Identification (1)

Applied Econometrics for Spatial Economics

Hans Koster

Professor of Urban Economics and Real Estate







- 2. Research design
- 3. Summary

- Academics usually aim to identify *causal* effects
- <u>Causal effects</u>: one process, *a cause*, contributes to the production of another process
 - the effect of a 'treatment' variable *x* on an outcome variable *y*



- 1. Introduction
- 2. Research design
- 3. Summary

VU

VRIJE UNIVERSITEIT AMSTERDAM

Interesting correlations



- Spatial correlation of cholera deaths in 1854
 - John Snow
 - Contaminated water...

- 1. Introduction
- Research design
 Summary

Spurious correlations



 $\rho = 0.90$



- 1. Introduction
- Research design
 Summary

Spurious correlations



 $\rho = 0.79$



3. Summary

This week

- Setting up a research project
- Discuss why RCT measures an average causal effect of a treatment
- Alternatives to RCTs
 - OLS with controls
 - IV
 - Quasi-experimental methods



More economic reasoning than pure econometrics!

- 1. Introduction
- 2. Research design
- 3. Summary

This week

- Learn to set up your own research project
- ... and think about identification issues

• Plan:

Lecture #1:Research designLecture #2:Randomised experiments, OLS, IVLecture #3:Quasi-experiments, standard errorsAssignment:Estimate gravity models of trade



- 1. Introduction
- 2. <u>Research design</u>
- 3. Summary

- 8 steps when undertaking research
- 1. Formulate your <u>hypotheses</u>
- 2. Determine the '<u>treatment</u>' variable(s) and the '<u>outcome</u>' variable(s)
- 3. Think of an <u>identification strategy</u> to identify causal effects
- 4. <u>Select samples</u>, discuss <u>measurement error</u> and provide <u>descriptives</u>
- 5. Determine <u>functional form</u> of variables of interest
- 6. Think of different issues in estimating <u>standard</u> <u>errors</u>
- 7. <u>Estimate model and interpret</u> the results
- 8. Provide <u>robustness</u> checks of the results



- 2. <u>Research design</u>
- 3. Summary

1. Formulate your <u>hypotheses</u>

Economic hypotheses

Based on economic theory

- Humans often use *reverse* causal reasoning
 - *"House prices have gone down the last years, but why?"*
 - Forward causal inference supplies answers
 - <u>Reverse causal inference supplies questions</u>



- **1. Introduction**
- 2. <u>Research design</u>
- 3. Summary

- 2. Determine the '<u>treatment</u>' variable(s) and the '<u>outcome</u>' variable(s)
- Define what variables are available in your data
- Focus on one (or a few) x variable(s) and one (or a few y variables
- Think about expected order of magnitude



- 1. Introduction
- 2. <u>Research design</u>
- 3. Summary

- 3. Think of an <u>identification strategy</u> to identify causal effects
- What is your <u>'treatment' group</u> and what is your <u>'control' group</u>?

- Discuss endogeneity issues
 - Might there be a selection effect?
 - What are potential unobserved factors? Are these correlated with the treatment status?
 - Reverse causality?
 - (Measurement error?)



- **1. Introduction**
- 2. <u>Research design</u>
- 3. Summary

- 3. Think of an <u>identification strategy</u> to identify causal effects
- Define the appropriate econometric methods
 - Discuss the identifying assumptions at length!



- 2. <u>Research design</u>
- 3. Summary

- 4. <u>Select samples</u>, discuss <u>measurement error</u> and provide <u>descriptives</u>
- Should you use the full dataset?
- Variance in *x* is necessary!





- 1. Introduction
- 2. <u>Research design</u>
- 3. Summary

- 4. <u>Select samples</u>, discuss <u>measurement error</u> and provide <u>descriptives</u>
- Data cleaning





- **1. Introduction**
- 2. <u>Research design</u>
- 3. Summary

- 4. <u>Select samples</u>, discuss <u>measurement error</u> and provide <u>descriptives</u>
- Measurement error is present in many datasets



- 2. <u>Research design</u>
- 3. Summary

- 4. <u>Select samples</u>, discuss <u>measurement error</u> and provide <u>descriptives</u>
- Random measurement error in y is not so much of a problem



•
$$y_i^* - u_i = \beta x_i + \epsilon_i \Rightarrow y_i^* = \beta x_i + (\epsilon_i + u_i)$$

- 2. <u>Research design</u>
- 3. Summary

- 4. <u>Select samples</u>, discuss <u>measurement error</u> and provide <u>descriptives</u>
- Random measurement error in x biases the effect towards zero



• $y_i = \beta(x_i + u_i) + \epsilon_i \rightarrow \beta \rightarrow 0$ if u_i is large

- **1. Introduction**
- 2. <u>Research design</u>
- 3. Summary

- 5. Determine <u>functional form</u> of variables of interest
- The specification of $f(\cdot)$ is referred to as the functional form

 $y_i = f(x_i, c_i, \epsilon_i)$

• Often a linear functional form is assumed: $y_i = \beta x_i + \gamma c_i + \epsilon_i$



- **1. Introduction**
- 2. <u>Research design</u>
- 3. Summary

- 5. Determine <u>functional form</u> of variables of interest
- Economists are often interested in elasticities
 - Elasticity is percentage change of *y* in response to a change in *x*
 - $\frac{\partial y}{\partial x} \frac{x}{y}$
- They therefore often estimate log-linear regressions:

• **because**
$$\beta = \frac{\partial \log y}{\partial \log x} = \frac{\partial y}{\partial x}$$



- 1. Introduction
- 2. <u>Research design</u>
- 3. Summary

- 5. Determine <u>functional form</u> of variables of interest
- When use logs?
 - Economic theory
 - Residuals have a skewed distribution
 - Heteroscedasticity
 - Different unit sizes



- **1. Introduction**
- 2. <u>Research design</u>
- 3. Summary

- 6. Think of different issues in estimating <u>standard</u> <u>errors</u>
- Whether β is statistically significant depends on standard error
 - The smaller the standard error, the more precise your conclusions are

- Issues to bear in mind...
 - Should you cluster your standard errors?
 - Is heteroscedasticity a problem?
 - Is there serial/spatial autocorrelation?



- 1. Introduction
- 2. <u>Research design</u>
- 3. Summary

- 7. <u>Estimate model and interpret</u> the results
- Use statistical software to estimate your model

- Usually we are interested in <u>marginal effects</u>
 - How much does *y* change (in units or %) when *x* change with one unit (or %)

•
$$\frac{\partial y}{\partial x}$$
 (in levels) or $\frac{\partial y}{\partial x} \frac{x}{y}$ (in %)



- 2. <u>Research design</u>
- 3. Summary

- 7. <u>Estimate model and interpret</u> the results
- Properly interpret β and its statistical significance
 - "When x increases by 1 (units) y increases by
 .. (units). This effect is statistically significant
 at the ...% level."



- 1. Introduction
- 2. <u>Research design</u>
- 3. Summary

- 7. <u>Estimate model and interpret</u> the results
- Statistical hypothesis testing is dependent on statistical significance
- <u>Economic significance ≠ statistical significance</u>
 - A large effect may be imprecise
 - A small, but stat. sign. effect may be irrelevant

 <u>Always discuss both economic and statistical</u> <u>significance</u>



See McCloskey and Ziliak (1996)

- 2. <u>Research design</u>
- 3. Summary

7. <u>Estimate model and interpret</u> the results

- Make sure how your variables are measured
 - logs, dummies, etc.

Specifications:	X	$\log x$
У	$y = \rho x + \eta$	$y = \rho \log x + \eta$
	$\hat{\rho} = \frac{\partial y}{\partial x}$	$\hat{\rho} = \frac{\partial y}{\partial \log x}$
	$x \uparrow 1 \rightarrow y \uparrow \hat{\rho}$	$x \uparrow 1\% \to y \uparrow \hat{\rho}/100$
log y	$\log y = \rho x + \eta$	$\log y = \rho \log x + \eta$
log y	$\log y = \rho x + \eta$ $\hat{\rho} = \frac{\partial \log y}{\partial x}$	$\log y = \rho \log x + \eta$ $\hat{\rho} = \frac{\partial \log y}{\partial \log x}$



- **1. Introduction**
- 2. <u>Research design</u>
- 3. Summary

7. <u>Estimate model and interpret</u> the results

- Note on <u>dummy variables</u>. Let's assume the model $\log y = \rho x + \epsilon$, with $x \in 0,1$.
 - One unit increase in *x* is *large*
 - Halvorsen & Palmquist: $x \uparrow 1 \rightarrow y \uparrow ((e^{\hat{\rho}} 1) * 100)\%$



- **1. Introduction**
- 2. <u>Research design</u>
- 3. Summary

8. Provide <u>robustness</u> checks of the results

• You make many somewhat arbitrary choices

- Test for sensitivity of your results with respect to these choices
 - ... sensitivity analysis



2. Research design

3. Summary

Today:

- Economists are generally interested in *causal* effects
- 8 steps when undertaking research
 - 1. Formulate your <u>hypotheses</u>
 - 2. Determine the '<u>treatment</u>' variable(s) and the '<u>outcome</u>' variable(s)
 - 3. Think of an <u>identification strategy</u> to identify causal effects
 - 4. <u>Select samples</u>, discuss <u>measurement error</u> and provide <u>descriptives</u>
 - 5. Determine <u>functional form</u> of variables of interest
 - 6. Think of different issues in estimating <u>standard errors</u>
 - 7. <u>Estimate model and interpret</u> the results
 - 8. Provide <u>robustness</u> checks of the results



- Research design
 <u>Summary</u>

Tomorrow:

- Randomised experiments
- OLS with controls
- Instrumental variables



Identification (1)

Applied Econometrics for Spatial Economics

Hans Koster

Professor of Urban Economics and Real Estate







Identification (2)

Applied Econometrics for Spatial Economics

Hans Koster

Professor of Urban Economics and Real Estate







- 1. Introduction
- 2. Randomised experiments
- 3. OLS with controls
- 4. IV
- 5. Summary

This week

- Learn to set up your own research project
- ... and think about identification issues

• Plan:

Research desig

Lecture #2: Lecture #3: Assignment:

Randomised experiments, OLS, IV

Quasi-experiments, standard errors Estimate gravity models of trade



- 1. Introduction
- 2. Randomised experiments
- 3. OLS with controls
- 4. IV
- 5. Summary

- 8 steps when undertaking research
- **1. Formulate your <u>hypotheses</u>**
- 2. Determine the '<u>treatment</u>' variable(s) and the '<u>outcome</u>' variable(s)
- 3. Think of an <u>identification strategy</u> to identify causal effects
- 4. <u>Select samples</u>, discuss <u>measurement error</u> and provide <u>descriptives</u>
- 5. Determine <u>functional form</u> of variables of interest
- 6. Think of different issues in estimating <u>standard</u> <u>errors</u>
- 7. <u>Estimate</u> model and <u>interpret</u> the results
- 8. Provide <u>robustness</u> checks of the results



- 1. Introduction
- 2. Randomised experiments
- 3. OLS with controls
- 4. IV
- 5. Summary

- In economics, identification of causal effects is of key importance
 - Step 3 is key → "think of an identification strategy to identify causal effects"

- Possible identification strategies
 - 1. <u>Randomised experiments</u>
 - 2. Exhaustive set of controls
 - 3. Instrumental variables (IV)
 - 4. Quasi-experiments (QE)
 - **Generation-discontinuity designs (RDD)**



- **1.** Introduction
- 2. Randomised experiments
- 3. OLS with controls
- 4. IV
- 5. Summary

 Let's first consider why randomised experiments are ideal in identifying causal effects

- We often have to rely on (imperfect) alternatives to RCTs
 - You have to think of <u>an identification strategy</u>



- **1.** Introduction
- 2. <u>Randomised experiments</u>
- 3. OLS with controls
- 4. IV
- 5. Summary

Ideal RCT

 We want to know the impact of housing subsidy on well-being

•
$$y_i = \begin{cases} y_{1i} & \text{if } x_i = 1 \\ y_{0i} & \text{if } x_i = 0 \end{cases}$$

•
$$y_i = y_{0i} + (y_{1i} - y_{0i}) \times x_i$$

- $(y_{1i} y_{0i})$ is a causal effect of housing subsidy
- You do not observe individuals twice


- **1.** Introduction
- 2. Randomised experiments
- 3. OLS with controls
- 4. IV
- 5. Summary

Ideal RCT

• The average effect is given by: $E[y_i|x_i = 1] - E[y_i|x_i = 0] =$ $E[y_{1i}|x_i = 1] - E[y_{0i}|x_i = 1] + ATT$ $E[y_{0i}|x_i = 1] - E[y_{0i}|x_i = 0] Selection bias$

 Selection: people who are eligible for housing subsidies have lower incomes and in turn lower well-being levels

Solution: <u>randomisation</u>



- **1. Introduction**
- 2. <u>Randomised experiments</u>
- 3. OLS with controls
- 4. IV
- 5. Summary

Ideal RCT

- The average effect is given by: $E[y_i|x_i = 1] - E[y_i|x_i = 0] =$ $E[y_{1i}|x_i = 1] - E[y_{0i}|x_i = 1] + ATT$ $E[y_{0i}|x_i = 1] - E[y_{0i}|x_i = 0] Selection bias$
- <u>Randomisation</u>: $E[y_{0i}|x_i = 1] = E[y_{0i}|x_i = 0]$
- Hence: $E[y_i|x_i = 1] - E[y_i|x_i = 0] =$ $E[y_{1i}|x_i = 1] - E[y_{1i}|x_i = 0]$
- RCTs solve the selection problem!

38

- **1. Introduction**
- 2. Randomised experiments
- 3. OLS with controls
- 4. IV
- 5. Summary

Ideal RCT

- RCTs are only informative on the mean effect of the treatment
 - Not always interesting

- Paradox: <u>the estimated effect may apply to no one</u>
 - Imagine: two subgroups for which one the treatment is effective while for the other not
 - ATT do not apply to either one of them



- **1. Introduction**
- 2. <u>Randomised experiments</u>
- 3. OLS with controls
- 4. IV
- 5. Summary

- Alternative interesting measures
 - Median treatment effect
 - Share of people that respond positively to treatment

- Analyse the ATT for different subgroups
 - But: requires a substantially large dataset
 - And information on individual characteristics



- **1. Introduction**
- 2. <u>Randomised experiments</u>
- 3. OLS with controls
- 4. IV
- 5. Summary

- A more philosophical critique on RCTs
 - We might find a causal effect of x on y, but <u>do</u> <u>not know why there is an effect</u>
- Theoretical models and reasoning are needed to explain why we would expect a causal effect
 - Deaton (2010)



- **1.** Introduction
- 2. <u>Randomised experiments</u>
- 3. OLS with controls
- 4. IV
- 5. Summary

Analogy of RCTs with regression

- **Recall:** $y_i = y_{0i} + (y_{1i} y_{0i}) \times x_i$
- $y_i = \alpha + \beta x_i + \epsilon_i$ $\alpha = E[y_{0i}]$ $\beta = y_{1i} - y_{0i}$ $\epsilon_i = y_{0i} - E[y_{0i}]$
- β is the treatment effect
- So:
 - $E[y_i|x_i = 1] E[y_i|x_i = 0] = \beta +$ ATT $E[\epsilon_i|x_i = 1] - E[\epsilon_i|x_i = 0]$ Selection bias
 - When there is a selection problem, ϵ_i is correlated with x_i



- **1. Introduction**
- 2. <u>Randomised experiments</u>
- 3. OLS with controls
- 4. IV
- 5. Summary

Analogy of RCTs with regression

- <u>Use control variables to 'balance' control and</u> <u>treatment groups</u>
- $y_i = \alpha + \beta x_i + \gamma c_i + \epsilon_i$
 - If x_i is uncorrelated to c_i , β will be identical
 - Standard errors go down

• *Conditional* on c_i , x_i should be uncorrelated to ϵ_i



- **1. Introduction**
- 2. <u>Randomised experiments</u>
- 3. OLS with controls
- 4. IV
- 5. Summary

- In most economic studies, RCTs are not applied
 - No experimental setting possible
 - Ethical reasons / fairness
 - Costly
 - Expected substantial heterogeneity in outcomes
 - Hard to measure long-run effects
 - Lab setting may bias outcomes
 - » Recall: biases in Stated Preference surveys



- 1. Introduction
- 2. Randomised experiments
- 3. <u>OLS with controls</u>
- 4. IV
- 5. Summary

- Possible identification strategies
 - L. Randomised experiments
 - 2. <u>Exhaustive set of controls</u>
 - 3. Instrumental variables (IV)
 - 4. Quasi-experiments (QE)
 - └→ Regression-discontinuity designs (RDD)



- **1. Introduction**
- 2. Randomised experiments
- 3. OLS with controls
- 4. IV
- 5. Summary

- <u>Use an exhaustive set of controls</u>
 - In some applications, you might know all explanatory variables

- For example, computers?
 - You aim to know the willingness to pay for a new processor
 - $price_i = \rho(processor \ quality_i) + (characteristics)'_i \gamma + \eta_i$

- Not all characteristics are available in the data
 - Houses, cars, etc.
 - Omitted variable bias...

46



- **1. Introduction**
- 2. Randomised experiments
- 3. OLS with controls
- 4. IV
- 5. Summary

- Use first-differencing or fixed effects to make this approach more convincing
 - Controls for all time-invariant factors
 - <u>Requires 'within' variation</u>



- **1. Introduction**
- 2. Randomised experiments
- 3. <u>OLS with controls</u>
- 4. IV
- 5. Summary

• <u>First-differencing</u> $\Delta y_{it} = \Delta \alpha_t + \beta \Delta x_{it} + \gamma \Delta c_{it} + \Delta \epsilon_i$ where $\Delta y_{it} = y_{it} - y_{it-1}$, etc.

- This controls for all time-invariant characteristics of i
 - Hence there should be variation in *x_{it}* over time



- **1.** Introduction
- 2. Randomised experiments
- 3. OLS with controls
- 4. IV
- 5. Summary

<u>Fixed effects</u>

 $y_{ig} = \bar{y}_i + \beta (x_{ig} - \bar{x}_g) + \gamma (c_{ig} - \bar{c}_g) + (\epsilon_{gt} - \bar{\epsilon}_g)$ $= \beta x_{ig} + \gamma c_{ig} + \mu_g + \epsilon_{ig}$

where μ_g is a fixed effect at the level of group g

Fixed effects vs. first-differencing

- Identical to first-differencing when having two observations per group
- Fixed effects is more efficient...
- ... but exogeneity assumption is stronger



- **1. Introduction**
- 2. Randomised experiments
- 3. OLS with controls
- 4. <u>IV</u>
- 5. Summary

Possible identification strategies

- 1. Randomised experiments
- 2. Exhaustive set of controls
- 3. Instrumental variables (IV)
- 4. Quasi-experiments (QE)
 - └→ Regression-discontinuity designs (RDD)



- **1. Introduction**
- Randomised experiments
 OLS with controls

- 4. <u>IV</u> 5. Summary

Find valid instrumental variables

•
$$x_i = \zeta + \eta z_i + \xi_i$$
 (1st stage)
 $y_i = \alpha + \beta \hat{x}_i + \epsilon_i$ (2nd stage)



- _ _ _
- 1. Introduction
- 2. Randomised experiments
- 3. OLS with controls
- 4. <u>IV</u>
- 5. Summary

- There are two conditions for valid instruments
- I. <u>Instrument relevance</u>: $cov[z_i, x_i] \neq 0$
 - μ should be statistically significant and strong
 - **Rule-of-thumb**: F > 10
 - Use Kleibergen-Paap *F*-statistic with multiple endogenous variables
- II. <u>Instrument exogeneity</u>: $cov[z_i, \epsilon_i] = 0$
 - Instrument should not be correlated to error term
 - **Instrument should only influence** *y* **via** *x*
 - Based on economic reasoning



- **1. Introduction**
- 2. Randomised experiments
- 3. OLS with controls
- 4. <u>IV</u>
- 5. Summary

- <u>Use exogenous sources of variation</u>
 - Use economic models to find valid instruments
 - Use national policies or natural shocks, etc. to construct instrument
 - Use historical/long-lagged instruments



- **1. Introduction**
- 2. Randomised experiments
- 3. OLS with controls
- 4. <u>IV</u>
- 5. Summary

4. IV

- Say that the instrument is only observed for a certain group, then IV identifies treatment effect for this group
- Different instruments may lead to different β

• Example: gender of children and housing demand



- **1. Introduction**
- 2. Randomised experiments
- 3. OLS with controls
- 4. IV
- 5. <u>Summary</u>

Today:

- We focus on step 3: "think of an identification strategy to identify causal effects"
- Examples of possible strategies
 - Randomised experiments
 - OLS with controls
 - Instrumental variables



- **1. Introduction**
- 2. Randomised experiments
- 3. OLS with controls
- 4. IV
- 5. <u>Summary</u>

Tomorrow:

- We continue to focus on step 3: "think of an identification strategy to identify causal effects"
 - Quasi-experiments including RDDs
- Issues with standard errors



Identification (2)

Applied Econometrics for Spatial Economics

Hans Koster

Professor of Urban Economics and Real Estate







Identification (3)

Applied Econometrics for Spatial Economics

Hans Koster

Professor of Urban Economics and Real Estate







- 1. Introduction
- 2. Quasi-experiments
- 3. RDD
- 4. Standard errors
- 5. Summary

This week

- Learn to set up your own research project
- ... and think about identification issues

• Plan:

: Research desi

Lecture #2: Lecture #3: Assignment:

Randomised experiments, OLS, IV Quasi-experiments, standard errors Estimate gravity models of trade



8 steps when undertaking research

- 1. Introduction
- 2. Quasi-experiments
- 3. RDD
- 4. Standard errors
- 5. Summary

1. Formulate your <u>hypotheses</u>

- 2. Determine the '<u>treatment</u>' variable(s) and the '<u>outcome</u>' variable(s)
- 3. Think of an <u>identification strategy</u> to identify causal effects
- 4. <u>Select samples</u>, discuss <u>measurement error</u> and provide <u>descriptives</u>
- 5. Determine <u>functional form</u> of variables of interest
- 6. Think of different issues in estimating <u>standard</u> <u>errors</u>
- 7. <u>Estimate</u> model and <u>interpret</u> the results
- 8. Provide <u>robustness</u> checks of the results

60

- **1. Introduction**
- 2. Quasi-experiments
- 3. RDD
- 4. Standard errors
- 5. Summary

2. Quasi-experiments

Possible identification strategies

- 1. Randomised experiments
- 2. Exhaustive set of controls
- **3. Instrumental variables (IV)**
- 4. Quasi-experiments (QE)
 - └ <u>Regression-discontinuity designs (RDD)</u>



- **1.** Introduction
- 2. Quasi-experiments
- 3. RDD
- 4. Standard errors
- 5. Summary

- Use <u>exogenous shocks</u> in the economy to identify causal effects
 - <u>'Quasi'-experiments</u> / natural experiments

- National policy changes, (arbitrary) policy rules, earthquakes, bombings
 - These shocks cannot be influenced by the individual decision makers
 - Recall: if shock is really random, selection effect is equal to zero
 - The <u>research context</u> indicates whether shock is indeed random



- **1.** Introduction
- 2. Quasi-experiments
- 3. <u>RDD</u>
- 4. Standard errors
- 5. Summary

- Regression-discontinuity design (RDD)
 - Quasi-experimental method

 Assume that we have a treatment effect that is dependent on r_i:

$$x_i = \begin{cases} 1 & \text{if } r_i \ge r_0 \\ 0 & \text{if } r_i < r_0 \end{cases}$$

• r_0 is some cutoff value

• This leads to a regression:

 $y_i = \alpha + \beta x_i + \gamma r_i + \epsilon_i$

• Note that x_i is a fully deterministic function of

 r_i

• Not perfectly collinear because r_i is continuous

- 1. Introduction
- 2. Quasi-experiments
- 3. <u>RDD</u>
- 4. Standard errors
- 5. Summary

- Example:
 - Students get a scholarship if they achieve a certain test-score
 - You aim to know the impact of the scholarship on job market outcomes
 - » e.g. wages



- **1. Introduction**
- Quasi-experiments
 <u>RDD</u>
- 4. Standard errors
- 5. Summary

Plot



Control for test scores and investigate the jump in treatment at r₀



- **1. Introduction**
- Quasi-experiments
 <u>RDD</u>
- 4. Standard errors
- 5. Summary

Plot



Control for test scores and investigate the jump in treatment at r₀



- **1.** Introduction
- 2. Quasi-experiments
- 3. <u>RDD</u>
- 4. Standard errors
- 5. Summary

- What if *x* is non-linearly related to *Y*
 - $y_i = \alpha + \beta x_i + f(r_i) + \epsilon_i$
 - **Specify** $f(r_i) = \gamma_1 r_i + \gamma_2 r_i^2 + \dots + \gamma_q r_i^q$
 - » *q*th-order polynomial
 - » <u>Can be estimated by OLS</u>





- **1. Introduction**
- 2. Quasi-experiments 3. <u>RDD</u>
- 4. Standard errors
- 5. Summary

To <u>check for nonlinearities</u> in a RDD is important





- 1. Introduction
- 2. Quasi-experiments
- 3. <u>RDD</u>
- 4. Standard errors
- 5. Summary

- To check for nonlinearities in a RDD is important
 - To reduce the possibility of mistakes, you may focus on observations 'close' to r_0
 - Reduces precision





- 1. Introduction
- 2. Quasi-experiments
- 3. <u>RDD</u>
- 4. Standard errors
- 5. Summary

- Two different versions
 - <u>Sharp RDD</u> \rightarrow Jump in treatment
 - <u>Fuzzy RDD</u> \rightarrow Jump in probability of treatment

Previous slides: sharp RDD

- Fuzzy RDDs are very common
 - Assignment is often 'fuzzy'



- **1.** Introduction
- Quasi-experiments
 <u>RDD</u>
- 4. Standard errors
- 5. Summary

Illustration of a fuzzy RDD





- 1. Introduction
- 2. Quasi-experiments
- 3. <u>RDD</u>
- 4. Standard errors
- 5. Summary

- Fuzzy RDD
 - Prob $[x_i = 1 | r_i] = \begin{cases} g_1(r_i) & \text{if } r_i \ge r_0 \\ g_0(r_i) & \text{if } r_i < r_0 \end{cases}$ where $g_1(r_i) \ne g_0(r_i)$
- Prob $[x_i = 1 | r_i] = g_0(r_i) + [g_1(r_i) g_0(r_i)]z_i$ • $z_i = \mathbb{I}(r_i \ge r_0)$
- Looks complicated it just means that treatment probability is discontinuous at some point

- This leads to a two-stage least squares estimator
 - **First stage** $\rightarrow x_i = \zeta + \eta z_i + g(r_i) + \xi_i$
 - **Second stage** $\rightarrow y_i = \alpha + \beta \hat{x}_i + f(r_i) + \epsilon_i$
- **1. Introduction**
- 2. Quasi-experiments
- 3. RDD
- 4. Standard errors
- 5. Summary

- This course has almost entirely focused on estimating average effects
- But how to assess statistical significance?
 - <u>Use correctly estimated standard errors</u>
 - May be very important!
- Some issues with standard errors
 - 1. Heteroscedasticity
 - 2. Clustering
 - 3. Serial correlation



- **1. Introduction**
- 2. Quasi-experiments
- 3. RDD
- 4. <u>Standard errors</u>
- 5. Summary

- 1. <u>Heteroscedasticity</u>
- Conditional variance of y_i given x_i changes with i





- **1. Introduction**
- 2. Quasi-experiments
- 3. RDD
- 4. Standard errors
- 5. Summary

- 1. <u>Heteroscedasticity</u>
- To estimate standard errors, we typically assume homoscedasticity
- Solution: use robust standard errors
 - In STATA, type *r* after REGRESS
 - This leads to consistent s.e.
- However, robust standard errors are biased
 - <u>Only a problem in small samples</u>



- **1. Introduction**
- 2. Quasi-experiments
- 3. RDD
- 4. Standard errors
- 5. Summary

- Imagine that you have some aggregate variable x_g
 - Let's say the municipal tax rate
- You aim to investigate the impact on individual outcomes, let's say well-being
- You may regress $y_g = \alpha + \beta x_g + \epsilon_i$
- However, the data is at the individual level *i*, so you aim to regress: $y_{ig} = \alpha + \beta x_g + \epsilon_{ig}$



- **1.** Introduction
- 2. Quasi-experiments
- 3. RDD
- 4. Standard errors
- 5. Summary

- **Issue** $y_{ig} = \alpha + \beta x_g + \epsilon_{ig}$
 - You basically multiply the size of the dataset leading to artificially low standard errors
 - The effective number of obs is much lower
 - Can make a big difference!
- More generally, to obtain consistent standard errors, you assume that $E[\epsilon_{ig}\epsilon_{jg}] = 0$
 - This is certainly not the case in the above example
- Solution: <u>cluster your standard errors at the</u> <u>appropriate level</u>



- **1.** Introduction
- 2. Quasi-experiments
- 3. RDD
- 4. Standard errors
- 5. Summary

- This is the formula for the standard error when $\mathbb{E}[\epsilon_{ig}\epsilon_{jg}] = 0$: $SE(\hat{\beta}) = \frac{\sigma_{\epsilon}}{\sqrt{N}} \frac{1}{\sigma_{r}}$
- Let's for simplicity assume that everyone in the municipality has the same well-being
 - Then the correct standard error is:

$$SE(\hat{\beta}) = \frac{\sigma_{\epsilon}}{\sqrt{G}} \frac{1}{\sigma_{x}}$$

- Say that you have 9,000 individuals but only 90 municipalities
 - Standard error is 10 times larger

- **1. Introduction**
- 2. Quasi-experiments
- 3. RDD
- 4. Standard errors
- 5. Summary

- Not always clear at what level you should cluster
 - ...when different variables are aggregated at different levels
 - Pragmatic approach: choose standard errors that lead to the most conservative conclusions (→ <u>highest standard errors</u>)
 - Use multi-way clustering
 - \rightarrow In REGHDFE command in Stata



 Note: <u>clustered standard errors are not correct</u> for a few number of clusters

- **1. Introduction**
- 2. Quasi-experiments
- 3. RDD
- 4. Standard errors
- 5. Summary

- 3. <u>Serial correlation</u>
- Time series specifications:
 - $\Delta y_{it} = \alpha_t + \rho \Delta x_{it} + \Delta X'_{it} \gamma + \Delta \eta_i$

- Same problem as previously: $\mathbb{E}[\epsilon_{it}\epsilon_{it-1}] \neq 0$
- Solution: cluster at individual level *i*?
 - **But:** $\mathbb{E}[\epsilon_{it}\epsilon_{it-1}] \neq \mathbb{E}[\epsilon_{it}\epsilon_{it-2}]$

- This issue is still under study!
 - Two-way clustering may be a solution



- **1. Introduction**
- 2. Quasi-experiments
- 3. RDD
- 4. Standard errors
- 5. Summary

- Statistical hypothesis testing is dependent on statistical significance
- Recall that <u>economic significance ≠ statistical</u> <u>significance</u>
 - A large effect may be imprecise
 - A small, but stat. sign. effect may be irrelevant

 <u>Always discuss both economic and statistical</u> <u>significance</u>



- **1. Introduction**
- 2. Quasi-experiments
- 3. RDD
- 4. Standard errors
- 5. <u>Summary</u>

Today:

- Several identification strategies have been discussed to measure causal effects
 - Quasi-experiments
 - \rightarrow RDD

- Some remarks on standard errors
 - Heteroscedasticity is not so important
 - Clustering is important
 - Economic vs. statistical significance



- **1.** Introduction
- 2. Quasi-experiments
- 3. RDD
- 4. Standard errors
- 5. <u>Summary</u>

- Friday, October 29, 12:15-14:15
- 3 questions on each of the topics
- Let's discuss a brief exam 'check-list'



- **1. Introduction**
- 2. Quasi-experiments
- 3. RDD
- 4. Standard errors
- 5. <u>Summary</u>

- Spatial econometrics
 - What is special about spatial data?
 - What is a spatial weight matrix?
 - How to measure spatial autocorrelation?
 - What is difference between global and local autocorrelation?
 - Understand the different spatial econometric models
 - What is the MAUP?



- **1.** Introduction
- 2. Quasi-experiments
- 3. RDD
- 4. Standard errors
- 5. Summary

- Discrete choice
 - When to use discrete choice methods?
 - Understand the random utility framework and why it forms the foundation of econometric applications
 - When to use binary discrete choice models?
 → LPM, Logit, Probit
 - Derive marginal effects
 - When to use multinomial discrete choice models?
 - → MNL with alternative-specific coefficients
 - → Conditional logit
 - \rightarrow Poisson
 - When to rely on RP and/or SP data?



- **1. Introduction**
- 2. Quasi-experiments
- 3. RDD
- 4. Standard errors
- 5. Summary

- Identification
 - What is a causal effect?
 - What are different steps in a research design?
 - Consider different identification strategies

5. Summary

- \rightarrow RCTs
- → Controls/fixed effects
- \rightarrow IV
- → Quasi-experiments
- \rightarrow RDD
- How to get correct standard errors
 - → Heteroscedasticity
 - → Clustering
 - \rightarrow Serial correlation
- What is the difference between economic and statistical significance?

- 1. Introduction
- 2. Quasi-experiments
- 3. RDD
- 4. Standard errors
- 5. <u>Summary</u>

• A concluding quote from *Mostly harmless* econometrics:

"Econometrics applied to coherent causal questions, regressions and 2SLS almost always make sense. Your standard errors probably won't be quite right, but they rarely are. Avoid embarrassment by being your own best skeptic, and especially, DON'T PANIC!"

- For any remaining questions, please drop me an email (<u>h.koster@vu.nl</u>) for an appointment
 - Don't wait until the day before the exam



Identification (3)

Applied Econometrics for Spatial Economics

Hans Koster

Professor of Urban Economics and Real Estate





