ECON3389 Econometric Methods

Module 7 Regression Discontinuity Design

Alberto Cappello

Department of Economics, Boston College

Spring 2025

Basic Idea of Causal Inference

- Social science (Economics) theories always ask causal questions.
- In general, a typical causal question is:

The effect of a treatment (D) on an outcome (Y).

- Outcome (Y): A variable that we are interested in.
- Treatment (D): A variable that has the (causal) effect on the outcome of our interest.
- A major problem of estimating the causal effect of treatment is the threat of selection bias.
- In many situations, individuals can select into treatment so those who get treatment could be very different from those who are untreated.
- The best way to deal with this problem is conducting a Randomized Experiment (RCT).

Experimental Idea

- In an RCT, researchers can eliminate selection bias by controlling the treatment assignment process.
- An RCT randomizes:
 - Treatment group: Receives a treatment.
 - Control group: Does not receive the treatment.
- Since we randomly assign treatment, the probability of getting treatment is unrelated to other confounding factors.
- But conducting an RCT is very expensive and may have ethical issues.

Causal Inference

- Instead of controlling the treatment assignment process, if researchers have detailed institutional knowledge of the treatment assignment process:
 - Then we could use this information to create an "experiment."
- Instrumental Variable (IV): Use IVs which are very much alike the endogenous variable but are exogenous (randomized) enough to proxy the treatment and control status.
- Regression Discontinuity Design (RDD): Another widely used method to make causal inference, which is considered more reliable and more robust.

Main Idea of Regression Discontinuity Design

- Regression Discontinuity Design (RDD) exploits the fact that:
 - Some rules are arbitrary and generate a discontinuity in treatment assignment.
- The treatment assignment is determined based on whether a unit exceeds some threshold on a variable (assignment variable, running variable, or forcing variable).
- Assume other factors do **NOT** change abruptly at the threshold.
- Then any change in the outcome of interest can be attributed to the assigned treatment.

A Motivating Example: Elite University

- Numerous studies have shown that graduates from more selective programs or schools earn more than others.
- Example: Students graduating from NJU averagely earn more than those graduating from other ordinary universities like NUFE
- But it is difficult to know whether the positive earnings premium is due to:
 - The true "causal" impact of human capital acquired in the academic program.
 - A spurious correlation linked to the fact that good students selected in these programs would have earned more no matter what (Selection Bias).
- OLS regression will not give us the right answer for the bias. (Because?)

A Motivating Example: Elite University

- But if we could know National College Entrance Exam Scores of all the students, then we can do something.
- Let us say that the entry cutoff for a score of the entrance exam is 100 for NJU.
 - Those with scores 95 or even 99 are unlikely to attend NJU, instead attend NUFE.
- Assume that those scoring 99 or 95 and those scoring 100 are essentially identical. The different scores can be attributed to some random events.
- **RD Strategy:** Compare the long-term outcomes (such as earnings in the labor market) for students scoring 100 (admitted to NJU) and those scoring 99 (admitted to NUFE).

Case Study: SAT Score

- In the United States, most schools used SAT (or ACT) scores in their admission process.
- For example, the flagship state university considered here uses a strict cutoff based on SAT score and high school GPA.
- For the sake of simplicity, we just focuses on the SAT score.
- The author is then able to match (using social security numbers) students applying to the flagship university in 1986-89 to their administrative earnings data for 1998 to 2005.

RD as Local Randomization

- RD provides "local" randomization if the following assumption holds:
 - ullet Agents have imperfect control over the assignment variable X.
- **Intuition:** The randomness guarantees that the potential outcome curves are smooth (e.g., continuous) around the cutoff point.
- There are no discrete jumps in outcomes at the threshold except due to the treatment.
- All observed and unobserved determinants of outcomes are smooth around the cutoff.

RDD and Potential Outcomes: Notations

- Treatment assignment variable (running variable): X_i
- Threshold (cutoff) for treatment assignment: c
- Treatment variable: D_i
- Treatment assignment rule:

$$D_i = \begin{cases} 1 & \text{if } X_i \ge c \\ 0 & \text{if } X_i < c \end{cases}$$

RDD and Potential Outcomes: Notations

Potential Outcomes:

- Potential outcome for an individual i with treatment: Y_{1i}
- Potential outcome for an individual i without treatment: Y_{0i}

Observed Outcomes:

- Y_{1i} if $D_i = 1 \ (X_i \ge c)$
- Y_{0i} if $D_i = 0$ $(X_i < c)$

Sharp RDD and Fuzzy RDD

 In general, depending on enforcement of treatment assignment, RDD can be categorized into two types:

Sharp RDD:

- Nobody below the cutoff gets the "treatment," everybody above the cutoff gets it.
- Everyone follows the treatment assignment rule (all are compliers).
- Local randomized experiment with perfect compliance around the cutoff.

Fuzzy RDD:

- The probability of getting the treatment jumps discontinuously at the cutoff (NOT a jump from 0 to 1).
- Not everyone follows the treatment assignment rule.
- Local randomized experiment with partial compliance around the cutoff.
- Using initial assignment as an instrument for actual treatment.

Identification for Sharp RDD: Continuity Assumption

- $E[Y_{1i}|X_i]$ and $E[Y_{0i}|X_i]$ are continuous at $X_i = c$.
- Assume potential outcomes do not change at the cutoff.
- This means that, except for treatment assignment, all other unobserved determinants of Y_i are continuous at cutoff c.
- This implies no other confounding factor affects outcomes at cutoff c.
- Any observed discontinuity in the outcome can be attributed to treatment assignment.

Continuity Assumption

- Continuity is a natural assumption but could be violated if:
 - There are differences between the individuals who are just below and above the cutoff that are NOT explained by the treatment.
 - 2 The same cutoff is used to assign some other treatment.
 - 3 Other factors also change at the cutoff.
- Individuals can fully manipulate the running variable in order to gain access to the treatment or to avoid it.

Sharp RDD Specification

A simple RD regression is:

$$Y_i = \alpha + \rho D_i + \gamma (X_i - c) + u_i$$

- where
 - Yi: Outcome variable.
 - D_i: Treatment variable (independent variable).
 - X_i: Running variable.
 - c: Value of the cutoff.
 - *u_i*: Error term including other factors.
- Questions: Which parameter do we care about the most?
- Functional form: The validity of RD estimates depends crucially on the function forms, which should provide an adequate representation of $\mathbb{E}[Y_{0i}|X_i]$ and $\mathbb{E}[Y_{1i}|X_i] \to \text{If not what looks like a jump may simply be a non-linear in <math>f(X_i)$ that the polynomials have not accounted for.

Sharp RDD Estimation

- There are two types of strategies for correctly specifying the functional form in an RDD:
 - Parametric/global method:
 - Use all available observations.
 - Estimate treatment effects based on a specific functional form for the relationship between the outcome and the assignment variable.
 - 2 Nonparametric/local method:
 - Use the observations around the cutoff.
 - Compare the outcome of treated and untreated observations that lie within a specific bandwidth.

Parametric Method

• Suppose that, in addition to the assignment mechanism above, potential outcomes can be described by some reasonably smooth function $f(X_i)$:

$$E[Y_{i0}|X_i] = \alpha + f(X_i)$$

- $Y_{i1} = Y_{i0} + \rho$
- Simply, we can construct RD estimates by fitting:

$$Y_i = \alpha + \rho D_i + f(X_i) + u_i$$

Estimation of RDD

More generally, we could estimate two separate regressions for each side respectively:

$$Y_i^b = \beta_b + f(X_i^b - c) + u_i^b$$

$$Y_i^a = \beta_a + g(X_i^a - c) + u_i^a$$

• Continuity Assumption: $f(\cdot)$ and $g(\cdot)$ are any continuous functions of $(X_i^{a,b}-c)$, and satisfy:

$$f(0) = g(0) = 0$$

- We estimate equations using only data above c and only data below c.
- Then, the treatment effect is:

$$\rho = \beta_b - \beta_a$$



Estimation of RDD

• Can do all in one step; just use all the data at once and estimate:

$$Y_i = \alpha + \rho D_i + f(X_i - c) + D_i \times h(X_i - c) + u_i$$

- Where D_i is a dummy variable for treated status.
- When $D_i = 0$, then:

$$Y_i = \alpha + f(X_i - c) + u_i$$

• When $D_i = 1$, let $g(X_i - c) = f(X_i - c) + h(X_i - c)$, then:

$$Y_i = \alpha + \rho + g(X_i - c) + u_i$$

ullet The treatment effect at c is ho



Nonparametric/Local Approach

Recall we can construct RD estimates by fitting:

$$Y_i = \alpha + \rho D_i + f(X_i) + u_i$$

- Nonparametric approach does **NOT** specify a particular functional form for the outcome and the assignment variable, thus $f(X_i)$.
- Instead, it uses only data within a small neighborhood (known as bandwidth) to estimate the discontinuity in outcomes at the cutoff:
 - Compare means in the two bins adjacent to the cutoff (treatment vs. control groups).
 - Use local linear regression (a formal nonparametric regression method).
 - boundary bias: However, comparing means in the two bins adjacent to the cutoff is generally biased in the neighborhood of the cutoff.
 - trade-off in the choice of the bandwidth between bias and precision. Larger bandwidth:
 - Get more precise treatment effect estimates since more data points are used in the regression.
 - But the linear specification is less likely to be accurate and the estimated treatment effect could be biased.



Nonparametric/Local Approach

- The standard solution to reduce boundary bias is to run local linear regression.
- It is a nonparametric method that is a linear smoother within a given bandwidth (window) of width h around the threshold.
- We estimate the following linear regression within a given window of width h around the cutoff:

$$Y_i = \alpha + \rho D_i + \beta_1 \tilde{X}_i + \beta_1^* D_i \tilde{X}_i + u_i$$

- Thus, the **bandwidth** is a key parameter:
 - Cross-Validation Procedure: Choose bandwidth h that produces the best fit for the relationship of outcome and assignment variable.
 - Usually, we would present the RD estimates by different choices of bandwidth.

