

CASE STUDY 1: VOCATIONAL TRAINING FOR DISADVANTAGED YOUTH

Impact evaluation methods



This case study is based on:

Orazio Attanasio, Adriana Kugler, and Costas Meghir: “Training Disadvantaged Youth in Latin America: Evidence from a Randomized Trial” *American Economic Journal: Applied Economics* 3 (July 2011): 188 – 220.

J-PAL thanks the authors for allowing us to use their paper as a teaching tool and for sharing their data.

KEY VOCABULARY	
Comparison Group	A group that is as similar as possible to the treatment group in order to be able to learn about the counterfactual. In an experimental design, the comparison group (also called the control group) is a group from the same population as the treatment group that, by random assignment, is not intended to receive the intervention.
Counterfactual	What would have happened to the participants of an intervention had they not received the intervention. The counterfactual can never be observed; it can only be inferred from a comparison group different from the treatment group.
Estimate	In statistics, a "best guess" about an unknown value in a population (such as the effect of a program on an outcome) according to a rule (known as the "estimator") and the values observed in a sample drawn from that population.
Impact	The impact of the intervention is the effect of the treatment on the whole population. The impact is estimated by measuring the differences in outcomes between the treatment group and its counterfactual, i.e., by measuring the difference in outcomes between treatment and comparison groups.
Omitted Variable Bias	Statistical bias that occurs when relevant (and often unobservable) variables/characteristics are left out of the regression analysis. When these variables are correlated with both the primary outcome and a variable of interest (e.g., participation in an intervention), their omission can lead to incorrectly attributing the measured impact solely to the program. For example, omitting socioeconomic status, which is correlated with test scores, could lead to overestimating the impact of a tutoring intervention on a group of wealthy students.
Treatment Group	The group that receives the intervention.
Selection Bias	<p>Selection bias is bias that occurs when the individuals who receive the program are systematically different from those who do not. For example, consider an elective after-school tutoring program. Is it effective at raising children's exam scores? If we compare those who take up the tutoring program to those who don't, we will get a biased estimate of the effect of the tutoring program, because those who chose to participate are likely different from those who don't (for example, those who took it up may be more motivated, or they may be weaker students). Randomization removes selection bias because it breaks the link between characteristics of the individual and their treatment status.</p> <p>Selection bias can occur in other ways in a randomized evaluation. For example, consider a situation where an intervention is making a phone call to a landline:</p> <ul style="list-style-type: none"> - Callers may be unable to reach certain participants (for example, participants in rural areas may have poor cell phone service and may be more likely to have landlines than those in urban areas). - Some participants may be less likely to pick up the phone depending on the time of day they are called (for example, calling a home phone during standard business hours).

LEARNING OBJECTIVES

To identify evaluation methods and how they estimate impact differently. To better understand issues of bias and causal inference and think through how to use evaluation methods to measure impact.

SUBJECTS COVERED

Causality, counterfactual, impact, comparison groups, selection bias, omitted variables, randomization, and balance.

INTRODUCTION

What is required in order for us to measure whether a program had impact and, if so, how much of an impact?

This is the same as asking whether changes in certain outcomes can be attributed directly to the intervention, which in turn requires ensuring that these measured outcome changes are not caused by other factors or events happening at the same time. Ideally, evaluators would do this by following the progress of a group of people as they participate in a program, measure any changes that occur, and then go back in time and measure the same group's progress without the program in place. This second set of outcomes is called the **counterfactual**. Since we cannot observe the true counterfactual, the best we can do is to approximate it by constructing (or mimicking) it.

The key challenge of an impact evaluation is constructing the counterfactual. We typically do this by selecting a group of people who resemble the participants as much as possible but who did not participate in the intervention. This group is called the **comparison group**. Because we want to be able to say that it was the intervention and not some other factor that caused the changes in outcomes, it is important that the comparison group and the participant group are, on average, as similar as possible so that we can attribute any differences in outcomes to the intervention. We then estimate the **impact** as the difference in outcomes observed at the end of the intervention between the comparison group and the **treatment group**.

An accurate impact estimate can only be attained if the comparison group is a good representation of the counterfactual, or what the treatment group would have looked like had the intervention not happened. If the comparison group poorly represents the counterfactual,

then the estimated impact will be **biased**. Therefore, the method used to select, construct, or estimate the comparison group is a key decision in the design of any impact evaluation.

This case study will explore different methods for measuring impact by looking at a training program for disadvantaged youth introduced in Colombia in 2005. We will show how different methods may produce different results.

THE “YOUTH IN ACTION” PROGRAM

All around the world, many young people struggle to find stable employment in both developed and developing countries. By the end of 2010, around 75.1 million young people worldwide were unemployed (ILO). Youth unemployment is commonly blamed on a lack of skills, especially in countries where education systems fail to equip young people with the skills they need to get a stable job.

In 2001, the Colombian government started a vocational training program for disadvantaged youth in its seven largest cities to tackle the problem of youth unemployment. The training program included three months of in-classroom training and three months of on-the-job training for people between the ages of 18 and 25¹, who were placed in the two lowest deciles of the income distribution. The classroom training was provided by private institutions selected through a competitive bidding process, while the on-the-job training was provided by legally registered companies operating in various sectors, including manufacturing, retail and trade, and services. Participating youth were given US\$2.20 per day to defray transportation and lunch costs; women with children under seven years of age were given US\$3.00.

What was the impact of this program? The intention of the program was to equip participating youth with skills valued by employers and the main outcome measure was employment rate. Asking if the training program “worked” is to ask if it increased the probability that participating youth would be employed following the program. The impact is the difference between the employment rate of those who participated in the program to what their employment rate would have been had they not participated in the training program.

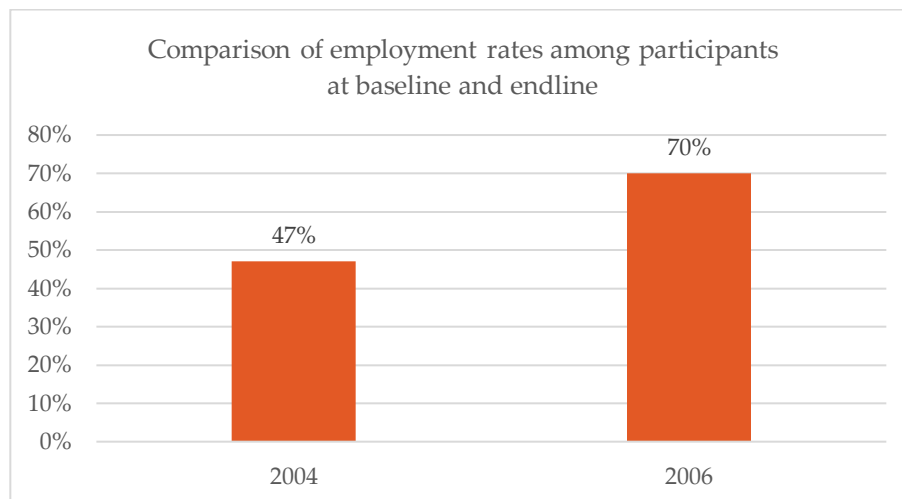
What comparison groups can we use? The following experts illustrate different methods of evaluating impact.

¹ While both men and women participated in the program, the sample of men in the evaluation showed some unbalance at baseline, so we present data only for women.

METHOD 1

Newspaper Article: Huge Gains for Women in Training Program

Statistics released today by a government agency indicate that the government-sponsored vocational training program “Youth in Action”, which has been running since 2001 in the seven largest cities of Colombia, increased the employment rate of participating women by 23 percentage points, a huge and important gain. At the beginning of the program, only 47 percent of participating women were employed, and when they were surveyed several months after completing the training program, they were 70 percent to have a job. These numbers provide evidence in support of vocational training programs, which governments all over the world have adopted to resolve the pressing problem of youth unemployment. Governments should take note of these results and start training programs or scale up existing ones.



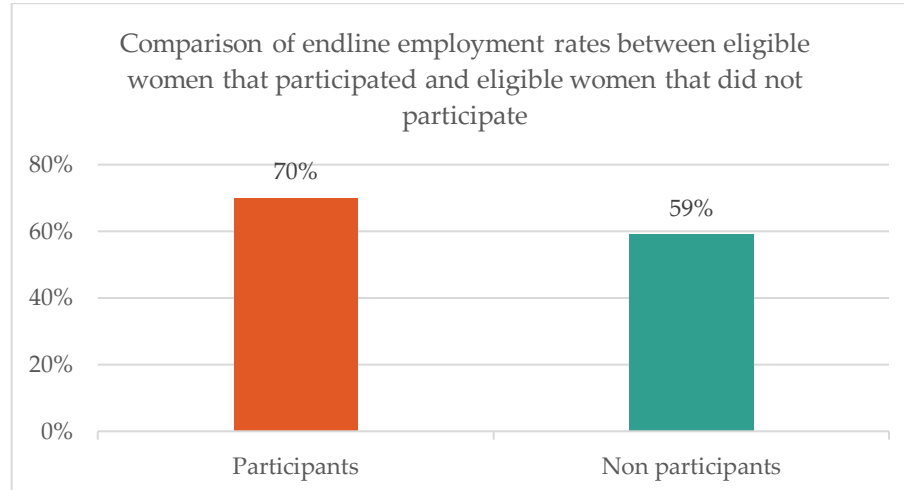
DISCUSSION TOPIC 1

- What type of evaluation method does this article imply?
- What represents the counterfactual?
- What are some potential problems with this type of evaluation method? List some confounding factors.

METHOD 2

Letter to the Editor: Youth in Action program—good but not great

Newspapers tend to exaggerate many claims and this is exactly what the article “Huge Gains for Women in Training Program” did last week when reporting about the impact of the government’s vocational training program. As an economist interested in labor markets, I have been following this training program since the government first announced it. Obviously, I hoped that the program would work and I am really happy to see positive results coming out from it. But the claims that the program had such a massive impact are very misleading. After all, many things could have happened to these women between the start and end of the training program. The Colombian economy has been experiencing healthy growth rates since 2002 and cities across the country have become safer. These confounding aspects could affect the results of the program’s evaluation, so we should get rid of these and focus instead on how women who participated in the training compare to women who did not participate in the training. I’ve gone ahead and collected data on women that were eligible to the program but did not participate. This shows that the program increased the employment rate of trained women by 11 percentage points, a far cry from the 23 percentage points increase claimed by the article, but still an increase nonetheless.



DISCUSSION TOPIC 2

1. What type of evaluation method is this letter employing?
2. What represents the counterfactual?

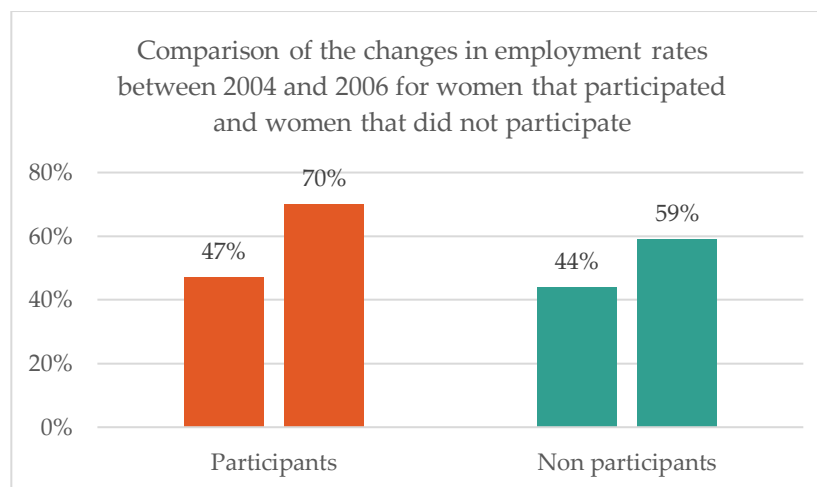
3. What are some potential problems with this type of evaluation method? List some confounding factors.

METHOD 3

Donor Report: Comparing apples to apples

The government's vocational training program has received a lot of press coverage recently. Some have claimed that the program has an enormous impact, while others argue that the impact is significantly more moderate. This report seeks to provide a more accurate measure of the impact of the program using a more appropriate method. Previous analyses have used the wrong metrics to calculate the training program's impact – possibly overestimating by how much the employment rate is actually increased by the program. For instance, if you compare the employment rates of those women who participated in the training program and those who did not, you might be introducing selection bias into the estimate. These two groups of women might be very different for many reasons beyond just participating or not in the training program.

What you need to do to get a more accurate estimate is to compare changes in the employment rates of the two groups. This way, we can see how fast employment rate changes for each group. When we repeat the analysis using this more appropriate outcome measure, we see that women participating in the program experienced an increase in their employment rate of 8 percentage point, showing that participating in a vocational training program does increase employment, but not by the magnitudes claimed by other analyses.



DISCUSSION TOPIC 3

1. What type of evaluation method is this report using?
2. What represents the counterfactual?
3. What are some potential problems with this type of evaluation method? List some confounding factors.

METHOD 4

Report: The numbers don't lie, unless your statisticians are asleep

Over the last few weeks, the public has received conflicting information about the impact of the Colombian government's vocational training program. Those who support the program assert that vocational training successfully equips young women with valuable skills, resulting in a substantially higher chance of being employed. Others, however, believe that this impact is grossly inflated and that actual gains are more modest, and perhaps driven by external factors and not the vocational training itself.

Unfortunately, both camps are using flawed instruments of analysis and the question of whether vocational training increases the chance of getting a job among women remains unanswered.

This report uses sophisticated statistical methods to measure the true impact of the vocational training program. We were concerned with other factors that might influence the results. As a result, we carried out a survey to collect information about age, marriage status, education levels, and the city where participants lived. To observe the possible bias caused by omitting key variables, we conducted one "naïve" analysis without controlling for these differences, and a separate analysis with controls. This helped us to obtain the true impact of the program.

TABLE 1
Probability of Employment

	(1)	(2)
Training	0.11 ** (0.022)	0.069** (0.022)
Age		0.004 (0.005)
Marriage		-0.066* (0.026)

Education Level		0.007 (0.006)
City		-0.036*** (0.005)
Constant	0.59** (0.02)	0.59 ** (0.14)
Observations	1,769	1,769

The results from column (1) suggest that the training program increased the probability of employment by 11 percent. If we look at column (2), which includes controls for confounding variables, the impact is diminished to 6.9 percent.

DISCUSSION TOPIC 4

1. What type of evaluation method is utilized in this report?
2. What represents the counterfactual?
3. What are some potential problems with this type of evaluation method? List some confounding factors.

METHOD 5:

Report: Understanding participation into “Youth in Action”

Two reports were recently published presenting different approaches to measure the impact of the program “Youth in Action”. Unfortunately, they miss a key point in the way participants were selected, which allows to estimate much more rigorously the impact of the program.

Each of the 114 training institutions participating in the program was instructed to select up to 50 percent more applicants than they could accommodate. Applicants were randomly selected from this list to fill over 26,000 total available slots, with the remaining people on the list forming the comparison group. This methodology allowed us to control for self-selection into the training program; all individuals in both the treatment and comparison groups had chosen to apply for the training and had been selected as suitable by the training institutions. Due to the large scale of the experiment, we surveyed a random sample of approximately 2,000 applicants from both the treatment and comparison groups. By simply comparing the employment rates of the two groups at endline, we found that **employment among women who participated in the training program was 5.4 percentage points higher than in the comparison group.**

DISCUSSION TOPIC 5

1. What type of evaluation method is used in this report?
2. What represents the counterfactual?
3. What are some potential problems with this type of evaluation method? List some confounding factors.

COMPARING ALL FIVE METHODS

Below are the impact estimates of the “Youth in Action” using the five different methods you have discussed in this case study.

Table 2: Comparing all five methods

Method	Estimated impact
Pre-Post	23 pp
Simple Difference	11 pp
Difference-in-Differences	8 pp
Multivariate Regression	6.9 pp
Randomized Evaluation	5.4 pp

As you can see, not all methods yield the same result. Hence, the choice of method is crucial.

REUSE AND CITATIONS

To request permission to reuse this case study or access the accompanying teachers' guide, please email training@povertyactionlab.org. Please do not reuse without permission. To reference this case study, please cite as:

Abdul Latif Jameel Poverty Action Lab (J-PAL). 2019. "Case Study: Why Randomize?: Vocational Training for Disadvantaged Youth." J-PAL Case Study. Last Modified 17 September, 2019.

	Method	Description	What assumptions are required, and how demanding are the assumptions?
Randomization	Randomized Evaluation/ Randomized Control Trial	Measure the differences in outcomes between randomly assigned program participants and non-participants after the program took effect.	<i>The outcome variable is only affected by program participation itself, not by assignment to participate in the program or by participation in the randomized evaluation itself.</i> Examples for such confounding effects could be information effects, spillovers, or experimenter effects. As with other methods, the sample size needs to be large enough so that the two groups are statistically comparable; the difference being that the sample size is chosen as part of the research design.
	Pre-Post	Measure the differences in outcomes for program participants before the program and after the program took effect.	<i>There are no other factors (including outside events, a drive to change by the participants themselves, altered economic conditions, etc.) that changed the measured outcome for participants over time besides the program.</i> In stable, static environments and over short time horizons, the assumption might hold, but it is not possible to verify that. Generally, a diff-in-diff or RDD design is preferred (see below).
	Simple Difference	Measure the differences in outcomes between program participants after the program took effect and another group who did not participate in the program.	<i>There are no differences in the outcomes of participants and non-participants except for program participation,</i> and both groups were equally likely to enter the program before it started. This is a demanding assumption. Non-participants may not fulfill the eligibility criteria, live in a different location, or simply see less value in the program (self-selection). Any such factors may be associated with differences in outcomes independent of program participation. Generally, a diff-in-diff or RDD design is preferred (see below).
Basic Non-Experimental Comparison Methods	Differences in Differences	Measure the differences in outcomes for program participants before and after the program <i>relative</i> to non-participants.	<i>Any other factors that may have affected the measured outcome over time are the same for participants and non-participants, so they would have had the same time trajectory absent the program.</i> Over short time horizons and with reasonably similar groups, this assumption may be plausible. A “placebo test” can also compare the time trends in the two groups before the program took place. However, as with “simple difference,” many factors that are associated with program participation may also be associated with outcome changes over time. For example, a person who expects a large improvement in the near future may not join the program (self-selection).
	Method	Description	What assumptions are required, and how demanding are the assumptions?
More Advanced Statistical Non-experimental	Multivariate Regression/OLS	The “simple difference” approach can be—and in practice almost always is—carried out using multivariate regression. Doing so allows accounting for other observable factors that might also affect the outcome, often called “control variables” or “covariates.” The regression filters out the effects of these covariates and measures differences in outcomes between participants and non-participants while holding the effect of the covariates constant.	Besides the effects of the control variables, <i>there are no other differences between participants and non-participants that affect the measured outcome.</i> This means that any unobservable or unmeasured factors that do affect the outcome must be the same for participants and non-participants. In addition, the control variables cannot in any way themselves be affected by the program. While the addition of covariates can alleviate some concerns with taking simple differences, limited available data in practice and unobservable factors mean that the method has similar issues as simple difference (e.g., self-selection).
	Statistical Matching	<u>Exact matching:</u> participants are matched to non-participants who are identical based on “matching variables” to measure differences in outcomes. <u>Propensity score matching</u> uses the control variables to predict a person’s likelihood to participate and uses this	Similar to multivariable regression: <i>there are no differences between participants and non-participants with the same matching variables that affect the measured outcome.</i> Unobservable differences are the main concern in exact matching. In propensity score matching, two individuals with the same score may be very different even along observable dimensions. Thus, the assumptions that need to hold in order to draw valid conclusions are quite demanding.

e r i m e n t a l m e t h o d s		predicted likelihood as the matching variable.	
	Regression Discontinuity Design (RDD)	In an RDD design, eligibility to participate is determined by a cutoff value in some order or ranking, such as income level. Participants on one side of the cutoff are compared to non-participants on the other side, and the eligibility criterion is included as a control variable (see above).	<i>Any difference between individuals below and above the cutoff (participants and non-participants) vanishes closer and closer to the cutoff point. A carefully considered regression discontinuity design can be effective. The design uses the “random” element that is introduced when two individuals who are similar to each other according to their ordering end up on different sides of the cutoff point. The design accounts for the continual differences between them using control variables. The assumption that these individuals are similar to each other can be tested with observables in the data. However, the design limits the comparability of participants further away from the cutoff.</i>
	Instrumental Variables	The design uses an “instrumental variable” that is a predictor for program participation. The method then compares individuals according to their predicted participation, rather than actual participation.	<i>The instrumental variable has no direct effect on the outcome variable. Its only effect is through an individual's participation in the program. A valid instrumental variable design requires an instrument that has no relationship with the outcome variable. The challenge is that most factors that affect participation in a program for otherwise similar individuals are also in some way directly related to the outcome variable. With more than one instrument, the assumption can be tested.</i>

