



Williams College ECON 523:

Program Evaluation for International Development

**Lecture 8: Regression Discontinuity Designs**

Professor: Pamela Jakiela

## Regression Discontinuity: Intuition

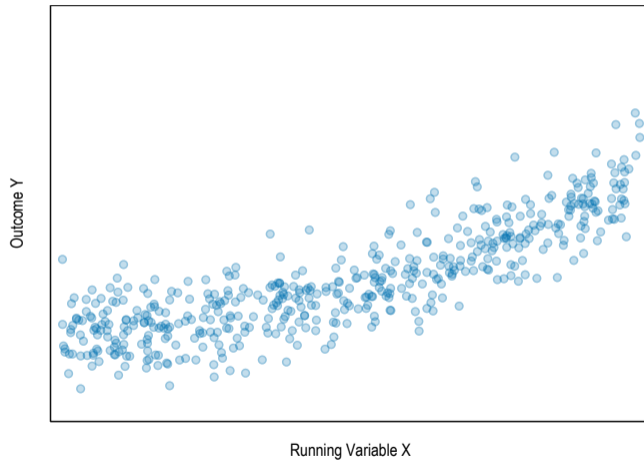
# What Is a Regression Discontinuity?

In a **regression discontinuity** (RD) design: treatment status is determined by a precise rule based on a threshold value of a continuous characteristic (the running or forcing variable)

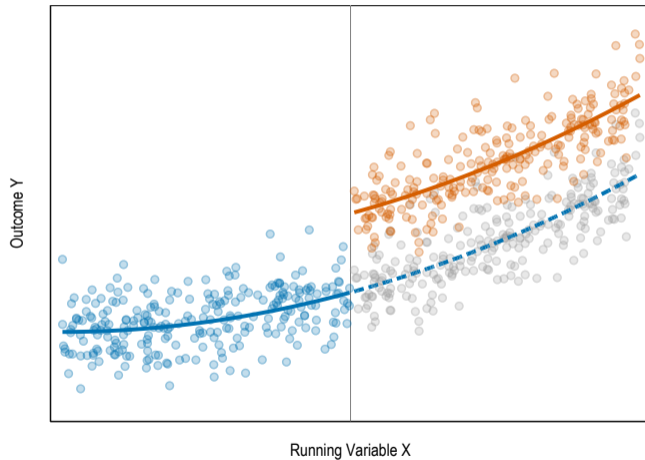
A few examples of potential RDs in development economics:

- **Academic test scores:** admission or scholarships awarded to those with a certain score (e.g. above the threshold for passing), or based on a precise percentile cutoff
- **Household income:** earned-income tax credit for those with incomes below cutoff
- **Poverty indices:** anti-poverty programs are often “proxy-means-tested”
- **Landholdings:** microfinance available to households with less than an acre of land
- **Date of birth:** eligibility for government programs such as preschool/kindergarten, medicare, vaccines, etc. often depends on whether you are above or below an age cutoff
- **Elections:** candidate who gets most (or more than fifty percent of the) votes wins

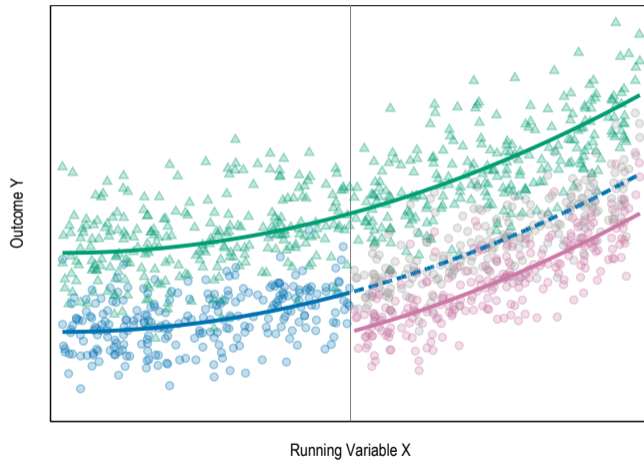
# $Y_0$ Is a Continuous Function of Running Variable $X$



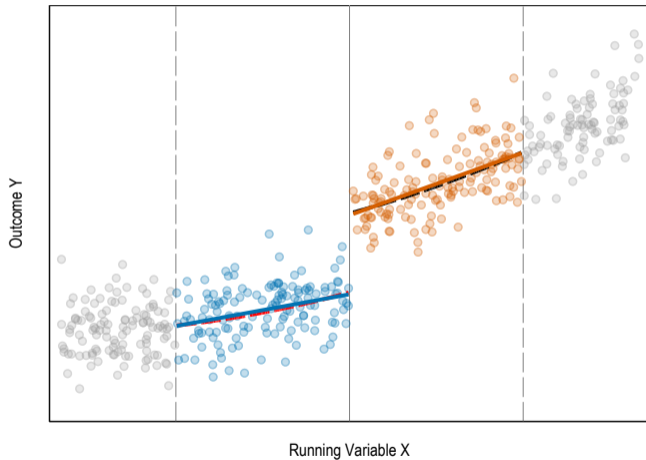
# A Jump in $f(X)$ at the Discontinuity Indicates a Treatment Effect



# Analogous to Diff-in-Diff, But Without a Comparison Group



# With a Narrow Bandwidth, $f(X)$ Should Be Approximately Linear



# When Can You Use RD?

RD designs are possible when eligibility is determined by some smooth continuous index, but participants don't have precise control over the index (so they can't select into treatment)

- Individuals with low vs. high values of the eligibility index may look very different, but differences should be quite small in the neighborhood right around the discontinuity
- Individuals might be able to influence the likelihood of being **near** the discontinuity, but those near the discontinuity can't make sure that they end up on one side vs. the other
  - ▶ Near discontinuity, treatment is “as good as random”

As in a randomized trial, you either have a valid discontinuity or you don't

# When Can You Use RD?

For a valid RD, the following must be true:

1. Treatment status determined by a continuous variable (treated  $\Leftrightarrow$  if above/below cutoff)
  - ▶ Eligibility index must be smooth, take many values around cutoff
  - ▶ Estimated treatment effect is limit as you approach cutoff
2. Individuals/units/etc. cannot **precisely** manipulate eligibility index
  - ▶ Precise manipulation  $\Rightarrow$  selection bias
3. No other programs/treatments can use the same discontinuity

## Implementing RD

# RD Specification

RD setup:

- Let  $X_i$  denote the running variable
- Observations with  $X_i$  above cutoff  $C$  are treated, those below are not
- Define a re-centered running variable  $Z_i = X_i - C$ 
  - ▶ Let  $T_i$  be an indicator for having  $Z_i \geq 0$ , and hence being treated

The most common empirical specification for a regression discontinuity design is:

$$Y_i = \alpha + \beta T_i + \theta_1 Z_i + \theta_2 (Z_i \times T_i) + \varepsilon_i$$

# RD Specification

RD setup:

- Let  $X_i$  denote the running variable
- Observations with  $X_i$  above cutoff  $C$  are treated, those below are not
- Define a re-centered running variable  $Z_i = X_i - C$ 
  - ▶ Let  $T_i$  be an indicator for having  $Z_i \geq 0$ , and hence being treated

The most common empirical specification for a regression discontinuity design is:

$$Y_i = \alpha + \beta T_i + \theta_1 Z_i + \theta_2 (Z_i \times T_i) + \varepsilon_i$$

Notice that we allow for separate slopes above and below the cutoff

## Alternative RD Specifications

Quadratic RD specification includes linear, squared distance from the cutoff (above and below)

$$Y_i = \alpha + \beta T_i + \theta_1 Z_i + \theta_2 (Z_i \times T_i) + \lambda_1 Z_i^2 + \lambda_2 (Z_i \times T_i)^2 + \varepsilon_i$$

Adding additional polynomial terms allows for better fit of  $Y = f(X)$

- Including higher order polynomial terms also leads to over-fitting, so we usually don't
- All continuous, differentiable functions are locally linear

Do make sure you allow  $f(X)$  to differ above vs. below the discontinuity

# Bandwidth Selection

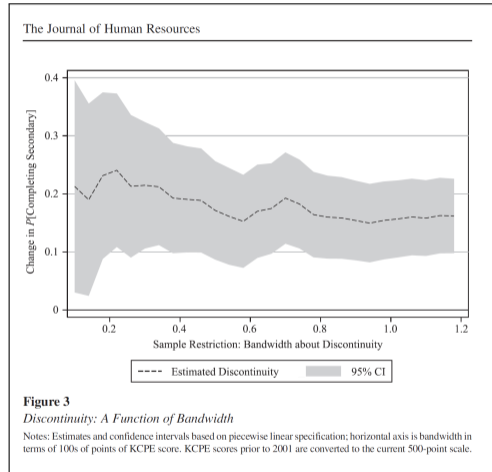
Bandwidth selection implies a sample restriction

- Example: `reg y treatment z_below_cutoff z_above_cutoff if abs(z) ≤ bw, r`
- What is `bw`? How wide of a window around the cutoff is included in the analysis?
  - ▶ Bias:  $f(X)$  is locally linear, but (probably) not linear; forcing linearity introduces bias
  - ▶ Variance: a narrow bandwidth means a small(er) sample, and hence a noisier estimate

Traditionally, most common approach to bandwidth selection is ad hoc: what looks good?

- Choose a range of data small enough such that relationship **looks** linear (in a graph)
- Compare results as you vary the bandwidth, from too small to too large

# Bandwidth Selection: Examples



Source: Ozier (2018)

# Bandwidth Selection: Examples

TABLE III  
ALTERNATIVE RD SPECIFICATIONS<sup>a</sup>

	Bandwidth				
	1 (1)	0.5 (2)	0.25 (3)	0.1 (4)	0.05 (5)
Panel A: Women					
<i>Polynomial order of control function</i>					
None	0.012** (0.006)	0.015** (0.006)	0.018*** (0.006)	0.025*** (0.007)	0.018* (0.010)
Linear	0.014** (0.007)	0.021*** (0.006)	0.025*** (0.007)	0.028** (0.012)	0.039** (0.019)
Quadratic	0.027*** (0.007)	0.030*** (0.007)	0.033*** (0.010)	0.032* (0.018)	0.051 (0.032)
Cubic	0.031*** (0.007)	0.026*** (0.010)	0.036** (0.015)	0.057** (0.028)	0.054 (0.042)
Quartic	0.030*** (0.009)	0.032** (0.012)	0.044** (0.017)	0.067** (0.033)	0.028 (0.056)
Observations	2628	2177	1049	489	257

Source: Myerston (2014)

# Optimal Bandwidth Selection

More recent approaches aim to optimize bias-variance tradeoff

- Imbens and Kalyanaraman (2012): bandwidth used in Meyersson (2014)
- Calonico, Catteneo, and Titiunik (2014)
- Calonico, Catteneo, Farrell, and Titiunik (2017): `rdrobust` command in Stata/R

# There's Always a Graph

Common elements of regression discontinuity graphs:

- The running variable is on the  $x$ -axis, and the cutoff is indicated
- Observations are binned and bin-level averages are shown in a scatter plot (**bin scatter**)
  - ▶ Bins either contain fixed  $N$  or fixed range of  $X$  values
- The continuous relationship between the running variable and the outcome is shown, always with separate graphs for above vs. below the cutoff for treatment eligibility
  - ▶ Sometimes a linear fit and sometimes a local-polynomial regression (1poly)
  - ▶ Sometimes with confidence intervals, and sometimes without

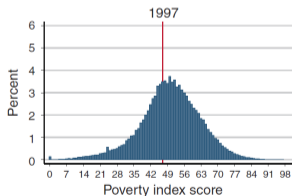
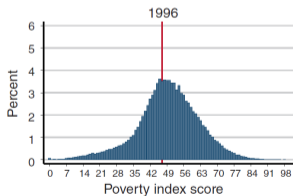
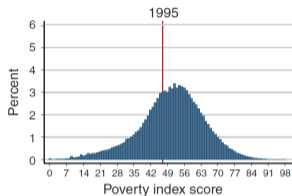
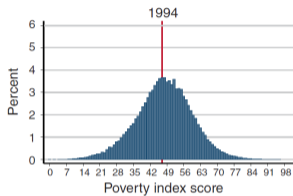
# Fuzzy RD

In some cases, being above the discontinuity increases treatment, but not from 0 to 1

- An RD with imperfect compliance is a **fuzzy RD**
- Being above the cutoff is a valid instrument for treatment (controlling for  $X$ )
  - ▶ First stage: RD regression of treatment on cutoff, slopes above and below
  - ▶ Second stage: RD regression of outcome on predicted treatment, again with slope controls

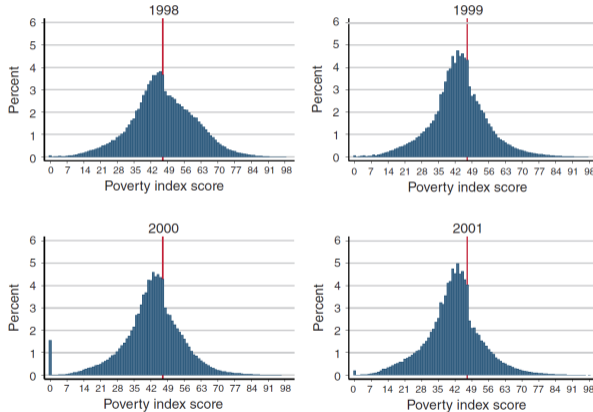
## Assessing the Validity of an RD

# Q: When Is an RD Not an RD?



Source: Conover and Comacho (2011)

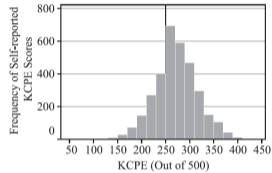
# A: After Potential Participants Learn the Allocation Rule



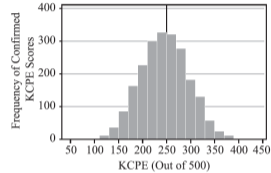
Source: Conover and Comacho (2011)

# A: Or, Potentially, When Scores Are Self-Reported

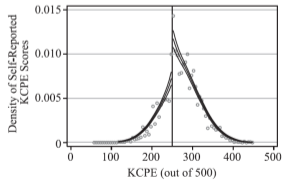
Panel A: KLPS2 Data,  $N=3305$



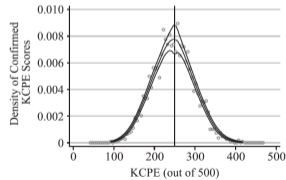
Panel B: KLPS2 Data,  $N=2167$ , Restricted to Confirmed First KCPE Scores



Panel C: Density Discontinuity  $p<0.001$



Panel D: Density Discontinuity  $p=0.953$



Source: Ozier (2018)

# Defending the Validity of an RD Design

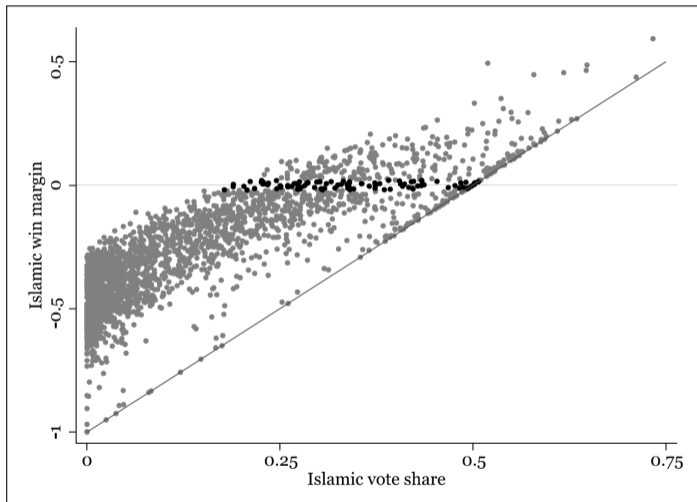
The assumption of no sorting/manipulation can't be tested directly

- Plot the distribution (McCrary (2007) also proposes a test)
- Plot the covariates: pre-treatment characteristics shouldn't jump with treatment

Two approaches empirical approaches:

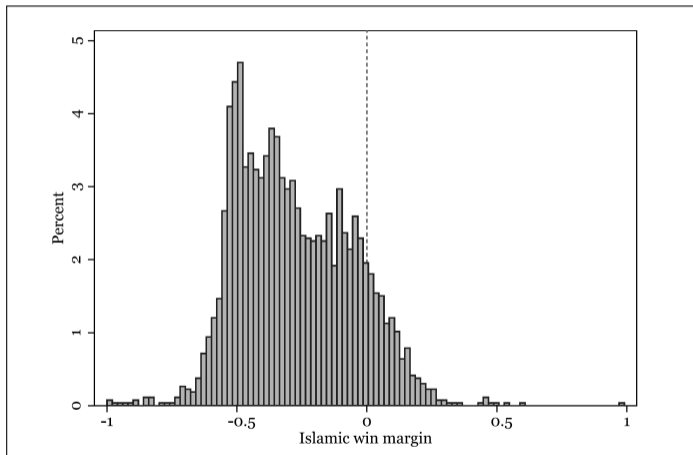
- Histogram (of either running variable or covariates)
- Local polynomial regression above vs. below (using `lpol` command in Stata)

# Checking for Manipulation



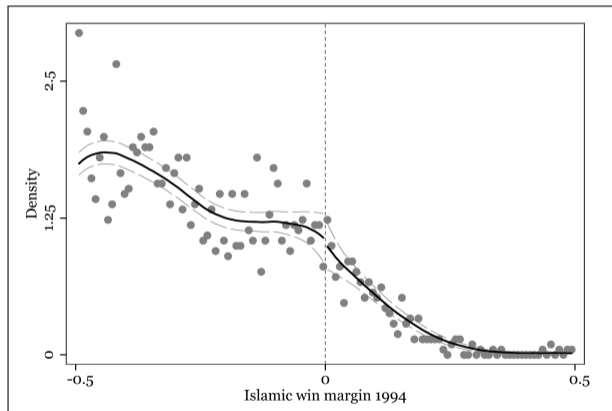
Source: Meyerson (2014)

# Checking for Manipulation: Histogram



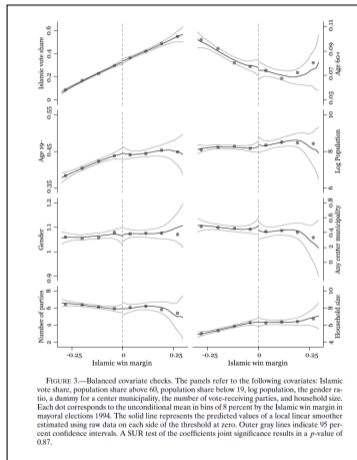
Source: Meyerson (2014)

# Checking for Manipulation: Density Plot



Source: Meyerson (2014)

# Checking Covariate Balance



Source: Meyerson (2014)

# RD Checklist

1. Do you have an RD?
  - 1.1 Program eligibility determined by continuous running variable
  - 1.2 Many observations close to the discontinuity
  - 1.3 No manipulation
2. Bin scatter showing relationship between running variable and treatment probability
3. Show that running variable and covariate densities look smooth at cutoff
4. Choose your bandwidth and run your regressions
  - 4.1 Report 2SLS for fuzzy RD designs
  - 4.2 Show alternative bandwidths, higher-order polynomials as robustness checks