

14.74

Lecture 9: The effect of school buildings on schooling: A “natural experiment”

Esther Duflo

March 2, 2011

Does top-down school construction lead to an increase in schooling? And does this lead to an increase in wages down the road? Or is school construction useless unless it is really desired by the community?

Last time we saw one tool for program evaluation: randomized evaluation. While randomized evaluations are very powerful when it is possible to do them, some time we want to evaluate policies that have already taken place, and it is not possible to randomize. We will see today another tool which comes handy in such cases, the method of “differences in differences”

1 School construction in Indonesia: The Set up

1.1 The INPRES school construction program

Second five year plan (1974-79)-Oil shock.

- A large program:
 - 61,807 primary schools constructed from to 1973/74 to 1978/79.
Number of schools multiplied by 2. 1 schools for every 500 children.
 - A *change* in policy: Before 1973, no construction, ban on recruiting for public service positions.

- A program meant to favor low-enrollment regions.

Allocation rule: number of schools constructed in a district proportional to the number of children (ages 7 to 12) *not enrolled in primary school*.

1.2 Data

SUPAS 95: A survey done in 1995: after the children educated in these schools have completed their schooling, and have started working.

- 150,000 men born 1950-1972
- Variables: education, year and region of birth, wages.

1.3 Sources of variation

Two factors affect the intensity of the program.

- *Year of birth* :
- *Region of birth* The government was targeting low enrollment regions \Rightarrow substantial variation in program intensity across districts.

2 The “Difference in differences” methodology

• Basic idea

Suppose that there are two regions in the data: a “high program” region, and a “low program” region.

Suppose that we have to age group of the individuals: “young people”, born after 1967 and who could fully benefit from the schools, and “old people” born before 1962, and who could not benefit at all from the schools.

So in total, we have four groups: YOUNG and High program, OLD and high program,....

Let us construct the average education for each of this group, and put them into a box. Use the stata handout, and the template (table 3).

- Calculate, D_{11} , the difference between the “HIGH” and “LOW” average among the young: what do we find and why?
- Calculate D_{21} the difference between the YOUNG and the OLD in the high program region: what do we find and why?
- Calculate, D_{12} the difference between the “HIGH” and “LOW” average among the old: how does it compare to D_{11} ? why?
- Calculate the difference $DD_a = D_{11} - D_{12}$. How do you interpret it?
- Calculate the difference $DD_b = D_{21} - D_{22}$. How does it compare to DD_a ? How do you interpret it?
- Could DD_a or DD_b be a good measure of the program?
 - Under what assumption?
 - Is assumption likely to be satisfied?

- **Control experiment**

We have a possibility to check that the assumption is not rejected in the available data.

Suppose we fill the same boxes, but we now compare the “OLD” to the “VERY OLD”. Neither of them benefited from the program: what do we expect to see if the assumption is satisfied?

What do we expect to see if the assumption is not satisfied?

Do it: what do we see?

3 Extending difference in differences

3.1 Using all the regional variation

There are 280 districts in Indonesia, and we know how many schools each district has received: grouping the region into two groups is throwing away some information!

Before, we had 2 regional group, and 2 age group, we formed 4 age-region group. Now we have 280 regional group, 2 age group, how many groups can we form? What are these groups?

First, we form the average for each group (we can use the stata command: collapse). See an extract of the data set in the handout. We will note S_{Yj} the average education of the young in any region j , and S_{Oj} the average education of the young in any region j .

What can we do next?

-Take the difference between young and old in all the regions

-Plot the differences against the number of school constructed per 1000 child during the INPRES program (see graph)

- What do we see? What does this suggest?

- Suppose we run the regression:

$$S_{Yj} - S_{Oj} = \alpha P_j + v_j \quad (1)$$

Where can you see the slope of this regression?

- See stata handout: what is the result of running this regression? What can we conclude?

-Under what assumption is this conclusion valid?

-Any suggestion to test this assumption?

-Do you see this test anywhere in the handout?

3.2 Using regional and age variation

The last generalization (after that, we are done!) is that we don't have only 3 age groups (young, old, and very old): we have 23 age groups (everybody born between 1950 and 1972).

How many groups can we now form?

Note $S_{j2}, S_{j3}, \dots, S_{jk}, \dots, S_{j24}$ the average education of people born in region j , and who were of age 2, 3, ... $k, \dots, 24$, when the program started.

Suppose we run the regression:

$$S_{j2} - S_{j24} = \alpha_2 P_j + v_{j2}$$

What is α_2 ?

Suppose we run the regression:

$$S_{j23} - S_{j24} = \alpha_{23} P_j + v_{j23}$$

What is α_{23} ? What should α_{23} be equal to?

In general, suppose that for all ages k we run the regression:

$$S_{jk} - S_{j24} = \alpha_k P_j + v_{jk}$$

For what values of k should we see a positive α_k ? (remember that children attend primary school until age 12). Should we see the coefficient be larger for younger children or older children?

Look at figure 2 in the handout: what does each dot represent? Do the dots have the expected pattern?