



**clear**

Regional Centers for  
Learning on Evaluation and Results

# Impact Evaluation: Scope and Limits in the Real World

*Claudia Maldonado*

*17th October 2017*

*NEC UNDP*

# Evaluation

Paraphrasing Wildavsky, evaluation is:  
a way of **speaking truth to power**

The use of the tools of social science to answer questions about attributes of public programs/intervention.

Which attributes?

# Learning Objectives

- Understanding of key concepts in impact evaluation literature
- Gain familiarity with experimental and quasiexperimental approaches
- Critically discuss scope and limits of impact evaluation in the real world

# Warm-up discussion

- What is evidence?
- What is impact?
- What is strong evidence of impact?
- Why is strong evidence expected from evaluation?

# Impact as casual effect

- Causal link = independent effect
- Average treatment effect (ATE)
- Impact evaluation implies the *identification* and *isolation* of the independent causal effect of an intervention on the treated population.

# Impact evaluation

... requires meeting the challenge of making **causal inferences** using the methodological tools of scientific inference.

# The Clinical Metaphor

- A program is a discrete and exogenous intervention ( $X_{t1}$ ) on a given problem/situation ( $Y_{t1}$ ) that seeks to alter the *status quo in some desired direction*. ( $Y_{t1} \neq Y_{t2}$ )
- If Y is rural poverty, then it is desired that ( $Y_{t1} > Y_{t2}$ )
- $X_{t1}$  = treatment       $Y_{t1}$  = problem



# Types of evaluation (questions)

- Is  $Y_{t1}(z_1, z_2, z_3, \dots, z_n)$  a public problem that is susceptible of public intervention?  
**(Diagnostics)**
- Is  $X_{t1}$  an adequate instrument to address  $Y_{t1}$ ?  
**(Design)  $X \in (Z)$**
- What is the actual change in  $Y_{t2}$  directly and exclusively attributable to  $X_{t1}$ ? **(Impact)**



# Another way of looking at impact

- The empirical verification of a program's theory of change
- The answer to a counterfactual
- **“What would have happened to the treated population in the absence of the program? (impact variable, i.e. rural poverty)”**

# What type of impact?

- Define impact variable (variable to be observed)
- Define aggregation level (individual, household, community?)
- Define measurement and means of verification (how will you know if and when change occurs?)

# Impact Evaluation: How to identify net impact?

- $\Delta$  in impact variable is NOT impact
- $\Delta=0$  does NOT mean there is no impact
- $\Delta=X$  does NOT reflect the *magnitude* or *direction* of impact

## Why?

# The problem of attribution

There are factors that affect the problem that coexist with the program, but are not related to the program:

- Time
- External shocks
- Demographic change
- Other programs
- Preferences of the population in the program
- Changes in relative prices
- Characteristics of the treated
- And a long.... very long etcétera.

# Counterfactual

What would have happened if you had chosen not to attend NEC 2017?

Well, we may NEVER know!!!!

# Because...in order to answer that question

- We must observe the universe in Yt1 without a program and that same universe in Yt1 with program, at the same time.
- If we observe Yt2 with program and compare it to Yt1 without program (before and after), we are going to observe the effects of the program with the effects of everything else (the world turning around).

# Common problems ....

- **Before and after comparisons**
- **Comparison of non comparable groups/units**
- **Bias** (selection, attrition, measurement error)
- **Observables**
- **Non observables**



# Exercise 1: the fallacy of attribution

- A community called Penkalingo has suffered from severe unemployment for the last 4 years. Last year, the President of Penkalingo launched a ambitious public works program in order to create 800,000 new temporary jobs, and cut unemployment by half.
- 2012-2016: unemployment 34% average
- 2017: unemployment rate 38% so far.

# The problem of attribution

- Since it is impossible to observe two state of things that are mutually exclusive at the same time....
- We need to build a counterfactual ....
- Using the experimental method ...

# Experimental Design

- Random assignment of treatment & control group.
- Average impact variable value is compared between treatment and control groups.
- The difference in means is attributable to the program
- This is a randomized controlled trial (RCT)

# Random selection and assignment

- What is it?
- What does it solve?
- What about internal/external validity?

# Causal inference

- If  $X=1$  (treatment), then  $Y_t$
- If  $X=0$  (control), then  $Y_t$  does not occur ( $Y_c$ )
- The difference between  $Y_t$  and  $Y_c$  is the average treatment effect (ATE) of the intervention  $X$ .

# Summary

- RCTs are an evaluation strategy for impact evaluation that credibly claims to rigorously make causal inferences about what works.
- RCTs answer the counterfactual question.
- RCTs require impact evaluation to be part of the design phase of an intervention, so that the logistical procedures to preserve methodological control are taking place (randomization, internal validity, group integrity, time horizon).
- Despite being considered the gold standard within the evidence based policy movement, they are less frequent than other evaluation strategies, and tend to answer very specific question about what works.
- What types of questions do impact evaluation with RCT NOT answer?



**clear**

Regional Centers for  
Learning on Evaluation and Results

# Quasiexperimental Methods

---



# Criticisms to “RCT”

- Ethical considerations
- Logistical considerations
- External validity problems
- Trade off between internal and external validity

# Plan B?

- The statistical (observational) construction of conditions that are close enough to experimental ones.
- PSM: Propensity score matching

# PSM step by step ...

- Construction of a statistical clone of the participants in the treatment, conditional on their observed characteristics ( $X$ )
- In other words, groups that may not be identical on “ $X$ ”, but are equally likely to have participated in the program, given “ $X$ ” (equivalently propense).

# And how do you know ...

- We have a group of participants/treated ( $T=1$ ) and group of not treated ( $T=0$ )
- We gather data on their individual characteristics  $I_1(x_1, x_2, x_3, x_4)$
- We calculate how each characteristic affects their probability of participating in the program/treatment and we get their overall propensity score, given  $x$ .
- We compare the averages of the treated and not treated in the impact variable: the difference is the treatment effect on the treated (TOT).

# How is it done?

$$P = \alpha + \beta_1 x_1 + \beta_2 x_2 + \beta_3 x_3 + \dots + \varepsilon \quad (\text{e.g. Probit model})$$

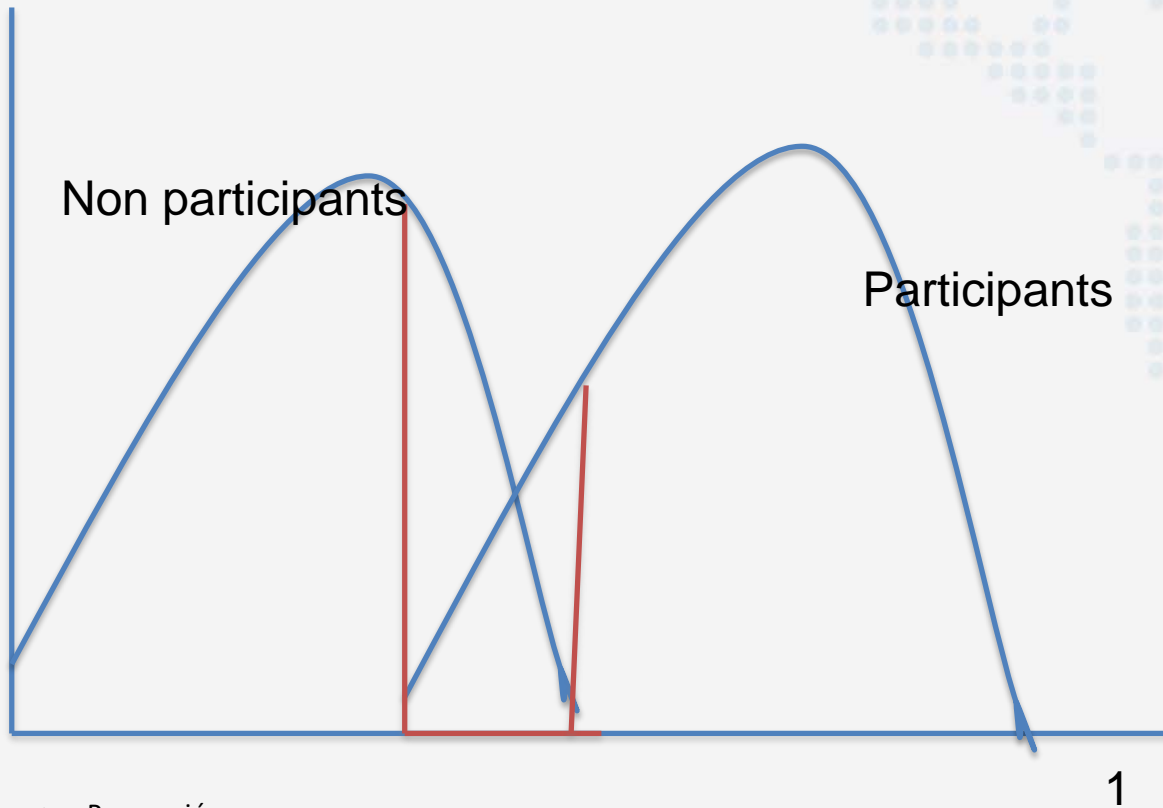
- $P$  = participation (discrete variable)
- $\beta$  = marginal contribution of a characteristic  $x$  on the probability of participating in the program

# Assumptions

- Conditional independence: no unobservable variables affect  $P$  (the probability of participating)
- Common support: some commonalities in propensity score in control and treatments group.
- Problem: only controls for bias due to observables

# Common support

- Density



- Propensión



# Formally...

$$TOT_{PSM} = E_{P(X) | T=1} \{ E[Y_T | T=1, P(X)] - E[Y_C | T=0, P(X)] \}.$$

Expected value of Y on the treated

Expected value of Y on the not treated  
(control)

*TOT: Treatment effect on the treated*

*psm: propensity score matching*

*E = expected value*

*If there is no « invisible hand » or unobserved variable that is affecting the probability of participation. Voilà!!!! There is our effect!*

*Las X no son afectadas por el programa (ej. Si es un programa de ingreso/empleo, no puedo usar esa variable X para el parametrico, porque se correlaciona con Y directamente).*

# Methods for matching

- Closest neighbour:
- Caliper : por radio
- Stratification
- Kernell (estimación no paramétrica de promedios ponderados de no participantes)

# Example: PROGRESA

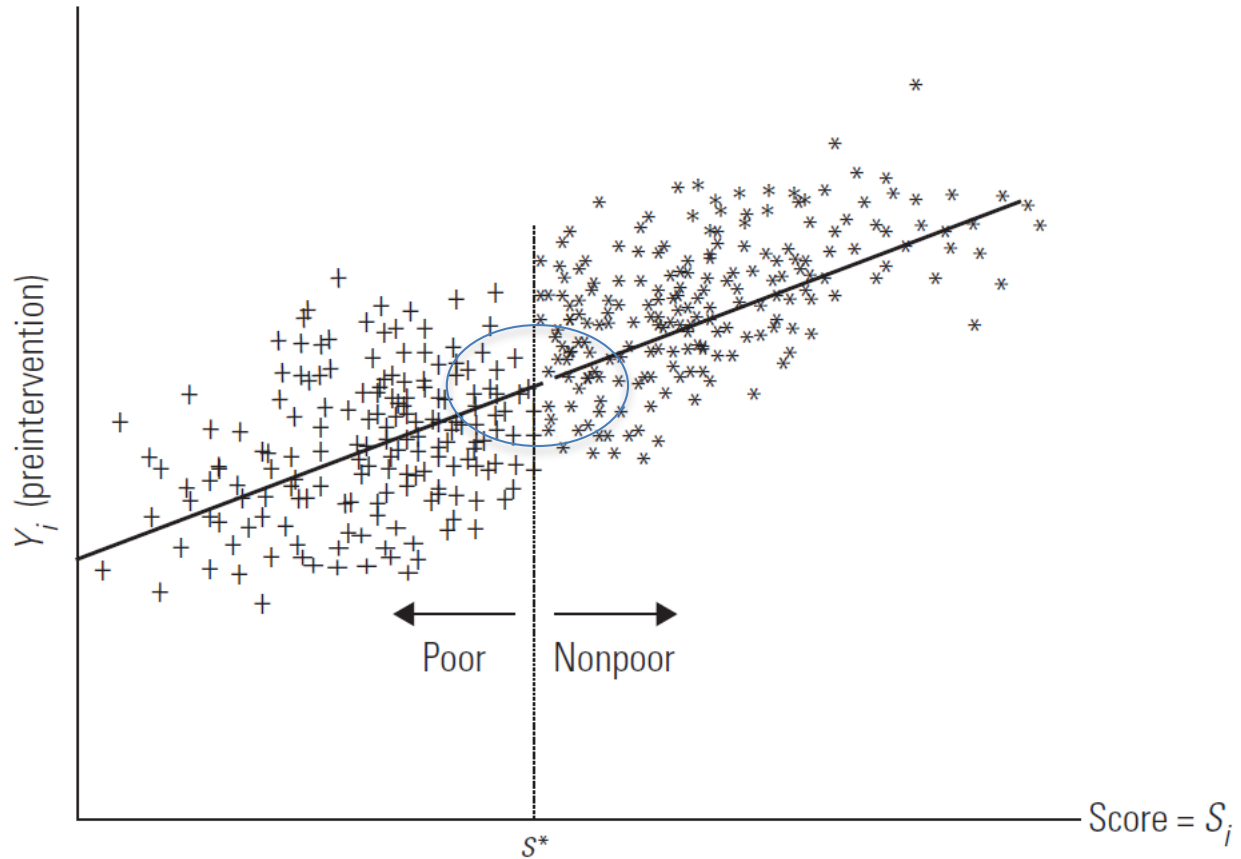
- Díaz y Handa: PSM is reliable when experimental results are compared to quasiexperimental ones. But if data sources (different instruments to measure income, for instance) vary too much, bias is a real problem.

# Regression discontinuity (RD)

- Non experimental method that uses exogenous rules of eligibility to identify participants and non participants.
- “Controls” for observed and unobserved heterogeneity.
- Compare people/groups before and after the eligibility threshold, assuming that they are fundamentally similar and with the same take up rate.

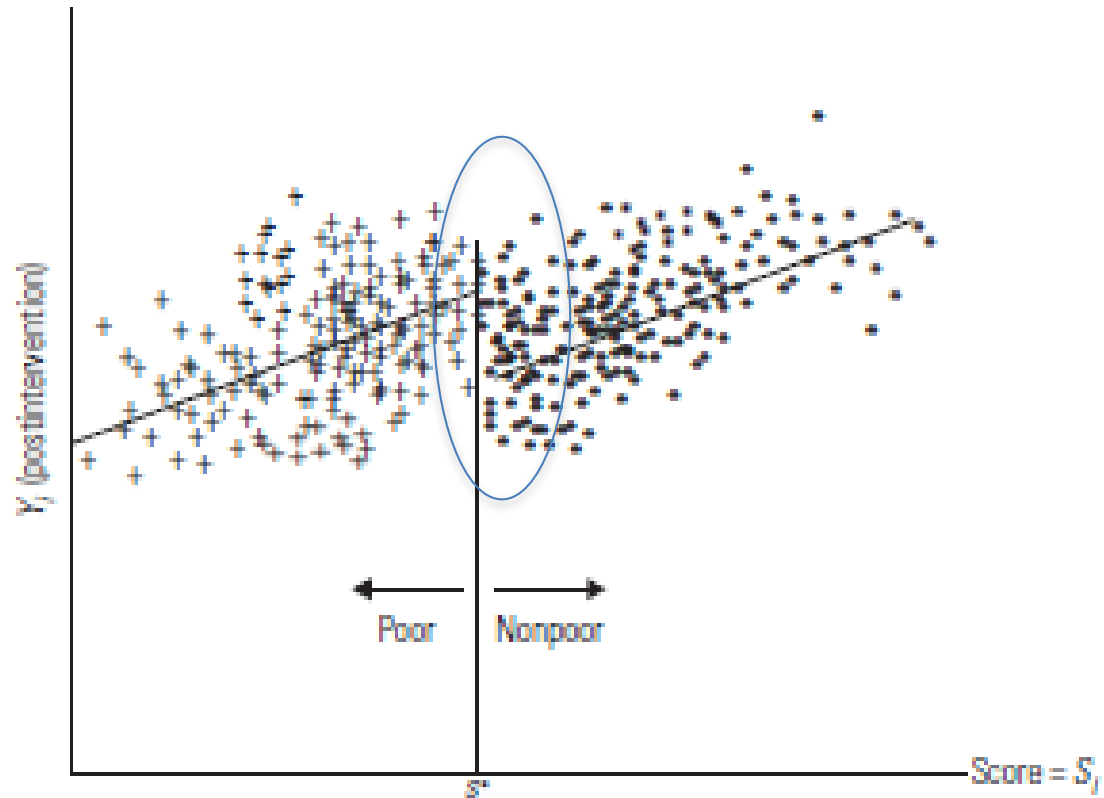
# Pre-intervention scenario

Figure 7.1 Outcomes Before Program Intervention



# Escenario post-intervención

Figure 7.2 Outcomes after Program Intervention

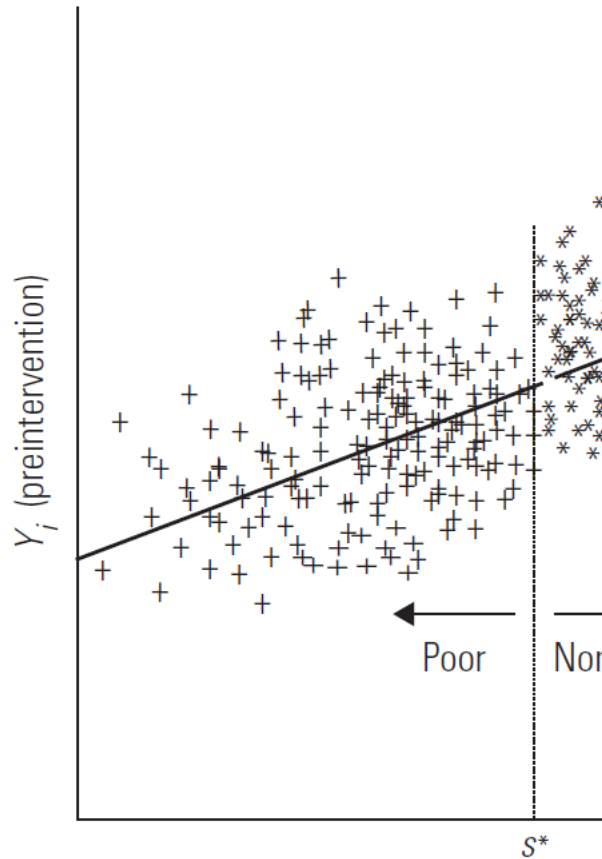


Source: Authors' representation.

# Regression Discontinuity

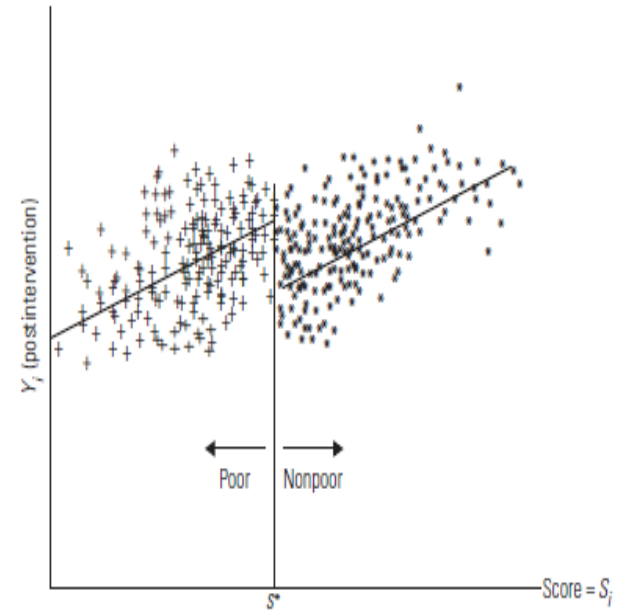
Figure 7.1 Outcomes before Program Intervention

• PRE



• POST

Figure 7.2 Outcomes after Program Intervention



Source: Authors' representation.



# Examples?

- We need a clearcut threshold that participants cannot manipulate as eligibility rule.
- If that is the case how different can people be around the cut point?

# Examples

- 65 and older
- Piso firme
- GRE scores for admissions
- Others?

# Advantages

- Identifies unbiased estimate around the cut point
- Uses program rules in its favour
- Does not imply rationing for control group
- BUT:
- The rule **MUST** be binding and exogenous.

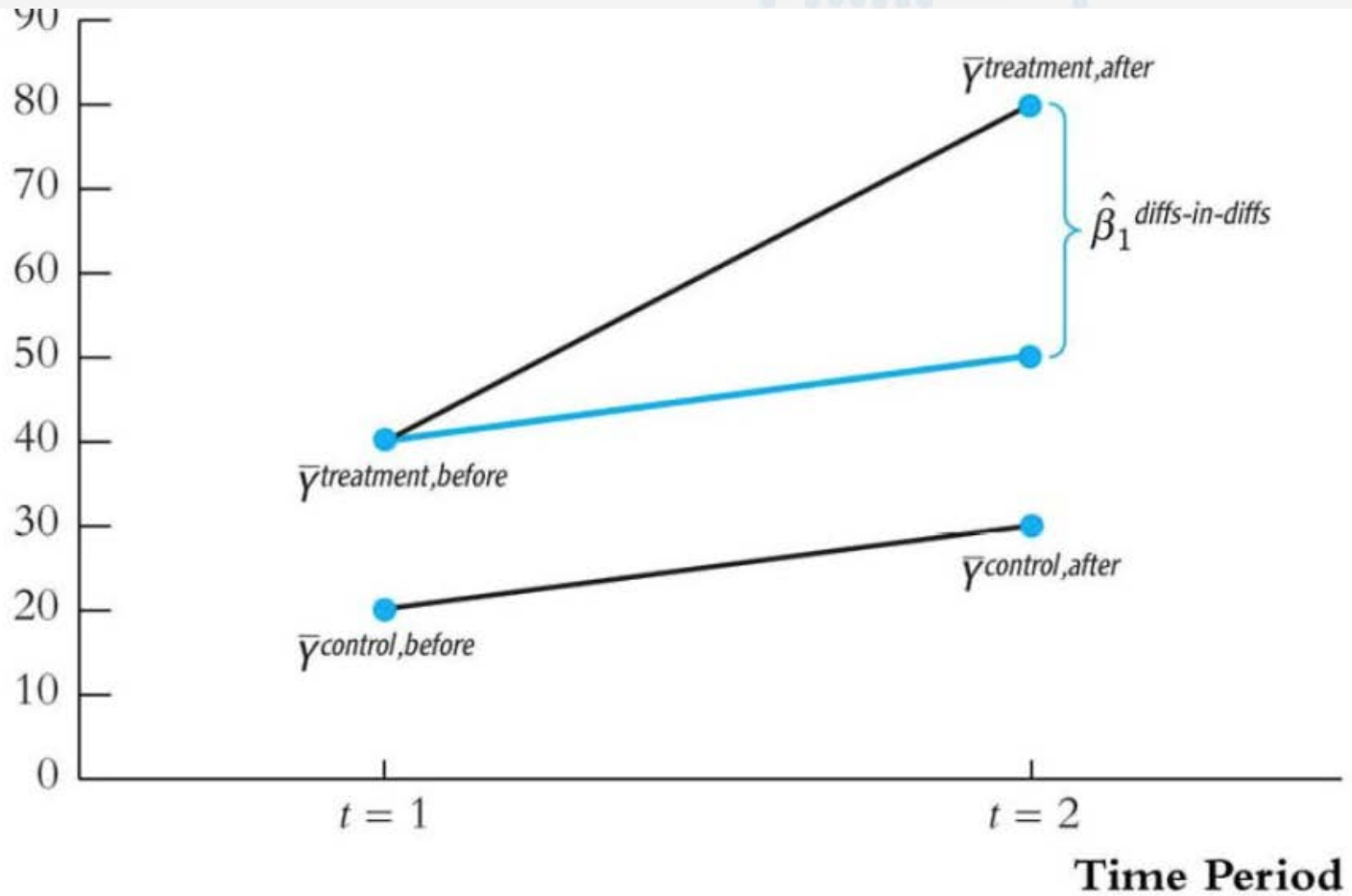
# Disadvantages

- Identifies average “local” impact, not always generalizable.
- Fewer observations around the cut point. (vs. experimental)
- Results can be less robust ( i.e. functional form or interactions).

# Using Diff in Diff as contrafactual: DD

- Comparison of groups (treatment vs control) according to their trayectories before and after the intervention.
- Baseline of both groups.
- Before and after analysis.
- The difference in the differences IS THE PROGRAM´S EFFECT.

# EXAMPLE



# DD

$$DD = E(Y_{1t} - Y_{0t} \mid T_1 = 1) - E(Y_{1c} - Y_{0c} \mid T_1 = 0)$$



# Assumptions

- Combined with PSM: time invariant unobservables
- *Solves bias for unobserved heterogeneity iff:*
- *Biases are time invariant*
- *Parallel trajectories*

# Ashenfelter's dip

- Effect of capacity building for labour inclusion
- Participants just suffered a temporal negative shock.
- Treatment groups will anyway show an increasing income over time...

How does this affect the estimate of impact?

# Beyond the black box

- Correlation is NOT causation!!!!
- In development is it as important to know **what works** as to know **how it works** or doesn't
- We need causal mechanisms and causal narratives ...

# Some examples

- Condoms and pregnancy prevention
- Retroviral and adherence rate
- Preventing infant mortality in indigenous areas
- Influencers and healthy habits

# Criticisms to EBPP

- *Boundary conditions and replicability.*
- Empirical record (CCTs)
- Pilots or prototypes?
- What does impact really mean?

# Qualitative response

- Old fashioned concept of causality from the C. XIX! (positivismo *naive*)
- Contingency and randomness are not the same
- Sampling (representative vs. intentional)
- *Process-tracing* and nested inference
- Interventions are open ended processes in themselves, not treatment like pills.

# Examples

- CCTs in indigenous areas
- No health impacts of proven nutritional supplement
- Others?



# Use of experimental impact evaluation

- Credibility and rigour widely recognized
- Legitimacy
- Evidence for scaling up (replicability)
- Knowledge accumulation on what works
- Global learning and diffusion of good practice
- New trend: delivery sciences

# Final Remarks?

- Gold standard of what?
- What about communication and decision making value?
- Sociology of knowledge matters ...
- There is no single right question to ask in an evaluation, there are reasonable and credible ways to answer a specific question, but what really matters is what you want the answer for...