

Regression Discontinuity Design:  
A test for manipulation of the running variable  
Research Module in Econometrics and Statistics  
WS 2019/20

Sofia Badini, Max Schäfer, Caroline Kraymer

University of Bonn

January 22, 2020

# Outline

1. Manipulation Test
2. Properties and Parameter Choice
3. Simulation Study
4. Real Data Application
5. Conclusions

## Manipulation Test

## Fundamental Problem of Causal Inference

- Identifying effect from cause requires us to ask:  
*What would have been if a certain state had not occurred?*
- Individual unit in one state only (treatment or control)

## Fundamental Problem of Causal Inference

- Identifying effect from cause requires us to ask:  
*What would have been if a certain state had not occurred?*
- Individual unit in one state only (treatment or control)
- Regression Discontinuity Design (RDD) is quasi-experimental research design creating local random assignment in treatment and control
- Treatment assignment changes discontinuously as function of running variable

## Fundamental Problem of Causal Inference

- Identifying effect from cause requires us to ask:  
*What would have been if a certain state had not occurred?*
- Individual unit in one state only (treatment or control)
- Regression Discontinuity Design (RDD) is quasi-experimental research design creating local random assignment in treatment and control
- Treatment assignment changes discontinuously as function of running variable
- In absence of manipulation, treatment and control units locally very similar
- Can recover ATE, weighted by probability to be close to cutoff

## Manipulation of Running Variable

- With perfect manipulation treatment and control units locally not similar anymore
- Counterfactual outcomes and treatment effect not identifiable
- Perfect manipulation if individuals (i) have precise control of running variable, (ii) know cutoff and (iii) benefit from treatment

## Manipulation of Running Variable

- With perfect manipulation treatment and control units locally not similar anymore
- Counterfactual outcomes and treatment effect not identifiable
- Perfect manipulation if individuals (i) have precise control of running variable, (ii) know cutoff and (iii) benefit from treatment
- Need to know assignment rule and context very well
- Nevertheless, testing for manipulation is inevitable



## A Manipulation Test (McCrary, 2008)

- With manipulation we expect many barely qualifying and few barely not qualifying for treatment
- Idea: test for discontinuity in density function of running variable at cutoff

## A Manipulation Test (McCrary, 2008)

- With manipulation we expect many barely qualifying and few barely not qualifying for treatment
- Idea: test for discontinuity in density function of running variable at cutoff
- Estimation at cutoff for data approaching cutoff from above and below separately:
  - Step 1: Bin data and create histogram
  - Step 2: Smooth with local linear regression

## Local Linear Density Estimation: First Step

- To deploy regression methods we need regressor and response variable
- Bin data, count observations in bins, center and normalize
  - Returns bin midpoints and bin heights, i.e. histogram

## Local Linear Density Estimation: First Step

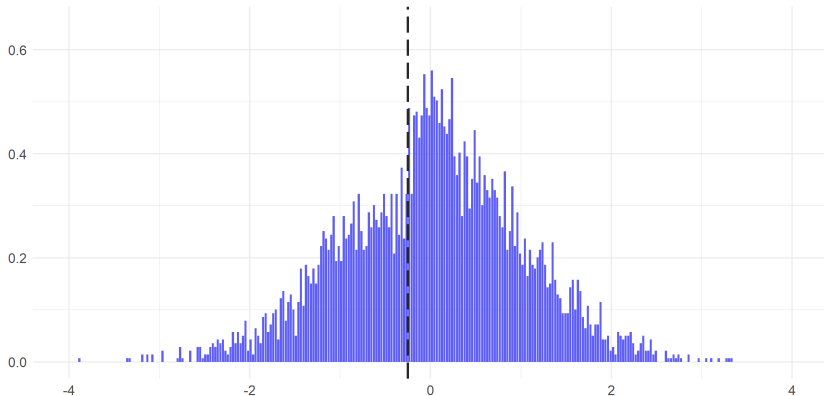
- To deploy regression methods we need regressor and response variable
- Bin data, count observations in bins, center and normalize
  - Returns bin midpoints and bin heights, i.e. histogram
- Binsize is first smoothing parameter we choose
  - Depends on number of observations and variance

## First Step: Histogram

### Local Linear Density Estimation

5000 observations in 260 bins

Cutoff at -0.25 and share of manipulators is 0.2



## Local Linear Density Estimation: Second Step

- Use non-parametric methods to smooth histogram
- Local linear regression essentially Weighted Least Squares
- At point of interest, fit linear model to weighted data
  - Weight increases the closer observation is to point of interest
  - Zero weight for data on other side of cutoff
- Second smoothing parameter is bandwidth, which decides range of data to use for estimation

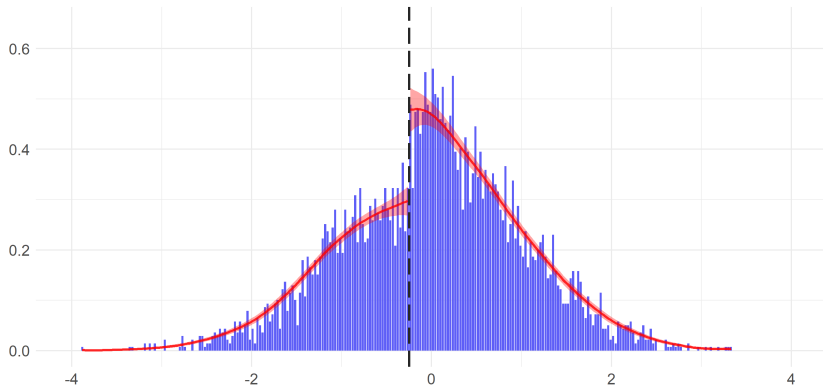
## Second Step: Local Linear Smoothing

### Local Linear Density Estimation

5000 observations in 260 bins

Cutoff at -0.25 and share of manipulators is 0.2

Bandwidth: 0.767



## Estimator and Wald Test

- Estimator is transformed difference of density estimates  $\hat{f}(r)$  approaching cutoff from above and below:

$$\hat{\theta} = \ln(\lim_{r \downarrow c} \hat{f}(r)) - \ln(\lim_{r \uparrow c} \hat{f}(r))$$



## Estimator and Wald Test

- Estimator is transformed difference of density estimates  $\hat{f}(r)$  approaching cutoff from above and below:

$$\hat{\theta} = \ln(\lim_{r \downarrow c} \hat{f}(r)) - \ln(\lim_{r \uparrow c} \hat{f}(r))$$

- Test  $H_0 : \theta = 0$  with Wald test:

$$\frac{\hat{\theta}}{\hat{\sigma}_{\theta}} \xrightarrow{d} \mathcal{N}(0, 1)$$

## Properties and Parameter Choice

## Properties of Estimator

- Asymptotic normality
- Consistency (with correct choice of bandwidth)

## Properties of Estimator

- Asymptotic normality
- Consistency (with correct choice of bandwidth)

Problems:

- Non-parametric setting: need large sample size
- Asymptotic bias with incorrect choice of bandwidth

## Parameter Choice

- Binsize: determines appearance of the first-step histogram
- Bandwidth: determines size of the local neighborhood used for estimation in step 2
- McCrary uses automatic selection procedures and subjectively adjusts afterwards

## Parameter Choice

- Binsize: determines appearance of the first-step histogram
- Bandwidth: determines size of the local neighborhood used for estimation in step 2
- McCrary uses automatic selection procedures and subjectively adjusts afterwards

Problem:

- Automatically selected bandwidth leads to asymptotic bias
  - Need to shrink (undersmooth) obtained bandwidth
  - Optimal amount of undersmoothing is unknown

## Simulation Study

## Data-Generating Process

Idea:

- Assess test's power when increasing discontinuity gap
- Assess validity of asymptotic normality



## Data-Generating Process

Idea:

- Assess test's power when increasing discontinuity gap
- Assess validity of asymptotic normality

Data-generating process:

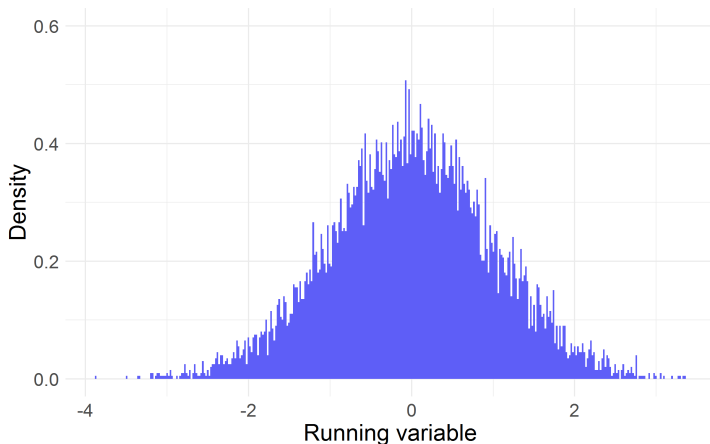
- Draw i.i.d. running variable values from fixed distribution
- Let individuals in fixed region below cutoff flip an unfair coin
- In case of “success”: assign additional value of running variable leading to a change in treatment status

## Data-Generating Process

Density of the running variable without manipulation

Standard normal distribution

10000 observations

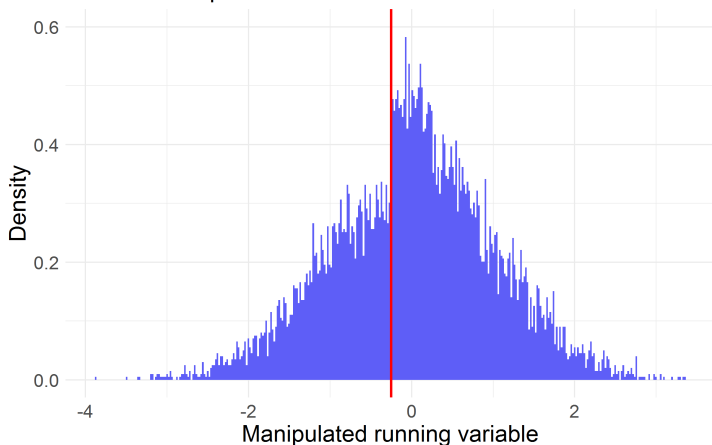


## Data-Generating Process

Density of the manipulated running variable

Cutoff -0.25

Share of manipulators 0.2



## Power of Test

- Investigate test's ability to detect presence of manipulation
- Sample non-manipulated running variable values from different probability distributions

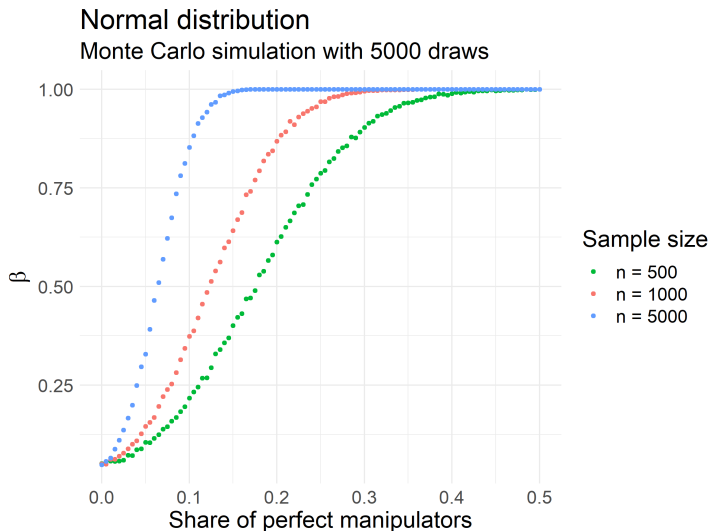
## Power of Test

- Investigate test's ability to detect presence of manipulation
- Sample non-manipulated running variable values from different probability distributions

Expect:

- Power approaches one when increasing manipulation
- Faster convergence for larger sample sizes

## Power of Test



## Asymptotic Normality

- Assess quality of test statistic's normal approximation
- Use quantile-quantile plots to plot standard normal quantiles against quantiles of test distribution

## Asymptotic Normality

- Assess quality of test statistic's normal approximation
- Use quantile-quantile plots to plot standard normal quantiles against quantiles of test distribution

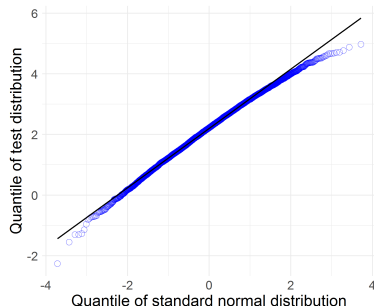
Expect:

- Normality becomes more valid with increasing sample size
- Normality becomes more valid when undersmoothing bandwidth

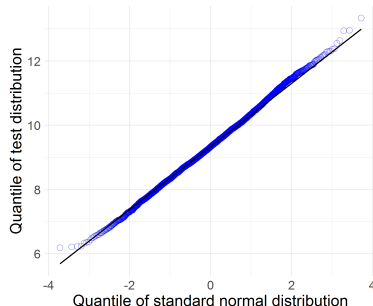


# Asymptotic Normality

Normal distribution  
500 observations, 5000 draws

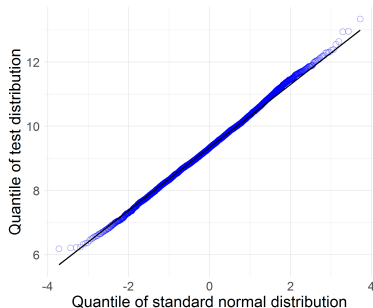


Normal distribution  
20000 observations, 5000 draws

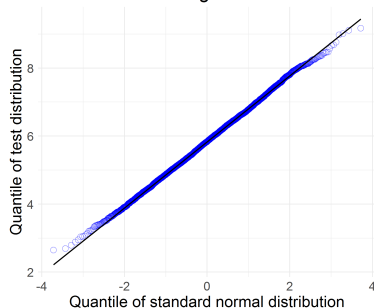


# Asymptotic Normality

Normal distribution  
20000 observations, 5000 draws



Normal distribution  
20000 observations, 5000 draws  
0.5 undersmoothing



## Real Data Application

## Manipulation of Social Program Eligibility in Colombia: the Case of SISBEN

- Camacho and Conover (2011), published in *American Economic Journal: Economic Policy*
- Manipulation of Census of the Poor, a score-based targeting system for social welfare programs in Colombia

## Manipulation of Social Program Eligibility in Colombia: the Case of SISBEN

- Camacho and Conover (2011), published in *American Economic Journal: Economic Policy*
- Manipulation of Census of the Poor, a score-based targeting system for social welfare programs in Colombia
- Large-scale corruption at local government level and mass manipulation of data at entry stage
- Estimated costs of corruption: roughly 7% of National Health and Social Security budget

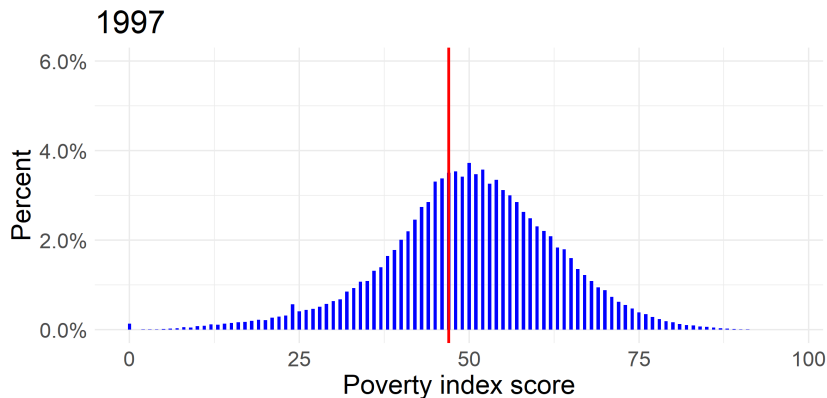
## The Census of the Poor (SISBEN)

- National survey managed by municipalities which collects comprehensive information on long-term quality of life
- Used to assign poverty index score, from 0 (poorest) to 100 (least poor), determining eligibility to social welfare programs

## The Census of the Poor (SISBEN)

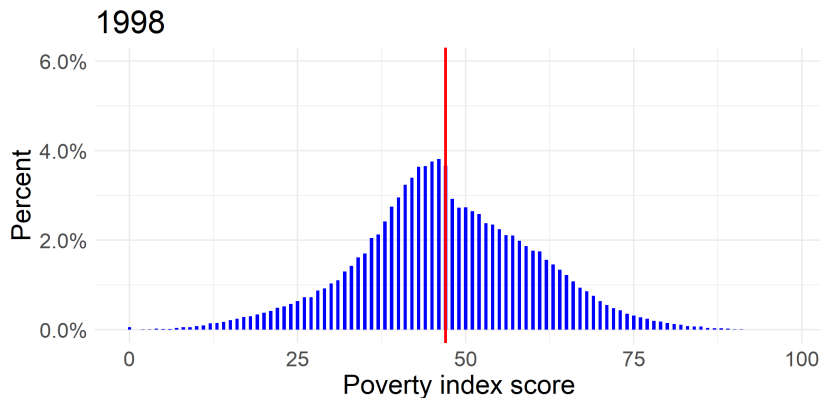
- National survey managed by municipalities which collects comprehensive information on long-term quality of life
- Used to assign poverty index score, from 0 (poorest) to 100 (least poor), determining eligibility to social welfare programs
- 1993-2004: first SISBEN, most common cutoff for urban families is 47
- Algorithm revealed to local officers in 1997

## Score Distribution Over Time

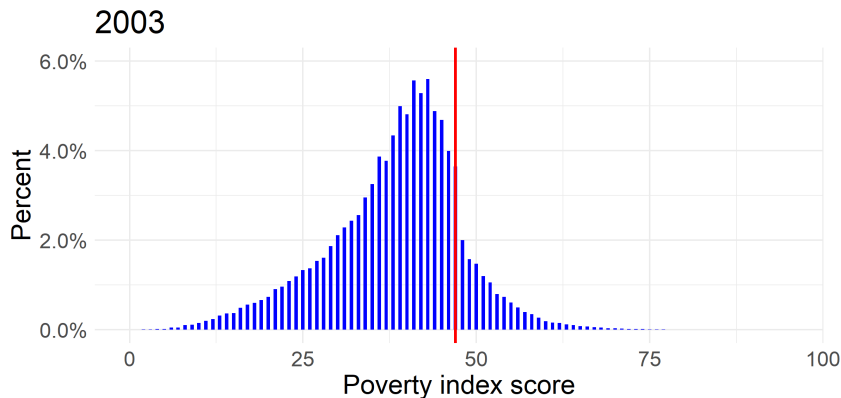




## Score Distribution Over Time



## Score Distribution Over Time



## Size of Discontinuity at Cutoff

McCRARY'S DENSITY TEST

Year	Pilot Bandwidth		50% Undersmoothing	
	Estimator	SE	Estimator	SE
1994	-0.010	[0.007]	0.014	[0.010]
1995	-0.026***	[0.004]	0.004	[0.006]
1996	-0.002	[0.006]	-0.001	[0.009]
1997	-0.006	[0.005]	-0.016**	[0.007]
1998	-0.256***	[0.005]	-0.219***	[0.007]
1999	-0.314***	[0.004]	-0.271***	[0.006]
2000	-0.374***	[0.003]	-0.313***	[0.004]
2001	-0.510***	[0.003]	-0.450***	[0.005]
2002	-0.571***	[0.004]	-0.486***	[0.006]
2003	-0.542***	[0.008]	-0.459***	[0.012]

\*\*\* Significant at 1 percent level. \*\* Significant at 5 percent level.

## Further Evidence of Manipulation

- A test rejection is **not** definitive proof of manipulation! Need for additional evidence
- Manipulation could have occurred in different ways at different stages of surveying process

## Further Evidence of Manipulation

- A test rejection is **not** definitive proof of manipulation! Need for additional evidence
- Manipulation could have occurred in different ways at different stages of surveying process
- Analysis of data suggests centralized manipulation rather than coordination at individual level
- Empirical findings motivated by theoretical framework: mayors attempt to influence their chances of reelection

## Conclusions

## Conclusions

- Manipulation test performs like it should
- Power of test increases when increasing manipulation or sample size
- Quality of normal approximation improves with larger number of observations
- Parameter choice depends on data at hand
- Further analysis of manipulation always needed

Thanks for your attention!