

Lecture 8: Regression Discontinuity, 2 of 2

March 5, 2025

Course Administration

1. Problem set 3 should be in
2. Quantitative summary due March 19
3. April 9 workshop instructions posted
4. Sign up for consultations
 - March 19 and 20
 - in lieu of class lecture 10
 - see link in lecture 10
 - let me know if more spots needed

Course Administration

1. Problem set 3 should be in
2. Quantitative summary due March 19
3. April 9 workshop instructions posted
4. Sign up for consultations
 - March 19 and 20
 - in lieu of class lecture 10
 - see link in lecture 10
 - let me know if more spots needed
5. Presentation dates assigned
 - April 16 and 23
 - see lectures tab for names
6. Please come see me about your replication paper
7. Any other issues?

Lecture 8: RD and RDK

Background on RD and RDK

1. RD Recap
2. Fuzzy RD
3. Regression Kink
4. How-to, Again

Manoli and Turner

1. Research question, endogeneity and data
2. Discontinuity, estimating equations, validity
3. Results

RD Recap

What do you need for a RD design?

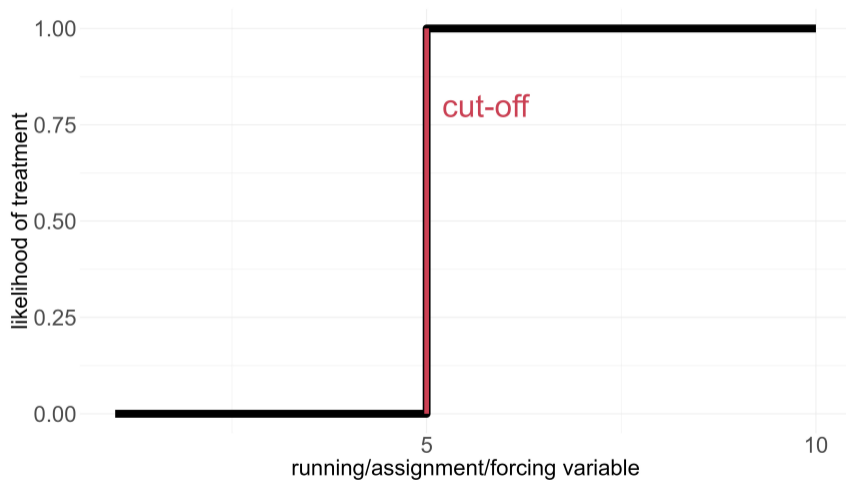
What do you need for a RD design?

- You need a discontinuity!

What do you need for a RD design?

- You need a discontinuity!
- Where you have treatment on one side and no (or little) treatment on the other

Sharp Discontinuity



We usually don't graph sharp discontinuities because they are boring!

We Can Analogize RD to an Experiment – Why?

We Can Analogize RD to an Experiment – Why?

Because

- close to the discontinuity
- treated and untreated units
- should be observationally equivalent
- pre-treatment

What kind of estimating equation do we use for a RD?

$$Y_i = \alpha_0 + \alpha_1 D + \alpha_2 f(X - c) + \alpha_3 D * f(X - c) + \alpha_4 Q + \epsilon$$

What kind of estimating equation do we use for a RD?

$$Y_i = \alpha_0 + \alpha_1 D + \alpha_2 f(X - c) + \alpha_3 D * f(X - c) + \alpha_4 Q + \epsilon$$

Suppose we use X such that $\underline{X} < X < \bar{X}$. How can we parameterize $f(X - c)$?

What kind of estimating equation do we use for a RD?

$$Y_i = \alpha_0 + \alpha_1 D + \alpha_2 f(X - c) + \alpha_3 D * f(X - c) + \alpha_4 Q + \epsilon$$

Suppose we use X such that $\underline{X} < X < \bar{X}$. How can we parameterize $f(X - c)$?

$$Y_i = \alpha_0 + \alpha_1 D + \alpha_2 (X - c) + \alpha_3 D * (X - c) + \alpha_4 Q + \epsilon \quad (1)$$

What kind of estimating equation do we use for a RD?

$$Y_i = \alpha_0 + \alpha_1 D + \alpha_2 f(X - c) + \alpha_3 D * f(X - c) + \alpha_4 Q + \epsilon$$

Suppose we use X such that $\underline{X} < X < \bar{X}$. How can we parameterize $f(X - c)$?

$$Y_i = \alpha_0 + \alpha_1 D + \alpha_2 (X - c) + \alpha_3 D * (X - c) + \alpha_4 Q + \epsilon \quad (1)$$

or

$$Y_i = \alpha_0 + \alpha_1 D + \alpha_2 (X - c) + \alpha_3 (X - c)^2 + \alpha_4 D * (X - c) + \alpha_5 D * (X - c)^2 + \alpha_6 Q + \epsilon \quad (2)$$

What kind of estimating equation do we use for a RD?

$$Y_i = \alpha_0 + \alpha_1 D + \alpha_2 f(X - c) + \alpha_3 D * f(X - c) + \alpha_4 Q + \epsilon$$

Suppose we use X such that $\underline{X} < X < \bar{X}$. How can we parameterize $f(X - c)$?

$$Y_i = \alpha_0 + \alpha_1 D + \alpha_2 (X - c) + \alpha_3 D * (X - c) + \alpha_4 Q + \epsilon \quad (1)$$

or

$$Y_i = \alpha_0 + \alpha_1 D + \alpha_2 (X - c) + \alpha_3 (X - c)^2 + \alpha_4 D * (X - c) + \alpha_5 D * (X - c)^2 + \alpha_6 Q + \epsilon \quad (2)$$

Or, make \underline{X} bigger and \bar{X} smaller so that the window around c is narrower.

Specification Matters Because Data Can Be Odd

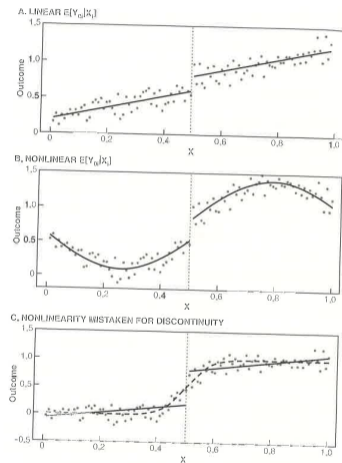


Figure 6.1.1 The sharp regression discontinuity design.

Source: Angrist and Pischke, *Mostly Harmless Econometrics*.

What are the testable implications for the validity of a RD design?

What are the testable implications for the validity of a RD design?

- Discontinuity of treatment
- Dontinuity of pre-determined observables
- No bunching of observations at discontinuity

RD: When it's Fuzzy

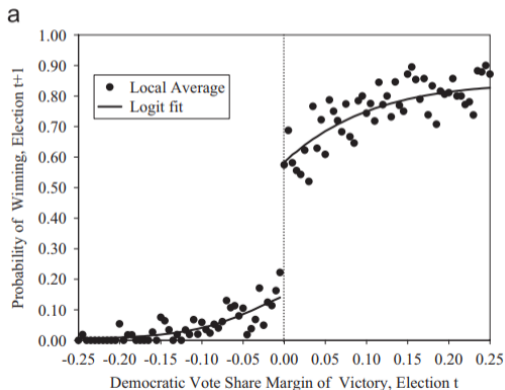
Sharp vs Fuzzy

- Sharp RD: perfect compliance with treatment at the cut-off
- Fuzzy RD: a higher likelihood of compliance with treatment at the cut-off

Sharp Discontinuity Outcome Examples

RQ: Impact of Incumbency on election at time $t + 1$

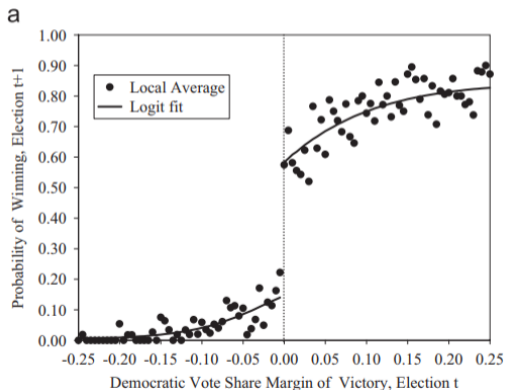
Treatment: Winning election at time t



Sharp Discontinuity Outcome Examples

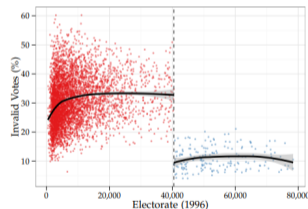
RQ: Impact of Incumbency on election at time $t + 1$

Treatment: Winning election at time t



RQ: What is the impact of electronic voting on invalid votes?

Treatment: Electronic voting, electorate $> 40,000$



In Brazil, no paper ballots with electorate $> 40,000$

Almond and 1917 Flu

A Fuzzy Treatment

- RQ: How does maternal health impact child's long-term health?
- Use random variation over time in maternal health

Almond and 1917 Flu

A Fuzzy Treatment

- RQ: How does maternal health impact child's long-term health?
- Use random variation over time in maternal health
- Treating is getting 1917 Spanish flu
- Why is this fuzzy?

Almond and 1917 Flu

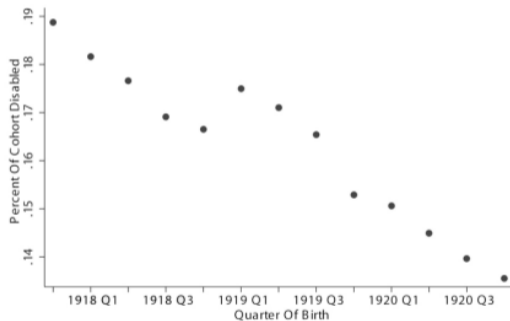
A Fuzzy Treatment

- RQ: How does maternal health impact child's long-term health?
- Use random variation over time in maternal health
- Treating is getting 1917 Spanish flu
- Why is this fuzzy?
- Hits kids born 1919Q1 to 1919Q4
- What do you expect for health outcome?

Almond and 1917 Flu

A Fuzzy Treatment

- RQ: How does maternal health impact child's long-term health?
- Use random variation over time in maternal health
- Treating is getting 1917 Spanish flu
- Why is this fuzzy?
- Hits kids born 1919Q1 to 1919Q4
- What do you expect for health outcome?



MHE Says That Fuzzy RD is IV

Recall two IV conditions

MHE Says That Fuzzy RD is IV

Recall two IV conditions

1. Z is correlated with endogenous variable

MHE Says That Fuzzy RD is IV

Recall two IV conditions

1. Z is correlated with endogenous variable
 - In RD case, this means that discontinuity is correlated with treatment
 - True by definition.

MHE Says That Fuzzy RD is IV

Recall two IV conditions

1. Z is correlated with endogenous variable
 - In RD case, this means that discontinuity is correlated with treatment
 - True by definition.
2. Z is uncorrelated with error – equivalent to saying that instrument affects outcome only through relationship with endogenous variable

MHE Says That Fuzzy RD is IV

Recall two IV conditions

1. Z is correlated with endogenous variable
 - In RD case, this means that discontinuity is correlated with treatment
 - True by definition.
2. Z is uncorrelated with error – equivalent to saying that instrument affects outcome only through relationship with endogenous variable
 - In RD case, this means that discontinuity matters to outcome only through relationship with treatment
 - May be plausible

MHE Says That Fuzzy RD is IV

Recall two IV conditions

1. Z is correlated with endogenous variable
 - In RD case, this means that discontinuity is correlated with treatment
 - True by definition.
2. Z is uncorrelated with error – equivalent to saying that instrument affects outcome only through relationship with endogenous variable
 - In RD case, this means that discontinuity matters to outcome only through relationship with treatment
 - May be plausible

So we can then think about fuzzy RD as a LATE, where the effect we estimate is for compliers.

RD: Regression Kink

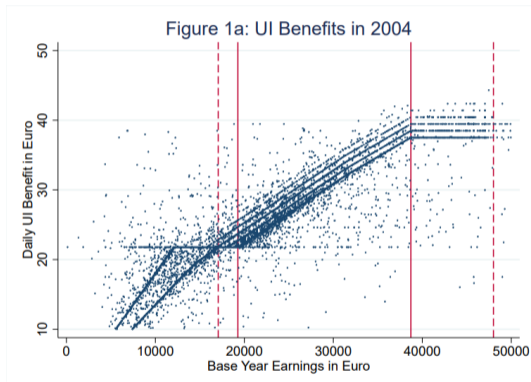
Basic Intuition

- RD uses a break in likelihood of treatment at cut-off
- RK uses a shift in likelihood of treatment at cut-off
 - can be sharp or fuzzy

Basic Intuition

- RD uses a break in likelihood of treatment at cut-off
- RK uses a shift in likelihood of treatment at cut-off
 - can be sharp or fuzzy
- As before, we are interested in change in Y given change in X , where the variation in X comes from the variation at the kink
- RK identification requirements from Card et al (Econometrica, 2015; pictures are from working paper version)
 1. “conditional on the unobservable determinants of the outcome variable, the density of the assignment variable is smooth (i.e., continuously differentiable) at the kink point in the policy rule” and
 2. “the treatment assignment rule is continuous at the kink point”
- Second one affects interpretation; not strictly necessary

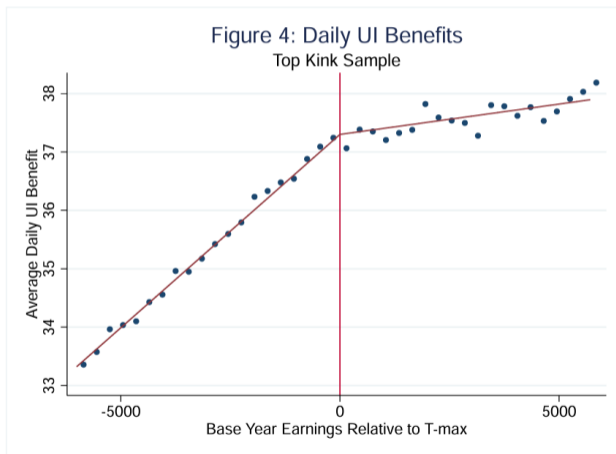
A Treatment Example



Source: Card et al (2015) in *Econometrica*; but all figures are from NBER working paper no. 18564.

- Figure 1 is unemployment insurance receipt by income in Austria
- Lines are not drawn – those are perfect compliers with the schedule
- Authors suggest there are various errors that lead them not to get everyone on the line
- Five lines in the middle portion of the figure are for number of dependents
- We'll focus on top kink

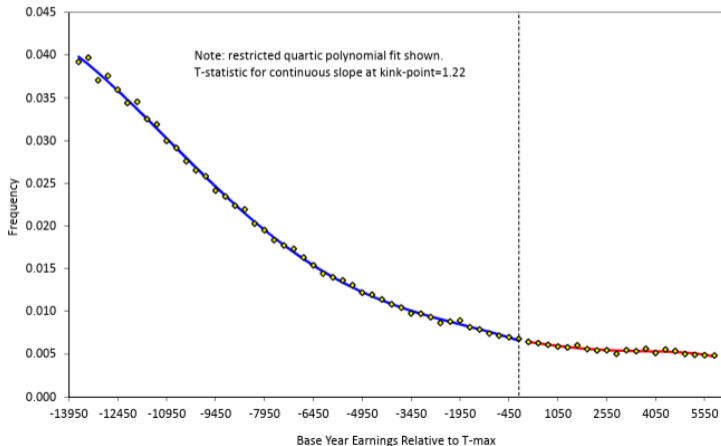
Figure 4 Zooms in on Parts of Figure 1(a) for Single Workers



Now formulated as $(X - c)$

Can You Precisely Manipulate the Running Variable?

Figure 2b: Density in Top Kink Sample



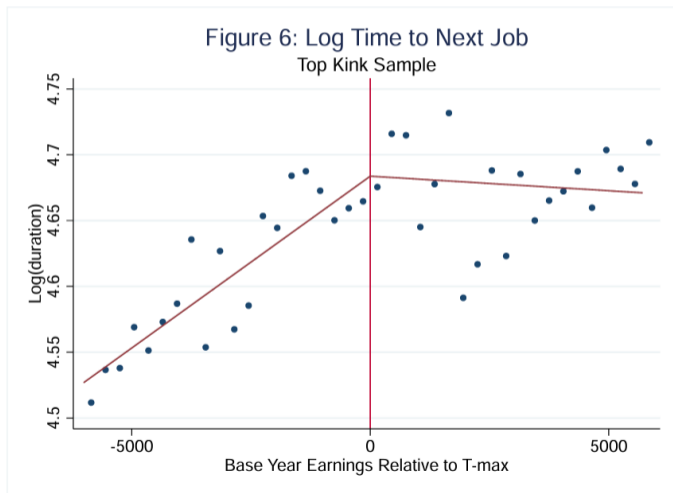
- Vertical axis is concentration of observations
- No bunching in the distribution of earnings at the kink points
- → cannot manipulate distance to the running variable

Outcomes: How Long Does it Take to Find a New Job?



- Figure 6 is the outcome
- How long does it take you to find a new job?

Outcomes: How Long Does it Take to Find a New Job?



- Figure 6 is the outcome
- How long does it take you to find a new job?
- As benefits plateau, people finding jobs more quickly

Formal statement of RK (from Card, Lee, Pei and Weber, 2012)

- Suppose we have a model $Y = \tau b(V) + g(V) + \epsilon$, where the kink is at $V = v_0$
- $b(V)$ is the assignment function: how much UI you get
- $g(V)$ is how the assignment variable otherwise impacts the outcome Y
- Assume
 1. “ $b(V)$ [is] a deterministic and continuous function of V , with a kink and $V = 0$ ”
 2. “ $g(\cdot)$ and $E(\epsilon|V = v)$ have derivatives that are continuous in V at $V = 0$ ”
- Then

$$\tau = \frac{\lim_{v_0 \rightarrow 0^+} \frac{dE(Y|V=v)}{dv} | V = v_0 - \lim_{v_0 \rightarrow 0^-} \frac{dE(Y|V=v)}{dv} | V = v_0}{\lim_{v_0 \rightarrow 0^+} b'(v_0) - \lim_{v_0 \rightarrow 0^-} b'(v_0)} \approx \frac{\Delta Y}{\Delta b(V)}$$

RD: How-to, Redux

How-to Steps

1. Find a discontinuity that's credible
2. Make a graph
 - treatment variable
 - outcome variable
 - number of observations
3. Do a RD regression
4. Tests for validity

Manoli and Turner: Income and College Attendance

Manoli and Turner on Income and College Attendance

1. Research question, endogeneity and data
2. Discontinuity, estimating equations, validity
3. Results

Research Question and Endogeneity

Research question

Research Question and Endogeneity

Research question

- How does family income affect college attendance?

Research Question and Endogeneity

Research question

- How does family income affect college attendance?

Endogeneity concerns

- Suppose we wanted to just do a simple OLS to answer the question that Manoli and Turner are interested in. What would we estimate?

$$\text{attendance}_i = \beta_0 + \beta_1 \text{income}_i + \beta_2 X_i + \epsilon_i$$

- Why is this a bad idea?

Research Question and Endogeneity

Research question

- How does family income affect college attendance?

Endogeneity concerns

- Suppose we wanted to just do a simple OLS to answer the question that Manoli and Turner are interested in. What would we estimate?

$$\text{attendance}_i = \beta_0 + \beta_1 \text{income}_i + \beta_2 X_i + \epsilon_i$$

- Why is this a bad idea?
- Family income correlated with lots of other things that affect college going: parental education, education habits, etc.

Manoli and Turner: Data

- What is the unit of observation?

Manoli and Turner: Data

- What is the unit of observation?
 - Person in a year

Manoli and Turner: Data

- What is the unit of observation?
 - Person in a year
- Data?

Manoli and Turner: Data

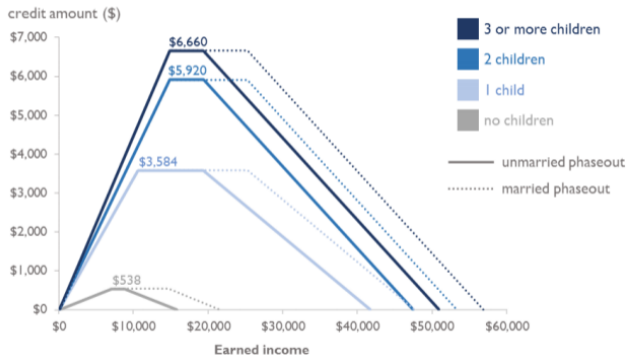
- What is the unit of observation?
 - Person in a year
- Data?
 - Tax data, but it seems like people are defined from social security records

Manoli and Turner: Data

- What is the unit of observation?
 - Person in a year
- Data?
 - Tax data, but it seems like people are defined from social security records
- There are ~ 4 million high school seniors/year in the US, so this is a selected sample
- See Appendix Table 4 for sample selection

Manoli and Turner: EITC

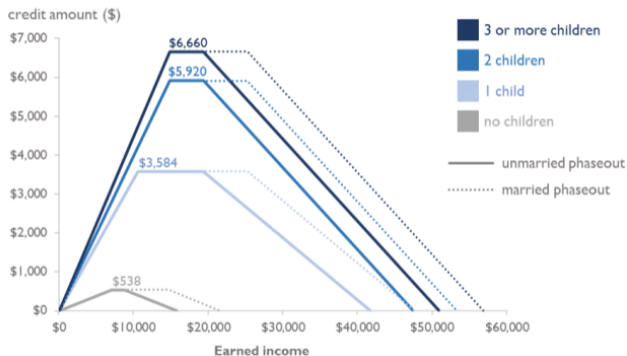
EITC Amount by Number of Qualifying Children, Marital Status, and Income, 2020



- What's the idea behind this program?

Manoli and Turner: EITC

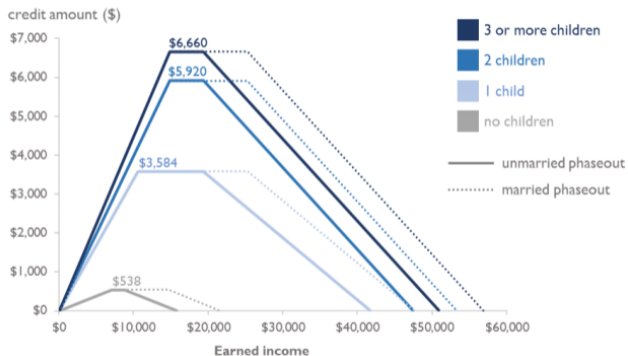
EITC Amount by Number of Qualifying Children, Marital Status, and Income, 2020



- What's the idea behind this program?
- Dates from Nixon Administration, and has been scaled up over time
- What does the pyramid shape of the benefit mean?

Manoli and Turner: EITC

EITC Amount by Number of Qualifying Children, Marital Status, and Income, 2020

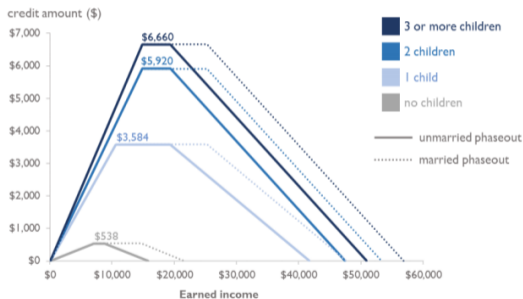


- What's the idea behind this program?
- Dates from Nixon Administration, and has been scaled up over time
- What does the pyramid shape of the benefit mean?
- Where might you want to lie about your income if you could? why?

Potential Kinks

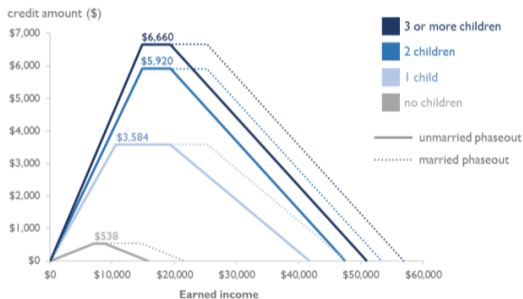
- Where are the potential kinks in this setup?

ETC Amount by Number of Qualifying Children, Marital Status, and Income, 2020



Potential Kinks

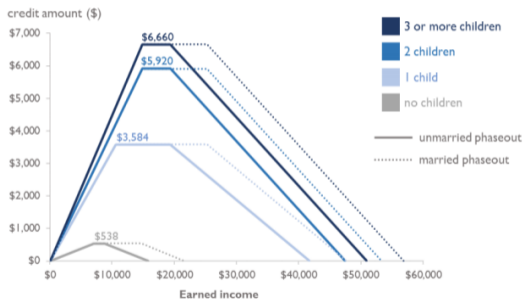
EITC Amount by Number of Qualifying Children, Marital Status, and Income, 2020



- Where are the potential kinks in this setup?
- Why is the first kink more attractive than the second one?

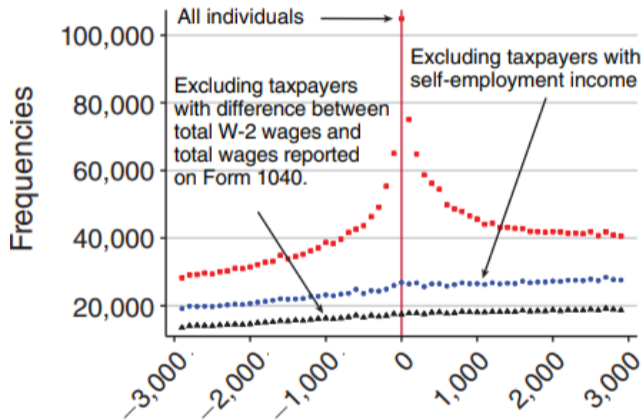
Potential Kinks

EITC Amount by Number of Qualifying Children, Marital Status, and Income, 2020



- Where are the potential kinks in this setup?
- Why is the first kink more attractive than the second one?
 - fn 16, p. 9 says “ We have explored using a regression kink design at the second and third EITC kink points using a restricted sample that has earned income equal to AGI. ... Furthermore, when restricting to this unusual sample, we find small but statistically significant evidence of sorting along the running variable (i.e. bunching at the kink point)”
 - What does “bunching” mean here?

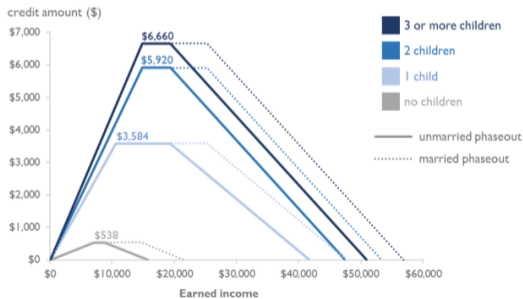
A Bad Kink for Some People



The Slope Comparison

On which kink in this picture are we focusing?

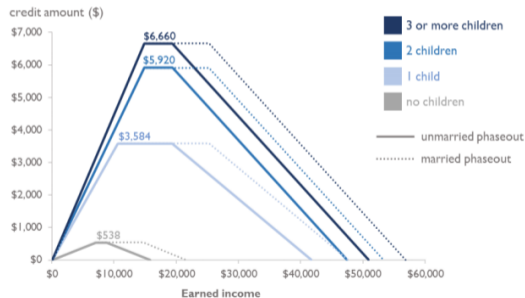
EITC Amount by Number of Qualifying Children, Marital Status, and Income, 2020



The Slope Comparison

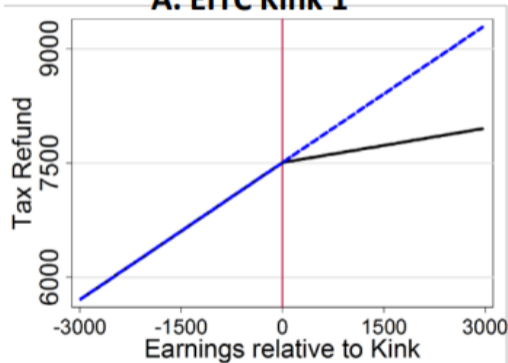
On which kink in this picture are we focusing?

EITC Amount by Number of Qualifying Children, Marital Status, and Income, 2020



Zooming in on the first kink

A. EITC Kink 1



Estimation

$$\text{enroll}_i = \beta \text{refund}_i + g(\text{kinkdist}_i) + \epsilon_i$$

- enroll_i is 0 or 1, depending on whether you enroll in higher education in the fall following the spring refund
- refund_i is the tax refund (not just EITC)
- kinkdist_i is income distance to the kink
- The idea is to find the $\hat{\beta} = \frac{\Delta \text{enrollment}}{\Delta \text{refund}}$

Estimation, Again

But more exactly, they estimate

$$\begin{aligned} \text{enroll}_i &= \alpha \text{kinkdist}_i + \delta_e D_i \text{kinkdist}_i + \alpha_2 X_i + \epsilon_i \\ \text{refund}_i &= \gamma \text{kinkdist}_i + \delta_r D_i \text{kinkdist}_i + \alpha_2 X_i + \epsilon_i \end{aligned}$$

- $D_i = 1$ after the kink
- Note that

$$\hat{\beta} = \frac{\hat{\delta}_e}{\hat{\delta}_r}$$

- Alternatively, we could instrument for refund_i in the below with $D_i \text{kinkdist}_i$:

$$\text{enroll}_i = \beta \text{refund}_i + \text{kinkdist}_i + \epsilon$$

Estimation Intuition

$$\beta = \frac{\text{enrollment rate of change, post cut-off} - \text{enrollment rate of change, pre cut-off}}{\text{refund rate of change, post cut-off} - \text{refund rate of change, pre cut-off}}$$

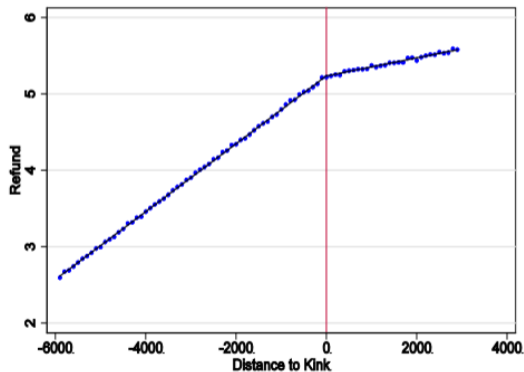
Estimation Intuition

$$\beta = \frac{\text{enrollment rate of change, post cut-off} - \text{enrollment rate of change, pre cut-off}}{\text{refund rate of change, post cut-off} - \text{refund rate of change, pre cut-off}}$$

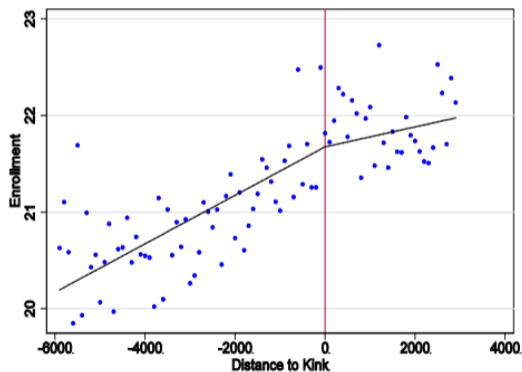
In a sharp RD, analogous denominator is 1!

Results In Pictures

A. Change in Tax Refund



B. Change in Enrollment



Results in Table

- How do we interpret the top two coefficients?

Table 2: RKD Estimates

	Full Sample		1 Child	
	First Stage Refund	Reduced Form Enrollment	First Stage Refund	Reduced Form Enrollment
Slope Change	-0.343 (0.002)	-0.147 (0.060)	-0.337 (0.001)	-0.152 (0.107)
Effect of \$1000 on Enrollment (IV)		0.430 (0.175)		0.450 (0.317)
N	1,427,447		465,745	

Notes: Each coefficient is estimated from a separate regression. Each regression includes dummy variables for filing status. Standard errors are clustered based on \$100 bins of earnings relative to the kink point.

Results in Table

- How do we interpret the top two coefficients?
 - An extra \$1,000 dollars yields a -0.343 change in the slope of the refund

Table 2: RKD Estimates

	Full Sample		1 Child	
	First Stage Refund	Reduced Form Enrollment	First Stage Refund	Reduced Form Enrollment
Slope Change	-0.343 (0.002)	-0.147 (0.060)	-0.337 (0.001)	-0.152 (0.107)
Effect of \$1000 on Enrollment (IV)		0.430 (0.175)		0.450 (0.317)
N	1,427,447		465,745	

Notes: Each coefficient is estimated from a separate regression. Each regression includes dummy variables for filing status. Standard errors are clustered based on \$100 bins of earnings relative to the kink point.

Results in Table

- How do we interpret the top two coefficients?
 - An extra \$1,000 dollars yields a -0.343 change in the slope of the refund
 - An -0.15 slope change in enrollment
- Is enrollment lower after the kink?

Table 2: RKD Estimates

	Full Sample		1 Child	
	First Stage Refund	Reduced Form Enrollment	First Stage Refund	Reduced Form Enrollment
Slope Change	-0.343 (0.002)	-0.147 (0.060)	-0.337 (0.001)	-0.152 (0.107)
Effect of \$1000 on Enrollment (IV)		0.430 (0.175)		0.450 (0.317)
N	1,427,447		465,745	

Notes: Each coefficient is estimated from a separate regression. Each regression includes dummy variables for filing status. Standard errors are clustered based on \$100 bins of earnings relative to the kink point.

Results in Table

- How do we interpret the top two coefficients?
 - An extra \$1,000 dollars yields a -0.343 change in the slope of the refund
 - An -0.15 slope change in enrollment
- Is enrollment lower after the kink? no!

Table 2: RKD Estimates

	Full Sample		1 Child	
	First Stage Refund	Reduced Form Enrollment	First Stage Refund	Reduced Form Enrollment
Slope Change	-0.343 (0.002)	-0.147 (0.060)	-0.337 (0.001)	-0.152 (0.107)
Effect of \$1000 on Enrollment (IV)		0.430 (0.175)		0.450 (0.317)
N	1,427,447		465,745	

Notes: Each coefficient is estimated from a separate regression. Each regression includes dummy variables for filing status. Standard errors are clustered based on \$100 bins of earnings relative to the kink point.

Results in Table

Table 2: RKD Estimates

	Full Sample		1 Child	
	First Stage Refund	Reduced Form Enrollment	First Stage Refund	Reduced Form Enrollment
Slope Change	-0.343 (0.002)	-0.147 (0.060)	-0.337 (0.001)	-0.152 (0.107)
Effect of \$1000 on Enrollment (IV)		0.430 (0.175)		0.450 (0.317)
N	1,427,447		465,745	

Notes: Each coefficient is estimated from a separate regression. Each regression includes dummy variables for filing status. Standard errors are clustered based on \$100 bins of earnings relative to the kink point.

- How do we interpret the top two coefficients?
 - An extra \$1,000 dollars yields a -0.343 change in the slope of the refund
 - An -0.15 slope change in enrollment
- Is enrollment lower after the kink? no!
- Remember this is a change in the rate of change
- IV result is 0.430 (about -0.147/-0.343)

Results in Table

Table 2: RKD Estimates

	Full Sample		1 Child	
	First Stage Refund	Reduced Form Enrollment	First Stage Refund	Reduced Form Enrollment
Slope Change	-0.343 (0.002)	-0.147 (0.060)	-0.337 (0.001)	-0.152 (0.107)
Effect of \$1000 on Enrollment (IV)		0.430 (0.175)		0.450 (0.317)
N	1,427,447		465,745	

Notes: Each coefficient is estimated from a separate regression. Each regression includes dummy variables for filing status. Standard errors are clustered based on \$100 bins of earnings relative to the kink point.

- How do we interpret the top two coefficients?
 - An extra \$1,000 dollars yields a -0.343 change in the slope of the refund
 - An -0.15 slope change in enrollment
- Is enrollment lower after the kink? no!
- Remember this is a change in the rate of change
- IV result is 0.430 (about -0.147/-0.343)
- Interpretation? An additional \$1k of refund yields 0.43 pp increase in enrollment

Validity Tests

- What are the two key underlying assumptions?

Validity Tests

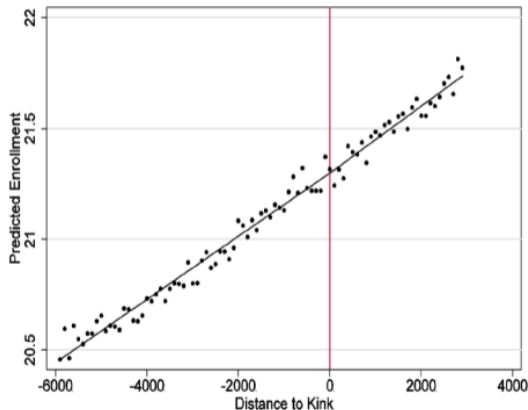
- What are the two key underlying assumptions? (listed p. 14)
 1. “other covariates do not change at the kink points” and
 2. “taxpayers do not sort along the tax schedule”
- And, implicitly, that their kink is important
- What kind of evidence can you show on this?

1. “Other Covariates Do Not Change at the Kink Points”

- Put all covariates together into “predicted enrollment”
- What is this “predicted enrollment”?

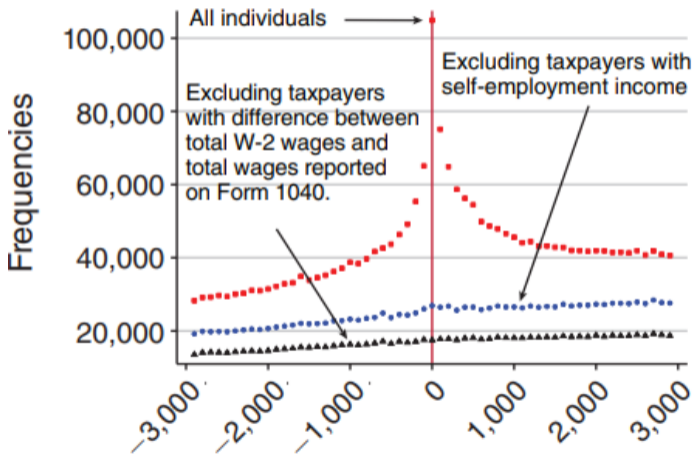
$$\text{enrollment}_i = \gamma X_i + \epsilon_i$$

- Make $\widehat{\text{enrollment}}_i = \hat{\gamma} X_i$
- A way to combining a bunch of covariates into one
- Footnote 19 says individual covariates could have a kink even if their aggregate does not. They say it's not true, but don't say more than that.



2. “Taxpayers Do Not Sort Along the Tax Schedule”

- How do we do this?
- Look at number of people around the kink



More Implicit: This is an Important Kink

- How can we tell if this kink is important?

More Implicit: This is an Important Kink

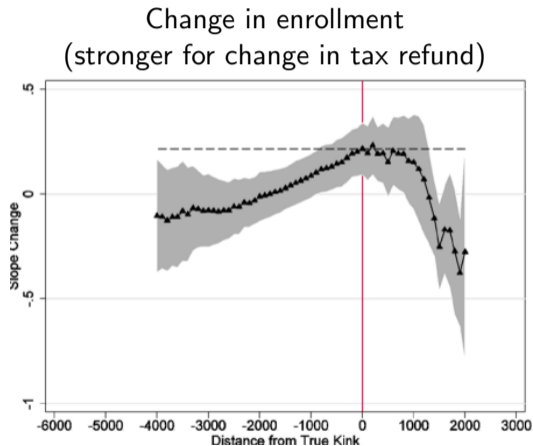
- How can we tell if this kink is important?
 - where should difference in slope be largest?

More Implicit: This is an Important Kink

- How can we tell if this kink is important?
 - where should difference in slope be largest?
 - at the true kink
 - the move the kink location around and look for differences

More Implicit: This is an Important Kink

- How can we tell if this kink is important?
 - where should difference in slope be largest?
 - at the true kink
 - the move the kink location around and look for differences



Next Lectures

- Next two classes on matching
- Matching 1
 - comment on my paper, and it's ok to complain
- Next class: quantitative summary due