

# Treatment Effects

Richard L. Sweeney

based on slides by Chris Conlon

Empirical Methods  
Spring 2019

# 1 Setup Conditional Independence

## 2 Matching

## 3 IV Basics Example: Dobbie et al

## 4 DiD

## 5 RDD

## 6 MTE

This lecture draw heavily upon

- 2012 AEA continuing education lectures by Imbens and Wooldridge (full materials available [here](#).)
- Abadie and Cattaneo (2018)

Setup

Conditional  
Independence

Matching

IV

Basics  
Example:  
Dobbie et al

DiD

RDD

MTE

References

# The Evaluation Problem

- The issue we are concerned about is identifying the effect of a policy or an investment or some individual action on one or more outcomes of interest
- This has become the workhorse approach of the applied microeconomics fields (Public, Labor, etc.)
- Examples may include:
  - The effect of taxes on labor supply
  - The effect of education on wages
  - The effect of incarceration on recidivism
  - The effect of competition between schools on schooling quality
  - The effect of price cap regulation on consumer welfare
  - The effect of indirect taxes on demand
  - The effects of environmental regulation on incomes
  - The effects of labor market regulation and minimum wages on wages and employment

Typically attributed to Rubin

- Observe  $N$  units, indexed by  $i$ , drawn randomly from a larger population
- Postulate two **potential outcomes** for each unit  $\{Y_i(1), Y_i(0)\}$  depending on whether they receive treatment or not.
- Observe additional *exogenous* covariates  $X_i$
- Consider a binary treatment  $W_i$  such that

$$Y_i \equiv Y_i(W_i) = \begin{cases} Y_i(0) & \text{if } W_i = 0 \\ Y_i(1) & \text{if } W_i = 1 \end{cases}$$

# SUTVA

- Note there is already an important assumption embedded in this setup, the stable unit treatment value assumption (**SUTVA**).
- Assume that the outcome, in either state for unit  $i$  does not depend on the assignment of other units.
- This is likely to fail in many important settings. Examples?

- ① Matching
  - ② Instrumental Variables
  - ③ Difference in Difference and Natural Experiments
  - ④ RCTs
  - ⑤ Structural Models
- Key distinction: the treatment effect of some program (a number) from understanding how and why things work (the mechanism).
  - Models let us link numbers to mechanisms.

# The Evaluation Problem

- Two major problems:
  - All individuals have different treatment effects (**heterogeneity**).
  - We don't actually observe any one person's treatment effect ! (Missing Data problem)
  - Individual treatment effects  $\tau_i = Y_{1i} - Y_{0i}$  are never observed (FPOCI)
- We need strong assumptions in order to recover  $f(\beta_i)$  from data.

Setup

Conditional  
Independence

Matching

IV

Basics

Example:  
Dobbie et al

DiD

RDD

MTE

References



What is hard here?

- Heterogeneous effect of  $\beta_i$  in population.
- Selection in treatment may be endogenous. That is  $W_i$  depends on  $Y_i(1), Y_i(0)$ .
- Fisher or Roy (1951) model:

$$Y_i = (Y_i(1) - Y_i(0))W_i + Y_i(0) = \alpha + \beta_i W_i + u_i$$

- Agents usually choose  $W_i$  with  $\beta_i$  or  $u_i$  in mind.
- Can't necessarily pool across individuals since  $\beta_i$  is not constant.

Setup

Conditional  
Independence

Matching

IV

Basics  
Example:  
Dobbie et al

DiD

RDD

MTE

References

# Structural vs. Reduced Form

Setup

Conditional  
Independence

Matching

IV

Basics  
Example:  
Dobbie et al

DiD

RDD

MTE

References

- Usually we are interested in one or two parameters of the distribution of  $\beta_i$  (such as the average treatment effect or average treatment on the treated).
- Most program evaluation approaches seek to identify one effect or the other effect. This leads to these as being described as **reduced form** or **quasi-experimental**.
- The **structural** approach attempts to recover the entire joint  $f(\beta_i, u_i)$  distribution but generally requires more assumptions, but then we can calculate whatever we need.
- Instead we often focus on simpler estimands.

## Common Objects of Interest

- Population average treatment effect (PATE)

$$\tau_P = E[Y_i(1) - Y_i(0)]$$

- Population average treatment effect for treated units (PATT)

$$\tau_{P,T} = E[Y_i(1) - Y_i(0) | W = 1]$$

- Sample average treatment effect (SATE)

$$\tau_S = \frac{1}{N} \sum_{i=1}^N (Y_i(1) - Y_i(0))$$

- Sample average treatment effect for treated units (SATT)

$$\tau_{S,T} = \frac{1}{N_T} \sum_{i \in W_i=1} (Y_i(1) - Y_i(0))$$

Setup

Conditional  
Independence

Matching

IV

Basics

Example:  
Dobbie et al

DiD

RDD

MTE

References

- Consider the *association*

$$\tau = E[Y|W = 1] - E[Y|W = 0]$$

- Then  $\tau = \tau_{ATE} + b_{ATE}$
- Where  $b$  is the *bias*

$$\begin{aligned} b_{ATE} = & (E[Y_1|W = 1] - E[Y_1|W = 0])Pr(W = 0) \\ & + (E[Y_0|W = 1] - E[Y_0|W = 1])Pr(W = 1) \end{aligned}$$

- So the bias disappears only if the potential outcomes are independent of treatment assignment.
- This is called **unconfoundedness**.

## Estimation under unconfoundedness

## Assumption: 1

$$(Y_i(0), Y_i(1)) \perp W_i | X_i$$

- Sometimes called “conditional independence assumption” or “selection on observables”.
- Can see this is implicit in the regression  $Y_i = \alpha + \tau W_i + X_i' \beta + \epsilon_i$  where  $\epsilon_i \perp X_i$  under the assumption of a constant treatment effect (otherwise this is not the same)

## Assumption 2 (Overlap)

$$0 < Pr(W_i = 1 | X_i) < 1$$

## How useful are these assumptions?

Imbens (2015) has a good discussion on this. Suggests following motivations:

- This is a natural starting point. Compare treatment and control units, after adjusting for observables. Need not be the last word!
- *All* comparisons involve comparing treated to untreated units. Absent RCT, its up to researcher to investigate which comparisons to emphasize
- Often specifying a model can clarify how sensible this is. Guido has a good example on costs in the paper.

# Under these assumptions, can we just use regression?

## Setup

Conditional  
Independence

## Matching

## IV

Basics  
Example:  
Dobbie et al

## DiD

## RDD

## MTE

## References

- Let  $\mu_w(x) = E[Y_i(w)|X_i = x]$
- A regression estimate of  $\tau$  is then

$$\hat{\tau}_{reg} = \frac{1}{N} \sum_i W_i(Y_i - \hat{\mu}_0(X_i)) + (1 - W_i)(\hat{\mu}_1(X_i) - Y_i)$$

- Typically estimate

$$Y_i = \alpha + \beta'X_i + \tau W_i + \epsilon_i$$

which assumes  $\mu_w(x) = \beta'x + \tau * w$

- Could easily also compute

$$\mu_w(x) = \alpha_w + \beta'_w x$$

- Key point is that this estimator can be viewed as a **missing data** problem, where predictions are computed using regression.

## When is this likely to be a problem?

Setup

Conditional  
Independence

Matching

IV

Basics

Example:  
Dobbie et al

DiD

RDD

MTE

References

- Note  $\mu_0(x)$  is used to predict the "missing" control outcomes for the treated observations.
- Want this prediction at the average treated covariates  $\bar{X}_T$
- With linear regression, our average control prediction for the treated observations is going to be  $\bar{Y}_C + \hat{\beta}'(\bar{X}_T - \bar{X}_C)$
- Ok if:
  - ①  $\mu()$  is properly specified
  - ② treated and control observations are similar (in  $X$ )
- First condition is untestable, but in practice predictions are often sensitive to functional form
- Leads to a big emphasis on covariate balance.



# Matching

- Regression imputes missing potential outcomes using regression.
- Matching imputes using the *realized* outcome of (nearly) identical units in the opposite assignment group.
- Remember, we're in a world where we've assumed unconfoundedness. Only challenge is that the treatment group and the control group don't have the same distribution of  $X$ 's.
- **Re-weight** the un-treated population so that it resembles the treated population.
- Once distribution of  $X_i$  is the same for both groups  $X_i|W_i \sim X_i$  then we assume all other differences are irrelevant and can just compare means.

## Matching

Let  $F^1(x)$  be the distribution of characteristics in the treatment group, we can define the ATE as

$$\begin{aligned} E[Y(1) - Y(0)|T = 1] &= E_{F^1(x)}[E(Y(1) - Y(0)|T = 1, X)] \\ &= E_{F^1(x)}[E(Y(1)|T = 1, X)] - E_{F^1(x)}[E(Y(0)|T = 1, X)] \text{ line} \end{aligned}$$

The first part we observe directly:

$$= E_{F^1(x)}[E(Y(1)|T = 1, X)]$$

But the counterfactual mean is not observed!

$$= E_{F^1(x)}[E(Y(0)|T = 1, X)]$$

But conditional independence does this for us:

$$E_{F^1(x)}[E(Y(0)|T = 1, X)] = E_{F^1(x)}[E(Y(0)|T = 0, X)]$$

# A Matching Example

Here is an example where I found that matching was helpful in my own work with Julie Mortimer:

- We ran a randomized experiment where we removed Snickers bars from around 60 vending machines in office buildings in downtown Chicago.
- We have a few possible control groups:
  - ① Same vending machine in other weeks (captures heterogeneous tastes in the cross section)
  - ② Other vending machines in the same week (might capture aggregate shocks, ad campaigns, etc.)
- We went with #1 as #2 was not particularly helpful.

# A Matching Example

Major problem was that there was a ton of heterogeneity in the overall level of (potential) weekly sales which we call  $M_t$ .

- Main source of heterogeneity is how many people are in the office that week, or how late they work.
- Based on total sales our average over treatment weeks was in the 74th percentile of all weeks.
- This was after removing a product, so we know sales should have gone down!
- How do we fix this without running the experiment for an entire year!

## Treatment Effects

Richard L. Sweeney

### Setup

Conditional Independence

### Matching

### IV

Basics

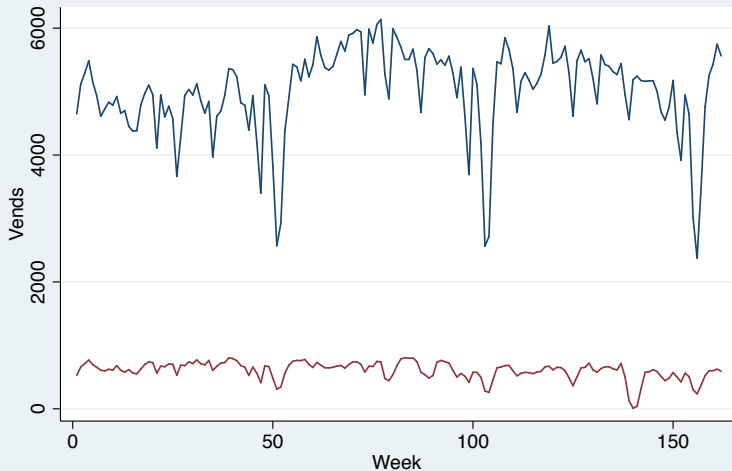
Example: Dobbie et al

### DiD

### RDD

### MTE

### References



# A Matching Example

Ideally we could just observe  $M_t$  directly and use that as our matching variable  $X$

- We didn't observe it directly and tried a few different measures:
  - Sales at the soda machine next to the snack machine
  - Sales of salty snacks at the same machine (not substitutes for candy bars).
  - We used k-NN with  $k = 4$  to select control weeks – notice we re-weight so that overall sales are approximately same (minus the removed product).
- We also tried a more structured approach:
  - Define controls weeks as valid IFF
  - Overall sales were weakly lower
  - Overall sales were not less than Overall Sales less expected sales less Snickers Sales.

Product	Control Mean	Control %ile	Treatment Mean	Treatment %ile	Mean Difference
<i>Vends</i>					
Peanut M&Ms	359.9	73.6	478.3*	99.4	118.4*
Twix Caramel	187.6	55.3	297.1*	100.0	109.5*
Assorted Chocolate	334.8	66.7	398.0*	95.0	63.2*
Assorted Energy	571.9	63.5	616.2	76.7	44.3
Zoo Animal Cracker	209.1	78.6	243.7*	98.1	34.6*
Salted Peanuts	187.9	70.4	216.3*	93.7	28.4
Choc Chip Famous Amos	171.6	71.7	193.1*	95.0	21.5*
Ruger Vanilla Wafer	107.3	59.7	127.9	78.6	20.6*
Assorted Candy	215.8	43.4	229.6	60.4	13.7
Assorted Potato Chips	279.6	64.2	292.4*	66.7	12.8
Assorted Pretzels	548.3	87.4	557.7*	88.7	9.4
Raisinets	133.3	66.0	139.4	74.2	6.1
Cheetos	262.2	60.1	260.5	58.2	-1.8
Grandmas Choc Chip	77.9	51.3	72.5	37.8	-5.4
Doritos	215.4	54.1	203.1	39.6	-12.3*
Assorted Cookie	180.3	61.0	162.4	48.4	-17.9
Skittles	100.1	62.9	75.1*	30.2	-25.1*
Assorted Salty Snack	1382.8	56.0	1276.2*	23.3	-106.7*
Snickers	323.4	50.3	2.0*	1.3	-321.4*
Total	5849.6	74.2	5841.3	73.0	-8.3

Notes: Control weeks are selected through the-neighbor matching using four control observations for each treatment week. Percentiles are relative to the full distribution of control weeks.

## How do you actually do this?

- One dimension is easy: just sort
- In multiple dimensions, there are a variety of built in nearest neighbor packages (Abadie Imbens (2006))
- What's nice about these is that the reasearcher only has to pick the number of matches (although the default tolerances not always innocuous)
- This is still cursed in that our nearest neighbors get further away as the dimension grows.
- Suppose instead we had a **sufficient statistic**

Setup

Conditional  
Independence

Matching

IV

Basics  
Example:  
Dobbie et al

DiD

RDD

MTE

References



## Propensity Score

- Rosenbaum and Rubin propose the **propensity score**

$$e(x) = Pr(W_i = 1|X_i) = E[W_i|X_i = x]$$

- They prove that under the assumption of unconfoundedness,

$$(Y_i(0), Y_i(1)) \perp W_i | e(X_i)$$

- So even if  $X$  is high dimensional, it is sufficient to condition on a scalar function
- Of course, the true propensity score is not known...

This suggests an attractive  
weighing

## Setup

Conditional  
Independence

## Matching

## IV

## Basics

Example:  
Dobbie et al

## DiD

## RDD

## MTE

## References

## 4.B.3 Propensity Score Estimators: Weighting

$$\mathbb{E} \left[ \frac{WY}{e(X)} \right] = \mathbb{E} \left[ \mathbb{E} \left[ \frac{WY_i(1)}{e(X)} \middle| X \right] \right] = \mathbb{E} \left[ \mathbb{E} \left[ \frac{e(X)Y_i(1)}{e(X)} \right] \right] = \mathbb{E}[Y_i(1)],$$

and similarly

$$\mathbb{E} \left[ \frac{(1-W)Y}{1-e(X)} \right] = \mathbb{E}[Y_i(0)],$$

implying

$$\tau_P = \mathbb{E} \left[ \frac{W \cdot Y}{e(X)} - \frac{(1-W) \cdot Y}{1-e(X)} \right].$$

With the propensity score known one can directly implement this estimator as

$$\tilde{\tau} = \frac{1}{N} \sum_{i=1}^N \left( \frac{W_i \cdot Y_i}{e(X_i)} - \frac{(1-W_i) \cdot Y_i}{1-e(X_i)} \right). \quad (3)$$

## Approaches now look similar

- One option is "inverse probability weighting"
- Nonparametrically estimate  $e(x)$ , then compute

$$\hat{\tau} = \sum_i^N \frac{W_i Y_i}{\hat{e}(X_i)} / \sum_i^N \frac{W_i}{\hat{e}(X_i)} - \sum_i^N \frac{(1 - W_i) Y_i}{1 - \hat{e}(X_i)} / \sum_i^N \frac{(1 - W_i)}{1 - \hat{e}(X_i)}$$

where this is slightly more complicated than just plugging in  $\hat{e}()$  because in your sample the weights won't necessarily sum to one (Hirano, Imbens and Ridder (2003))

- Alternatively we could flexibly estimate  $\mu_w$  then plug in these predictions for each observation manually.
- With discrete covariates, these will be equivalent!
- Otherwise their finite sample properties will vary depending on the smoothness of the regression and propensity score functions.

## What about matching on the (estimated) propensity score?

- VERY widely used approach
- Large sample properties not known
- "Why Propensity Scores Should Not Be Used for Matching" (King and Nielsen, Forthcoming)
- Show this performs poorly in simulations compared to matching on X's directly.
- One alternative from the same author's: Coarsened Exact Matching
  - Available in R and Stata from [Gary King's website](#)
  - The idea: temporarily coarsen each variable into substantively meaningful groups, exact match on these coarsened data, and then retain only the original (uncoarsened) values of the matched data.

# CEM has many uses

- Linh To's JMP:
- Question: Is there a signal value to parental leave?
- Theory: many PBNE's. In practice depends on pooling.
- Setting: Extension of leave in Denmark.
- Look for response among three types of women:
  - ① pool, pool
  - ② pool, separate
  - ③ separate, separate
- Convincing RD: restrict to sample already pregnant when law announced
- Challenge: Only see mothers in one group or the other
- Solution: Match each pre period mother using their closest post-period counterpart, and assign her to that post-group.

Setup

Conditional  
Independence

Matching

IV

Basics

Example:  
Dobbie et al

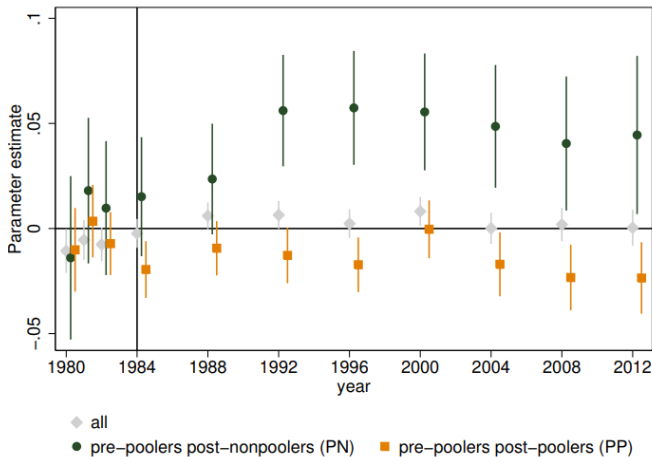
DiD

RDD

MTE

References

(a) Log wages



## What can ML add here?

- Estimating the propensity score is a pure **prediction** problem. We don't care what causes someone to be treated in this setup
- This is a natural place for ML (decision trees, random forests).
- What should we use to predict?

## Some recent ML proposals I

Belloni, Chernozhukov, Fernández, and Hansen (2013)

- "double selection" procedure
- use LASSO to select  $X$  which predict  $Y$ , and another LASSO to find  $X$  that predict  $W$
- then do OLS on the union of the two sets of covariates
- show this performs better than simple regularized regression of outcome on treatment and covariates in one step

Setup

Conditional  
Independence

Matching

IV

Basics

Example:  
Dobbie et al

DiD

RDD

MTE

References



## Some recent ML proposals II

Athey, Imbens, and Wager (2016)

“Approximate Residual Balancing: De-Biased Inference of Average Treatment Effects in High Dimensions)”

- Idea: In order to predict the counterfactual outcomes that the treatment group would have had in the absence of the treatment, it is necessary to extrapolate from control
- This is confounded by imbalance.
- AIW construct weights so these samples are equivalent, and run penalized regression to compute  $\tau$

Setup

Conditional  
Independence

Matching

IV

Basics

Example:  
Dobbie et al

DiD

RDD

MTE

References

# Assessing Unconfoundedness

- This assumption is fundamentally untestable
- However people have proposed a number of tests which, if failed, might be *inconsistent* with unconfoundedness.
- One option is to look for an "effect" on an untreated group.
- Imagine you had one sample of "eligible" units, some who were treated and some who weren't. And another sample of "ineligible" units, all of whom are also untreated by construction.
- You could estimate a difference in outcomes within the two untreated groups. If eligible but untreated units look different than ineligible, that should be worrisome.
- Imbens lecture does this with the Lalonde data and the CPS.
- Another natural approach is to use "psuedo outcomes", like lagged Y.

Setup

Conditional  
Independence

Matching

IV

Basics

Example:  
Dobbie et al

DiD

RDD

MTE

References

# Assessing Overlap

- Obviously want to start with a summary table comparing the means of your treatment and control groups.
- What's a big difference? t-stats reflective of sample size
- Instead report the normalized difference in covariates. According to Imbens, a an average difference bigger than 0.25 standard deviations is worrisome.
- Another alternative is to plot the propensity score for the two groups.

Setup

Conditional  
Independence

Matching

IV

Basics

Example:  
Dobbie et al

DiD

RDD

MTE

References

## Matching wrapup

- Even under unconfoundedness, very important to ensure overlap
- Restrict your sample so that its balanced, using exact matching if low dimensional, coarse or propensity score otherwise
- Assess unconfoundedness using a psuedo-outcome if possible
- Run regression on your matched sample

# Instrumental Variables

See Guido Imben's [NBER Slides](#).

## How Close to ATE?

Angrist and Imbens give some idea how close to the ATE the LATE is:

$$\hat{\beta}_1^{TSLS} \rightarrow^p \frac{E[\beta_{1i}\pi_{1i}]}{E[\pi_{1i}]} = LATE$$

$$LATE = ATE + \frac{Cov(\beta_{1i}, \pi_{1i})}{E[\pi_{1i}]}$$

- Weighted average for people with large  $\pi_{1i}$ .
- Late is treatment effect for those whose probability of treatment is most influenced by  $Z_i$ .
- If you always (never) get treated you don't show up in LATE.

## How Close to ATE?

- With different instruments you get different  $\pi_{1i}$  and TSLS estimators!
- Even with two valid  $Z_1, Z_2$ 
  - Can be influential for different members of the population.
  - Using  $Z_1$ , TSLS will estimate the treatment effect for people whose probability of treatment  $X$  is most influenced by  $Z_1$
  - The LATE for  $Z_1$  might differ from the LATE for  $Z_2$
  - A J-statistic might reject even if both  $Z_1$  and  $Z_2$  are exogenous! (Why?).

Setup

Conditional  
Independence

Matching

IV

Basics

Example:  
Dobbie et al

DiD

RDD

MTE

References

## Example: Cardiac Catheterization

- $Y_i$  = survival time (days) for AMI patients
- $X_i$  = whether patient received cardiac catheterization (or not) (intensive treatment)
- $Z_i$  = differential distance to CC hospital

$$\begin{aligned} SurvivalDays_i &= \beta_0 + \beta_{1i}CardCath_i + u_i \\ CardCath_i &= \pi_0 + \pi_{1i}Distance_i + v_i \end{aligned}$$

- For whom does distance have the great effect on probability of treatment?
- For those patients what is their  $\beta_{1i}$ ?



## Example: Cardiac Catheterization

- IV estimates causal effect for patients whose value of  $X_i$  is most heavily influenced by  $Z_i$ 
  - Patients with small positive benefit from CC in the expert judgement of EMT will receive CC if trip to CC hospital is short (**compliers**)
  - Patients that need CC to survive will always get it (**always-takers**)
  - Patients for which CC would be unnecessarily risky or harmful will not receive it (**never-takers**)
  - Patients for who would have gotten CC if they lived further from CC hospital (hopefully don't see) (**defiers**)
- We mostly weight towards the people with small positive benefits.

# Local Average Treatment Effect

So how is this useful?

- It shows why IV can be meaningless when effects are heterogeneous.
- It shows that if the monotonicity assumption can be justified, IV estimates the effect for a particular subset of the population.
- In general the estimates are specific to that instrument and are not generalisable to other contexts.
- As an example consider two alternative policies that can increase participation in higher education.
  - Free tuition is randomly allocated to young people to attend college ( $Z_1 = 1$  means that the subsidy is available).
  - The possibility of a competitive scholarship is available for free tuition ( $Z_1 = 1$  means that the individual is allowed to compete for the scholarship).

Setup

Conditional  
Independence

Matching

IV

Basics

Example:  
Dobbie et al

DiD

RDD

MTE

References

# Local Average Treatment Effect

- Suppose the aim is to use these two policies to estimate the returns to college education. In this case, the pair  $\{Y^1, Y^0\}$  are log earnings, the treatment is going to college, and the instrument is one of the two randomly allocated programs.
- First, we need to assume that no one who intended to go to college will be discouraged from doing so as a result of the policy (monotonicity).
- This could fail as a result of a General Equilibrium response of the policy; for example, if it is perceived that the returns to college decline as a result of the increased supply, those with better outside opportunities may drop out.

Setup

Conditional  
Independence

Matching

IV

Basics

Example:  
Dobbie et al

DiD

RDD

MTE

References

## Local Average Treatment Effect

- Now compare the two instruments.
- The subsidy is likely to draw poorer liquidity constrained students into college but not necessarily those with the highest returns.
- The scholarship is likely to draw in the best students, who may also have higher returns.
- It is not a priori possible to believe that the two policies will identify the same parameter, or that one experiment will allow us to learn about the returns for a broader/different group of individuals.

Setup

Conditional  
Independence

Matching

IV

Basics

Example:  
Dobbie et al

DiD

RDD

MTE

References

# Local Average Treatment Effect

Finally, we need to understand what monotonicity means in terms of restrictions on economic theory.

- To quote from Vytlacil (2002) *Econometrica*:  
*“The LATE assumptions are not weaker than the assumptions of a latent index model, but instead impose the same restrictions on the counterfactual data as the classical selection model if one does not impose parametric functional form or distributional assumptions on the latter.”*
- This is important because it shows that the LATE assumptions are equivalent to whatever economic modeling assumptions are required to justify the standard Heckman selection model and has no claim to greater generality.
- On the other hand there are no magical solutions to identifying effects when endogeneity/selection is present; this problem is exacerbated when the effects are heterogeneous and individuals select into treatment on the basis of the returns.

## Example: Pretrial Detention

- In US, innocent until proven guilty.
- Some defendants are detained prior to trial.
- Extreme cases are obvious, but lots of discretion in the middle.
- What are the impacts on:
  - time served
  - future crime
  - rehabilitation in to workforce

# Example: Pretrial Detention

## The Effects of Pretrial Detention on Conviction, Future Crime, and Employment: Evidence from Randomly Assigned Judges†

By WILL DOBBIE, JACOB GOLDIN, AND CRYSTAL S. YANG\*

*Over 20 percent of prison and jail inmates in the United States are currently awaiting trial, but little is known about the impact of pretrial detention on defendants. This paper uses the detention tendencies of quasi-randomly assigned bail judges to estimate the causal effects of pretrial detention on subsequent defendant outcomes. Using data from administrative court and tax records, we find that pretrial detention significantly increases the probability of conviction, primarily through an increase in guilty pleas. Pretrial detention has no net effect on future crime, but decreases formal sector employment and the receipt of employment- and tax-related government benefits. These results are consistent with (i) pretrial detention weakening defendants' bargaining positions during plea negotiations and (ii) a criminal conviction lowering defendants' prospects in the formal labor market. (JEL J23, J31, J65, K41, K42)*

Means for detained vs released  
defendants

## Setup

Conditional  
Independence

## Matching

## IV

## Basics

Example:  
Dobbie et al

## DiD

## RDD

## MTE

## References

*Panel E. Outcomes*

Any guilty offense	0.578	0.486
Guilty plea	0.441	0.207
Any incarceration	0.300	0.145
Failure to appear in court	0.121	0.179
Rearrest in 0–2 years	0.462	0.398
Earnings (\$ thousands) in 1–2 years	5.224	7.911
Employed in 1–2 years	0.378	0.509
Any income in 1–2 years	0.458	0.522
Earnings (\$ thousands) in 3–4 years	5.887	8.381
Employed in 3–4 years	0.378	0.483
Any income in 3–4 years	0.461	0.508
Observations	186,938	234,127

*Notes:* This table reports descriptive statistics for the sample of defendants from Philadelphia and Miami-Dade counties. Data from Philadelphia are from 2007–2014 and data from Miami-Dade are from 2006–2014. Information on ethnicity, gender, age, and criminal outcomes is derived from court records. Information on earnings, employment, and income is derived from the IRS data and is only available for



# First stage: Judges Matter

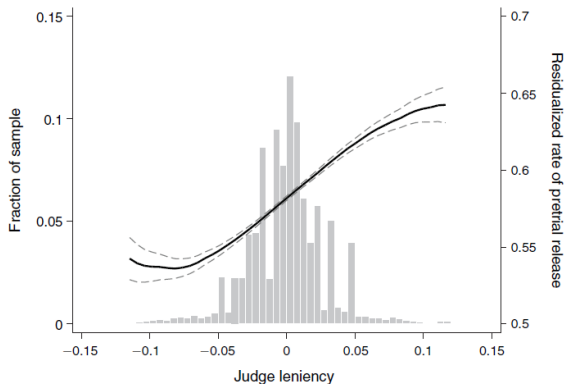


FIGURE 1. DISTRIBUTION OF JUDGE LENIENCY MEASURE AND FIRST STAGE

*Note:* This figure reports the distribution of the judge leniency measure that is estimated using data from other cases assigned to a bail judge in the same year following the procedure described in Section III.

## Is assignment random?

TABLE 3—TEST OF RANDOMIZATION

	Pretrial release (1)	Judge leniency (2)
Male	−0.11781 (0.00716)	0.00007 (0.00015)
Black	−0.03941 (0.00362)	0.00003 (0.00017)
Age at bail decision	−0.01287 (0.00236)	−0.00005 (0.00006)
Prior offense in past year	−0.15492 (0.00739)	0.00019 (0.00012)
Number of offenses	−0.02409 (0.00120)	0.00000 (0.00002)
Felony offense	−0.25575 (0.01821)	0.00005 (0.00010)
Any drug offense	0.12528 (0.00909)	0.00013 (0.00019)
Any DUI offense	0.10966 (0.01679)	0.00019 (0.00024)

Setup

Conditional  
Independence

Matching

IV

Basics

Example:  
Dobbie et al

DiD

RDD

MTE

References

TABLE 4—PRETRIAL RELEASE AND CRIMINAL OUTCOMES

	Detained mean	OLS results			2SLS results	
	(1)	(2)	(3)	(4)	(5)	(6)
<i>Panel A. Case outcomes</i>						
Any guilty offense	0.578 (0.494)	−0.072 (0.014)	−0.057 (0.009)	−0.046 (0.007)	−0.123 (0.047)	−0.140 (0.042)
Guilty plea	0.441 (0.497)	−0.188 (0.008)	−0.099 (0.010)	−0.082 (0.007)	−0.095 (0.056)	−0.108 (0.052)
Any incarceration	0.300 (0.458)	−0.161 (0.012)	−0.104 (0.006)	−0.110 (0.007)	0.006 (0.029)	−0.012 (0.030)
<i>Panel B. Court process outcomes</i>						
Failure to appear in court	0.121 (0.326)	0.063 (0.004)	0.010 (0.008)	0.021 (0.007)	0.158 (0.046)	0.156 (0.046)
Absconded	0.002 (0.045)	0.005 (0.000)	0.002 (0.000)	0.002 (0.000)	0.005 (0.004)	0.005 (0.004)
<i>Panel C. Future crime</i>						
Rearrest in 0–2 years	0.462 (0.499)	−0.050 (0.011)	−0.015 (0.006)	0.016 (0.005)	0.024 (0.061)	0.015 (0.063)
Rearrest prior to disposition	0.155 (0.362)	0.051 (0.008)	0.066 (0.007)	0.100 (0.007)	0.192 (0.038)	0.189 (0.042)
Rearrest after disposition	0.343 (0.475)	−0.075 (0.006)	−0.049 (0.002)	−0.041 (0.003)	−0.114 (0.057)	−0.121 (0.055)
Court × time fixed effects	—	Yes	Yes	Yes	Yes	Yes
Baseline controls	—	No	Yes	Yes	No	Yes
Complier weights	—	No	No	Yes	No	No
Observations	186,938	421,065	421,065	421,065	421,065	421,065

*Notes* This table reports OLS and two-stage least squares results of the impact of pre-trial release. The regressions are estimated on the sample as described in the notes to Table 1. The dependent variable is listed in each

# Interpretation: Who is marginal here?

# Interpretation: Who is marginal here?

- Instrument isn't binary here
- Thought experiment is the same though: identify which defendants get out under the most lenient judge minus those that get out under the strictest judge

Table C.1: Sample Share by Compliance Type

Model Specification:	Local Linear Model			Linear Model		
Leniency Cutoff:	1%	1.5%	2%	1%	1.5%	2%
Compliers	0.13	0.13	0.13	0.11	0.10	0.09
Never Takers	0.36	0.36	0.36	0.39	0.39	0.40
Always Takers	0.51	0.51	0.51	0.50	0.51	0.51

# Who are the compliers?

## Who are the compliers?

- Follow strategy of Dahl et al (QJE 2014)
- Estimate complier share by subgroup

Table C.1: Sample Share by Compliance Type

Model Specification:	Local Linear Model			Linear Model		
Leniency Cutoff:	1%	1.5%	2%	1%	1.5%	2%
Compliers	0.13	0.13	0.13	0.11	0.10	0.09
Never Takers	0.36	0.36	0.36	0.39	0.39	0.40
Always Takers	0.51	0.51	0.51	0.50	0.51	0.51

# Interpretation: Who is marginal here?



# Interpretation: Who is marginal here?

- Instrument isn't binary here
- Thought experiment is the same though: identify which defendants get out under the most lenient judge minus those that get out under the strictest judge

Table C.2: Characteristics of Marginal Defendants

	$P[X = x]$	$P[X = x   \text{complier}]$	$\frac{P[X=x \text{complier}]}{P[X=x]}$
White	0.402 (0.001)	0.375 (0.017)	0.931 (0.042)
Non-White	0.598 (0.001)	0.624 (0.017)	1.047 (0.028)
Drug	0.274 (0.001)	0.301 (0.015)	1.099 (0.054)
Non-Drug	0.726 (0.001)	0.699 (0.015)	0.963 (0.020)
Violent	0.173 (0.001)	0.010 (0.012)	0.058 (0.068)
Non-Violent	0.827 (0.001)	0.990 (0.012)	1.197 (0.014)
Felony	0.459 (0.001)	0.318 (0.016)	0.692 (0.036)
Misdemeanor	0.541 (0.001)	0.682 (0.016)	1.261 (0.030)
Prior Last Year	0.269 (0.001)	0.310 (0.013)	1.154 (0.049)

Further approaches to evaluation of  
program effects:

## Difference in Differences

## Setup

Conditional  
Independence

## Matching

## IV

Basics  
Example:  
Dobbie et al

## DiD

## RDD

## MTE

## References

- Sometimes we may feel we can impose more structure on the problem.
- Suppose in particular that we can write the outcome equation as

$$Y_{it} = \alpha_i + d_t + \beta_i T_{it} + u_{it}$$

- In the above we have now introduced a time dimension  $t = \{1, 2\}$ .
- Now suppose that  $T_{i1} = 0$  for all  $i$  and  $T_{i2} = 1$  for a well defined group of individuals in our population.
- This framework allows us to identify the ATT effect under the assumption that the growth of the outcome in the non-treatment state is independent of treatment allocation:

$$E[Y_{i2}^0 - Y_{i1}^0 | T] = E[Y_{i2}^0 - Y_{i1}^0]$$

# Before and After

An even simpler estimator is the **before and after** or **event study**.

- We look an outcome before or after an event
  - A news event: the announcement of a merger or stock split.
  - A tax change, a new law, etc.

$$\begin{aligned} E[Y_{i2} - Y_{i1} | T_{i2} = 1] &= E[Y_{i2}^1 - Y_{i1}^1 | T_{i2} = 1] \\ &= d_2 - d_1 + E[\beta_i | T_{i2} = 1] \end{aligned}$$

- Except under strong conditions  $d_2 = d_1$  we shouldn't believe the results of the before and after estimator.
- Main Problem: we attribute changes to treatment that might have happened anyway **trend**.
- e.g: Cigarette consumption drops 4% after a tax hike. (But it dropped 3% the previous four years).
- Also worry about: **anticipation**, **gradual rollout**, etc.

# Difference in Differences

Let's try and estimate  $d_2 - d_1$  directly and then difference it out. Here we use **parallel trends**:

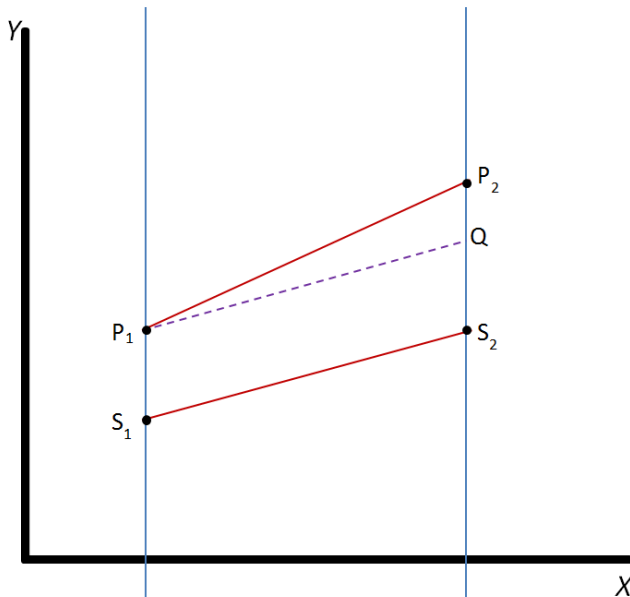
$$\begin{aligned}E[Y_{i2}^0 - Y_{i1}^0 | T_{i2} = 1] &= E[Y_{i2}^0 - Y_{i1}^0 | T_{i2} = 0] \\E[Y_{i2} - Y_{i1} | T_{i2} = 0] &= d_2 - d_1\end{aligned}$$

We now obtain an estimator for ATT:

$$E[\beta_i | T_{i2} = 1] = E[Y_{i2} - Y_{i1} | T_{i2} = 1] - E[Y_{i2} - Y_{i1} | T_{i2} = 0]$$

which can be estimated by the difference in the growth between the treatment and the control group.

# Parallel Trends



## Difference in Differences

Now consider the following problem:

- Suppose we wish to evaluate a training program for those with low earnings. Let the threshold for eligibility be  $B$ .
- We have a panel of individuals and those with low earnings qualify for training, forming the treatment group.
- Those with higher earnings form the control group.
- Now the low earning group is low for two reasons
  - ① They have low permanent earnings ( $\alpha_i$  is low) - this is accounted for by diff in diffs.
  - ② They have a negative transitory shock ( $u_{i1}$  is low) - this is not accounted for by diff in diffs.

Setup

Conditional  
Independence

Matching

IV

Basics  
Example:  
Dobbie et al

DiD

RDD

MTE

References

## Difference in Differences

- #2 above violates the assumption

$$E[Y_{i2}^0 - Y_{i1}^0 | T] = E[Y_{i2}^0 - Y_{i1}^0].$$

- To see why note that those participating into the program are such that  $Y_{i0}^0 < B$ . Assume for simplicity that the shocks  $u$  are *iid*. Hence  $u_{i1} < B - \alpha_i - d_1$ . This implies:

$$E[Y_{i2}^0 - Y_{i1}^0 | T = 1] = d_2 = d_1 - E[u_{i1} | u_{i1} < B - \alpha_i - d_1]$$

For the control group:

$$E[Y_{i2}^0 - Y_{i1}^0 | T = 1] = d_2 = d_1 - E[u_{i1} | u_{i1} > B - \alpha_i - d_1]$$

- Hence

$$\begin{aligned} &E[Y_{i2}^0 - Y_{i1}^0 | T = 1] - E[Y_{i2}^0 - Y_{i1}^0 | T = 0] = \\ &E[u_{i1} | u_{i1} > B - \alpha_i - d_1] - E[u_{i1} | u_{i1} < B - \alpha_i - d_1] > 0 \end{aligned}$$

- This is effectively regression to the mean: those unlucky enough to have a bad shock recover and show greater growth relative to those with a good shock. The nature of the bias depends on the stochastic properties of the shocks

## Difference in Differences

Ashefelter (1978) was one of the first  
to consider difference in differences to evaluate training programs.

TABLE 1.—MEAN EARNINGS PRIOR, DURING, AND SUBSEQUENT TO TRAINING FOR 1964 MDTA CLASSROOM TRAINEES AND A COMPARISON GROUP

	White Males		Black Males		White Females		Black Females	
	Trainees	Comparison Group	Trainees	Comparison Group	Trainees	Comparison Group	Trainees	Comparison Group
1959	\$1,443	\$2,588	\$ 904	\$1,438	\$ 635	\$ 987	\$ 384	\$ 635
1960	1,533	2,699	976	1,521	687	1,076	440	635
1961	1,572	2,782	1,017	1,573	719	1,163	471	719
1962	1,843	2,963	1,211	1,742	813	1,308	566	813
1963	1,810	3,108	1,182	1,896	748	1,433	531	904
1964	1,551	3,275	1,273	2,121	838	1,580	688	1,017
1965	2,923	3,458	2,327	2,338	1,747	1,698	1,441	1,163
1966	3,750	4,351	2,983	2,919	2,024	1,990	1,794	1,441
1967	3,964	4,430	3,048	3,097	2,244	2,144	1,977	1,698
1968	4,401	4,955	3,409	3,487	2,398	2,339	2,160	1,990
1969	\$4,717	\$5,033	\$3,714	\$3,681	\$2,646	\$2,444	\$2,457	\$2,160
Number of Observations	7,326	40,921	2,133	6,472	2,730	28,142	1,356	5,104



## Difference in Differences

Ashenfelter (1978) reports the following results.

TABLE 2.—CRUDE ESTIMATES (AND ESTIMATED STANDARD ERRORS), ASSUMING  $B=0$  AND  $\beta_j'=0$  FOR  $j>1$ , OF THE EFFECT OF TRAINING ON EARNINGS DURING AND AFTER TRAINING, WHITE MALE MDTA 1964 CLASSROOM TRAINEES

Effect in (value of $t$ )	Value of Effects for		
	$t-s=1963$	$t-s=1962$	$t-s=1961$
1962	—	—	91 (13)
1963	—	-179 (14)	-88 (17)
1964	-426 (16)	-605 (18)	-514 (20)
1965	763 (20)	584 (22)	675 (23)
1966	697 (25)	518 (27)	609 (28)
1967	833 (28)	655 (30)	746 (31)
1968	745 (34)	566 (35)	657 (36)
1969	984 (37)	805 (39)	896 (40)

# Difference in Differences

- The assumption on growth of the non-treatment outcome being independent of assignment to treatment may be violated, but it may still be true conditional on  $X$ .

- Consider the assumption

$$E[Y_{i2}^0 - Y_{i1}^0 | X, T] = E[Y_{i2}^0 - Y_{i1}^0 | X]$$

- This is just matching assumption on a redefined variable, namely the growth in the outcomes. In its simplest form the approach is implemented by running the regression

$$Y_{it} = \alpha_i + d_t + \beta_i T_{it} + \gamma'_t X_i + u_{it}$$

which allows for differential trends in the non-treatment growth depending on  $X_i$ . More generally one can implement propensity score matching on the growth of outcome variable when panel data is available.

# Difference in Differences with Repeated Cross Sections

## Setup

Conditional  
Independence

## Matching

## IV

Basics  
Example:  
Dobbie et al

## DiD

## RDD

## MTE

## References

- Suppose we do not have available panel data but just a random sample from the relevant population in a pre-treatment and a post-treatment period. We can still use difference in differences.
- First consider a simple case where  $E[Y_{i2}^0 - Y_{i1}^0 | T] = E[Y_{i2}^0 - Y_{i1}^0]$ .
- We need to modify slightly the assumption to  $E[Y_{i2}^0 | \text{Group receiving training}] - E[Y_{i1}^0 | \text{Group receiving training in the next period}] = E[Y_{i2}^0 - Y_{i1}^0]$

which requires, in addition to the original independence assumption that conditioned on particular individuals that population we will be sampling from does not change composition.

- We can then obtain immediately an estimator for ATT as

$$E[\beta_i | T_{i2} = 1]$$

Difference in Differences with  
Repeated Cross Sections

## Setup

Conditional  
Independence

## Matching

## IV

## Basics

Example:  
Dobbie et al

## DiD

## RDD

## MTE

## References

- More generally we need an assumption of conditional independence of the form

$$\begin{aligned} E[Y_{i2}^0 | X, \text{Group receiving training}] - E[Y_{i1}^0 | X, \text{Group receiving training next period}] \\ = E[Y_{i2}^0 | X] - E[Y_{i1}^0 | X] \end{aligned}$$

- Under this assumption (and some auxiliary parametric assumptions) we can obtain an estimate of the effect of treatment on the treated by the regression

$$Y_{it} = \alpha_g + d_t + \beta T_{it} + \gamma' X_{it} + u_{it}$$

Difference in Differences with  
Repeated Cross Sections

## Setup

Conditional  
Independence

## Matching

## IV

## Basics

Example:  
Dobbie et al

## DiD

## RDD

## MTE

## References

- More generally we can first run the regression

$$Y_{it} = \alpha_g + d_t + \beta(X_{it})T_{it} + \gamma'X_{it} + u_{it}$$

where  $\alpha_g$  is a dummy for the treatment or comparison group, and  $\beta(X_{it})$  can be parameterized as  $\beta(X_{it}) = \beta'X_{it}$ . The ATT can then be estimated as the average of  $\beta'X_{it}$  over the (empirical) distribution of  $X$ .

- A non parametric alternative is offered by Blundell, Dias, Meghir and van Reenen (2004).

Difference in Differences and  
Selection on Unobservables

## Setup

Conditional  
Independence

## Matching

## IV

## Basics

Example:  
Dobbie et al

## DiD

## RDD

## MTE

## References

- Suppose we relax the assumption of *no selection* on unobservables.
- Instead we can start by assuming that

$$E[Y_{i2}^0|X, Z] - E[Y_{i1}^0|X, Z] = E[Y_{i2}^0|X] - E[Y_{i1}^0|X]$$

where  $Z$  is an instrument which determines training eligibility say but does not determine outcomes in the non-training state. Take  $Z$  as binary (1,0).

- Non-Compliance: not all members of the eligible group ( $Z = 1$ ) will take up training and some of those ineligible ( $Z = 0$ ) may obtain training by other means.
- A difference in differences approach based on grouping by  $Z$  will estimate the impact of being allocated to the eligible group, but not the impact of training itself.

# Difference in Differences and Selection on Unobservables

Setup

Conditional  
Independence

Matching

IV

Basics  
Example:  
Dobbie et al

DiD

RDD

MTE

References

- Now suppose we still wish to estimate the impact of training on those being trained (rather than just the effect of being eligible)
- This becomes an IV problem and following up from the discussion of LATE we need stronger assumptions
  - Independence: for  $Z = a$ ,  $\{Y_{i2}^0 - Y_{i1}^0, Y_{i2}^1 - Y_{i1}^1, T(Z = a)\}$  is independent of  $Z$ .
  - Monotonicity  $T_i(1) \geq T_i(0) \forall i$
- In this case LATE is defined by

$$[E(\Delta Y|Z = 1) - E(\Delta Y|Z = 0)]/[Pr(T(1) = 1) - Pr(T(0) = 1)]$$

assuming that the probability of training in the first period is zero.

# Regression Discontinuity Design

- Another popular research design is the **Regression Discontinuity Design**.
- In some sense this is a special case of IV regression. (RDD estimates a LATE).
- Most of this is taken from the JEL Paper by Lee and Lemieux (2010).

Setup

Conditional  
Independence

Matching

IV

Basics  
Example:  
Dobbie et al

DiD

RDD

MTE

References

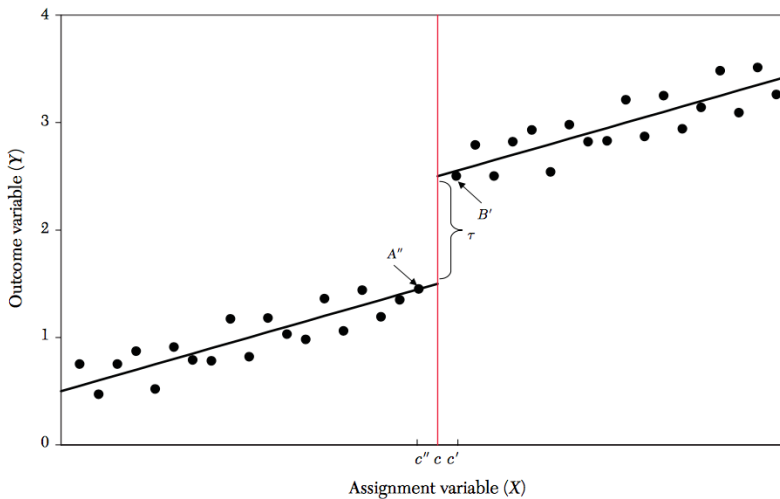


- We have a **running or forcing variable**  $x$  such that

$$\lim_{x \rightarrow c^+} P(T_i | X_i = x) \neq \lim_{x \rightarrow c^-} P(T_i | X_i = x)$$

- The idea is that there is a **discontinuous jump** in the **probability of being treated**.
- For now we focus on the **sharp discontinuity**:  
 $P(T_i | X_i \geq c) = 1$  and  $P(T_i | X_i < c) = 0$
- There is no single  $x$  for which we observe treatment and control. (Compare to Propensity Score!).
- The most important assumption is that of **no manipulability**  $\tau_i \perp D_i$  in some neighborhood of  $c$ .
- Example: a social program is available to people who earned less than \$25,000.
  - If we could compare people earning \$24,999 to people earning \$25,001 we would have as-if random assignment. (MAYBE)
  - But we might not have that many people...

# RDD: In Pictures



## RDD: Sharp RD Case

RDD uses a set of assumptions distinct from our LATE/IV assumptions. Instead it depends on **continuity**.

- We need that  $E[Y^{(1)}|X]$  and  $E[Y^{(0)}|X]$  both be continuous at  $X = c$ .
- People just to the left of  $c$  are a valid control for those just to the right of  $c$ .
- **This is not a testable assumption** → draw pictures!
- We could run the regression where  $D_i = \mathbf{1}[X_i > c]$ .

$$Y_i = \beta_0 + \tau D_i + X_i \beta + \epsilon_i$$

- This puts a lot of restrictions (linearity) on the relationship between  $Y$  and  $X$ .
- Also (without additional assumptions) we only learn about  $\tau_i$  at the point  $X = c$ .

First thing to relax is assumption of linearity.

$$Y_i = f(x_i) + \tau D_i + \epsilon_i$$

This is known as **partially linear model**.

- Two options for  $f(x_i)$ :
  - ① Kernels: Local Linear Regression
  - ② Polynomials:
$$Y_i = \beta_0 + \beta_1 x_i + \beta_2 x_i^2 + \cdots + \beta_p x_i^p + \tau D_i + \epsilon_i.$$
    - Actually, people suggest different polynomials on each side of cutoff! (Interact everything with  $D_i$ ).
- Same objective. Want to flexibly capture what happens on both sides of cutoff.
- Otherwise risk confusing nonlinearity with discontinuity!

# RDD: Kernel Boundary Problem

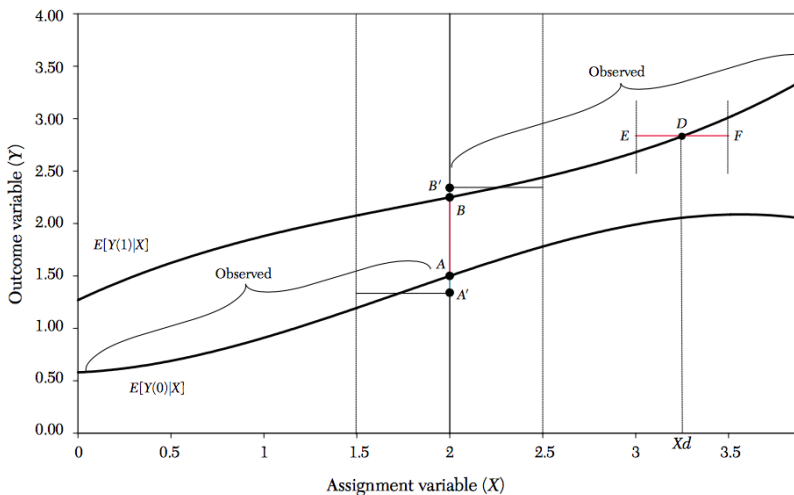


Figure 2. Nonlinear RD

RDD: Polynomial Implementation  
Details

## Setup

Conditional  
Independence

## Matching

## IV

Basics  
Example:  
Dobbie et al

## DiD

## RDD

## MTE

## References

To make life easier:

- replace  $\tilde{x}_i = x_i - c$ .
- Estimate coefficients  $\beta: (1, \tilde{x}, \tilde{x}^2, \dots, \tilde{x}^p)$  and  $\tilde{\beta}: (D_i, D_i\tilde{x}, D_i\tilde{x}^2, \dots, D_i\tilde{x}^p)$ .
- Now treatment effect at  $c$  just the coefficient on  $D_i$ . (We can ignore the interaction terms).
- If we want treatment effect at  $x_i > c$  then we have to account for interactions.
  - Identification away from  $c$  is somewhat dubious.
- Lee and Lemieux (2010) suggest estimating a coefficient on a dummy for each bin in the polynomial regression  $\sum_k \phi_k B_k$ .
  - Add polynomials until you can satisfy the test that the joint hypothesis test that  $\phi_1 = \dots \phi_k = 0$ .
  - There are better ways to choose polynomial order...

## RDD: Checklist

Most RDD papers follow the same formula (so should yours)

- Plot of  $P(D|X)$  so that we can see the discontinuity
- Plot of  $E[Y|X]$  so that we see discontinuity there also
- Plot of  $E[W|X]$  so that we don't see a discontinuity in controls.
- Density of  $X$  (check for manipulation).
- Show robustness to different “windows”
- The OLS RDD estimates
- The Local Linear RDD estimates
- The polynomial (from each side) RDD estimates
- An f-test of “bins” showing that the polynomial is flexible enough.

Read Lee and Lemieux (2010) before you get started.

## Application: Lee (2008)

Looked at incumbency advantage in the US House of Representatives

- Running variable was vote share in previous election
  - Problem of naive approach: good candidates get lots of votes!
  - Compare outcomes of districts with barely  $D$  to barely  $R$ .
- First we plot bin-scatter plots and quartic (from each side) polynomials.
- Discussion about how to choose bin-scatter bandwidth (CV).



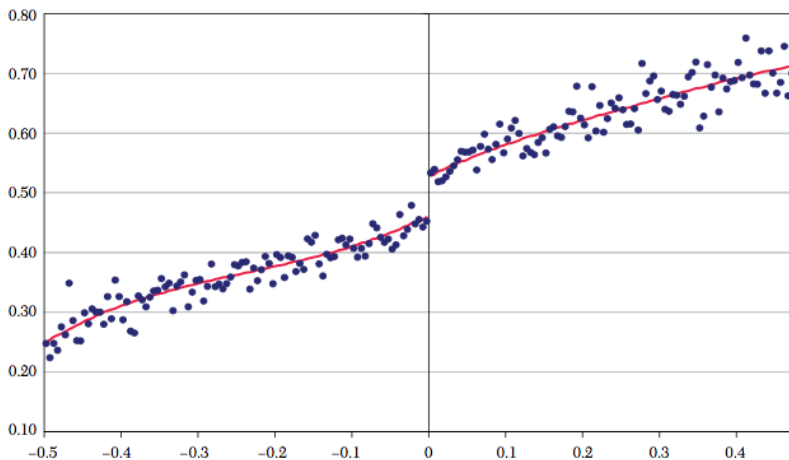


Figure 8. Share of Vote in Next Election, Bandwidth of 0.005 (200 bins)

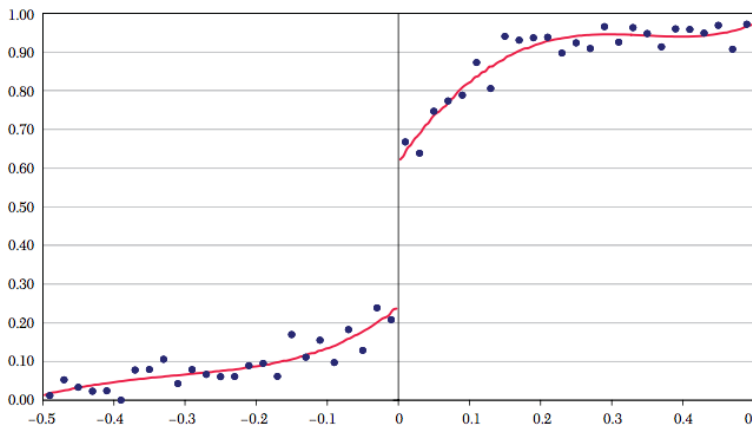


Figure 9. Winning the Next Election, Bandwidth of 0.02 (50 bins)

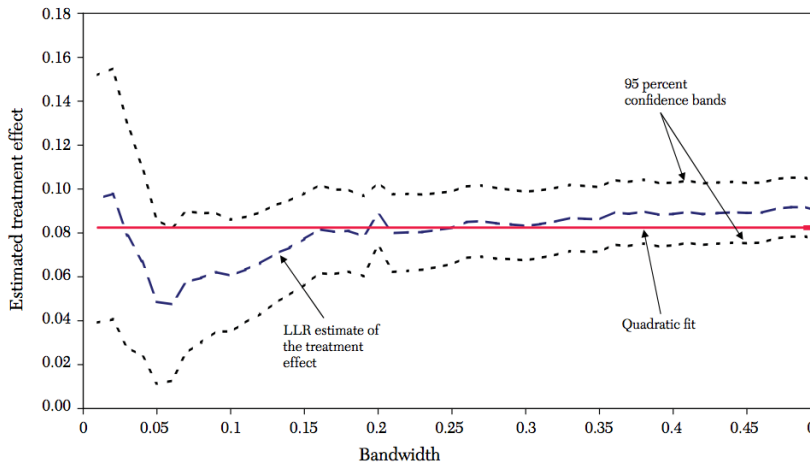


Figure 18. Local Linear Regression with Varying Bandwidth: Share of Vote at Next Election

## Other Examples

### Luca on Yelp

- Have data on restaurant revenues and yelp ratings.
- Yelp produces a yelp score (weighted average rating) to two decimals ie: 4.32.
- Score gets rounded to nearest half star
- Compare 4.24 to 4.26 to see the impact of an extra half star.
- Now there are multiple discontinuities: Pool them?  
Estimate multiple effects?

An important extension in the **Fuzzy RD**. Back to where we started:

$$\lim_{x \rightarrow c^+} P(T_i | X_i = x) \neq \lim_{x \rightarrow c^-} P(T_i | X_i = x)$$

- We need a discontinuous jump in probability of treatment, but it doesn't need to be  $0 \rightarrow 1$ .

$$\tau_i(c) = \frac{\lim_{x \rightarrow c^+} P(Y_i | X_i = x) - \lim_{x \rightarrow c^-} P(Y_i | X_i = x)}{\lim_{x \rightarrow c^+} P(T_i | X_i = x) - \lim_{x \rightarrow c^-} P(T_i | X_i = x)}$$

- Under sharp RD everyone was a **complier**, now we have some **always takers** and some **never takers** too.
- Now we are estimating the treatment effect only for the population of compliers at  $x = c$ .
- This should start to look familiar. We are going to do IV!

## Related Idea: Kinks

A related idea is that of **kinks**.

- Instead of a discontinuous jump in the outcome there is a discontinuous jump in  $\beta_i$  on  $x_i$ .
- Often things like tax schedules or government benefits have a kinked pattern.

# One quantity to rule them all: MTE

Heckman and Vytlacil provide a unifying non-parametric framework to categorize treatment effects. Their approach is known as the **marginal treatment effect** or MTE

- The MTE isn't a number it is a **function**.
- All of the other objects (LATE, ATE, ATT, etc.) can be written as integrals (weighted averages) of the MTE.
- The idea is to bridge the treatment effect parameters (stuff we get from running regressions) and the structural parameters: features of  $f(\beta_i)$ .

## One quantity to rule them all: MTE

## Setup

Conditional  
Independence

## Matching

## IV

## Basics

Example:  
Dobbie et al

## DiD

## RDD

## MTE

## References

- Consider a treatment effect  $\beta_i = Y_i(1) - Y_i(0)$ .
- Think about a single-index such that  $T_i = 1(v_i \leq Z'_i\gamma)$ .
- Think about the person for whom  $v_i = Z'_i\gamma$  (just barely untreated).

$$\Delta^{MTE}(X_i, v_i) = E[\beta_i | X_i, v_i = Z'_i\gamma]$$

- MTE is average impact of receiving a treatment for everyone with the same  $Z'\gamma$ .
- For any single index model we can rewrite

$$T_i = 1(v_i \leq Z'_i\gamma) = 1(u_{is} \leq F(Z'_i\gamma)) \text{ for } u_s \in [0, 1]$$

- $F$  is just the cdf of  $v_i$
- Now we can write  $P(Z) = Pr(T = 1|Z) = F(Z'\gamma)$ .



Now we can write,

$$Y_0 = \gamma'_0 X + U_0$$

$$Y_1 = \gamma'_1 X + U_1$$

Setup

Conditional  
Independence

Matching

IV

Basics

Example:  
Dobbie et al

DiD

RDD

MTE

References

$P(T = 1|Z) = P(Z)$  works as our instrument with two assumptions:

- ①  $(U_0, U_1, u_s) \perp P(Z)|X$ . (Exogeneity)
- ② Conditional on  $X$  there is enough variation in  $Z$  for  $P(Z)$  to take on all values  $\in (0, 1)$ .
  - This is much stronger than typical **relevance** condition. Much more like the **special regressor** method we will discuss next time.

## MTE: Derivation

Now we can write,

$$\begin{aligned} Y &= \gamma'_0 X + T(\gamma_1 - \gamma_0)'X + U_0 + T(U_1 - U_0) \\ E[Y|X, P(Z) = p] &= \gamma'_0 X + p(\gamma_1 - \gamma_0)'X + E[T(U_1 - U_0)|X, P(Z) = p] \end{aligned}$$

Observe  $T = 1$  over the interval  $u_s = [0, p]$  and zero for higher values of  $u_s$ . Let  $U_1 - U_0 \equiv \eta$ .

$$\begin{aligned} E[T(U_1 - U_0)|P(Z) = p, X] &= \int_{-\infty}^{\infty} \int_0^p (U_1 - U_0) f((U_1 - U_0)|U_s = u_s) du_s d\eta \\ E[T(\eta)|P(Z) = p, X] &= \int_{-\infty}^{\infty} \int_0^p \eta f(\eta|U_s = u_s) d\eta du_s \end{aligned}$$

$$\begin{aligned} \Delta^{MTE}(p) &= \frac{\partial E[Y|X, P(Z) = p]}{\partial p} = (\gamma_1 - \gamma_0)'X + \int_{-\infty}^{\infty} \eta f(\eta|U_s = p) d\eta \\ &= (\gamma_1 - \gamma_0)'X + E[\eta|u_s = p] \end{aligned}$$

What is  $E[\eta|u_s = p]$ ? The expected unobserved gain from treatment of those people who are on the treatment/no-treatment margin  $P(Z) = p$ .

## How to Estimate an MTE

## Easy

- 1 Estimate  $P(Z) = Pr(T = 1|Z)$  nonparametrically (include exogenous part of  $X$  in  $Z$ ).
- 2 Nonparametric regression of  $Y$  on  $X$  and  $P(Z)$  (polynomials?)
- 3 Differentiate w.r.t.  $P(Z)$
- 4 plot it for all values of  $P(Z) = p$ .

So long as  $P(Z)$  covers  $(0, 1)$  then we can trace out the full distribution of  $\Delta^{MTE}(p)$ .

## Everything is an MTE

Calculate the outcome given  $(X, Z)$  (actually  $X$  and  $P(Z) = p$ ).

- ATE : This one is obvious. We treat everyone!

$$\int_{-\infty}^{\infty} \Delta^{MTE}(p) = (\gamma_1 - \gamma_0)'X + \underbrace{\int_{-\infty}^{\infty} E(\eta|u_s) du_s}_0$$

- LATE: Fix an  $X$  and  $P(Z)$  varies from  $b(X)$  to  $a(X)$  and we integrated over the area between (compliers).

$$LATE(X) = \int_{-\infty}^{\infty} \Delta^{MTE}(p) = (\gamma_1 - \gamma_0)'X + \frac{1}{a(X) - b(X)} \int_{b(X)}^{a(X)} E(\eta|u_s) du_s$$

- ATT

$$TT(X) = \int_{-\infty}^{\infty} \Delta^{MTE}(p) \frac{Pr(P(Z|X) > p)}{E[P(Z|X)]} dp$$

- Weights for IV and OLS are a bit more complicated. See the Heckman and Vytlacil paper(s).

# Carneiro, Heckman and Vytlacil (AER 2010)

Setup

Conditional  
Independence

Matching

IV

Basics

Example:  
Dobbie et al

DiD

RDD

MTE

References

- Estimate returns to college (including heterogeneity of returns).
- NLSY 1979
- $Y = \log(wage)$
- Covariates  $X$ : Experience (years), Ability (AFQT Score), Mother's Education, Cohort Dummies, State Unemployment, MSA level average wage.
- Instruments  $Z$ : College in MSA at age 14, average earnings in MSA at 17 (opportunity cost), avg unemployment rate in state.

# Carneiro, Heckman and Vytlacil

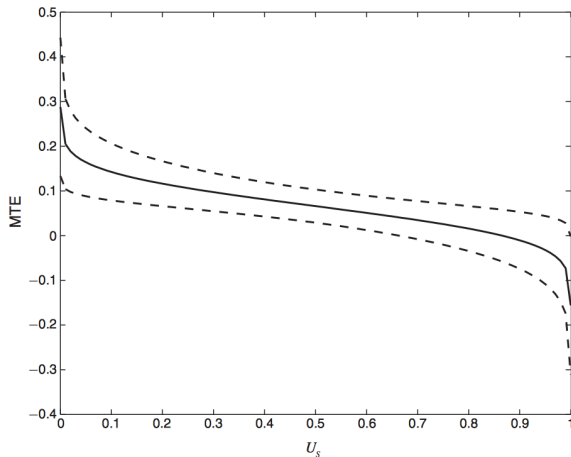


FIGURE 1. MTE ESTIMATED FROM A NORMAL SELECTION MODEL

*Notes:* To estimate the function plotted here, we estimate a parametric normal selection model by maximum likelihood. The figure is computed using the following formula:

$$\Delta^{\text{MTE}}(\mathbf{x}, u_s) = \mu_1(\mathbf{x}) - \mu_0(\mathbf{x}) - (\sigma_{1V} - \sigma_{0V})\Phi^{-1}(u_s),$$

## Carneiro, Heckman and Vytlacil

Setup  
Conditional  
Independence

Matching

IV

Basics  
Example:  
Dobbie et al

DiD

RDD

MTE

References

TABLE 4—TEST OF LINEARITY OF  $E(Y|\mathbf{X}, P = p)$  USING POLYNOMIALS IN  $P$ ; AND  
TEST OF EQUALITY OF LATES OVER DIFFERENT INTERVALS ( $H_0: LATE^j(U_S^L, U_S^H) - LATE^{j+1}(U_S^{L+j+1}, U_S^{H+j+1}) = 0$ )Panel A. Test of linearity of  $E(Y|\mathbf{X}, P = p)$  using models with different orders of polynomials in  $P^a$ 

Degree of polynomial for model	2	3	4	5
$p$ -value of joint test of nonlinear terms	0.035	0.049	0.086	0.122
Adjusted critical value	0.057			
Outcome of test	Reject			

Panel B. Test of equality of LATES ( $H_0: LATE^j(U_S^L, U_S^H) - LATE^{j+1}(U_S^{L+j+1}, U_S^{H+j+1}) = 0$ )<sup>b</sup>

Ranges of $U_S$ for $LATE^j$	(0, 0.04)	(0.08, 0.12)	(0.16, 0.20)	(0.24, 0.28)	(0.32, 0.36)	(0.40, 0.44)
Ranges of $U_S$ for $LATE^{j+1}$	(0.08, 0.12)	(0.16, 0.20)	(0.24, 0.28)	(0.32, 0.36)	(0.40, 0.44)	(0.48, 0.52)
Difference in LATES	0.0689	0.0629	0.0577	0.0531	0.0492	0.0459
$p$ -value	0.0240	0.0280	0.0280	0.0320	0.0320	0.0520
Ranges of $U_S$ for $LATE^j$	(0.48, 0.52)	(0.56, 0.60)	(0.64, 0.68)	(0.72, 0.76)	(0.80, 0.84)	(0.88, 0.92)
Ranges of $U_S$ for $LATE^{j+1}$	(0.56, 0.60)	(0.64, 0.68)	(0.72, 0.76)	(0.80, 0.84)	(0.88, 0.92)	(0.96, 1)
Difference in LATES	0.0431	0.0408	0.0385	0.0364	0.0339	0.0311
$p$ -value	0.0520	0.0760	0.0960	0.1320	0.1800	0.2400
Joint $p$ -value	0.0520					

## Carneiro, Heckman and Vytlacil

TABLE 5—RETURNS TO A YEAR OF COLLEGE

Model	Normal	Semiparametric
$ATE = E(\beta)$	0.0670 (0.0378)	Not identified
$TT = E(\beta S = 1)$	0.1433 (0.0346)	Not identified
$TUT = E(\beta S = 0)$	-0.0066 (0.0707)	Not identified
MPRTE		
Policy perturbation	Metric	
$Z_{\alpha}^k = Z^k + \alpha$	$ \mathbf{Z}\gamma - V  < e$	0.0662 (0.0373)
		0.0802 (0.0424)
$P_{\alpha} = P + \alpha$	$ P - U  < e$	0.0637 (0.0379)
		0.0865 (0.0455)
$P_{\alpha} = (1 + \alpha)P$	$ \frac{P}{U} - 1  < e$	0.0363 (0.0569)
		0.0148 (0.0589)
Linear IV (Using $P(\mathbf{Z})$ as the instrument)		0.0951 (0.0386)
OLS		0.0836 (0.0068)

*Notes:* This table presents estimates of various returns to college, for the semiparametric and the normal selection models: average treatment effect (ATE), treatment on the treated (TT), treatment on the untreated (TUT), and different versions of the marginal policy relevant treatment effect (MPRTE). The linear IV estimate uses  $P$  as the instrument. Standard errors are



# Carneiro, Heckman and Vytlacil

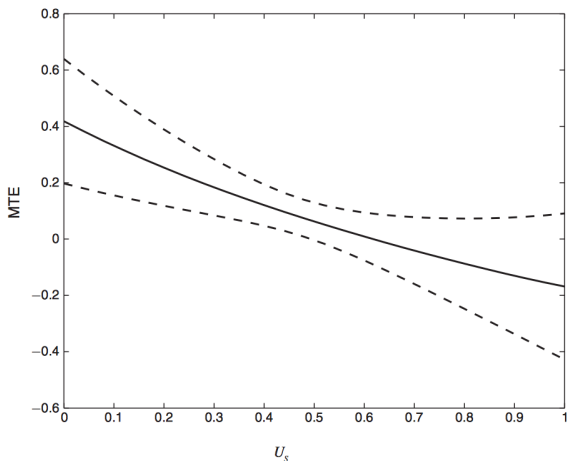


FIGURE 4.  $E(Y_1 - Y_0 | \mathbf{X}, U_s)$  WITH 90 PERCENT CONFIDENCE INTERVAL—  
LOCALLY QUADRATIC REGRESSION ESTIMATES

*Notes:* To estimate the function plotted here, we first use a partially linear regression of log wages on polynomials in  $\mathbf{X}$ , interactions of polynomials in  $\mathbf{X}$  and  $P$ , and  $K(P)$ , a locally quadratic function of  $P$  (where  $P$  is the predicted probability of attending college), with a bandwidth of 0.32;  $\mathbf{X}$  includes experience, current average earnings in the  $\theta^4 / 100$

# Carneiro, Heckman and Vytlacil

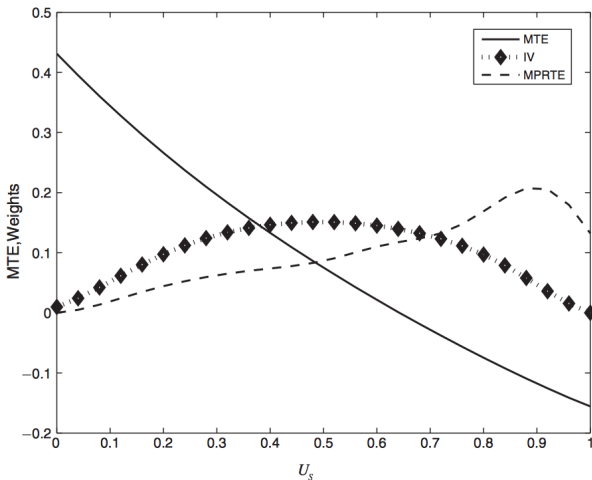


FIGURE 6. WEIGHTS FOR IV AND MPRT

*Note:* The scale of the y-axis is the scale of the MTE, not the scale of the weights, which are scaled to fit the picture.

# Diversion Example

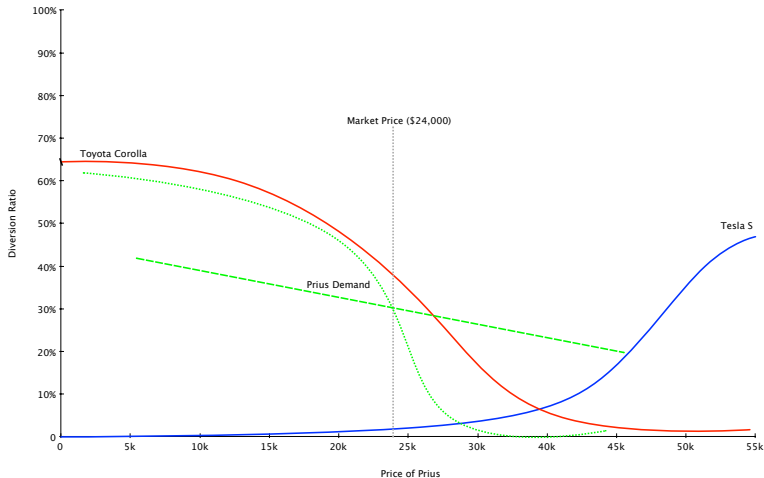
I have done some work trying to bring these methods into merger analysis.

- Key quantity: **Diversion Ratio** as I raise my price, how much do people switch to a particular competitor's product

$$D_{jk}(p_j, p_{-j}) = \left| \frac{\partial q_k}{\partial p_j}(p_j, p_{-j}) / \frac{\partial q_j}{\partial p_j}(p_j, p_{-j}) \right|$$

- We hold  $p_{-j}$  fixed and trace out  $D_{jk}(p_j)$ .
- The **treatment** is leaving good  $j$ .
- The  $Y_i$  is increased sales of good  $k$ .
- The  $Z_i$  is the price of good  $j$ .
- The key is that all changes in sales of  $k$  come through people leaving good  $j$  (no direct effects).

# Diversion for Prius (FAKE!)



## Diversion Example

## Setup

Conditional  
Independence

## Matching

## IV

## Basics

Example:  
Dobbie et al

## DiD

## RDD

## MTE

## References

$$\widehat{D_{jk}^{LATE}} = \frac{1}{\Delta q_j} \int_{p_j^0}^{p_j^0 + \Delta p_j} \underbrace{\frac{\partial q_k(p_j, p_{-j}^0)}{\partial q_j}}_{\equiv D_{jk}(p_j, p_{-j}^0)} \left| \frac{\partial q_j(p_j, p_{-j}^0)}{\partial p_j} \right| dp_j$$

- $D_{jk}(p_j, p_{-j}^0)$  is the MTE.
- Weights  $w(p_j) = \frac{1}{\Delta q_j} \frac{\partial q_j(p_j, p_{-j}^0)}{\partial p_j}$  correspond to the lost sales of  $j$  at a particular  $p_j$  as a fraction of all lost sales.
- When is  $LATE \approx ATE$ ?
  - Demand for Prius is steep: everyone leaves right away
  - $D_{j,k}(p_j)$  is relatively flat.
  - We might want to think about raising the price to choke price (or eliminating the product from the consumers choice set) same as treating everyone!

Abadie, Alberto and Matias D. Cattaneo. 2018. “Econometric Methods for Program Evaluation.” *Annual Review of Economics* 10 (1):465–503. URL <https://doi.org/10.1146/annurev-economics-080217-053402>.

Athey, Susan, Guido W Imbens, Stefan Wager et al. 2016. “Efficient inference of average treatment effects in high dimensions via approximate residual balancing.” Tech. rep.

Belloni, Alexandre, Victor Chernozhukov, Ivan Fernández-Val, and Christian Hansen. 2015. “Program evaluation with high-dimensional data.” Tech. rep., cemmap working paper, Centre for Microdata Methods and Practice.

Dobbie, Will, Jacob Goldin, and Crystal S. Yang. 2018. “The Effects of Pretrial Detention on Conviction, Future Crime, and Employment: Evidence from Randomly Assigned Judges.” *American Economic Review* 108 (2):201–40. URL <http://www.aeaweb.org/articles?id=10.1257/aer.20161503>.

Imbens, Guido W. 2015. "Matching methods in practice: Three examples." *Journal of Human Resources* 50 (2):373–419.

King, Gary and Richard Nielsen. Forthcoming. "Why Propensity Scores Should Not Be Used for Matching." *Political Analysis*