$(4) \quad \Delta_{SR} = 1(\beta_0 + \beta_1(\text{proxy}(Y_{1b} - Y_{0b})) + \beta_X X + \varepsilon > 0),$$

where ε has a logistic distribution and X is a vector of observable characteristics with corresponding coefficients β_X . This is, of course, a standard logit model, which means that we can identify the coefficients only up to scale; we report estimates of mean derivatives below.

We employ (4) rather than (3) in this case because the proxy variables we examine, such as labor market outcomes and the types of services received, do not contain any measurement or estimation error (unlike our econometric impact estimates). In addition, because we measure these variables directly, rather than predicting them as a linear combination of the X , we can include conditioning variables X . These conditioning variables soak up residual variance and thus make our estimates more precise. They may also proxy, in part, for $B_{1b} - B_{0b}$, $E(Y_{1a} - Y_{0a})$ and $E(B_{1a} - B_{0a})$, thus clarifying the interpretation of our estimates.

4.0 Econometric impact estimators using randomized treatment information

This section describes our two preferred econometric estimators (those that utilize the randomized assignment of the JNS study) for obtaining impact estimates that vary among participants.

4.1 Experimental impacts at the subgroup level

The first method we employ for generating impact estimates that vary among participants takes advantage of the experimental data and the fact that random assignment remains valid for subgroups defined on characteristics measured at or before random assignment, as discussed in, e.g. Heckman (1997b).

We estimate regressions of the form

$$(5) \quad Y_i = \beta_0 + \beta_D D_i + \beta_X X_i + \beta_I D_i X_i + \varepsilon_i,$$

where Y_i is some outcome measure, D_i denotes assignment to the experimental treatment group, X_i denotes a vector of characteristics measured at or before random assignment and $D_i X_i$ represent interactions between the characteristics and the treatment indicator.

It is these terms that yield variation in predicted impacts among individuals at the subgroup level. The predicted impacts based on (5) for the treatment group members are given by

$$(6) \quad \overline{Y}_{1i} - Y_{0i} = \hat{\beta}_D + \hat{\beta}_I X_i.$$

Though quite straightforward conceptually, our experimental subgroup impact estimates do raise a few important issues, which we now discuss. The first issue concerns the choice of variables to interact with the treatment indicator. We address this

issue by presenting two sets of estimates based on vectors of characteristics selected in very different ways. One set of estimates simply borrows the vector of characteristics employed by Heckman, Heinrich and Smith (2002) in their analysis of the JTPA data. The notes to Table 3 list these variables. The second set of estimates utilizes a set of characteristics selected using the somewhat unsavory method of stepwise regression. While economists typically shun stepwise procedures as atheoretic, for our purposes here that bug becomes a feature, as it makes the variable selection procedure completely mechanical. Thus, we can be assured of not having stacked the deck in one direction or another. In both cases, we restrict our attention to main effects in order to keep the problem manageable.

We implement the stepwise procedure using essentially all of the variables from the BIF including variables measuring participant demographics, site, receipt of means-tested monetary and in-kind transfers, labor force status and work history. We include a missing indicator for each variable (to avoid losing a large fraction of the sample due to item non-response), and interact both the variables and the missing indicators with the treatment group indicator. The stepwise procedure has to keep or drop each variable along with the missing indicator and interactions with the treatment indicator as a group. The stepwise procedure, which we perform separately for each of the four demographic groups, iteratively drops variables with coefficients not statistically different from zero in a regression with self-reported earnings in the 18 months after random assignment as the dependent variable.⁵

⁵ We employ the “step up” stepwise procedure as it has more power than the “step down” and “single step” procedures. See Dunnett and Tamhane (1992) and Lui (1997) for details. We set the p-value for choosing variables in the final specification at 0.05.

The second issue concerns the amount of subgroup variation in impacts in the NJS data within the four demographic groups – adult males and females ages 22 and older and male and female out-of-school youth ages 16-21 – for which both we and the official reports conduct separate analyses. Although the NJS impact estimates differ substantially between youth and adults (and between male and female youth when considering the full samples), the experimental evaluation reports – see Exhibits 4.15, 5.14, 6.6 and 6.5 in Bloom et al. (1993) for the 18 month impacts and Exhibits 5.8, 5.9, 5.19 and 5.20 in Orr, et al. (1994) for the 30 month impacts – do not reveal a huge amount of statistically significant variation in impacts among subgroups defined by the observables available on the BIF. If the impact does in fact vary a lot among individuals, but not in a way that is correlated with the characteristics we use in our model, then we may reach the wrong conclusions about participants’ ability to construct consistent estimates of earnings impacts. This case has more than academic interest given that Heckman, Smith and Clements (1997) bound the variance of the impacts in the JTPA data away from zero for adult women; their lower bound on the standard deviation of the impacts equals \$674.50 with a standard error of \$137.53 (see their Table 3).⁶ In addition to simply keeping it in mind, we attempt to address this concern in part by examining the quantile treatment effect estimates described in the next section, which do vary a lot among participants, and by looking, in other work, at data from other experimental evaluations with more in the way of subgroup variation in impacts.

The third issue concerns an additional assumption that we must make in order to interpret our results in the way that we have described here. A simple example illustrates

⁶ Our subgroup impacts have standard deviations that range from \$840 to \$2600 depending on the demographic group and set of covariates. The quantile treatment effects have lower standard deviations; they range between \$257 and \$477.

the need for this assumption. Suppose participants care only about earnings impacts, and give a positive survey evaluation only if they receive a positive earnings impact. Now consider two groups. In group one, just 10 percent of the individuals receive a \$1000 impact while the rest receive no impact, so that the mean group one impact is \$100 and the fraction giving positive evaluations is 0.1. In contrast, 20 percent of group two individuals receive a \$400 impact while the rest receive no impact, so that the mean group two impact is \$80 and the fraction giving positive evaluations is 0.2. This example shows that subgroup mean impacts could vary inversely with the fraction receiving a positive impact. In interpreting our results below, we assume that this does not occur in our data. Put differently, we assume that mean impacts and the fraction with a positive impact have positive covariance across subgroups.

4.2 Quantile treatment effects

The second econometric method we use to derive individual level treatment effect estimates relies on an additional non-experimental assumption. In particular, we make the assumption of a perfect positive rank correlation between the outcomes in the treated and untreated states described in Heckman, Smith and Clements (1997). Intuitively, we assume that the expected counterfactual for an individual at a given percentile of the treatment group outcome distribution consists of the mean outcome at the same percentile of the control group outcome distribution. One way to think about this assumption is that expected labor market outcomes depend on a single factor, so that individuals who do well in the treatment state also do well in the control state. This represents a very different view of the world than, for example, the classic model of Roy (1951), but may

represent a reasonable approximation for treatments, such as those offered by JTPA, that have small average impacts relative to the mean outcomes in question.

Using this method, we estimate the impact for treated individual “ i ” with an outcome at percentile “ j ” of the treatment group outcome distribution as

$$(7) \quad \hat{Y}_{i1} - Y_{i0} = \hat{Y}_1^{(j)} - \hat{Y}_0^{(j)},$$

where the superscript “ (j) ” denotes the percentile. This estimator underlies the literature on quantile treatment effects, as in Abadie, Angrist and Imbens (2002) and Bitler, Gelbach and Hoynes (2004), with the difference that rather than interpreting the estimates as the effect of treatment on the quantiles of the outcome distribution, we make the additional rank correlation assumption. As discussed in Heckman, Smith and Clements (1997), the rank correlation assumption pins down the joint distribution of outcomes, which in turn pins down which quantile of the control outcome distribution provides the counterfactual for each quantile of the treatment outcome distribution and allows us to assign impact estimates to specific individuals.

5.0 The relationship between econometric impact estimates and participant evaluations

5.1 Bivariate relationships

We begin our analysis of the data from the NJS with simple bivariate relationships between mean experimental impacts and fractions of treated participants giving a positive evaluation of JTPA. This analysis, presented in Table 1, extends that presented in Table 8.11 of Heckman and Smith (1998). It represents a very basic application of the methodology outlined in Section 4.1.

Throughout our discussion of results, we refer to adult males, adult females, youth males and youth females as “the four basic (demographic) groups” of the experiment. The first four rows of Table 1 present statistics for each of these groups. The first column presents the fraction of treated individuals in each group that give JTPA a positive evaluation in the NJS survey. The remaining eight columns report experimental impact estimates (for each group) for eight different earnings and employment outcomes. The first two outcomes are self-reported earnings, and any self-reported employment, in the first 18 months after random assignment, which corresponds roughly to the whole period prior to the survey response for most sample members. The second two outcomes instead focus on respondent’s status right around the time of the survey question: These are self-reported earnings, and any employment, in month 18 after random assignment. The remaining four outcomes repeat the first four, but use the quarterly data from the matched UI earnings records instead of respondents’ self-reported outcome data. We include both sets because they appear substantially different at both the individual and aggregate levels – see the discussions in Smith (1997a,b) and Kornfeld and Bloom (1999). Throughout our discussion of results, we refer to these collectively as the “eight outcomes.” Entries in the last row of each outcome column in Table 1 are the correlation (and p-value against zero correlation) between the fractions of treated individuals giving a positive evaluation (the first column) and the impact estimates in that column.

Table 1 reveals that male youths, who have the worst experimental impact estimates for all eight outcomes, also give the second highest fraction of positive evaluations. Female youths, who usually have the second-worst impact estimates for the eight outcomes, also have the highest fraction of positive self-reported evaluations.

Consistent with this basic pattern, the correlations reported in the last row end up negative six out of eight times, though we can never reject the null of a zero correlation (which is no surprise, since each correlation is based on four observations—one for each group). At this crude level, then, we find little evidence of a positive association between participant evaluations and experimental impact estimates.

Given that the survey question wording focuses more on employment than earnings, one might expect a more consistent picture in employment outcomes alone. Yet there is no clear pattern of correlation signs across the four alternative employment impact measures (two are positive and two are negative). More broadly, if we make the rank correlation assumption described in Section 4.2, the experimental impacts show that the program improved the employment situation of at most a few percent of the respondents, yet well over half gave positive self-reports in all four groups.

Tables 2A, 2B, 2C and 2D (one table for each of the four basic demographic groups) show summaries of a similar analyses for a variety of subgroups. As before, columns correspond to the eight outcomes. Row labels indicate a variable used to create finer subgroups, and each row contains correlation coefficients (and their p-values) between positive self-evaluation fractions and mean estimated experimental impacts across those finer subgroups (similar to the final row of Table 1). The variables we use are race/ethnicity, years of schooling categories, marital status, time since last employment categories, site and age categories (except for youth groups, in which age is already limited to individuals from 16 to 21 in any case).

The bottom of each table presents the number of positive and negative correlations in the table and, for each sign type, the number of statistically significant

coefficients at the five and ten percent levels. These counts of “significant votes” provide a useful summary of the 48 (or 40 for the youth) entries in each table. For adults, for instance, random variation would lead to an expected 4 or 5 significant votes at the ten percent level, and 2 or 3 at the five percent level. Of course, this ignores the lack of independence among the estimates used to create each table, nor does it attempt to weight or value the different estimates based on their precision. The results in Tables 2A to 2D reinforce the impression made by Table 1. No clear patterns emerge in terms of coefficient signs, and the number of statistically significant correlations resembles what one would expect if the population coefficients all equal zero.

5.2 Regression results for experimental subgroup estimates

We now turn to evidence based on regressions of experimental subgroup impact estimates for the eight outcomes on participant evaluations. In terms of our earlier discussion, we report estimates of β_1 from equation (3), where the dependent variable is an experimental impact estimate based on subgroup variation in impacts (estimated from equations (5) and (6)). Table 3 shows the results. Each entry in the table shows β_1 and its standard error from a different regression. Rows show all regression results using one of the eight impact estimates as the dependent variable. Columns are grouped into four pairs of two columns. Each pair of columns shows regression results in one of the four basic demographic groups. Within each column pair, the two columns show results based on the two different sets of covariates used to estimate equation (5). The columns headed by (1) contain the estimates based on the covariates from Heckman, Heinrich and Smith (2002), while the columns headed by (2) contain the estimates using the covariate set

chosen by the stepwise procedure described earlier. The final two rows of the table summarize the evidence in each column; in particular, they give the numbers of positive and negative estimates and, within each category, the number of statistically significant estimates at the five and ten percent levels.

The regression evidence in Table 3 suggests little, if any, relationship between the experimental impact estimates based on subgroup variation and the self-reported evaluations, echoing the results of Tables 1 and 2. While the estimates lean negative in the aggregate, only a handful of the estimates reach conventional levels of statistical significance (and not all of those fall on the negative side of the ledger). Either the participants do not weigh labor market impacts very heavily in their response, or else their impact estimates (or ours) do not do a very good job of capturing the actual impacts.

5.3 Results based on quantile treatment effect estimates

This section presents evidence on the relationship between impact estimates constructed under the perfect positive rank correlation assumption described in Section 4.2. We focus on one particular outcome in this analysis: Self-reported earnings over the whole eighteen months since randomized assignment. Figures 1A to 1D begin with graphical presentation of the evidence for each of the four basic demographic groups. The horizontal axis in each figure corresponds to percentiles of the untreated outcome distribution. The solid line in each graph presents impact estimates at every fifth percentile (5, 10, 15, ..., 95) constructed as in equation (7). The broken line in each graph represents an estimate of the fraction with a positive self-reported evaluation at every fifth percentile. For percentile “j”, this estimate consists of the fraction of the

treatment group sample members in the interval between percentile “j-2.5” and percentile “j+2.5” with a positive self-reported evaluation. If the assumptions underlying the percentile difference estimator hold, if participants care enough about labor market outcomes in answering the survey question, and if participants consistently estimate their own impacts, then the two lines should move together in the figures.

Several features of the figures merit notice. First, in the lower percentiles in each figure the econometric impact estimate equals zero. This results from the fact that the lowest percentiles in both the treated and untreated outcome distributions have zero earnings in the 18 months after random assignment; the difference between the two then equals zero as well. Surprisingly, a substantial fraction (over half in all four demographic groups) of treatment group members at these percentiles respond positively to the survey evaluation question, even though they have zero earnings in the 18 months after random assignment. This could mean that respondents view the question as asking about longer term labor market impacts beyond the 18-month window of the measured outcomes; but it also could mean that respondents are acting as lay theorists (more on this shortly).

Second, the fraction with a positive self-reported evaluation has remarkably little variation across percentiles of the outcome distribution. For all four demographic groups, it remains within a band from 0.6 to 0.8. For the adults, the mean increases with the percentile; for the youth, the data fail to reveal a clear pattern.

Third, no obvious relationship between the two variables emerges from the figures except for adult females: for them, both variables increase with the percentile of the outcome distribution. More specifically, for adult women, both variables have a

higher level for percentiles where the impact estimate exceeds zero. Within the two intervals defined by this point, both variables remain more or less constant.

Table 4 presents some of the numbers underlying the figures. In particular, the first five rows present the values for the 5th, 25th, 50th, 75th and 95th percentiles. The last two rows of the table give the correlation between the quantile treatment effects and the fraction with a positive self-reported evaluation for each group (and the corresponding p-value from a test of the null that the correlation equals zero) along with the estimated coefficient from a regression of the quantile treatment effects on the fraction with a positive self-reported evaluation (and the corresponding standard error). The correlation and regression estimates quantify and confirm what the figures indicate: a strong positive relationship for adult women, a weak and statistically insignificant positive relationship for adult men, and moderately strong and negative relationship for male youth and a similar, but not statistically significant, relationship for female youth. Although we find a bit more here than in the estimates that rely on subgroup variation, once again the data do not suggest a strong, consistent relationship between the econometric impact estimates and the self-reported evaluations.

To summarize this section, we find little evidence of positive relationships between treated participants' evaluations of JTPA's effects on them and individual experimental impact estimates based on subgroup variation or the perfect rank correlation assumption of the quantile estimator. At the same time, average evaluations are positive in almost all subgroups, and across distributions of outcomes, and show remarkably little variance across these; this is even true for individuals with no earnings over the eighteen months following randomized assignment. These results motivate the next section.

6.0 Relationship between positive self-evaluation and proxies for impacts

6.1 Motivation and caveats

The largely negative results of section 5 could of course occur if our econometric impact estimates lack the precision (and/or consistency, in the case of the quantile estimator) required to correlate well with participant evaluations. [*Jeff to elaborate here with citations and text.*] Beyond this, the results are consistent with our subjective rationality interpretation: Participants may care mostly about effects of JTPA not captured by labor market outcomes over specified time periods. Yet the results also strongly suggest lay scientists at work. For instance, consider the finding that participant evaluations are both positive and vary remarkably little across groups, subgroups and quantiles of impact distributions. This suggests that participants' evaluations may be largely theory-driven inferences based on shared folk theories. In particular, we suspect that participants may share the theory that impacts are monotone increasing in inputs (the expense or resource-intensiveness of program services received). To explore this, we estimate relationships between participant evaluations and services received by participants, and expect relatively large positive effects for relatively expensive (resource-intensive) services.

An obvious interpretive caveat is that different program services may themselves have different subjective or direct costs and/or benefits not captured by labor market outcomes. For example, classroom training may be more fun (or more tedious) than, say, job search assistance. Alternatively, classroom training at a distant community college imposes higher direct costs than job search assistance provided at a local neighborhood location. Thus, the subjective rationality interpretation also allows for significant

relationships between participant evaluations and service types—though it makes no obvious prediction about the direction of those relationships.

Our results so far also suggest that some participants may act as lay empiricists, making judgments based on proxy variables that correlate only weakly with true impacts, and perhaps with insufficient notice of potential confounds. If so, their evaluations can be both inconsistent and full of nuisance variance, undermining any relationship between them and consistent impact estimates. The proxies we examine are actual labor market outcomes (employment and earnings) and simple before-after differences in those outcomes. If respondents really do know the impacts, then such proxies should have little explanatory power...*[Jeff to add citations and/or discussion elaborating here]*.

With both these big picture interpretations and caveats in mind, we now turn to estimation of equation (4) in Section 3.3, using inputs (program service types), outcomes and/or before-after outcome differences as regressors.

6.2 Results with service types

Table 5 presents logit estimates of equation (4) that include not one but two measures of the service type received by JTPA treatment group members. The two measures of service receipt derive from self-reports collected in the NJS follow-up surveys and from administrative data from the individual sites participating in the experiment.⁷ As shown in Smith and Whalley (2005), these two measures differ substantially; as a result we do not run into collinearity problems when including them both. The two data sources code

⁷ In fact, two versions of the administrative data on service receipt exist, one created by MDRC and one created by Abt Associates. Both rely on the original MIS files from the 16 sites in the experiment. Our experience with both files, described in detail in Smith and Whalley (2005), leads us to employ the Abt version in this paper.

the service types somewhat differently; for comparability and ease of interpretation, we employ just five service types: CT-OS, OJT/WE (which is almost all OJT), JSA, ABE and “other”. We code a dummy variable for each service type in each data source indicating whether or not the respondent received it; a respondent who received more than one service type in a given data source gets coded based on the training type they receive in their first spell. We expect positive evaluations to be relatively more likely when relatively expensive service types, such as CT-OS or OJT, are received.

The logit models presented in Table 5 also include a variety of background variables. These variables play two roles. First, we expect them to pick up parts of the overall impact of participation unrelated to the labor market outcomes we examine. For example, the site dummies will pick up differences in the friendliness and efficiency of site operation as perceived by the respondents. The variable “work for pay”, which is an indicator variable for whether or not the respondent has ever worked for pay, relates to the opportunity cost of participation, as does the variable for having a young child. The AFDC receipt at random assignment variable captures variation in the cost of classroom training due to the availability of an income source not tied to employment. In order to avoid losing a large fraction of the sample due to item non-response, we recode missing values to zero and include indicator variables for missing values of each variable.

The columns in Table 5 give results for the four basic demographic groups. The table presents mean derivatives, estimated standard errors for the mean derivatives in parentheses, and the p-value of a test of the null hypothesis that the mean derivative equals zero in square brackets. Table 6 summarizes the results in Table 5 by presenting

test statistics and p-values from tests of the joint null that the mean derivatives for groups of related covariates (e.g. all of the self-reported training type variables) equal zero.

Consider the variables other than the service type variables first. Although they are not shown in the table, the site variables have a strong effect on the probability of a positive self-reported evaluation. The magnitudes vary a lot as well; for example, for adult males the coefficients on the site dummies range from -0.257 to 0.093. Moreover, Table 6 shows that these variables are strongly statistically significant as a group. Respondents may take account of non-pecuniary aspects of their JTPA experience, such as staff friendless and efficiency, or the attractiveness and ease of access of the JTPA office and the local service providers. Variation of local environmental conditions across sites, such as hiring opportunities, also might affect respondents' evaluations through an influence on outcomes. Although there might also be site differences in program impacts, this seems less likely given the findings in Section 5.2.

With the exception of age for adults, race for youth, and age and education for female youth, the other demographic variables play surprisingly little role in determining the probability of a positive self-reported evaluation. Among adults, age has a strong negative effect on the probability of a positive self-evaluation, while black male youth and Hispanic male and female youth have higher probabilities of a positive response. The limited role played by background characteristics in the analysis surprised us.

In contrast to these background characteristics, service type plays a major role in determining individual self-reported evaluations. Table 6 shows that, taken together, both the self-reported and administrative training type variables achieve high levels of statistical significance for adults, and the administrative measures do so for youth.

Smith and Whalley (2005) show that the self-reported and administrative measures of receipt of classroom training in occupational skill tend to agree; as such, we can (as a crude approximation) simply add their coefficients. Doing so reveals that CT-OS has a large positive effect on the probability of a positive self-reported evaluation for all four groups. Subsidized on-the-job training reports often do not coincide in the two data sources; here we find that self-reported OJT has a strong (and usually statistically significant) positive effect, as does administratively reported OJT! Recall that CT-OS and OJT generally represent the largest resource investment in the JTPA participant.

In contrast, job search assistance (JSA) is the cheapest of the services and, as expected, elicits less of a positive effect. This service type also tends to get reported differently in the two data sources. Here, except for self-reported JSA for adult males and administratively reported JSA for adult females, we find modest and statistically insignificant effects. Adult basic education (ABE) and “other” services do not yield precise estimates, except for “other” services in the administrative data, which has, somewhat puzzlingly, a negative and statistically significant effect for adult women and a positive and statistically effect for female youth. These two service types are, in general, least frequently assigned, which may help explain the imprecision of the estimates.

Keep in mind the caveat that service types might matter to participant evaluations under the subjective rationality interpretation of these results. Still, service types matter in the direction expected on the basis of the lay theorist notion: Relatively more expensive services are associated with relatively higher likelihoods of a positive evaluation.

6.3 Results with labor market outcomes

Tables 7 and 8 report results of various estimates of logit equation (4) that include the same background variables as the models of Table 5, but add various versions of Y_1 , the labor market outcome in the treated state. Acting as lay empiricists, respondents may be relatively more likely to infer a positive program impact if they have done relatively well in the labor market over the period between random assignment and the survey, or if they are doing well around the time of the survey.

The top panel of Table 7 shows results of a specification where we divide self-reported earnings in the 18 months after random assignment into five categories: zero, and four quartiles of the distribution of positive earnings, and then include dummy variables for four of these five categories, with the highest category as the excluded category. The second panel of Table 7 reports results of the same specification, but using UI administrative earnings (in the six calendar quarters after random assignment) instead of self-reported earnings. The last two lines of Table 7 report results of a specification that includes a dummy variable for any employment in the 18 months after random assignment, first based on self-reported employment and then on UI employment data.

The broad picture from Table 7 is that labor market outcomes appear to predict self-reported evaluations. This is particularly true for adults. For both adult males and adult females, and for both earnings measures, all of the estimated coefficients are negative (as expected when compared to the highest earnings quintile) and most are statistically significant. Also, broadly speaking, the estimated coefficients decrease as earnings increase, as expected. The coefficients also turn out largely negative (indicating a positive relationship) for youth, but are rather imprecisely estimated. The self-reported employment measure also has a strong positive and statistically significant relationship to

the self-reported evaluations, but the UI employment outcome measure does not. This latter finding may result from measurement error in UI employment (due to its omission of government jobs and informal jobs).

Table 8 summarizes the evidence in Table 7, as well as evidence from alternative specifications not fully reported here for reasons of space. As in Table 6, the summary takes the form of chi-square statistics, and their p-values, for tests of the null hypothesis that a coefficient (or all coefficients) on a specific labor market outcome measure (or vector of outcome measures) equal zero. These results show that relationships tend to be statistically stronger for adults than for youth, and stronger for measures based on self-reported data than UI data. Also, earnings measures tend to yield more statistically significant relationships than employment measures, especially for outcomes at or just around the time of the survey.

Overall, the evidence that participants use outcomes as proxies for impacts is strong. This is consistent with lay science, but with real science too: Outcomes may be correlated with actual impacts. Indeed, some fraction of the treated group would have zero counterfactual (untreated) earnings outcomes after the treatment period; in the case where they also have zero earnings after treatment, outcomes and impacts perfectly coincide. We cannot rule out this possibility by an appeal to our section 5 results since the measurement error in our own impact estimates may be quite large.

6.4 Results with before-after comparisons of labor market outcomes

In this section, we explore the relationship between participant evaluations and before-after comparisons of employment and earnings. The cognitive appeal and simplicity of

before-after comparisons as an estimator of impacts is undeniable. Moreover, despite their simplicity, before-after comparisons are consistent impact estimates in the absence of confounds, that is, if there is no change in any outcome-relevant factor over the period between the two measurements (initial outcomes will then consistently estimate the final outcome that would have occurred in the absence of treatment). Unfortunately, there is a dramatic confound over the NJS experiment period. Heckman and Smith (1999) show that, because of “Ashenfelter’s dip” in earnings in the pre-program period, before-after impact estimates tend to have a strong upward bias. Remembering that lay empiricists may fail to correct for nonsalient confounds, participants making judgments on the basis of before-after comparisons may well fail to appreciate a confounding factor like Ashenfelter’s dip and, as a result, would also produce upward-biased impact estimates.

Tables 9 and 10 present estimates of logit models with participant evaluations as the dependent variable and three different measures of before-after earnings changes as independent variables, along with all of the variables in Column 1 of Table 5. The first measure, for which the estimates appear in Table 9, consists of the difference in average monthly earnings between the 12 months before random assignment and the 18 months after random assignment. We can use only the 12 months before random assignment due to the limitations of the survey data on pre-random assignment earnings for the treatment group.⁸ The second and third measures, for which the estimates appear in Table 10, rely on the UI earnings data. The first measure, in the top panel of Table 10, consists of mean monthly earnings in the six calendar quarters after random assignment minus mean monthly earnings in the six calendar quarters before random assignment. The second measure, denoted UI(2) in the lower panel of Table 10, consists of the difference in mean

⁸ In particular, we only have the response to a question about earnings in the previous year from the BIF.

monthly earnings between just the sixth quarter after and the sixth quarter before random assignment. In each case, we let the data speak to the functional form by including indicator variables for quintiles of the before-after difference.

We find strong evidence that before-after differences in labor market outcomes predict self-reported impacts. For the self-reported measure, the relationship is clearest for the adult females and the male youth, where the estimated coefficients increase monotonically (or almost so) and are statistically and substantively significant for the upper quintiles. Even stronger findings appear for the UI earnings difference measure in the top panel of Table 10, with large, and almost always statistically significant, coefficients for all four groups for the two upper quintiles. For male youth, the key difference seems to be between the lowest quintile and the other four; the four coefficients are all relatively large, all about equal and all statistically significant. For the other three groups, there is a general pattern of increasing coefficients as one moves down the table. The results for the UI(2) measure in the bottom panel of Table 10 are weaker, in both a substantive and a statistical sense, than those in the top panel; this suggests that respondents use outcomes over the entire pre- and post-random assignment periods in constructing their implicit before-after estimates of program impact.

Given that the second of the two survey questions that compose our self-reported evaluation measure asks directly about finding a job, in Table 11 we consider its relationship to before-after employment changes. We coded an employment status difference variable based on employment at the date of random assignment and 18 months after random assignment. This yields four patterns. We include dummy variables for three of the four patterns, with employed at both points in time as the

omitted pattern. The findings here are, perhaps, less strong than expected. In general, relative to the always employed, those who are never employed or who lose a job tend to have less positive self-reported evaluations. For the adults, those who gain a job tend to be somewhat more positive. However, only a handful of the differences achieve statistical significance. Measurement error in the “after” employment status may account for our weak results. By looking at employment around the time of the survey, we have given the respondents plenty of time to lose jobs that JTPA helped them find and, in the case of dropouts, to find jobs without the help of JTPA.

Finally, Table 12 presents the results of chi-squared tests for the joint significance of the before-after difference variables considered in Tables 9, 10 and 11. The test statistics and p-values in this table confirm that respondents’ self-reported evaluations depend (in a statistical sense) on these before-after differences. Indeed, the joint tests for the employment change variable look stronger than the individual t-tests, suggesting that our omitted group lies in the middle of the categories in terms of its effect on the self-evaluation measure. Overall, the findings in this section lend support to the view that respondents implicitly or explicitly use natural and cognitively simple (but nonetheless quite biased) before-after comparisons in constructing their self-reported evaluations.

7.0 Results with performance measures

In this section we present results on the relationship between participant self-reported evaluations and performance measures based on program outcomes commonly used in employment and training programs both in the U.S. and elsewhere. Performance standards systems attempt to provide information on the impacts of programs quickly and

at low cost by relying on crude proxies. In this sense, as the reader will quickly discern, some of the performance measures are fairly closely related to the outcome proxies examined in Section 6.2. In the JTPA program, the performance measures had real effects on site budgets; sites that did well on them received budgetary rewards while sites that did exceptionally poorly could receive “technical assistance” and also experienced the threat of formal reorganization. See, e.g., Heckman, Heinrich and Smith (2002), Heckman and Heinrich (2005) and Barnow and Smith (2004) for more detailed descriptions of the performance standards systems in JTPA, WIA and other programs, for evidence that the usual performance measures have little, if any, relationship to the actual impacts of the program, for evidence of strategic behavior by program staff in response to the incentives provided by the JTPA performance standards system, and for additional pointers to the literature.

The performance measures we examine here are subsets of those included in the JTPA and WIA performance standards systems. From the JTPA system we consider employment status at termination from the program, wages at termination from the program (which is defined only for those employed at termination), employment at “follow-up”, which is 13 weeks after termination, and weekly earnings (not including zeros) at follow-up. We use self-reported information to construct the JTPA performance measures, as was done in that program. From the WIA system we consider employment at termination, employment at six months after termination conditional on employment at termination (this measure aims to count “retention”, although it does not require the individual to stay at the same job), and the difference in quarterly earnings between the two calendar quarters after termination and the two quarters prior to random assignment.

Note that the earnings gain measure, which was an innovation in the WIA system relative to JTPA, will contain the pre-program dip in mean earnings discussed in, e.g. Heckman and Smith (1999), and therefore will invariably suggest positive earnings impacts whether the program works or not. We follow the WIA program in relying mainly on the UI earnings records in constructing the WIA performance measures for our sample.

The top panel of Table 13 presents results based on estimating logit models with self-reported evaluations as the dependent variable and one of the performance measures, again including all of the variables in Table 5 as covariates. The bottom panel reports results from similar models using the WIA performance measures.

Three patterns emerge from the findings in Table 13. First, among the JTPA measures, employment at termination and employment at follow-up are significantly related, in both senses, to self-reported evaluations. The estimated mean change in the probability of a positive self-evaluation due to employment at termination ranges from 0.08 for male youth to 0.13 for adult females. Employment at follow-up shows a similarly strong relationship. Second, the WIA measures, other than the earnings change measure, show little in the way of a consistent relationship with the self-reported evaluation measure. In particular, using the UI earnings data to measure employment at termination rather than survey data on employment spells adds enough measurement error to yield a much weaker relationship for all four groups and one that is statistically significant only for adult males. Third, the relationship between self-reported evaluations and the performance measures appears stronger for women than for men.

Overall, Table 13 yields a mixed picture. Some performance measures based on labor market outcomes have substantively and statistically significant relationships with

self-reported evaluations but even these account for only a modest fraction of the variance. Thus, the self-reported evaluations capture something related to, but very different from, the performance measures.

8.0 Conclusions

Broadly speaking, and putting aside the material in Section 7 regarding the performance standards, we have two main findings. The first is that self-reported evaluations by treatment group members from the JTPA experimental evaluation have, in general, little if any relationship to either experimental impact estimates at the subgroup level or to what we regard as relatively plausible econometric impact estimates based on percentile differences. The second is that the self-reported evaluation measures do have consistent relationships with crude proxies for impacts, such as measures of service type (a proxy for resources expended on the participant), labor market outcome levels (which measure impacts only if the counterfactual state consists of no employment or earnings, which it does not for the vast majority of our sample), and before-after comparisons.

Taken together, these two findings provide strong support for the view that respondents avoid the cognitive burden associated with trying to construct (implicitly or explicitly) the counterfactual outcome they would have experienced had they been in the control group and thus excluded from JTPA. Instead, they appear to act as lay scientists, using readily available proxies and simple heuristics to conclude, for example, that if they are employed at the time of the survey or if their earnings have risen relative to the period prior to random assignment, that the program probably helped them find a job or get a

better job. At the same time, our evidence does not rule out the view that respondents consider factors in their answers not captured in our experimental and econometric impact estimates, such as expected impacts in later periods or subjective and/or direct costs and benefits associated with the services they received. The proxy variables still leave much variation in the self-evaluation measure to be explained by other factors.

However, our “lay science” interpretation of our results is in fact borrowed from a large literature in social psychology on the fallibility of self-reports. The “study skills” experiment of Conway and Ross (1984) is the most parallel study we know of. Conway and Ross recruited subjects from one large introductory psychology course who expressed interest in taking a three-week study skills class, and randomly assigned them to either the class (treatment) or a waiting list (control). Both groups gave self-reports on their own study skill proficiency both before and after the three-week class. Since subjects came from one course, and the experiment took place between the midterm and final in that course, comparable performance measures (in the form of grades on the midterm and final in the same class) were available to Conway and Ross, as were overall semester grades collected from registrar records.

The objective measures confirmed what professional evaluators of such classes, e.g. Gibbs (1981), have found: The class had no significant effect on outcome measures (grades in the course or on the final, or for the semester overall). Yet treatment subjects reported significantly greater improvement in their study skills, and expected significantly better grades, than did control subjects. Curiously, the former effect seems mainly traceable to a biased recall of pre-class study skill level. In the “after” surveys, Conway and Ross (1984) also asked subjects to recall what their “before” evaluation of

their study skills had been. Although treatment and control subjects did not actually differ in their before and after self-reports of their study skills, treatment subjects recall of their “before” study skills were significantly lower than they had actually reported in “before” surveys, while control subjects showed no such bias. Conway and Ross interpret these results as showing that their subjects act as lay theorists. Subjects have a theory that a study skills class will improve study skills. Treatment and control subjects’ “after” evaluations of their study skills differ little; put differently, there is no evidence that treatment subjects are exaggerating their “after” skill. If their memory of their “before” study skill level is hazy, they may instead opt to theoretically infer positive change from the class, and then postdict their before study skill level on the basis of that theoretical prediction. This kind of process wholly explains Conway and Ross’ findings.

The literature on using surveys to measure expectations, as discussed in Manski (2004), provides some hope that more sophisticated survey questions might do a better job of measuring the underlying objects of interest. Depending on whether and how expectations and evaluations are normally constructed by human minds, recall processes could be a crucial part of both of these. It is worth mentioning that some forms of biased recall of “before” judgments (similar but not identical to the kind studied by Conway and Ross 1984) are known to be mitigated if the subject’s actual “before” judgment is known and the subject is given a real incentive to correctly recall it (Aderman and Brehm 1976). Put differently, the experimental economist’s normal emphasis on incentive-compatible mechanisms for truthful reports might be helpful here, and such methods are increasingly incorporated into various surveys.

Overall, we conclude that participant self-evaluations of the type analyzed here are very poor substitutes for rigorous experimental or non-experimental estimates of program impact. Although our paper is (to our knowledge) the first to seriously study what these questions actually measure in social experiments, the discussion above shows that our findings, and our interpretations of them, are far from unique: Very similar results can be found in many laboratory experiments conducted by social psychologists and, indeed, we have drawn our “lay science” interpretation from this literature. We can think of no better summary warning than the one Ross (1989, p. 354) gave in his musings on the implications of Conway and Ross (1984):

“The biased retrospections evidenced by participants in improvement programs may cause them to remain in worthless programs and not to search for more effective treatments. Indeed, when self-reports are a primary indicant of improvement, a conspiracy of ignorance may emerge in which both the helper and the helped erroneously believe in the achievement of their common goal.”

References

Abbring, Japp, James Heckman and Edward Vytlacil. 2005. "Econometric Evaluation of Social Programs" in James Heckman and Edward Leamer (eds.) *Handbook of Econometrics, Volume 6*. Amsterdam: North-Holland: forthcoming.

Abadie, Alberto, Joshua Angrist and Guido Imbens. 2002. "Instrumental Variables Estimates of the Effect of Subsidized Training on the Quantiles of Trainee Earnings." *Econometrica*. 70(1): 91-117.

Aderman, David and Sharon S. Brehm. 1976. On the recall of initial attitudes following counterattitudinal advocacy: An experimental reexamination. *Personality and Social Psychology Bulletin* 19:59-62.

Angrist, Joshua. 2004. "Treatment Effect Heterogeneity in Theory and Practice." *Economic Journal* 114(494): C52-C83.

Barnow, Burt and Jeffrey Smith. 2004. "Performance Management of U.S. Job Training Programs: Lessons from the Job Training Partnership Act." *Public Finance and Management* 4(3): 247-287.

Bitler, Marianne, Jonah Gelbach and Hilary Hoynes. 2004. "What Mean Impacts Miss: Distributional Effects of Welfare Reform Experiments." Unpublished manuscript, University of Maryland.

Bloom, Howard, Larry Orr, George Cave, Stephen Bell and Fred Doolittle. 1993. *The National JTPA Study: Title II-A Impacts on Earnings and Employment at 18 Months*. Bethesda, MD: Abt Associates.

Bloom, Howard, Larry Orr, Stephen Bell, George Cave, Fred Doolittle, Winston Lin, and Johannes Bos. 1997. "The Benefits and Costs of JTPA Title II-A Programs: Findings from the National Job Training Partnership Act Study." *Journal of Human Resources* 32: 549-576.

Carniero, Pedro, James Heckman and Edward Vytlacil. 2005. "Understanding What Instrumental Variables Estimate: Estimating Marginal and Average Returns to Education." Unpublished manuscript, University of Chicago.

Conway, Michael, and Michael Ross. 1984. Getting What You Want by Revising What You Had. *Journal of Personality and Social Psychology* 47:738-748.

Devine, Theresa, and Heckman, James. 1996. "The Structure and Consequences of Eligibility Rules for a Social Program" in Solomon Polachek (ed.) *Research in Labor Economics Volume 15*. Greenwich, CT: JAI Press. 111-170.

- Dominitz, Jeff and Charles Manski. 1996. "Eliciting Student Expectations of the Returns to Schooling." *Journal of Human Resources*. 31(1): 1-26.
- Doolittle, Fred, and Linda Traeger. 1990. *Implementing the National JTPA Study*. New York, NY: Manpower Demonstration Research Corporation.
- Dunnett, Charles and Ajit Tamhane. 1992. "A Step-Up Multiple Test Procedure." *Journal of the American Statistical Association* 87(417): 162-170.
- Greenberg, David. 1997. "The Leisure Bias in Cost-Benefit Analyses of Employment and Training Programs." *Journal of Human Resources* 32(2): 413-439.
- Greenberg, David and Daniel Shroder. 2004. *Digest of the Social Experiments, Third Edition*. Washington, DC: Urban Institute Press.
- Heckman, James. 1997a. "Instrumental Variables: A Study of the Implicit Behavioral Assumptions Used in Making Program Evaluations." *Journal of Human Resources* 32(3): 441-452.
- Heckman, James. 1997b. "Randomization as an Instrumental Variable." *Review of Economics and Statistics* 78(2): 336-341.
- Heckman, James and Carolyn Heinrich, eds. 2005. *Performance Standards in a Government Bureaucracy*. W.E. Upjohn Institute for Employment Research.
- Heckman, James, Carolyn Heinrich, and Jeffrey Smith. 2002. "The Performance of Performance Standards." *Journal of Human Resources* 36(4): 778-811.
- Heckman, James, Neil Hohmann, Jeffrey Smith, with the assistance of Michael Khoo. 2000. "Substitution and Drop Out Bias in Social Experiments: A Study of an Influential Social Experiment." *Quarterly Journal of Economics* 115(2): 651-694.
- Heckman, James, Robert LaLonde, and Jeffrey Smith. 1999. "The Economics and Econometrics of Active Labor Market Programs" in Orley Ashenfelter and David Card (eds.), *Handbook of Labor Economics, Volume 3A*. Amsterdam: North-Holland, 1865-2097.
- Heckman, James and Jeffrey Smith. 1998. "Evaluating the Welfare State" in Steiner Strom (ed.), *Econometrics and Economic Theory in the 20th Century: The Ragnar Frisch Centennial*. Cambridge University Press for Econometric Society Monograph Series, 241-318.
- Heckman, James, and Jeffrey Smith. 1999. "The Pre-Program Earnings Dip and the Determinants of Participation in a Social Program: Implications for Simple Program Evaluation Strategies." *Economic Journal* 109(457): 313-348.

Heckman, James, and Jeffrey Smith. 2000. "The Sensitivity of Experimental Impact Estimates: Evidence from the National JTPA Study" in David Blanchflower and Richard Freeman (eds.), *Youth Employment and Joblessness in Advanced Countries*, Chicago: University of Chicago Press for NBER, 331-356.

Heckman, James, and Jeffrey Smith. 2004. "The Determinants of Participation in a Social Program: Evidence from the Job Training Partnership Act." *Journal of Labor Economics* 22(4): 243-298.

Heckman, James, Jeffrey Smith, with the assistance of Nancy Clements. 1997. "Making the Most Out of Programme Evaluations and Social Experiments: Accounting for Heterogeneity in Programme Impacts." *Review of Economic Studies* 64(4): 487-535.

Heckman, James and Edward Vytlacil. 2005. "Structural Equations, Treatment Effects and Econometric Policy Evaluation. *Econometrica*, forthcoming.

Heinrich, Carolyn, Gerald Marschke and Annie Zhang. 1998. "Using Administrative Data to Estimate the Cost-Effectiveness of Social Program Services." Technical report. The University of Chicago.

Jacob, Brian and Lars Lofgren. 2005. "Principals as Agents: Subjective Performance Measurement in Education." NBER Working Paper No. 11463.

Kornfeld, Robert and Howard Bloom. 1999. "Measuring Program Impacts on Earnings and Employment: Do Unemployment Insurance Wage Reports From Employers Agree With Surveys of Individuals?" *Journal of Labor Economics* 17(1): 168-197.

Kemple, James, Fred Doolittle, and John Wallace. 1993. *The National JTPA Study: Site Characteristics and Participation Patterns*. New York, NY: Manpower Demonstration Research Corporation.

Liu, Wei. 1997. "Stepwise Tests When the Test Statistics Are Independent." *Australian Journal of Statistics* 39(2): 169-177.

Manski, Charles. 1990. "The Use of Intentions Data to Predict Behavior: A Best-Case Analysis." *Journal of the American Statistical Association* 85(412): 934-940.

Manski, Charles. 2004. "Measuring Expectations." *Econometrica* 72(5): 1329-1376.

Nisbett, Richard E. and Nancy Bellows. 1977. Verbal Reports about Causal Influences on Social Judgments: Private Access vs. Public Theories. *Journal of Personality and Social Psychology* 35:613-624.

Nisbett, Richard E. and Lee Ross. 1980. *Human Inference: Strategies and Shortcomings of Social Judgment*. Englewood Cliffs, N.J: Prentice-Hall.

- Nisbett, Richard E. and Timothy D. Wilson. 1977. Telling More Than We Can Know: Verbal Reports on Mental Processes. *Psychological Review* 84(3): 231-259.
- Orr, Larry, Howard Bloom, Stephen Bell, Winston Lin, George Cave and Fred Doolittle. 1994. *The National JTPA Study: Impacts, Benefits and Costs of Title II-A*. Bethesda, MD: Abt Associates.
- Orr, Larry, Howard Bloom, Stephen Bell, Fred Doolittle, Winston Lin, and George Cave. 1996. *Does Training Work for the Disadvantaged? Evidence from the National JTPA Study*. Washington, DC: Urban Institute Press.
- Pearson, R. W., Michael Ross, and Robyn Dawes. 1992. Personal Recall and the Limits of Retrospective Questions in Surveys. In J. M. Tanur (Ed.), *Questions about questions*. New York: Russell Sage.
- Philipson, Tomas and Larry Hedges. 1998. "Subject Evaluation in Social Experiments." *Econometrica* 66(2): 381-408.
- Polanyi, Michael. 1964. *Personal Knowledge: Toward a Post-Critical Philosophy*. New York: Harper.
- Ross, Michael. 1989. Relation of Implicit Theories to the Construction of Personal Histories. *Psychological Review* 96(2):341-357.
- Roy, A. D. 1951. "Some Thoughts on the Distribution of Earnings." *Oxford Economic Papers* 3: 135-146.
- Simon, Herbert A. 1956. A Comparison of Game Theory and Learning Theory. *Psychometrika* 21:267-272.
- Smith, Jeffrey. 1997a. "Measuring Earnings Dynamics Among the Poor: Evidence from Two Samples of JTPA Eligibles." Unpublished manuscript, University of Western Ontario.
- Smith, Jeffrey. 1997b. "Measuring Earnings Levels Among the Poor: Evidence from Two Samples of JTPA Eligibles." Unpublished manuscript, University of Western Ontario.
- Smith, Jeffrey and Alexander Whalley. 2005. "How Well Do We Measure Public Job Training?" Unpublished manuscript, University of Maryland.
- Wood, Michelle. 2005. "National JTPA Study – SDA Unit Costs." Abt Associates Memo to Jerry Marsky [sic] and Larry Orr.

Appendix 1: Data Appendix

1. Sample Selection Criteria for the Samples Used

Our data set combines self-reported information from the Background Information Form, completed at or near the time of random assignment and the First Follow-Up Survey, collected around 18 months after random assignment with administrative data on quarterly earnings from matched UI wage records.

The full experimental sample contains 6639 observations in the control group and 13972 observations in the treatment group. If we restrict our sample to only those with valid self-reported earnings for the 18 months after random assignment we lose 2080 observations from the control group and 4329 observations from the treatment group. If we instead restrict the sample to only those with valid UI earnings over the six quarters after random assignment we lose 122 observations from the control group and 232 observations from the treatment group. We only require sample members to have valid values for earnings for the analysis in question; that is, we use all available observations for a given dependent variable. The analyses presented in Tables 5 to 13 require only the data from the experimental treatment group.

Our self-reported earnings data consists of the self-reported data used in Bloom et al. (1993), the official 18-month impact report. The data we use include the recoded values for outliers (which were examined individually and by hand by staff of Abt Associates) but do not include the imputed values based on the matched UI earnings records that they employed in some of their analyses. This earnings variable is not available on the public use CD but is available from the authors by request.

The matched administrative data from UI records consists of earnings in each calendar quarter. As a result, for some sample members, the six calendar quarters after the calendar quarter of random assignment (the period used in some of our dependent variables from the UI data) will cover a slightly different set of months than the 18 months after the month after random assignment (the period covered in some of our dependent variables from the self-reported data).

We do not drop observations with missing values of covariates from the sample for any of our analyses; instead we include dummy variables for those with missing values of the covariates used in each analysis. If we had instead listwise deleted observations from the sample having any missing value for the covariates we would lose 18327 observations out of the 20601 observations in the full experimental sample.

2. Variable Definitions

Predicted impact: This consists of the experimentally estimated predicted impact of the program for an individual based on either the individual's measured characteristics or the individual's quantile in the untreated outcome distribution.

Percent positive self-evaluation: This is the mean of a binary indicator for a positive participant self-evaluation. It is defined only for individuals in the treatment group.

Earnings one: This is total earnings over the 18 months after random assignment based on the self-reported earnings data.

Employment one: This is a binary variable indicating any employment over the 18 months after random assignment using self-reported earnings data. The variable equals one if self-reported earnings over the 18 months after random assignment are positive and zero otherwise.

Earnings two: This is total earnings in the 18th month after random assignment based on the self-reported earnings data.

Employment two: This is a binary variable indicating employment in month 18 after random assignment based on the self-reported earnings data. The variable equals one if self-reported earnings in the 18th month after random assignment are positive and zero otherwise.

Earnings three: This is total earnings in the six calendar quarters after the calendar quarter of random assignment based on the matched UI administrative earnings data.

Employment three: This is a binary variable indicating any employment over the six calendar quarters after the calendar quarter of random assignment based on the matched UI administrative earnings data. This variable equals one if UI earnings over the six calendar quarters after the calendar quarter of random assignment are positive and zero otherwise.

Earnings four: This is total earnings in month 18 after random assignment based on the matched UI administrative earnings data.

Employment four: This is a binary variable indicating any employment in the sixth calendar quarter after the calendar quarter of random assignment based on the matched UI administrative earnings data. This variable equals one if UI earnings in the sixth calendar quarter after random assignment are positive and zero otherwise.

EXHIBIT 1: JTPA Self-Evaluation Survey Questions

(D7)

According to (LOCAL JTPA PROGRAM NAME) records, you applied to enter (LOCAL JTPA PROGRAM NAME) in (MONTH/YEAR OF RANDOM ASSIGNMENT). Did you participate in the program after you applied?

YES (SKIP TO D9)

NO (GO TO D8)

(D9)

Do you think that the training or other assistance that you got from the program helped you get a job or perform better on the job?

YES

NO

Source: JTPA First Follow-Up Study Survey Instrument

TABLE 1: Bivariate Results for the relationship between Experimental Impacts and Positive Self-Evaluation By Demographic Group

	Percentage Positive Self- Evaluation	Earnings One	Employ One	Earnings Two	Employ Two	Earnings Three	Employ Three	Earnings Four	Employ Four
Adult Males	0.63 (0.01)	538.20 (379.22)	0.03 (0.01)	23.58 (28.55)	0.02 (0.02)	-36.42 (293.50)	0.00 (0.01)	-24.10 (65.69)	-0.03 (0.02)
Adult Females	0.65 (0.01)	750.87 (236.17)	0.03 (0.01)	56.79 (18.34)	0.04 (0.14)	594.08 (195.48)	0.04 (0.01)	131.24 (44.18)	0.03 (0.01)
Male Youth	0.67 (0.02)	-777.33 (463.33)	0.01 (0.01)	-82.93 (37.00)	-0.03 (0.02)	-381.03 (328.19)	-0.02 (0.02)	-128.07 (73.54)	-0.03 (0.02)
Female Youth	0.72 (0.01)	-44.89 (295.12)	0.04 (0.02)	8.38 (29.87)	-0.00 (0.02)	-233.74 (227.97)	0.01 (0.02)	-13.84 (50.74)	0.00 (0.02)
Correlation with Positive Self-Evaluation	--	-0.4620 [0.538]	0.5510 [0.449]	-0.2239 [0.776]	-0.4553 [0.545]	-0.4381 [0.562]	-0.1486 [0.851]	-0.1858 [0.814]	0.1426 [0.857]

Notes: Source: Authors' calculations using the NJS data. Values in the table are means for Positive Self-Evaluation, and experimental impacts for the eight outcomes. The values in parentheses are standard errors and the values in square brackets are p-values. Percentage Positive Self-Assessment is calculated as the mean of the binary indicator positive self-assessment variable for those who self-report participating and are in the treatment group. Earnings one and employment one are earnings and any employment over the 18-months after random assignment using self-reported earnings data. Earnings two and employment two are earnings and employment in month 18 after random assignment using self-reported earnings data. Earnings three and employment three are earnings and any employment over the 18-months after random assignment using UI-reported earnings data. Earnings four and employment four are earnings and employment in month 18 after random assignment using UI-reported earnings data. Those with missing outcomes are dropped from the estimate for that outcome only.

TABLE 2A: Bivariate results for the Correlation between Experimental Impacts and Self-Evaluation for Eight Outcomes, Adult Males

	Earnings One	Employment One	Earnings Two	Employment Two	Earnings Three	Employment Three	Earnings Four	Employment Four
Race	0.1742 [0.826]	-0.3556 [0.644]	-0.2184 [0.782]	-0.0346 [0.965]	0.7508 [0.249]	0.6956 [0.304]	0.7023 [0.298]	0.6447 [0.355]
Age Category	0.9974 [0.046]	-0.9989 [0.030]	0.9870 [0.103]	-0.9137 [0.266]	0.8198 [0.388]	-0.0455 [0.971]	0.7573 [0.453]	0.0169 [0.989]
Education Category	0.4984 [0.393]	-0.7003 [0.188]	0.7613 [0.135]	-0.3996 [0.505]	0.9077 [0.033]	0.1717 [0.783]	0.9798 [0.003]	-0.8810 [0.048]
Marital Status	-0.5476 [0.631]	-0.9999 [0.007]	-0.9946 [0.066]	-0.9404 [0.229]	0.9939 [0.070]	-0.2063 [0.868]	0.7701 [0.440]	0.6103 [0.582]
Employ Category	-0.1606 [0.897]	-0.8638 [0.336]	-0.1177 [0.925]	-0.5880 [0.600]	-0.6909 [0.514]	0.3855 [0.748]	-0.3794 [0.752]	-0.9451 [0.212]
Site	0.3380 [0.200]	0.1495 [0.581]	0.1682 [0.533]	0.4170 [0.108]	-0.2132 [0.428]	-0.1497 [0.580]	-0.1015 [0.709]	0.0265 [0.923]

Positive Correlations

Overall: 24 of 48 (50 %); significant at 0.10: 4 of 48 (8 %); significant at 0.05: 3 of 48 (6 %)

Negative Correlations

Overall: 24 of 48 (50 %); significant at 0.10: 4 of 48 (8 %); significant at 0.05: 3 of 48 (6 %)

Notes: Source: Authors' calculations using the NJS data. Values in the table are the correlation between the mean of Positive Self-Evaluation, and the experimental impacts by subgroup. The values in square brackets are p-values. Percentage Positive Self-Evaluation is calculated as the mean of the binary

indicator positive self-evaluation variable for those who self-report participating and are in the treatment group. Earnings one and employment one are earnings and any employment over the 18-months after random assignment using self-reported earnings data. Earnings two and employment two are earnings and employment in month 18 after random assignment using self-reported earnings data. Earnings three and employment three are earnings and any employment over the 18-months after random assignment using UI-reported earnings data. Earnings four and employment four are earnings and employment in month 18 after random assignment using UI-reported earnings data. Those with missing outcomes are dropped from the estimate for that outcome only. The categories are defined as the following. Race: White, Black, Hispanic and Other. Age: less than 19 years, 19-21 years, 22-25 years, 26-34 years and 35+ years. Education: under 10 years, 10-11 years, 12 years, 13-15 years and 16+ years. Marital Status: single, married, and divorced/widowed/separated. Employment Status: out of labor force, unemployed, and employed. Site: sixteen site categories.

TABLE 2B: Bivariate results for the Correlation between Experimental Impacts and Self-Evaluation for Eight Outcomes, Adult Females

	Earnings One	Employment One	Earnings Two	Employment Two	Earnings Three	Employment Three	Earnings Four	Employment Four
Race	0.3681 [0.632]	0.4127 [0.587]	-0.2490 [0.751]	-0.2490 [0.751]	0.3618 [0.638]	0.2877 [0.712]	0.4188 [0.581]	0.6203 [0.380]
Age Category	-0.8947 [0.295]	0.1282 [0.918]	-0.9136 [0.267]	-0.9681 [0.161]	0.5726 [0.612]	-0.4926 [0.672]	-0.9490 [0.204]	-0.9597 [0.181]
Education Category	-0.6549 [0.230]	-0.6872 [0.200]	-0.5260 [0.363]	0.1167 [0.852]	-0.1673 [0.788]	0.6492 [0.236]	-0.6820 [0.205]	-0.2442 [0.692]
Marital Status	0.0802 [0.949]	0.8084 [0.401]	0.4329 [0.715]	0.5232 [0.650]	0.4284 [0.718]	0.3944 [0.742]	-0.7971 [0.413]	0.5074 [0.661]
Employ Category	0.9296 [0.240]	0.7264 [0.482]	0.6294 [0.567]	0.2549 [0.836]	0.9978 [0.042]	0.8663 [0.333]	0.2949 [0.809]	0.8199 [0.388]
Site	-0.0745 [0.784]	-0.2296 [0.392]	-0.0628 [0.817]	0.1812 [0.502]	-0.0753 [0.782]	0.1143 [0.674]	-0.0250 [0.927]	0.1923 [0.476]
Positive Correlations								
Overall: 28 of 48 (58 %); significant at 0.10: 0 of 48 (0 %); significant at 0.05: 0 of 48 (0 %)								
Negative Correlations								
Overall: 20 of 48 (42 %); significant at 0.10: 0 of 48 (0 %); significant at 0.05: 0 of 48 (0 %)								

Notes: Source: Authors' calculations using the NJS data. Values in the table are the correlation between the mean of Positive Self-Evaluation, and the experimental impacts by subgroup. The values in square brackets are p-values. Percentage Positive Self-Evaluation is calculated as the mean of the binary indicator positive self-evaluation variable for those who self-report participating and are in the treatment group. Earnings one and employment one are earnings and any employment over the 18-months after random assignment using self-reported earnings data. Earnings two and employment two are earnings and employment in month 18 after random assignment using self-reported earnings data. Earnings three and employment three are earnings and any employment over the 18-months after random assignment using UI-reported earnings data. Earnings four and employment four are earnings and employment in month 18 after random assignment using UI-reported earnings data. Those with missing outcomes are dropped from the estimate for that outcome only. The categories are defined as the following. Race: White, Black, Hispanic and Other. Age: less than 19 years, 19-21 years, 22-25 years, 26-34 years and 35+ years. Education: under 10 years, 10-11 years, 12 years, 13-15 years and 16+ years. Marital Status: single, married, and divorced/widowed/separated. Employment Status: out of labor force, unemployed, and employed. Site: sixteen site categories.

TABLE 2C: Bivariate results for the Correlation between Experimental Impacts and Self-Evaluation for Eight Outcomes, Male Youths

	Earnings One	Employment One	Earnings Two	Employment Two	Earnings Three	Employment Three	Earnings Four	Employment Four
Race	0.471 [0.528]	0.2314 [0.769]	-0.0844 [0.916]	-0.3902 [0.610]	0.2621 [0.738]	0.1283 [0.872]	0.2866 [0.713]	0.2097 [0.790]
Education Category	-0.4412 [0.559]	0.2749 [0.725]	-0.1031 [0.897]	0.8667 [0.133]	0.4519 [0.548]	-0.8642 [0.136]	0.6511 [0.349]	0.9301 [0.070]
Marital Status	-0.8985 [0.289]	-0.8039 [0.406]	-0.9816 [0.122]	-0.2344 [0.849]	0.9262 [0.246]	-0.9926 [0.077]	0.9954 [0.061]	0.9671 [0.164]
Employ Category	0.8770 [0.319]	0.9606 [0.179]	0.7889 [0.421]	0.8352 [0.371]	0.9985 [0.035]	0.8704 [0.328]	0.9865 [0.105]	0.7949 [0.415]
Site	0.3258 [0.236]	-0.1278 [0.637]	0.1120 [0.680]	0.6023 [0.014]	0.2063 [0.443]	0.1062 [0.696]	-0.1386 [0.609]	-0.4438 [0.085]

Positive Correlations

Overall: 27 of 40 (68 %); significant at 0.10: 4 of 40 (10 %); significant at 0.05: 2 of 40 (5 %)

Negative Correlations

Overall: 13 of 40 (32 %); significant at 0.10: 2 of 40 (5 %); significant at 0.05: 0 of 40 (0 %)

Notes: Source: Authors' calculations using the NJS data. Values in the table are the correlation between the mean of Positive Self-Evaluation, and the experimental impacts by subgroup. The values in square brackets are p-values. Percentage Positive Self-Evaluation is calculated as the mean of the binary indicator positive self-evaluation variable for those who self-report participating and are in the treatment group. Earnings one and employment one are earnings and any employment over the 18-months after random assignment using self-reported earnings data. Earnings two and employment two are earnings and employment in month 18 after random assignment using self-reported earnings data. Earnings three and employment three are earnings and any employment

over the 18-months after random assignment using UI-reported earnings data. Earnings four and employment four are earnings and employment in month 18 after random assignment using UI-reported earnings data. Those with missing outcomes are dropped from the estimate for that outcome only. The categories are defined as the following. Race: White, Black, Hispanic and Other. Education: under 10 years, 10-11 years, 12 years, 13-15 years and 16+ years. Marital Status: single, married, and divorced/widowed/separated. Employment Status: out of labor force, unemployed, and employed. Site: sixteen site categories.

TABLE 2D: Bivariate results for the Correlation between Experimental Impacts and Self-Evaluation for Eight Outcomes, Female Youths

	Earnings One	Employment One	Earnings Two	Employment Two	Earnings Three	Employment Three	Earnings Four	Employment Four
Race	-0.7560 [0.244]	-0.2111 [0.789]	-0.7309 [0.269]	-0.5835 [0.417]	-0.5849 [0.415]	0.4595 [0.541]	-0.6132 [0.387]	0.7253 [0.274]
Education Category	-0.9988 [0.000]	-0.9995 [0.000]	-0.9914 [0.001]	-0.9896 [0.001]	-0.9459 [0.015]	-0.9950 [0.000]	-0.9616 [0.009]	-0.9702 [0.006]
Marital Status	-0.4687 [0.690]	-0.1860 [0.881]	0.6774 [0.526]	0.3977 [0.740]	0.8910 [0.300]	-0.1495 [0.905]	0.9998 [0.012]	0.3252 [0.789]
Employ Category	0.8395 [0.366]	-0.5065 [0.662]	-0.8400 [0.365]	-0.9703 [0.156]	0.2660 [0.829]	0.9998 [0.012]	-0.3951 [0.741]	0.1339 [0.915]
Site	0.2620 [0.346]	0.3277 [0.233]	0.2164 [0.439]	0.3086 [0.263]	0.2643 [0.341]	0.1690 [0.547]	0.2554 [0.358]	0.2112 [0.450]

Positive Correlations

Overall: 19 of 40 (47 %); significant at 0.10: 2 of 40 (5 %); significant at 0.05: 2 of 40 (5 %)

Negative Correlations

Overall: 21 of 40 (53 %); significant at 0.10: 8 of 40 (20 %); significant at 0.05: 8 of 40 (20 %)

Notes: Source: Authors' calculations using the NJS data. Values in the table are the correlation between the mean of Positive Self-Evaluation, and the experimental impacts by subgroup. The values in square brackets are p-values. Percentage Positive Self-Evaluation is calculated as the mean of the binary indicator positive self-evaluation variable for those who self-report participating and are in the treatment group. Earnings one and employment one are earnings and any employment over the 18-months after random assignment using self-reported earnings data. Earnings two and employment two are earnings and employment in month 18 after random assignment using self-reported earnings data. Earnings three and employment three are earnings and any employment over the 18-months after random assignment using UI-reported earnings data. Earnings four and employment four are earnings and employment in month 18

after random assignment using UI-reported earnings data. Those with missing outcomes are dropped from the estimate for that outcome only. The categories are defined as the following. Race: White, Black, Hispanic and Other. Education: under 10 years, 10-11 years, 12 years, 13-15 years and 16+ years. Marital Status: single, married, and divorced/widowed/separated. Employment Status: out of labor force, unemployed, and employed. Site: sixteen site categories.

TABLE 3: Regression results for the relationship between Predicted Impacts and Positive Self-Evaluation for Eight Outcomes, By Demographic Group

	Adult Males		Adult Females		Male Youths		Female Youths	
	(1)	(2)	(1)	(2)	(1)	(2)	(1)	(2)
Earnings over 18 Months	-121.04 (134.85)	48.66 (83.41)	45.71 (85.10)	-16.86 (57.75)	-21.95 (244.01)	273.06 (214.97)	-208.51 (89.36)	24.82 (97.87)
Any Employment During 18 Months	-0.009 (0.003)	-0.002 (0.04)	0.001 (0.023)	0.005 (0.004)	-0.003 (0.003)	0.010 (0.006)	-0.005 (0.002)	-0.003 (0.008)
Earnings in Month 18	-21.20 (20.65)	-6.05 (7.73)	2.37 (4.95)	0.97 (3.80)	35.53 (31.26)	-6.67 (16.58)	-2.32 (9.49)	1.54 (10.37)
Employment in Month 18	-0.005 (0.003)	-0.002 (0.004)	-0.004 (0.004)	0.002 (0.003)	-0.003 (0.014)	-0.001 (0.011)	-0.002 (0.001)	-0.003 (0.011)
Earnings (UI) over 6 Quarters	-63.67 (94.48)	-61.17 (85.90)	-85.43 (50.92)	14.51 (35.79)	-71.24 (133.03)	271.47 (134.34)	-103.53 (68.76)	78.86 (68.61)
Any Employment (UI) During 6 Quarters	-0.003 (0.002)	-0.002 (0.003)	0.001 (0.003)	0.002 (0.003)	-0.007 (0.007)	0.000 (0.006)	0.003 (0.013)	0.007 (0.006)
Earnings (UI) in Quarter 6	-22.56 (23.25)	-14.52 (17.18)	0.73 (10.96)	-10.17 (7.65)	-0.16 (18.90)	80.59 (31.78)	2.60 (13.53)	-13.14 (12.60)
Employment (UI) in Quarter 6	-0.004 (0.003)	0.000 (0.004)	0.000 (0.003)	-0.004 (0.003)	0.008 (0.005)	0.019 (0.010)	-0.001 (0.002)	0.007 (0.009)
Positive (overall / 0.10 / 0.05)	0/0/0	1/0/0	5/0/0	5/0/0	2/0/0	5/3/2	2/0/0	5/0/0
Negative (overall / 0.10 / 0.05)	8/1/1	6/0/0	2/1/0	3/0/0	6/0/0	2/0/0	5/2/2	3/0/0

Notes: Source: Authors' calculations using the NJS data. Each cell in the table is a coefficient estimate from the regression of the estimated impacts for an individual (based on their X) as the dependent variable and self-evaluation as the independent variable. The population used is the treatment sample. The values

in parentheses are the heteroskedastic-consistent standard errors. The values in the bottom two rows are the counts of the number of cells in the column above which are positive or negative, and counts of those that are significantly different from zero at the 10% and 5% levels. Specification (1) selects the set of X's used to predict the impacts for each individual by a stepwise procedure. Specification (2) uses the specification of X's used in Heckman, Heinrich and Smith (2003) (HHS) to predict the impacts for each individual. The HHS set of X's contains: are race, age, education, marital status, employment status, AFDC receipt, receipt of food stamps and site.

TABLE 4: Relationship between Quantile Treatment Effects for 18-Month Earnings and the Percent with Positive Self-Evaluation, By Demographic Group

	Adult Males		Adult Females		Male Youths		Female Youths	
	Quantile Treatment Effects	Percentage Positive Self-Evaluation	Quantile Treatment Effects	Percentage Positive Self-Evaluation	Quantile Treatment Effects	Percentage Positive Self-Evaluation	Quantile Treatment Effects	Percentage Positive Self-Evaluation
5 th	0 (0.90)	0.51 (0.03)	0 (0.38)	0.57 (0.02)	0 (1.15)	0.56 (0.07)	0 (1.07)	0.68 (0.04)
25 th	1233 (452)	0.66 (0.05)	501 (193)	0.63 (0.04)	-516 (515)	0.62 (0.08)	402 (193)	0.72 (0.06)
50 th	825 (608)	0.56 (0.05)	747 (416)	0.63 (0.04)	-1161 (681)	0.83 (0.06)	-39 (371)	0.71 (0.07)
75 th	8 (590)	0.72 (0.05)	938 (383)	0.78 (0.04)	-1261 (701)	0.68 (0.08)	-479 (566)	0.83 (0.06)
95 th	1589 (1323)	0.65 (0.05)	1910 (740)	0.70 (0.04)	-887 (1959)	0.81 (0.07)	-53 (1012)	0.64 (0.07)
Correlation with Percentage Positive Self-Evaluation	0.0760 [0.750]	--	0.7652 [0.000]	--	-0.4527 [0.045]	--	-0.4209 [0.065]	--
Coefficient on Percentage Positive Self-Evaluation	511 (1686)	--	5489 (1204)	--	-2232 (909)	--	-1576 (931)	--

Notes: Source: Authors' calculations using the NJS data. The values in the left column of the upper panel for each demographic group are quantile treatment effects estimates with standard errors in parentheses for five quantiles. The values in the right column of the upper panel for each demographic group are the means of the binary positive self-evaluation indicator variable for each quantile of the outcome distribution for those in the treatment group. The first row of the lower panel contains the correlation between the treatment effect estimates and the percentage positive self-evaluation by quantile (where one observation is one of the 20 quantiles) and the p-value for the correlation is in square brackets. The second row of the lower panel contains the coefficient of the regression with percentage positive self-evaluation as the independent variable and the treatment effect estimate as the dependant variable (where one observation is one of the 20 quantiles). The hetero-skedastic consistent standard errors for these estimates appear in parentheses.

TABLE 5: Logit Estimates of the Determinants of Positive Self-Evaluation,
By Demographic Group

	Adult Males	Adult Females	Male Youths	Female Youths
Age: 19-21 Years	--	--	0.026 (0.039) [0.505]	-0.036 (0.031) [0.255]
Age: 26-34 Years	-0.020 (0.029) [0.500]	-0.023 (0.026) [0.383]	--	--
Age: 35+ years	-0.071 (0.034) [0.036]	-0.122 (0.032) [0.000]	--	--
Marital Status: Married	-0.035 (0.032) [0.283]	0.039 (0.030) [0.188]	-0.009 (0.065) [0.889]	-0.013 (0.048) [0.794]
Marital Status: Divorced/Widowed/ Separated	-0.032 (0.032) [0.327]	-0.009 (0.025) [0.706]	0.204 (0.119) [0.086]	-0.145 (0.053) [0.006]
Education: 10-11 Years	0.049 (0.035) [0.155]	-0.022 (0.032) [0.487]	0.040 (0.044) [0.364]	0.083 (0.035) [0.018]
Education: 12 Years	0.048 (0.032) [0.137]	-0.051 (0.029) [0.081]	-0.003 (0.049) [0.944]	0.082 (0.038) [0.030]
Education: 13-15 Years	-0.005 (0.041) [0.900]	-0.035 (0.037) [0.337]	-0.075 (0.090) [0.406]	0.102 (0.059) [0.066]
Education: 16+ Years	0.022 (0.057) [0.705]	0.065 (0.058) [0.266]	--	--
Race: Black	0.020 (0.034) [0.543]	-0.016 (0.030) [0.598]	0.122 (0.045) [0.007]	-0.035 (0.043) [0.425]
Race: Hispanic	0.048 (0.047) [0.313]	-0.001 (0.042) [0.979]	0.159 (0.052) [0.002]	0.103 (0.048) [0.031]
Race: Other	0.035 (0.075) [0.646]	-0.054 (0.072) [0.458]	0.039 (0.130) [0.764]	0.017 (0.106) [0.872]
English Language	0.051	0.106	0.135	-0.058

	(0.070)	(0.059)	(0.107)	(0.158)
	[0.471]	[0.073]	[0.210]	[0.715]
AFDC Receipt	0.039	-0.015	0.085	-0.007
	(0.042)	(0.024)	(0.052)	(0.037)
	[0.352]	[0.521]	[0.097]	[0.850]
Work for Pay	0.037	0.039	-0.059	0.015
	(0.045)	(0.028)	(0.053)	(0.036)
	[0.415]	[0.164]	[0.264]	[0.697]
Child less than Six	-0.019	0.003	-0.146	0.003
	(0.033)	(0.023)	(0.073)	(0.035)
	[0.571]	[0.900]	[0.044]	[0.940]
Self-Report Training: CT-OS	0.135	0.122	0.084	0.010
	(0.028)	(0.023)	(0.040)	(0.034)
	[0.000]	[0.000]	[0.038]	[0.764]
Self-Report Training: OJT/WE	0.153	0.115	0.109	0.081
	(0.040)	(0.034)	(0.062)	(0.055)
	[0.000]	[0.001]	[0.076]	[0.138]
Self-Report Training: JSA	0.085	0.022	0.000	0.029
	(0.042)	(0.038)	(0.088)	(0.066)
	[0.044]	[0.570]	[0.998]	[0.661]
Self-Report Training: ABE	0.046	0.019	0.075	0.095
	(0.049)	(0.038)	(0.049)	(0.040)
	[0.347]	[0.614]	[0.129]	[0.018]
Self-Report Training: Other	0.153	0.048	0.079	0.090
	(0.053)	(0.048)	(0.082)	(0.061)
	[0.004]	[0.313]	[0.334]	[0.140]
Administrative- Report Training: CT-OS	0.039	-0.005	-0.023	0.154
	(0.058)	(0.046)	(0.104)	(0.052)
	[0.493]	[0.916]	[0.824]	[0.003]
Administrative- Report Training: OJT/WE	0.055	0.088	-0.084	0.067
	(0.058)	(0.048)	(0.113)	(0.067)
	[0.340]	[0.066]	[0.459]	[0.313]
Administrative- Report Training: JSA	0.036	-0.045	0.077	0.062
	(0.058)	(0.052)	(0.099)	(0.063)
	[0.532]	[0.387]	[0.433]	[0.326]
Administrative- Report Training: ABE	0.108	-0.122	-0.184	-0.064
	(0.066)	(0.068)	(0.108)	(0.078)
	[0.102]	[0.073]	[0.089]	[0.409]
Administrative- Report Training: Other	0.050	-0.121	-0.148	0.123
	(0.063)	(0.057)	(0.110)	(0.060)
	[0.425]	[0.033]	[0.178]	[0.040]

Notes: Source: Authors' calculations using the NJS data. Columns two through five of the table report the results from a logit model where the binary positive self-evaluation variable is the dependant variable and the categorical variables listed in column one are the independent variables. The values in the table are mean numerical derivatives, with the standard errors in parentheses and p-values in square brackets. The population for these regressions is the treatment sample. Indicator variables for missing values for the independent variables are also included in the regression. The omitted age category for adults is age 22-25

years and is age less than 19 for youths. The omitted marital status is single, the omitted education category is less than 10 years, the omitted racial group is white, and the omitted training type for both self-report and administrative report is no training for all demographic groups.

TABLE 6: Test Statistics from Logit Models of the Determinants of Positive Self-Evaluation, By Demographic Group

	Adult Males	Adult Females	Male Youths	Female Youths
Site	65.30 [0.000] R ^{*2} =3.77	73.06 [0.000] R ^{*2} =3.18	29.66 [0.009] R ^{*2} =3.90	61.12 [0.000] R ^{*2} =6.25
Age Category	5.09 [0.078] R ^{*2} =0.27	21.71 [0.000] R ^{*2} =0.83	0.45 [0.504] R ^{*2} =0.05	1.26 [0.262] R ^{*2} =0.12
Marital Status	1.44 [0.487] R ^{*2} =0.10	3.68 [0.159] R ^{*2} =0.52	1.44 [0.488] R ^{*2} =0.52	8.55 [0.014] R ^{*2} =0.29
Education Category	4.50 [0.343] R ^{*2} =0.40	6.47 [0.167] R ^{*2} =0.27	2.46 [0.482] R ^{*2} =0.28	6.37 [0.095] R ^{*2} =0.71
Race	1.20 [0.753] R ^{*2} =0.07	0.80 [0.849] R ^{*2} =0.03	10.49 [0.015] R ^{*2} =1.31	6.36 [0.095] R ^{*2} =0.59
English Language	0.91 [0.633] R ^{*2} =0.05	3.16 [0.206] R ^{*2} =0.12	1.12 [0.290] R ^{*2} =0.13	0.15 [0.703] R ^{*2} =0.01
Other Individual Characteristics	6.13 [0.294] R ^{*2} =0.33	3.16 [0.675] R ^{*2} =0.12	6.72 [0.242] R ^{*2} =0.79	1.23 [0.942] R ^{*2} =0.12
Self-Reported Training Type	30.21 [0.000] R ^{*2} =1.77	30.66 [0.000] R ^{*2} =1.43	6.23 [0.284] R ^{*2} =1.23	7.05 [0.217] R ^{*2} =1.37

Administrative Reported Training Type	30.67 [0.000] R ^{*2} =1.76	53.40 [0.000] R ^{*2} =2.31	21.55 [0.002] R ^{*2} =2.79	25.61 [0.000] R ^{*2} =2.39
---------------------------------------	---	---	---	---

Notes: Source: Authors' calculations using the NJS data. Columns two through five of the table report the results from a logit model where the binary positive self-evaluation variable is the dependent variable and the categorical variables summarized in column one are the independent variables. The values in the table are χ^2 -statistics for joint tests that all of the coefficients equal zero for a given group of variables, with the p-values in square brackets. The value for R^{*2} are the partial r-squared times 100 for the group of covariates in the row. The population for these regressions is the treatment sample. The variables in 'Other Individual Characteristics' are AFDC receipt, child less than six indicator, and worked for pay indicator. Indicator variables for missing values for the independent variables are also included in the regressions.

TABLE 7: Logit Estimates of the Relationship between Outcomes and Positive Self-Evaluation: Four Outcomes, By Demographic Group

	Adult Males	Adult Females	Male Youths	Female Youths
Earnings over 18 Months = 0	-0.171 (0.042) [0.000]	-0.184 (0.038) [0.000]	-0.151 (0.084) [0.073]	-0.065 (0.064) [0.312]
Earnings over 18 Months Bottom Quartile	-0.096 (0.040) [0.015]	-0.094 (0.035) [0.007]	-0.158 (0.062) [0.011]	-0.028 (0.053) [0.593]
Earnings over 18 Months Lower Middle Quartile	-0.038 (0.035) [0.269]	-0.133 (0.034) [0.000]	-0.001 (0.056) [0.992]	0.068 (0.049) [0.172]
Earnings over 18 Months Upper Middle Quartile	-0.083 (0.033) [0.013]	-0.033 (0.034) [0.337]	-0.052 (0.056) [0.358]	0.053 (0.053) [0.319]
Earnings over 18 Months = 0 (UI)	-0.076 (0.041) [0.062]	-0.027 (0.034) [0.423]	-0.063 (0.073) [0.387]	-0.068 (0.064) [0.290]
Earnings over 18 Months Bottom Quartile (UI)	-0.080 (0.039) [0.038]	-0.112 (0.034) [0.001]	-0.120 (0.061) [0.049]	-0.081 (0.057) [0.153]
Earnings over 18 Months Lower Middle Quartile (UI)	-0.068 (0.035) [0.049]	-0.067 (0.032) [0.035]	-0.050 (0.056) [0.368]	-0.089 (0.057) [0.113]
Earnings over 18 Months Upper Middle Quartile (UI)	-0.026 (0.032) [0.420]	-0.035 (0.030) [0.249]	-0.030 (0.059) [0.605]	0.014 (0.054) [0.795]
Any Employment over 18 Months	0.122 (0.038) [0.001]	0.105 (0.027) [0.000]	0.084 (0.071) [0.237]	0.082 (0.044) [0.060]
Any Employment over 18 Months (UI)	0.036 (0.035) [0.308]	-0.029 (0.025) [0.253]	0.000 (0.054) [1.000]	0.006 (0.039) [0.886]

Notes: Source: Authors' calculations using the NJS data. Columns two through five of this table report the results from logit regressions where the binary positive self-evaluation variable is the dependant variable and the categorical variables listed in column one of Table 5 are the independent variables, in addition an outcome variable is included in each regression. The values in the table are mean numerical derivatives,

with the standard errors in parentheses and p-values in square brackets. For earnings outcomes the continuous variables are entered as four categorical variables: zero earnings, an indicator for being in the lowest quartile of the non-zero earnings distribution, lower middle quartile of the non-zero earnings distribution, upper middle quartile of the non-zero earnings distribution. The omitted category is for those with earnings in the highest quartile of the non-zero earnings distribution. For the employment outcomes a binary variable is included indicating whether the respondent was employed or not. Each set of cells in the table is the result for a different specification where the outcome to be included as an independent variable is different. The sets of cells are defined as two groups of four and two groups of two depending on how the outcome enters the regression. The population for these regressions is the treatment sample. Indicator variables for missing values for the independent variables are also included in the regression.

TABLE 8: Test Statistics from Logit Models of the Relationship between Outcomes and Positive Self-Evaluation,
By Demographic Group

	Adult Males	Adult Females	Male Youths	Female Youths
Earnings over 18 Months	20.37 [0.001] R ^{*2} =1.06	34.33 [0.000] R ^{*2} =1.38	15.96 [0.007] R ^{*2} =1.98	12.18 [0.032] R ^{*2} =1.10
Any Employment during 18 Months	11.39 [0.003] R ^{*2} =0.61	15.60 [0.000] R ^{*2} =0.63	5.58 [0.062] R ^{*2} =0.68	4.51 [0.105] R ^{*2} =0.39
Earnings in Month 18	15.75 [0.008] R ^{*2} =0.82	30.55 [0.000] R ^{*2} =1.22	8.02 [0.155] R ^{*2} =1.00	1.32 [0.933] R ^{*2} =0.12
Employment in Month 18	6.22 [0.045] R ^{*2} =0.32	12.36 [0.002] R ^{*2} =0.51	5.67 [0.059] R ^{*2} =0.67	0.99 [0.609] R ^{*2} =0.09
Earnings over 6 Quarters (UI)	7.37 [0.195] R ^{*2} =0.51	13.92 [0.016] R ^{*2} =0.85	4.98 [0.418] R ^{*2} =1.96	10.06 [0.074] R ^{*2} =0.41
Any Employment During 6 Quarters (UI)	1.17 [0.556] R ^{*2} =0.29	1.27 [0.529] R ^{*2} =0.22	0.39 [0.825] R ^{*2} =1.00	2.41 [0.300] R ^{*2} =0.13
Earnings in Quarter 6 (UI)	12.23 [0.032] R ^{*2} =0.38	12.49 [0.029] R ^{*2} =0.54	3.53 [0.473] R ^{*2} =0.62	7.14 [0.211] R ^{*2} =0.92
Employment in Quarter 6 (UI)	1.08 [0.583]	0.51 [0.776]	0.28 [0.595]	0.68 [0.710]

Earnings in the Month of the Survey	R ^{*2} =0.06 9.85 [0.080]	R ^{*2} =0.05 21.78 [0.001]	R ^{*2} =0.05 15.57 [0.008]	R ^{*2} =0.24 4.65 [0.460]
Employment in the Month of the Survey	R ^{*2} =0.64 5.73 [0.057]	R ^{*2} =0.49 5.59 [0.061]	R ^{*2} =0.44 8.33 [0.016]	R ^{*2} =0.62 1.50 [0.473]
Earnings in the Quarter of the Survey (UI)	R ^{*2} =0.05 10.46 [0.063]	R ^{*2} =0.02 9.00 [0.109]	R ^{*2} =0.19 11.74 [0.039]	R ^{*2} =0.05 3.53 [0.618]
Employment in the Quarter of the Survey (UI)	R ^{*2} =0.55 3.45 [0.178]	R ^{*2} =0.35 0.24 [0.889]	R ^{*2} =1.43 10.59 [0.005]	R ^{*2} =0.30 0.82 [0.664]
	R ^{*2} =0.19	R ^{*2} =0.01	R ^{*2} =1.29	R ^{*2} =0.07

Notes: Source: Authors' calculations using the NJS data. Columns two through five of this table report the results from logit models where the binary positive self-evaluation variable is the dependant variable and the categorical variables listed in column one of Table 5 are the independent variables, in addition an outcome variable is included in each regression. Each cell in the table is the result for a different specification where the outcome to be included as an independent variable is different. The values in the table are χ^2 -Statistics for joint tests that all of the coefficients are zero for a given outcome, with the p-values in square brackets. The value for R^{*2} are the partial r-squared times 100 for the group of covariates in the row. For earnings outcomes the continuous variables are entered as four categorical variables: zero earnings, an indicator for being in the lowest quartile of the non-zero earnings distribution, lower middle quartile of the non-zero earnings distribution, upper middle quartile of the non-zero earnings distribution. The omitted category is for those with earnings in the highest quartile of the non-zero earnings distribution. For the employment outcomes a binary variable is included indicating whether the respondent was employed or not. The population for these regressions is the treatment sample. Indicator variables for missing values for the independent variables are also included in the regression.

TABLE 9: Logit Estimates of the Relationship between Before-After Self-Reported Earnings Changes and Positive Self-Evaluation, By Demographic Group

	Adult Males	Adult Females	Male Youths	Female Youths
Before-After Self Reported Earnings 2 nd Quintile	-0.025 (0.041) [0.540]	-0.008 (0.035) [0.825]	0.030 (0.057) [0.600]	-0.077 (0.059) [0.190]
Before-After Self Reported Earnings 3 rd Quintile	0.033 (0.037) [0.375]	0.055 (0.031) [0.077]	0.107 (0.050) [0.031]	0.015 (0.047) [0.752]
Before-After Self Reported Earnings 4 th Quintile	0.044 (0.038) [0.250]	0.067 (0.031) [0.029]	0.115 (0.049) [0.019]	0.053 (0.045) [0.239]
Before-After Self Reported Earnings 5 th Quintile	0.020 (0.034) [0.547]	0.109 (0.027) [0.000]	0.132 (0.046) [0.004]	0.035 (0.040) [0.379]

Notes: Source: Authors' calculations using the NJS data. "Before-After Self-Reported Earnings" consists of monthly self-reported earnings over the 18 months after random assignment minus monthly self-reported earnings in the 12 months prior to random assignment. The estimates come from logit models with an indicator for a positive self-evaluation as the dependent variable and the before-after earnings change variable and the categorical variables listed in column one of Table 5 as independent variables. The values in the table are mean numerical derivatives, with standard errors in parentheses and p-values in square brackets. The before-after earnings changes enter in the form of indicator variables for being in the 2nd, 3rd, 4th, and 5th quintiles of the before-after earnings change distribution. The omitted category is the 1st quintile of the distribution. The population for these regressions is the treatment group. Indicator variables for missing values for the independent variables are also included in the regression.

TABLE 10: Logit Estimates of the Relationship between Before-After UI Earnings Changes and Positive Self-Evaluation, By Demographic Group

	Adult Males	Adult Females	Male Youths	Female Youths
Before-After UI Reported Earnings 2 nd Quintile	0.008 (0.036) [0.825]	0.006 (0.031) [0.853]	0.164 (0.044) [0.000]	0.029 (0.044) [0.494]
Before-After UI Reported Earnings 3 rd Quintile	0.053 (0.035) [0.124]	0.032 (0.030) [0.275]	0.127 (0.044) [0.004]	-0.021 (0.043) [0.623]
Before-After UI Reported Earnings 4 th Quintile	0.071 (0.034) [0.037]	0.083 (0.028) [0.003]	0.135 (0.044) [0.002]	0.076 (0.039) [0.050]
Before-After UI Reported Earnings 5 th Quintile	0.093 (0.033) [0.005]	0.107 (0.027) [0.000]	0.143 (0.044) [0.001]	0.135 (0.035) [0.000]
Before-After UI (2) Reported Earnings 2 nd Quintile	0.015 (0.051) [0.765]	-0.075 (0.051) [0.143]	0.015 (0.059) [0.792]	0.156 (0.112) [0.163]
Before-After UI (2) Reported Earnings 3 rd Quintile	-0.032 (0.038) [0.407]	-0.060 (0.035) [0.083]	0.021 (0.054) [0.697]	0.024 (0.044) [0.579]
Before-After UI (2) Reported Earnings 4 th Quintile	0.022 (0.030) [0.473]	-0.001 (0.026) [0.960]	0.045 (0.045) [0.326]	0.012 (0.035) [0.723]
Before-After UI (2) Reported Earnings 5 th Quintile	0.065 (0.030) [0.027]	0.085 (0.025) [0.001]	0.073 (0.045) [0.102]	0.092 (0.034) [0.006]

Notes: Source: Authors' Calculations using the NJS data. "Before-After UI Reported Earnings" consist of monthly UI earnings in the six quarters after random assignment minus monthly UI earnings in the 18 months before random assignment. "Before-After UI (2) Reported Earnings" consist of monthly UI earnings in the 6th quarter after random assignment minus monthly UI earnings in the 6th quarter before random assignment. The estimates come from logit models with an indicator for a positive self-evaluation as the dependent variable and the before-after earnings change variable and the categorical variables listed

in column one of Table 5 as independent variables. The values in the table are mean numerical derivatives, with standard errors in parentheses and p-values in square brackets. The before-after earnings changes enter in the form of indicator variables for being in the 2nd, 3rd, 4th, and 5th quintiles of the before-after earnings change distribution. The omitted category is the 1st quintile of the distribution. The population for these regressions is the treatment group. Indicator variables for missing values for the independent variables are also included in the regression.

TABLE 11: Logit Estimates of the Relationship between Before-After Employment Status Changes and Positive Self-Evaluation, By Demographic Group

	Adult Males	Adult Females	Male Youths	Female Youths
Employed Before & Not Employed After	-0.059 (0.102) [0.567]	-0.077 (0.093) [0.408]	-0.502 (0.230) [0.029]	-0.282 (0.139) [0.043]
Not Employed Before & Employed After	0.040 (0.028) [0.161]	0.014 (0.024) [0.570]	-0.112 (0.045) [0.013]	-0.013 (0.034) [0.696]
Always Not Employed	-0.123 (0.060) [0.038]	-0.042 (0.042) [0.319]	0.071 (0.142) [0.616]	-0.073 (0.064) [0.249]

Notes: Source: Authors' calculations using the NJS data. Employment status changes are based on changes in self-reported employment status measured at the date of random assignment and 18 months after random assignment. The omitted category is always employed. The estimates come from logit models with an indicator for a positive self-evaluation as the dependent variable and the before-after employment change variable and the categorical variables listed in column one of Table 5 as independent variables. The values in the table are mean numerical derivatives, with standard errors in parentheses and p-values in square brackets. The population for these regressions is the treatment group. Indicator variables for missing values for the independent variables are also included in the regression.

TABLE 12: Test Statistics from Logit Models of the Relationship between Before-After Estimates and Positive Self-Evaluation, By Demographic Group

	Adult Males	Adult Females	Male Youths	Female Youths
Before-After Self Reported Earnings	3.95 [0.4130] R ^{*2} =0.20	25.18 [0.0000] R ^{*2} =0.93	12.18 [0.0160] R ^{*2} =1.30	7.54 [0.1100] R ^{*2} =0.68
Before-After UI Reported Earnings	11.25 [0.0239] R ^{*2} =0.54	23.20 [0.0001] R ^{*2} =0.81	18.66 [0.0009] R ^{*2} =1.75	23.06 [0.0001] R ^{*2} =1.69
Before-After UI (2) Reported Earnings	7.14 [0.1286] R ^{*2} =0.36	23.04 [0.0001] R ^{*2} =0.84	2.95 [0.5659] R ^{*2} =0.33	9.74 [0.0451] R ^{*2} =0.73
Before-After Employment Status Changes	10.46 [0.0150] R ^{*2} =1.08	2.90 [0.4070] R ^{*2} =0.75	11.15 [0.0109] R ^{*2} =2.89	5.23 [0.1555] R ^{*2} =0.94

Notes: Source: Authors' Calculations using the NJS data. The values in the table are χ^2 -Statistics for joint tests of the null hypothesis that all of the coefficients equal zero for a given outcome, with the p-values in square brackets. The tests correspond to the estimates presented in Tables 9, 10 and 11. The value for R^{*2} are the partial r-squared times 100 for the group of covariates in the row. See the notes for those tables for further details.

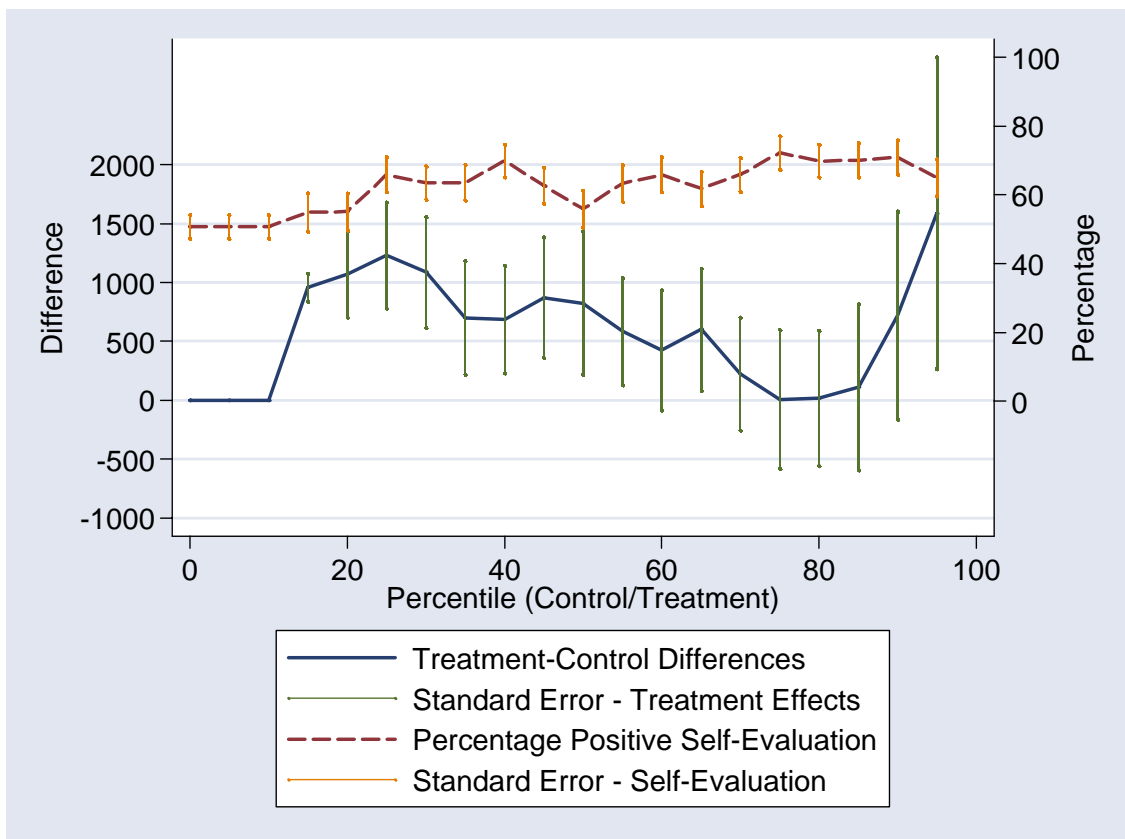
TABLE 13: Logit Estimates of the Relationship between Positive Self-Evaluation and Performance Standards, By Demographic Group

	Adult Males	Adult Females	Male Youths	Female Youths
<u>A. JTPA:</u>				
Employment at Termination	0.115 (0.034) [0.001] n=1507	0.126 (0.029) [0.000] n=1882	0.083 (0.047) [0.078] n=699	0.091 (0.039) [0.020] n= 890
Wages at Termination	0.012 (0.009) [0.173] n=617	0.028 (0.010) [0.003] n=873	0.026 (0.022) [0.232] n=280	0.021 (0.021) [0.316] n=319
Employment at Follow-up	0.097 (0.035) [0.005] n=1507	0.137 (0.029) [0.000] n=1882	0.043 (0.049) [0.380] n=699	0.076 (0.040) [0.056] n=890
Weekly Earnings at Follow-up	0.007 (0.015) [0.623] n=617	0.040 (0.015) [0.007] n=883	0.051 (0.033) [0.120] n=302	0.065 (0.030) [0.028] n=336
<u>B. WIA:</u>				
Employment at Termination	0.052 (0.036) [0.151] n=1155	0.048 (0.032) [0.122] n=1536	-0.073 (0.054) [0.179] n=528	0.035 (0.051) [0.484] n=705
Employment at 6-Months	0.030 (0.046) [0.511] n=659	-0.069 (0.044) [0.123] n=778	0.131 (0.065) [0.045] n=291	0.143 (0.060) [0.018] n=293

	0.000	0.001	0.002	0.004
Earnings Gain at	(0.000)	(0.000)	(0.001)	(0.001)
6-Months	[0.340]	[0.022]	[0.079]	[0.003]
	n=566	n=634	n=240	n=204

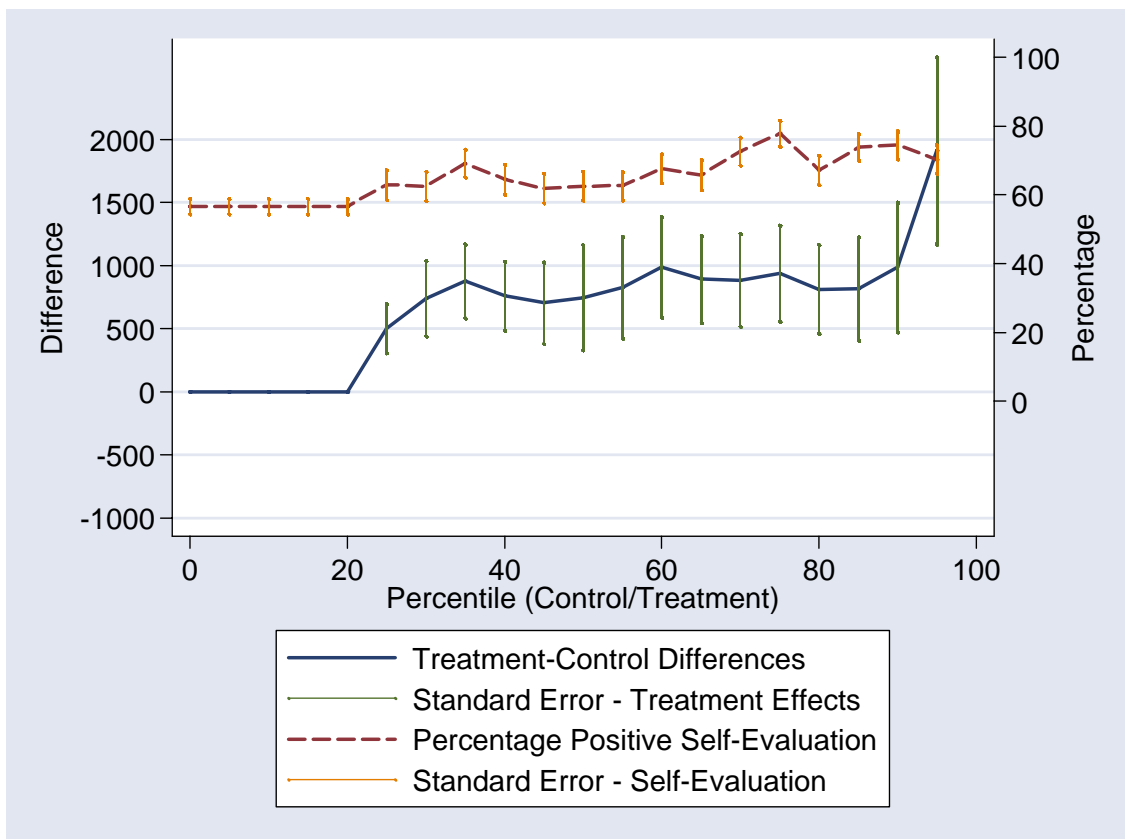
Notes: Source: Authors' calculations using the NJS data. The JTPA performance measures consist of (1) employment at JTPA termination date; (2) employment at follow-up, which is 13 weeks after termination in JTPA; (3) wage per hour at termination date (conditional on employment) in dollars; and (4) the average total weekly earnings at follow up (conditional on employment). Our construction of all of the JTPA performance measures relies on self-reported data. The WIA performance measures consist of (1) employment at exit, which we calculate as non-zero UI earnings in the calendar quarter of termination (conditional on non-employment at the date of random assignment based on self-reported labor force status); (2) employment at six months after termination, which we calculate as non-zero UI earnings in the third calendar quarter after termination (conditional on employment in the first quarter after termination); (3) earnings differences (conditional on employment in the first quarter after termination), which we calculate as the sum of UI earnings in the second and third calendar quarters after program termination minus the sum of earnings in the two calendar quarters prior to random assignment. The estimates in the table correspond to logit models with an indicator for a positive self-evaluation as the dependent variable and one of the performance measures as the only independent variable. The models also include all of the variables listed in Table 5 as additional covariates. The values in the table are mean numerical derivatives, with standard errors in parentheses and p-values in square brackets. We multiply the values for the earnings-based performance measures by 100 for ease of presentation. The final row in each cell gives the sample size for the sample used to produce each estimate. Before deleting observations with missing values of the performance measures, the treatment group samples contain 3067 adult males, 3922 adult females, 1308 male youths, and 1711 female youths. The population for these regressions is the treatment group.

FIGURE 1A: Quantile Treatment Effects and Percentage Reporting Positive Self-Evaluation, Adult Males



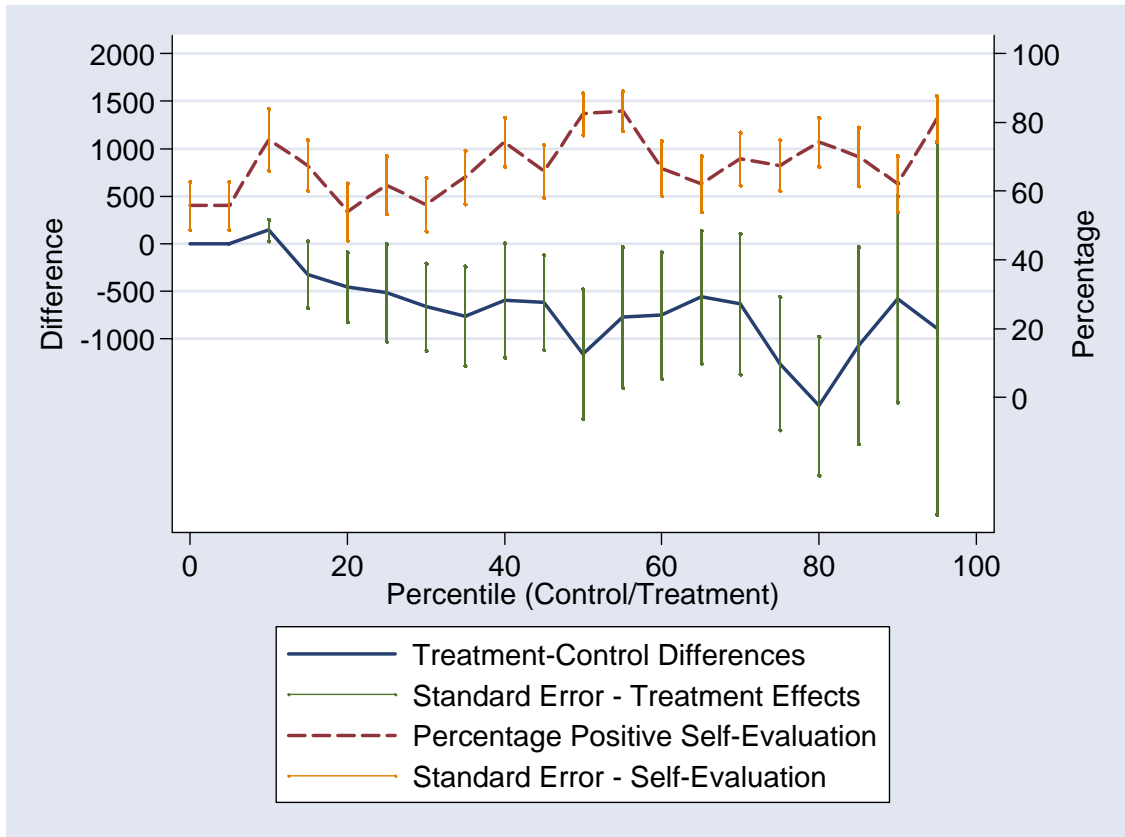
Notes: Source: Authors' Calculations using the JTPA data. The outcome used here is self-reported earnings over the 18months after random assignment.

FIGURE 1B: Quantile Treatment Effects and Percentage Reporting Positive Self-Evaluation, Adult Females



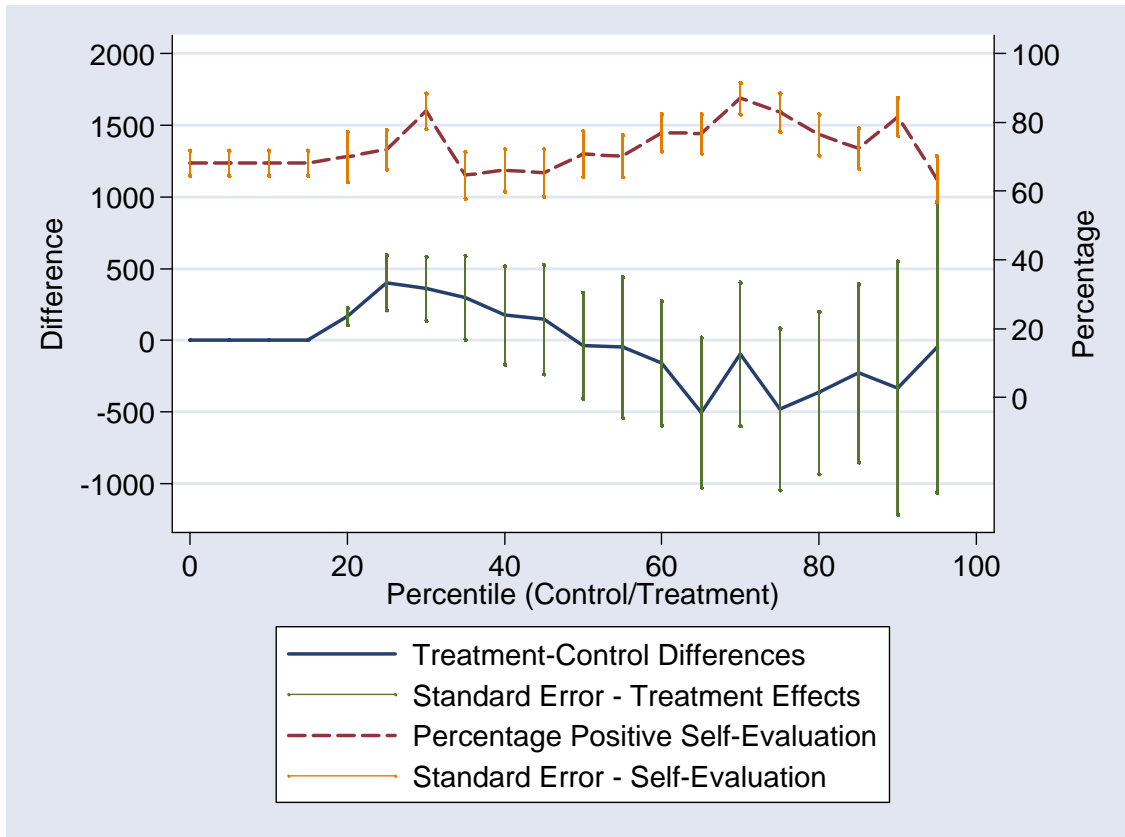
Notes: Source: Authors' Calculations using the JTPA data. The outcome used here is self-reported earnings over the 18months after random assignment.

FIGURE 1C: Quantile Treatment Effects and Percentage Reporting Positive Self-Evaluation, Male Youth



Notes: Source: Authors' Calculations using the JTPA data. The outcome used here is self-reported earnings over the 18months after random assignment.

FIGURE 1D: Quantile Treatment Effects and Percentage Reporting Positive Self-Evaluation, Female Youth



Notes: Source: Authors' Calculations using the JTPA data. The outcome used here is self-reported earnings over the 18months after random assignment.