

# **Economics 172: Problem Set #2**

Due on October 20, 2025 at 11:00pm

**Zachary Brandt**  
[zbrandt@berkeley.edu](mailto:zbrandt@berkeley.edu)

## Problem 1: RCT

This question uses an adapted dataset based on Muralidharan, Singh, and Ganimian's (2019) paper *Disrupting Education? Experimental Evidence on Technology-Aided Instruction in India*. The paper is available online and the replication data is available on ICPSR. Download the adapted dataset from bCourses.

This project evaluated the impact of a center-based and technology-aided after-school educational program on math and Hindi performance among middle schoolers living in low-income neighborhoods in urban India. The technology-based curriculum was designed to be high-quality, adaptive, and engaging. Approximately 600 middle schoolers were recruited to participate in the study. Half of these recruited students were randomly allocated by lottery to receive a voucher to participate in the program (treatment group), and half were not (control group). For the purposes of this question, you can assume that 100% of those assigned to treatment participated in the program and that 0% of those assigned to control participated in the program.

Below is a list of the variables included in the dataset, with a brief description of each. Note that BL refers to versions of variables collected at baseline (collected before the program began) while EL refers to variables collected at endline (collected at the conclusion of the program).

Variable	Description
student_id	Identification numbers that uniquely identify students
student_age	Age of the student (collected at baseline)
student_female	Indicator (0/1) for whether the student is female
student_grade	Grade of the student (collected at baseline)
treatment	Indicator (0/1) for treatment status of the student
BL_math_percent, EL_math_percent	Math score in percent correct
BL_hindi_percent, EL_hindi_percent	Hindi score in percent correct
BL_ses_index, EL_ses_index	Household wealth index

1. Import the data and generate a table of summary statistics. What is the range of ages and grade levels within this sample? What are average scores for math and Hindi at baseline? At endline?
2. Next you'll check whether the treatment and control groups are balanced in terms of each of the following variables: (1) age, (2) sex, (3) household wealth index at baseline, (4) math scores at baseline, and (5) Hindi scores at baseline. (Hint: To do this, you'll run a separate regression for each of these five variables. Use `stargazer` or other preferable command to output the tables.) Before you code up any regressions, write down the regression you plan to run for at least one of the variables. Which parameter represents the coefficient of interest? What do you expect the estimate of this parameter to be (positive, negative, zero) and why?
3. Use the `lm` command to run the regressions in R. Are there significant differences between the treatment and control groups for any of these five variables? What is the purpose of this exercise, and what are you able to conclude?
4. Next, you'll estimate the impact of the treatment on math and Hindi scores at endline.
  - (a) First, write the two regressions you will run to estimate these treatment effects. In words, what will each of these parameters capture?
  - (b) Run the two regressions using the `lm` command and use the `stargazer` command to produce tables containing the results of these two regressions. Interpret your results. What is the effect of the treatment on each of math and Hindi scores? Are these estimated treatment effects statistically significant? Explain.

## Problem 2: Diff-in-Diff

We often study the effect of interventions on those who actually received treatment. However, it is possible that economic interventions can create positive externalities for non-recipient households, for instance, through market interactions or financial transfers between households. In this part of the problem set, we will estimate the spillover effects of a fictitious cash transfer program in a rural Tanzanian setting.

Starting in 2015, imagine that thousands of households in rural northern Tanzania received large cash transfers of roughly US \$1,000 from the international non-governmental organization GiveDirectly (GD). These cash transfers were targeted to relatively poor households within certain villages but not others. The cash grants were unconditional, and the amount is equivalent to more than 50% of total annual income for many recipient households. GD expanded cash transfers in some villages and not others in a non-randomized fashion. Then in 2018, approximately three years after the distribution of the cash transfers, the households were surveyed and measures were collected on many economic, social and life outcomes, including household per capita income.

The goal of this question is to estimate the impact of these transfers on the households deemed to be too wealthy to be eligible for GD transfers. Comparisons will be made across these ineligible households in the "cash" villages (locations where poor households received cash transfers from GD) versus in control villages (where poor household did not receive GD transfers). Any effects on these ineligible households can be considered spillovers since none of them received cash transfers through the program.

1. Discuss the econometric assumptions needed to make a difference-in-differences (DD) approach appropriate in the case of understanding the spillover impacts of cash transfers on household incomes.
2. Why is it important that we have access to data from both the baseline (pre-intervention) and endline (post-intervention) survey rounds?
3. How does the lack of randomization in the allocation of cash transfers across villages affect the estimation of treatment effects, and how might it lead to omitted variable bias?
4. Please download the data to be used in the analysis. The dataset `Ec172_Fall125_PS2Q2.data.csv` is a partial extract of project data (although the actual data has been modified in various ways). Each observation (row) in the dataset represents one household in one time period.

All households in this problem set dataset are **ineligible** for the GD cash transfers (i.e., they are too wealthy to receive cash transfers). There are two observations for each household, one from the baseline survey (time=0) and one from the endline survey (time=1). Using the `lm.cluster` command, determine the average difference at baseline (time=0) between the households in villages where eligible households received cash transfers (`cash=1`) versus the households in control villages that did not receive transfers (`cash=0`) for each of the following two characteristics: household income per capita (`income_pc`) and if the household has an elderly member (`elder=1`). Importantly, these characteristics were collected by the survey team in the baseline survey before the cash transfers were sent out. (Consider these two characteristics one by one, that is, you should run two separate regressions here.) Make sure to account for the correlation among households in the same village by "clustering" your standard errors by village (`village_code`).

**Hint:** To cluster SEs we can use the `lm.cluster` command or the `feols` command. To use the `lm.cluster` command in R, first install the packages `miceadds` and `sandwich`, and then load them into your library. The `lm.cluster` command functions just like `lm`, but with the added option to specify a clustering variable.

For the `feols` command, install and load the `fixest` package. Clustering is specified at the end of the command using `, cluster = village_code`. The syntax of `feols` is otherwise similar to `lm`, except

that you do not include `formula =` before the regression equation.

- (a) Report the regression output for the two regressions mentioned above, and interpret the coefficients. Please also discuss the standard errors and t-statistics. Taken together, are ineligible households similar in cash villages versus control villages along these two dimensions? How does this finding inform the discussion of the validity of the difference-in-differences (DD) approach that you laid out above, if at all?
- (b) Carry out a difference-in-differences (DD) analysis estimating the spillover impact of GD cash transfers among ineligible households, where the outcome of interest is per capita income (`income_pc`). You will need to use data from both time periods. Once again, make sure to account for the correlation among households in the same village by clustering by village (`village_code`).

(Hint: recall that a DD analysis requires the inclusion of separate explanatory variables for post (`time`), treatment (`cash`), and an interaction between the two variables. Note you will need to create a new variable in your dataset for this interaction term.)

What is the spillover impact among ineligible households when poor households in their village receive large cash transfers? Is this effect significantly different from zero at 95% confidence?

- (c) Discuss at least two reasons why cash transfers received by other households in one's village might generate spillover effects over time even among non-recipient households. Given these likely mechanisms, does the magnitude of the estimated effect in part (b) seem plausible? If not, why not?