# THE JOURNAL OF

# EDUCATIONAL PSYCHOLOGY

Volume 51

December 1960

Number 6

# REGRESSION-DISCONTINUITY ANALYSIS:

#### AN ALTERNATIVE TO THE EX POST FACTO EXPERIMENT

DONALD L. THISTLETHWAITE AND

DONALD T. CAMPBELL

Northwestern University

National Merit Scholarship Corporation

While the term "ex post facto experiment" could refer to any analysis of records which provides a quasi-experimental test of a causal hypothesis, as described by Chapin (1938) and Greenwood (1945), it has come to indicate more specifically the mode of analysis in which two groups—an experimental and a control group—are selected through matching to yield a quasiexperimental comparison. studies the groups are presumed, as a result of matching, to have been equivalent prior to the exposure of the experimental group to some potentially change inducing event (the "experimental treatment"). If the groups differ on subsequent measures and if there are no plausible rival hypotheses which might account for the differences, it is inferred that the experi-

1 This study is a part of the research program of the National Merit Scholarship Corporation. This research was supported by the National Science Foundation, the Old Dominion Foundation, and by Ford Foundation grants to the National Merit Scholarship Corporation. The participation of the second author was made possible through the Northwestern University Carnegie Corporation Project in Psychology-Education. The mode of analysis illustrated in Figure 1 of this paper was first suggested by the second author, and will be presented in a chapter entitled "Experimental Designs in Research on Teaching" in the forthcoming NEA, AERA Handbook of Research on Teaching to be published by Rand McNally and edited by N. L. Gage.

mental treatment has caused the observed differences.

This paper has three purposes: first, it presents an alternative mode of analysis, called regression-discontinuity analysis, which we believe can be more confidently interpreted than the ex post facto design; second, it compares the results obtained when both modes of analysis are applied to the same data; and, third, it qualifies interpretations of the ex post facto study recently reported in this journal (Thistlethwaite, 1959).

Two groups of near-winners in a national scholarship competition were matched on several background variables in the previous study in order to study the motivational effect of public recognition. The results suggested that such recognition tends to increase the favorableness of attitudes toward intellectualism, the number of students planning to seek the MD or PhD degree, the number planning to become or scientific college teachers searchers, and the number who succeed in obtaining scholarships from other scholarship granting agencies. The regression-discontinuity analysis to be presented here confirms the effects upon success in winning scholarships from other donors but negates the inference of effects upon attitudes and is equivocal regarding career plans.

# **London School of Economics**

# **Equation Chapter 1 Section 1**The Regression Discontinuity Design

Sometimes whether something happens to you or not depends on your 'score' on a particular variable. You get a scholarship if you get above a certain mark in an exam, you get given remedial education if you get below a certain level, a policy is implemented if it gets more than 50% of the vote in a ballot, your sentence for a criminal offence is higher if you are above a certain age (an 'adult') etc etc.

All of these examples are candidates for an application of the regression discontinuity design. The essential element of a regression discontinuity design is that the probability of assignment to treatment depends in a discontinuous way on some *observable* variable W.

The simplest (and most common) form of the RDD has assignment to treatment being based on W being above some critical value  $w_0$  - I will use this case in what follows.

Note that, in some sense, the method of assignment to treatment is the very opposite here to that in random assignment – it is a deterministic function of some observable variable. But, it turns out that, in a sense I will explain, assignment to treatment is as 'good as random' in the neighbourhood of  $w_0$ , the discontinuity.

The basic RDD estimator can be understood very simply. Suppose we consider individuals with W in the interval  $[w_0 - \delta, w_0]$ . These are all in the control group so the outcome we will observe for these is  $y_0$ . Suppose that the expected value of  $y_0$  given W can be written as:

$$E(y|W,X=0) = g_0(W) \tag{1}$$

Take a first-order Taylor series expansion of this about the point  $W = w_0$ . We can then write:

$$g_0(w_0 - \delta) \approx g_0(w_0) - \delta g_0'(w_0)$$
(2)

Hence for the group of people with W in the interval  $[w_0 - \delta, w_0]$  we will have approximately that:

$$E(y|w_0 - \delta \le W < w_0) \approx g_0(w_0) - g_0'(w_0) E(\delta|w_0 - \delta \le W < w_0)$$
(3)

Now do the same exercise for individuals with W in the interval  $[w_0, w_0 + \delta]$ . These are all in the treatment group so the outcome we will observe for these is  $y_1$ . Suppose that the expected value of  $y_0$  given W can be written as:

$$E(y|W,X=1) = g_1(W)$$
(4)

Take a first-order Taylor series expansion of this about the point  $W = w_0$ . We can then write:

$$g_1(w_0 + \delta) \approx g_1(w_0) + \delta g_1'(w_0)$$
 (5)

Hence for the group of people with W in the interval  $[w_0, w_0 + \delta]$  we will have approximately that:

$$E(y|w_0 \le W \le w_0 + \delta) \approx g_1(w_0) + g_1'(w_0)E(\delta|w_0 \le W \le w_0 + \delta)$$

$$\tag{6}$$

# Economics 696F: Lecture Note 13

# Regression Discontinuity Design

# **Examples:**

- Thistlewaite and Campbell (1960): scholarship and career choice
- van der Klaauw (1997): financial aid and enrollment in college
- Angrist and Lavy (1997): class size and test scores
- Black (1999): school district and house prices

# Basic Setup:

 $Y_i(0), Y_i(1)$ : potential outcomes

 $T_i = 0, 1$ : treatment

Observed outcome:  $Y_i := T_i Y_i(1) + (1 - T_i) Y_i(0)$ .

 $Z_i$ : observed variable, scalar and continuous

"Sharp design":

$$T_i = 1(Z_i \ge z_0),$$

where  $z_0$  is fixed and known to the data analyst.

"Fuzzy Design":  $Pr(T_i = 1 | Z_i = z)$  has a discontinuity at  $z = z_0$ .

Assumption RD: The following limits exist:

$$T^+ := \lim_{z \downarrow z_0} \Pr(T = 1 | Z = z),$$

$$T^- := \lim_{z \uparrow z_0} \Pr(T = 1 | Z = z),$$

and  $T^+ \neq T^-$ .

Note that for the sharp design,  $T^+ = 1$  and  $T^- = 0$ .

We will focus on the sharp design; for extensions to the fuzzy design, see Hahn-Todd-van der Klaauw.



Available online at www.sciencedirect.com



Journal of Econometrics 142 (2008) 615-635



www.elsevier.com/locate/jeconom

# Regression discontinuity designs: A guide to practice

Guido W. Imbens<sup>a</sup>, Thomas Lemieux<sup>b,\*</sup>

<sup>a</sup>Department of Economics, Harvard University and NBER, M-24 Littauer Center, Cambridge, MA 02138, USA
<sup>b</sup>Department of Economics, University of British Columbia and NBER, 997-1873 East Mall, Vancouver, BC, V6T 1Z1, Canada

Available online 21 May 2007

#### **Abstract**

In regression discontinuity (RD) designs for evaluating causal effects of interventions, assignment to a treatment is determined at least partly by the value of an observed covariate lying on either side of a fixed threshold. These designs were first introduced in the evaluation literature by Thistlewaite and Campbell [1960. Regression-discontinuity analysis: an alternative to the ex-post Facto experiment. Journal of Educational Psychology 51, 309–317] With the exception of a few unpublished theoretical papers, these methods did not attract much attention in the economics literature until recently. Starting in the late 1990s, there has been a large number of studies in economics applying and extending RD methods. In this paper we review some of the practical and theoretical issues in implementation of RD methods. © 2007 Elsevier B.V. All rights reserved.

JEL classification: C14; C21

Keywords: Regression discontinuity; Treatment effects; Nonparametric estimation

#### 1. Introduction

Since the late 1990s there has been a large number of studies in economics applying and extending regression discontinuity (RD) methods, including Van Der Klaauw (2002), Black (1999), Angrist and Lavy (1999), Lee (2007), Chay and Greenstone (2005), DiNardo and Lee (2004), Chay et al. (2005), and Card et al. (2006). Key theoretical and conceptual contributions include the interpretation of estimates for fuzzy regression discontinuity (FRD) designs allowing for general heterogeneity of treatment effects (Hahn et al., 2001, HTV from hereon), adaptive estimation methods (Sun, 2005), specific methods for choosing bandwidths (Ludwig and Miller, 2005), and various tests for discontinuities in means and distributions of non-affected variables (Lee, 2007; McCrary, 2007).

In this paper, we review some of the practical issues in implementation of RD methods. There is relatively little novel in this discussion. Our general goal is instead to address practical issues in implementing RD designs and review some of the new theoretical developments.

After reviewing some basic concepts in Section 2, the paper focuses on five specific issues in the implementation of RD designs. In Section 3 we stress graphical analyses as powerful methods for illustrating

E-mail addresses: imbens@harvard.edu (G.W. Imbens), tlemieux@interchange.ubc.ca (T. Lemieux).

<sup>\*</sup>Corresponding author. Tel.: +16048222092; fax: +16048225915.

the design. In Section 4 we discuss estimation and suggest using local linear regression methods using only the observations close to the discontinuity point. In Section 5 we propose choosing the bandwidth using cross-validation. In Section 6 we provide a simple plug-in estimator for the asymptotic variance and a second estimator that exploits the link with instrumental variable methods derived by HTV. In Section 7 we discuss a number of specification tests and sensitivity analyses based on tests for (a) discontinuities in the average values for covariates, (b) discontinuities in the conditional density of the forcing variable, as suggested by McCrary, and (c) discontinuities in the average outcome at other values of the forcing variable.

### 2. Sharp and FRD designs

#### 2.1. Basics

Our discussion will frame the RD design in the context of the modern literature on causal effects and treatment effects, using the Rubin Causal Model (RCM) set up with potential outcomes (Rubin, 1974; Holland, 1986; Imbens and Rubin, 2007), rather than the regression framework that was originally used in this literature. For a general discussion of the RCM and its use in the economic literature, see the survey by Imbens and Wooldridge (2007).

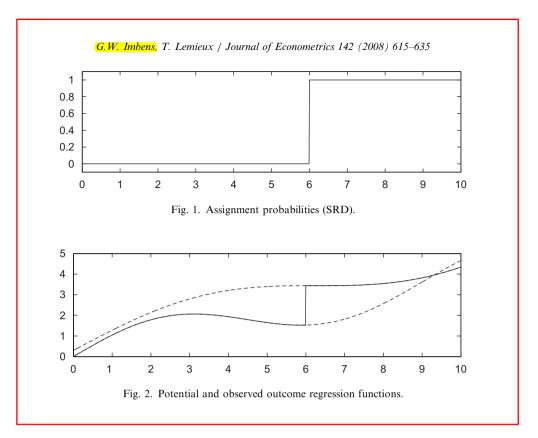
In the basic setting for the RCM (and for the RD design), researchers are interested in the causal effect of a binary intervention or treatment. Units, which may be individuals, firms, countries, or other entities, are either exposed or not exposed to a treatment. The effect of the treatment is potentially heterogenous across units. Let  $Y_i(0)$  and  $Y_i(1)$  denote the pair of potential outcomes for unit i:  $Y_i(0)$  is the outcome without exposure to the treatment and  $Y_i(1)$  is the outcome given exposure to the treatment. Interest is in some comparison of  $Y_i(0)$  and  $Y_i(1)$ . Typically, including in this discussion, we focus on differences  $Y_i(1) - Y_i(0)$ . The fundamental problem of causal inference is that we never observe the pair  $Y_i(0)$  and  $Y_i(1)$  together. We therefore typically focus on average effects of the treatment, that is, averages of  $Y_i(1) - Y_i(0)$  over (sub)populations, rather than on unit-level effects. For unit i we observe the outcome corresponding to the treatment received. Let  $W_i \in \{0,1\}$  denote the treatment received, with  $W_i = 0$  if unit i was not exposed to the treatment, and  $W_i = 1$  otherwise. The outcome observed can then be written as

$$Y_i = (1 - W_i) \cdot Y_i(0) + W_i \cdot Y_i(1) = \begin{cases} Y_i(0) & \text{if } W_i = 0, \\ Y_i(1) & \text{if } W_i = 1. \end{cases}$$

In addition to the assignment  $W_i$  and the outcome  $Y_i$ , we may observe a vector of covariates or pretreatment variables denoted by  $(X_i, Z_i)$ , where  $X_i$  is a scalar and  $Z_i$  is an M-vector. A key characteristic of  $X_i$  and  $Z_i$  is that they are known not to have been affected by the treatment. Both  $X_i$  and  $Z_i$  are covariates, with a special role played by  $X_i$  in the RD design. For each unit we observe the quadruple  $(Y_i, W_i, X_i, Z_i)$ . We assume that we observe this quadruple for a random sample from some well-defined population.

The basic idea behind the RD design is that assignment to the treatment is determined, either completely or partly, by the value of a predictor (the covariate  $X_i$ ) being on either side of a fixed threshold. This predictor may itself be associated with the potential outcomes, but this association is assumed to be smooth, and so any discontinuity of the conditional distribution (or of a feature of this conditional distribution such as the conditional expectation) of the outcome as a function of this covariate at the cutoff value is interpreted as evidence of a causal effect of the treatment.

The design often arises from administrative decisions, where the incentives for units to participate in a program are partly limited for reasons of resource constraints, and clear transparent rules rather than discretion by administrators are used for the allocation of these incentives. Examples of such settings abound. For example, Hahn et al. (1999) study the effect of an anti-discrimination law that only applies to firms with at least 15 employees. In another example, Matsudaira (2007) studies the effect of a remedial summer school program that is mandatory for students who score less than some cutoff level on a test (see also Jacob and Lefgren, 2004). Access to public goods such as libraries or museums is often eased by lower prices for individuals depending on an age cutoff value (senior citizen discounts and discounts for children under some age limit). Similarly, eligibility for medical services through medicare is restricted by age (Card et al., 2004).



### 2.2. The sharp regression discontinuity design

It is useful to distinguish between two general settings, the sharp and the fuzzy regression discontinuity (SRD and FRD from hereon) designs (e.g., Trochim, 1984, 2001; HTV). In the SRD design the assignment  $W_i$  is a deterministic function of one of the covariates, the forcing (or treatment-determining) variable  $X^1$ :

$$W_i = 1\{X_i \geqslant c\}.$$

All units with a covariate value of at least c are assigned to the treatment group (and participation is mandatory for these individuals), and all units with a covariate value less than c are assigned to the control group (members of this group are not eligible for the treatment). In the SRD design we look at the discontinuity in the conditional expectation of the outcome given the covariate to uncover an average causal effect of the treatment:

$$\lim_{x \downarrow c} \mathbb{E}[Y_i | X_i = x] - \lim_{x \uparrow c} \mathbb{E}[Y_i | X_i = x],$$

which is interpreted as the average causal effect of the treatment at the discontinuity point

$$\tau_{\text{SRD}} = \mathbb{E}[Y_i(1) - Y_i(0)|X_i = c]. \tag{2.1}$$

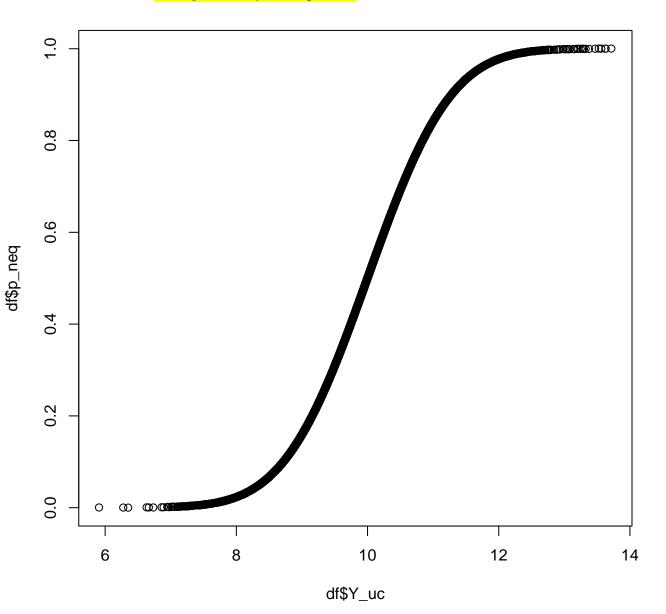
Figs. 1 and 2 illustrate the identification strategy in the SRD setup. Based on artificial population values, we present in Fig. 1 the conditional probability of receiving the treatment, Pr(W = 1|X = x) against the covariate x. At x = 6 the probability jumps from 0 to 1. In Fig. 2, three conditional expectations are plotted. The two continuous lines (partly dashed, partly solid) in the figure are the conditional expectations of the two potential outcomes given the covariate,  $\mu_w(x) = \mathbb{E}[Y(w)|X = x]$ , for w = 0, 1. These two conditional expectations are continuous functions of the covariate. Note that we can only estimate  $\mu_0(x)$  for x < c and  $\mu_1(x)$  for  $x \ge c$ .

where  $X_1 \subset X$ , and X is the covariate space.

<sup>&</sup>lt;sup>1</sup>Here we take  $X_i$  to be a scalar. More generally, the assignment can be a function of a vector of covariates. Formally, we can write this as the treatment indicator being an indicator for the vector  $X_i$  being an element of a subset of the covariate space, or

 $W_i = 1\{X_i \in \mathbb{X}_1\},\,$ 

RQ2, probability of assignment



Journal of Educational Statistics Spring 1977, Volume 2, Number 1, Pp. 1-26

# ASSIGNMENT TO TREATMENT GROUP ON THE BASIS OF A COVARIATE

#### Donald B. Rubin

#### Educational Testing Service

Key words: Non-Randomized Studies; Observational Studies; Covariance Adjustment; Causal Inference; Experimental Design; Treatment Assignment; Average Treatment Effects

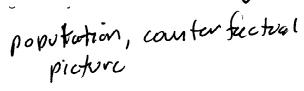
#### ABSTRACT

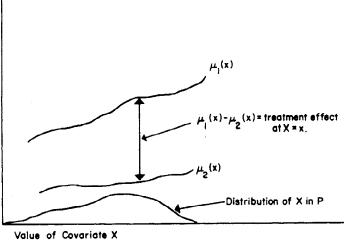
When assignment to treatment group is made solely on the basis of the value of a covariate, X, effort should be concentrated on estimating the conditional expectations of the dependent variable Y given X in the treatment and control groups. One then averages the difference between these conditional expectations over the distribution of X in the relevant population. There is no need for concern about "other" sources of bias, e.g., unreliability of X, unmeasured background variables. If the conditional expectations are parallel and linear, the proper regression adjustment is the simple covariance adjustment. However, since the quality of the resulting estimates may be sensitive to the adequacy of the underlying model, it is wise to search for nonparallelism and nonlinearity in these conditional expectations. Blocking on the values of X is also appropriate, although the quality of the resulting estimates may be sensitive to the coarseness of the blocking employed. In order for these techniques to be useful in practice, there must be either substantial overlap in the distribution of X in the treatment groups or strong prior information.

#### 1. INTRODUCTION

In some studies, the experimental units are divided

# Rubin (1977) Assignment on Covariate





Expected value of Y given X=x in Population P

pick-a-point

FIG. 1

The Treatment Effect in Population P:

$$\tau = \frac{\text{Ave}}{x \in P} \left[ \mu_1(x) - \mu_2(x) \right]$$

probabistic assignment HW2

for assignment on X jorobabilistic or not

Result 4: If  $\mu_1(x)$  and  $\mu_2(x)$  are both linear in x and parallel, then the simple analysis of covariance estimator

$$\overline{y}_1 - \overline{y}_2 - (\overline{x}_1 - \overline{x}_2) \hat{\beta}$$
 (8)

where 
$$\hat{\beta} = \frac{\sum_{\Sigma} \sum_{\Sigma} (y_{ij} - \overline{y_i})(x_{ij} - \overline{x_i})}{\sum_{\Sigma} \sum_{\Sigma} (x_{ij} - \overline{x_i})^2}$$

$$= \frac{\sum_{i=1}^{2} \sum_{j=1}^{n_i} (x_{ij} - \overline{x_i})^2}{\sum_{i=1}^{2} \sum_{j=1}^{n_i} (x_{ij} - \overline{x_i})^2}$$

is unbiased for  $\tau$  .

Subpopulation Px calcs (e.g. treatment exposure)

Belson ex. (aneone using control group slope) (cohrun
reading)

Data example (p.16) 700

# Rubin's Thm, HW illustration

```
#### aside for week 5, assignment based on covariate
> # note looking ahead week 5: what if we did ancova using Y uc as covariate, assignment based on
> Gnegancova = lm(Y utGneg ~ Gneg + Y uc)
> summary(Gnegancova)
Call:
lm(formula = Y utGneg ~ Gneg + Y uc)
Residuals:
    Min 10 Median 30
                                      Max
-2.01254 -0.14568 0.01800 0.12904 2.13266
Coefficients:
           Estimate Std. Error t value Pr(>|t|)
(Intercept) -1.50115 0.87795 -1.71 0.0905.
Gneg 1.95684 0.18445 10.61 <2e-16 ***
Y uc
     1.16248 0.09449 12.30 <2e-16 ***
Signif. codes: 0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
Residual standard error: 0.7054 on 97 degrees of freedom
Multiple R-Squared: 0.8813, Adjusted R-squared: 0.8789
F-statistic: 360.1 on 2 and 97 DF, p-value: < 2.2e-16
> # gets back to ACE = 2 coeff of Gneg!!
```

# The Regression-Discontinuity Design

The regression-discontinuity design. What a terrible name! In everyday language both parts of the term have connotations that are primarily negative. To most people "regression" implies a reversion backwards or a return to some earlier, more primitive state while "discontinuity" suggests an unnatural jump or shift in what might otherwise be a smoother, more continuous process. To a research methodologist, however, the term regression-discontinuity (hereafter labeled "RD") carries no such negative meaning. Instead, the RD design is seen as a useful method for determining whether a program or treatment is effective.

The label "RD design" actually refers to a set of design variations. In its simplest most traditional form, the RD design is a pretest-posttest program-comparison group strategy. The unique characteristic which sets RD designs apart from other pre-post group designs is the method by which research participants are assigned to conditions. In RD designs, participants are assigned to program or comparison groups solely on the basis of a cutoff score on a pre-program measure. Thus the RD design is distinguished from randomized experiments (or randomized clinical trials) and from other quasi-experimental strategies by its unique method of assignment. This cutoff criterion implies the major advantage of RD designs — they are appropriate when we wish to target a program or treatment to those who most need or deserve it. Thus, unlike its randomized or quasi-experimental alternatives, the RD design does not require us to assign potentially needy individuals to a no-program comparison group in order to evaluate the effectiveness of a program.

The RD design has not been used frequently in social research. The most common implementation has been in compensatory education evaluation where school children who obtain scores which fall below some predetermined cutoff value on an achievement test are assigned to remedial training designed to improve their performance. The low frequency of use may be attributable to several factors. Certainly, the design is a relative latecomer. Its first major field tests did not occur until the mid-1970s when it was incorporated into the nationwide evaluation system for compensatory education programs funded under Title I of the Elementary and Secondary Education Act (ESEA) of 1965. In many situations, the design has not been used because one or more key criteria were absent. For instance, RD designs force administrators to assign participants to conditions solely on the basis of quantitative indicators thereby often impalatably restricting the degree to which judgment, discretion or favoritism may be used. Perhaps the most telling reason for the lack of wider adoption of the RD design is that at first glance the design doesn't seem to make sense. In most research, we wish to have comparison groups that are equivalent to program groups on pre-program indicators so that post-program differences may be attributed to the program itself. But because of the cutoff criterion in RD designs, program and comparison groups are deliberately and maximally different on pre-program characteristics, an apparently insensible

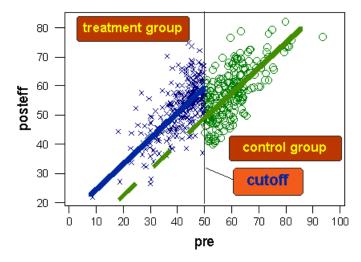


Figure 2. Regression-Discontinuity Design with Ten-point Treatment Effect.

Figure 2 is identical to Figure 1 except that all points to the left of the cutoff (i.e., the treatment group) have been raised by 10 points on the posttest. The dashed line in Figure 2 shows what we would expect the treated group's regression line to look like if the program had no effect (as was the case in Figure 1).

It is sometimes difficult to see the forest for the trees in these types of bivariate plots. So, let's remove the individual data points and look only at the regression lines. The plot of regression lines for the treatment effect case of Figure 2 is shown in Figure 3.

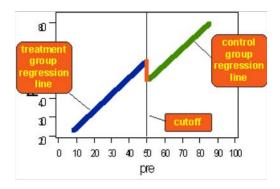


Figure 3. Regression lines for the data shown in Figure 2.

On the basis of Figure 3, we can now see how the RD design got its name - - a program effect is suggested when we observe a "jump" or discontinuity in the regression lines at the cutoff point. This is illustrated in Figure 4.

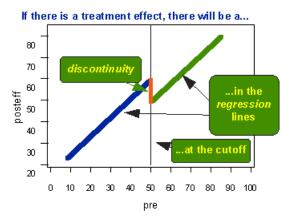


Figure 4. How the Regression-Discontinuity Design got its name.

# The Logic of the RD Design

The discussion above indicates what the key feature of the RD design is: assignment based on a cutoff value on a pre-program measure. The cutoff rule for the simple two-group case is essentially:

- all persons on one side of the cutoff are assigned to one group...
- all persons on the other side of the cutoff are assigned to the other
- need a continuous quantitative pre-program measure

Selection of the Cutoff. The choice of cutoff value is usually based on one of two factors. It can be made solely on the basis of the program resources that are available. For instance, if a program only has the capability of handling 25 persons and 70 people apply, one can choose a cutoff point that distinguishes the 25 most needy persons from the rest. Alternatively, the cutoff can be chosen on substantive grounds. If the pre-program assignment measure is an indication of severity of illness measured on a 1 to 7 scale and physicians or other experts believe that all patients scoring 5 or more are critical and fit well the criteria defined for program participants then a cutoff value of 5 may be used.

Interpretation of Results. In order to interpret the results of an RD design, one must know the nature of the assignment variable, who received the program and the nature of the outcome measure. Without this information, there is no distinct outcome pattern which directly indicates whether an effect is positive or negative.

To illustrate this, we can construct a new hypothetical example of an RD design. Let us assume that a hospital administrator would like to improve the quality of patient care through the institution of an intensive quality of care training program for staff. Because of financial constraints, the program is too costly to implement for all employees and so instead it will be administered to the entire staff from specifically targeted units or wards which seem most in need of improving quality of care. Two general measures of quality of care are available. The first is an aggregate rating of quality of care based on observation and rating by an administrative staff member and will be labeled here the QOC rating. The second is the ratio of the number of recorded patient complaints

controlled non-vandow Stat 209 wm Trochin हाराम हो हिल्लाचार्याची होते हो स्टाइ होते हैं Cornell onequivalent Groups Design ] [ The Regression-Discontinuity Design ] [ Other Quasi-Experimental Designs ] Sharp cut Hospital care examples If there is a treatment effect, there will be a... 70 Complaint Rate 80 Quality of Care RΠ 70 discontinuity 50 80 40 ...in the regression 30 50 30 lines 40 50 60 at the cutoff 30 Quality of Care Quality of Care 20 0 10 20 30 40 50 60 70 Complaint Rate Quality of Care 80 Threats 50 40 null effect dose-resp If the true pre-post relationship is not linear... dosc-respons 30 70 80 90 100 10 20 30 60 70 80 90 100 Complaint Rate Complaint Rate 70 60 posteff treatment group 50 70 40 60 30 50 20 10 20 30 40 50 60 70 80 90 100 40 control group pre 30 20 50 60 70 80 30 Data Example pre compansatury assignment XIT 7×10 (Title I) 80 The regression equation is 70 posteff = 49.8 + 0.824\*precut + 9.89\*group - 0.0196\*linint 60 t-ratio Stdev Predictor Coef 0.000 49.7508 0.6957 71.52 Constant 0.000 0.05889 13.99 precut 0.82371 40 0.000 9.8939 0.9528 10.38 group 0.813 -0.01963 0.08284 30 linint R-sq(adj) = 47.2%

note: probabilistic assingment preferred Hw2#6

40 50 60 70 80 90 100

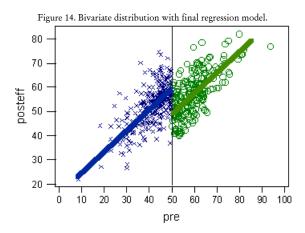
30

Figure 13. Regression results for final model.

Predictor	Coef	Stdev	t-ratio	p	
Constant	49.8421	0.5786	86.14	0.000	ancova
precut	0.81379	0.04138	19.67	0.000	aricova
group	9.8875	0.9515	10.39	0.000	

We see in these results that the treatment effect and SE are almost identical to the previous model and that the treatment effect estimate is an unbiased estimate of the true effect of 10 points. We can also see that all of the terms in the final model are statistically significant, suggesting that they are needed to model the data and should not be eliminated.

So, what does our model look like visually? Figure 14 shows the original bivariate distribution with the fitted regression model.



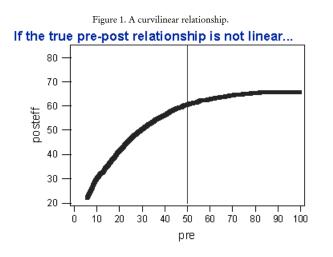
Clearly, the model fits well, both statistically and visually.

Copyright ©2006, William M.K. Trochim, All Rights Reserved

complicating the analysis somewhat.

# The Curvilinearity Problem

The major problem in analyzing data from the RD design is model misspecification. As will be shown below, when you misspecify the statistical model, you are likely to get biased estimates of the treatment effect. To introduce this idea, let's begin by considering what happens if the data (i.e., the bivariate pre-post relationship) are curvilinear and we fit a straight-line model to the data.



value of overlap; not step-function assignment, prob

Figure 1 shows a simple curvilinear relationship. If the curved line in Figure 1 describes the pre-post relationship, then we need to take this into account in our statistical model. Notice that, although there is a cutoff value at 50 in the figure, there is no jump or discontinuity in the line at the cutoff. This indicates that there is no effect of the treatment.

# Package 'rdd'

January 27, 2015



11

Maintainer Drew Dimmery <drewd@nyu.edu></drewd@nyu.edu>
Author Drew Dimmery
Version 0.56
<b>License</b> Apache License (== 2.0)
Title Regression Discontinuity Estimation
Description This package provides the tools to undertake estimation in Regression Discontinuity Designs. Both sharp and fuzzy designs are supported. Estimation is accomplished using local linear regression.  A provided function will utilize Imbens-Kalyanaraman optimal bandwidth calculation. A function is also included to test the assumption of no-sorting effects.
Type Package
<b>Date</b> 2013-10-11
<b>Depends</b> R (>= 2.15.0), sandwich, lmtest, AER, Formula
Collate 'kernelwts.R' 'DCdensity.R' 'IKbandwidth.R' 'RDestimate.R' 'plot.RD.R' 'summary.RD.R' 'rdd-package.R' 'print.RD.R'
NeedsCompilation no
Repository CRAN
<b>Date/Publication</b> 2013-10-12 00:25:44
R topics documented:
rdd-package DCdensity IKbandwidth kernelwts plot.RD

Index

6 RDestimate

Plot of the Regression Discontinuity

# Description

Plot the relationship between the running variable and the outcome

# Usage

```
## S3 method for class 'RD'
plot(x, gran = 400, bins = 100, which = 1,
  range, ...)
```

# Arguments

x	rd object, typically the result of RDestimate
gran	the granularity of the plot. This specifies the number of points to either side of the cutpoint for which the estimate is calculated.
bins	if the dependent variable is binary, include the number of bins within which to average
which	identifies which of the available plots to display. For a sharp design, the only possibility is 1, the plot of the running variable against the outcome variable. For a fuzzy design, an additional plot, 2, may also be displayed, showing the relationship between the running variable and the treatment variable. Both plots may be displayed with which=c(1,2).
range	the range of values of the running variable for which to plot. This should be a vector of length two of the format $c(\min,\max)$ . To plot from the minimum to the maximum value, simply enter $c("\min","\max")$ .
	unused

# Author(s)

Drew Dimmery <<drewd@nyu.edu>>

RDestimate

Regression Discontinuity Estimation

# Description

RDestimate supports both sharp and fuzzy RDD utilizing the **AER** package for 2SLS regression under the fuzzy design. Local linear regressions are performed to either side of the cutpoint using the Imbens-Kalyanaraman optimal bandwidth calculation, IKbandwidth.

RDestimate 7

## Usage

```
RDestimate(formula, data, subset = NULL, cutpoint = NULL,
bw = NULL, kernel = "triangular", se.type = "HC1",
cluster = NULL, verbose = FALSE, model = FALSE,
frame = FALSE)
```

## **Arguments**

formula the formula of the RDD. This is supplied in the format of y ~ x for a simple

sharp RDD, or  $y \sim x \mid c1 + c2$  for a sharp RDD with two covariates. Fuzzy RDD may be specified as  $y \sim x + z$  where x is the running variable, and z is the endogenous treatment variable. Covariates are then included in the same

manner as in a sharp RDD.

data an optional data frame

subset an optional vector specifying a subset of observations to be used

cutpoint the cutpoint. If omitted, it is assumed to be 0.

bw the bandwidth. If omitted, it is calculated using the Imbens-Kalyanaraman

method.

kernel a string specifying the kernel to be used in the local linear fitting. "triangular"

kernel is the default and is the "correct" theoretical kernel to be used for edge estimation as in RDD (Lee and Lemieux 2010). Other options are "rectangular",

"epanechnikov", "quartic", "triweight", "tricube", "gaussian" and "cosine".

se. type this specifies the robust SE calculation method to use. Options are, as in vcovHC,

"HC3", "const", "HC", "HC0", "HC1", "HC2", "HC4", "HC4m", "HC5". This op-

tion is overriden by cluster.

cluster an optional vector specifying clusters within which the errors are assumed to be

correlated. This will result in reporting cluster robust SEs. This option overrides anything specified in se. type. It is suggested that data with a discrete running variable be clustered by each unique value of the running variable (Lee and Card

2008).

verbose will provide some additional information printed to the terminal.

frame logical. If TRUE, the data frame used in model fitting will be returned.

model logical. If TRUE, the model object will be returned.

#### Value

RDestimate returns an object of class "RD". The functions summary and plot are used to obtain and print a summary and plot of the estimated regression discontinuity. The object of class RD is a list containing the following components:

type a string denoting either "sharp" or "fuzzy" RDD.

est the estimate of the discontinuity in the outcome under a sharp design, or the

Wald estimator in the fuzzy design

se the standard error z the z statistic

RDestimate 9

ci	the matrix of the 95 c("CI Lower Bound", "CI Upper Bound") for each corresponding bandwidth						
bw	numeric vector of each bandwidth used in estimation						
obs	vector of the number of observations within the corresponding bandwidth						
call	the matched call						
na.action	the observations removed from fitting due to missingness						
model	(if requested) For a sharp design, a list of the lm objects is returned. For a fuzzy design, a list of lists is returned, each with two elements: firststage, the first stage lm object, and iv, the ivreg object. A model is returned for each corresponding bandwidth.						

(if requested) Returns the model frame used in fitting.

#### Author(s)

frame

Drew Dimmery <<drewd@nyu.edu>>

#### References

```
Lee, David and Thomas Lemieux. (2010) "Regression Discontinuity Designs in Economics," Journal of Economic Literature. 48(2): 281-355. http://www.aeaweb.org/articles.php?doi=10. 1257/jel.48.2.281
```

Imbens, Guido and Thomas Lemieux. (2010) "Regression discontinuity designs: A guide to practice," *Journal of Econometrics*. 142(2): 615-635. http://dx.doi.org/10.1016/j.jeconom. 2007.05.001

Lee, David and David Card. (2010) "Regression discontinuity inference with specification error," *Journal of Econometrics*. 142(2): 655-674. http://dx.doi.org/10.1016/j.jeconom.2007.05.003

Angrist, Joshua and Jorn-Steffen Pischke. (2009) *Mostly Harmless Econometrics*. Princeton: Princeton University Press.

#### See Also

```
summary.RD, plot.RD, DCdensity IKbandwidth, kernelwts, vcovHC, ivreg, lm
```

# **Examples**

```
x<-runif(1000,-1,1)
cov<-rnorm(1000)
y<-3+2*x+3*cov+10*(x>=0)+rnorm(1000)
RDestimate(y~x)
# Efficiency gains can be made by including covariates
RDestimate(y~x|cov)
```

# Package 'rddtools'

August 29, 2016
Version 0.4.0
Title Toolbox for Regression Discontinuity Design ('RDD')
Description Set of functions for Regression Discontinuity Design ('RDD'), for data visualisation, estimation and testing,
Maintainer Bastiaan Quast <bquast@gmail.com></bquast@gmail.com>
Imports KernSmooth, ggplot2, rdd, sandwich, lmtest, Formula, locpol, methods
Depends AER, np
Suggests stats4, car, knitr, testthat
License GPL (>= 2)
<pre>URL https://github.com/bquast/RDDtools</pre>
<pre>BugReports https://github.com/bquast/RDDtools/issues</pre>
VignetteBuilder knitr
NeedsCompilation no
Author Matthieu Stigler [aut], Bastiaan Quast [aut, cre]
Repository CRAN
<b>Date/Publication</b> 2015-07-27 13:32:08
R topics documented:
as.npregbw
clusterInf
covarTest_dis
dens test
gen_mc_ik
house

x > 0TRUE

x

cov

9.99862

1.99487

2.97875

73.70

17.62

88.53

<2e-16 \*\*\*

<2e-16 \*\*\*

<2e-16 \*\*\*

0.13567

0.11323

0.03365

# STAT 209 Weels 5 Regression Discontinuity Designs

```
Example from Package 'rdd'
    Maintainer Drew Dimmery <drewd@nyu.edu> Title Regression Discontinuity Estimation
    Description This package provides the tools to undertake estimation in Regression Discontinuity
    Designs. Both sharp and fuzzy designs are supported. Estimation is accomplished using local linear
    R version 3.0.1 (2013-05-16) -- "Good Sport"
    > install.packages("rdd")
                                          install and load
    > library(rdd)
                      ## Artificial data example from p.9 rdd manual
    > x<-runif(1000,-1,1)
                             selection var
    > cov<-rnorm(1000) # extra auxiliary variable
    > y<-3+2*x+3*cov+10*(x>=0)+rnorm(1000) form outcome
    # example builds in a treatment effect of 10 points for those selected on x (x>0)
    # story? students with high (or higher) ability selected for enriched instruction
    # run the rdd function just using X
    > RDestimate(y~x)
    Call: RDestimate(formula = y \sim x)
    Coefficients:
local MIE LATE
                 Half-BW Double-BW
        9.821
                   9.599
                              10.005
    > summary(RDestimate(y~x))
    Call: RDestimate(formula = y \sim x)
    Type: sharp
    Estimates:
               Bandwidth Observations Estimate
                                                   Std. Error
                                                                z value
    LATE
               0.7613
                            721
                                          9.821
                                                   0.5264
                                                                18.65
                                                                          1.157e-77
               0.3807
    Half-BW
                            348
                                          9.599
                                                    0.7705
                                                                12.46
                                                                          1.258e-35
    Double-BW 1.5227
                           1000
                                                    0.4134
                                                                24.20
                                                                         2.001e-129
    F-statistics:
                       Num. DoF
                                  Denom. DoF
    LATE
                785.6
                                  717
                       3
                                              0
    Half-BW
                341.9
                       3
                                  344
                                              0
    Double-BW 1260.4 3
                                  996
    > plot(RDestimate(y~x)) -
    #compare with our simple ancova approach
    > rubin = lm(y \sim (x>0) + x)
    > summary(rubin)
    Call: lm(formula = y \sim (x > 0) + x)
    Coefficients:
                Estimate Std. Error t value Pr(>|t|)
    (Intercept)
                  3.1168
                              0.2227
                                     13.998 < 2e-16 ***
    x > 0TRUE
                 10.0795
                              0.4038
                                      24.960 < 2e-16 ***
                  1.9951
                              0.3371
                                       5.919 4.45e-09 ***
    Residual standard error: 3.046 on 997 degrees of freedo
    Multiple R-squared: 0.8017,
                                    Adjusted R-squared: 0.
    F-statistic: 2015 on 2 and 997 DF, p-value: < 2.2e-16
use extra into
                      (W=COV)
    > summary(RDestimate(y~x | cov))
    Call: RDestimate(formula = y ~ x | cov) Type: sharp
    Estimates:
                                                                                                           X
                                                                -1.0
                                                                                        0.0
                                                                                                    0.5
               Bandwidth
                          Observations Estimate
                                                   Std. Error
                                                                z value
                                                                         Pr(>|z|)
    LATE
               0.7613
                           721
                                          9.954
                                                   0.1831
                                                                54.36
                                                                          0.000e+00
    Half-BW
               0.3807
                           348
                                         10.005
                                                   0.2837
                                                                35.27
                                                                         1.621e-272
                                                                                     ***
    Double-BW 1.5227
                          1000
                                          9.989
                                                   0.1389
                                                                71.94
                                                                          0.000e+00
    > rubincov = lm(y \sim (x>0) + x + cov)
                                                                Note: assignment var x and auxilliary info
    > summary(rubincov)
                                                                cov generated uncorrelated above.
    Call: lm(formula = y \sim (x > 0) + x + cov)
    Coefficients:
                                                                Otherwise including cov very risky
                Estimate Std. Error t value Pr(>|t|)
                                                                UNWISE
    (Intercept) 3.04203
                            0.07480
                                       40.67
                                               <2e-16 ***
```

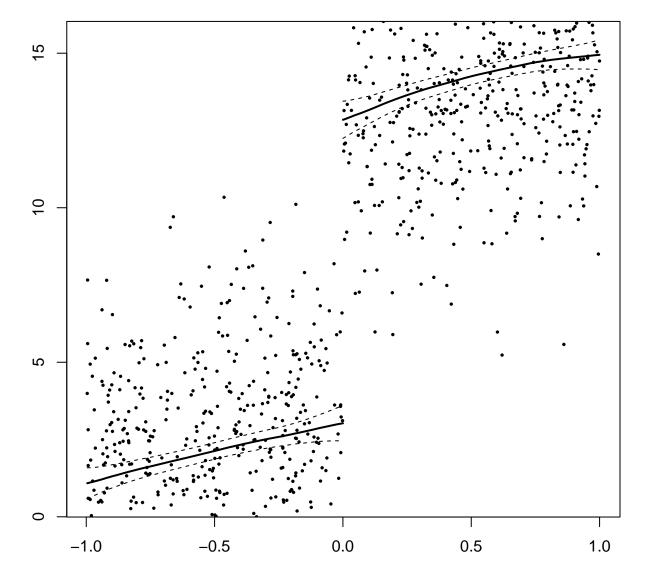
```
Example from Package 'rdd'
Maintainer Drew Dimmery <drewd@nyu.edu> Title Regression Discontinuity Estimation
Description This package provides the tools to undertake estimation in Regression Discontinuity
Designs. Both sharp and fuzzy designs are supported. Estimation is accomplished using local linear
R version 3.0.1 (2013-05-16) -- "Good Sport"
> install.packages("rdd")
> library(rdd)
                 ## Artificial data example from p.9 rdd manual
> x<-runif(1000,-1,1)
> cov<-rnorm(1000) # extra auxiliary variable
> y<-3+2*x+3*cov+10*(x>=0)+rnorm(1000)
# example builds in a treatment effect of 10 points for those selected on x (x>0)
# story? students with high (or higher) ability selected for enriched instruction
# run the rdd function just using X
> RDestimate(y~x)
Call: RDestimate(formula = y \sim x)
Coefficients:
     LATE
             Half-BW Double-BW
    9.821
              9.599
                     10.005
> summary(RDestimate(y~x))
Call: RDestimate(formula = y \sim x)
Type: sharp
Estimates:
           Bandwidth Observations Estimate Std. Error z value Pr(>|z|)
T.ATE
           0.7613
                      721
                                   9.821
                                              0.5264 18.65
                                                                    1.157e-77
                                                                               ***
           0.3807
                                              0.7705
                                                         12.46
Half-RW
                       348
                                    9.599
                                                                    1.258e-35
                                                                              ***
Double-BW 1.5227
                      1000
                                   10.005
                                              0.4134
                                                         24.20
                                                                   2.001e-129 ***
F-statistics:
                  Num. DoF Denom. DoF
                                         р
LATE
           785.6
                             717
                                         0
Half-BW
           341.9
                             344
                                         0
Double-BW 1260.4
                             996
                                         O
> plot(RDestimate(y~x))
#compare with our simple ancova approach
> rubin = lm(y \sim (x>0) + x)
> summary(rubin)
Call: lm(formula = y \sim (x > 0) + x)
Coefficients:
            Estimate Std. Error t value Pr(>|t|)
(Intercept)
              3.1168
                         0.2227 13.998 < 2e-16 ***
                         0.4038 24.960 < 2e-16 ***
x > 0TRUE
            10.0795
              1.9951
                                 5.919 4.45e-09 ***
                        0.3371
Residual standard error: 3.046 on 997 degrees of freedom
Multiple R-squared: 0.8017, Adjusted R-squared: 0.8013
F-statistic: 2015 on 2 and 997 DF, p-value: < 2.2e-16
> summary(RDestimate(y~x | cov))
Call: RDestimate(formula = y ~ x | cov) Type: sharp
Estimates:
           Bandwidth Observations Estimate Std. Error z value Pr(>|z|)
LATE
           0.7613
                      721
                                    9.954
                                              0.1831
                                                          54.36
                                                                    0.000e+00
Half-BW
           0.3807
                       348
                                    10.005
                                              0.2837
                                                          35.27
                                                                   1.621e-272 ***
Double-BW 1.5227
                      1000
                                     9.989
                                              0.1389
                                                          71.94
                                                                    0.000e+00 ***
> rubincov = lm(y \sim (x>0) + x + cov)
> summary(rubincov)
Call: lm(formula = y \sim (x > 0) + x + cov)
Coefficients:
           Estimate Std. Error t value Pr(>|t|)
(Intercept) 3.04203
                        0.07480
                                  40.67
                                          <2e-16 ***
x > 0TRUE
            9.99862
                       0.13567
                                  73.70
                                          <2e-16 ***
                                          <2e-16 ***
```

х

1.99487

0.11323

17.62



# Package 'rdrobust'

May 8, 2016

Type Package

<b>Title</b> Robust Data-Driven Statistical Inference in Regression-Discontinuity Designs
Version 0.90
<b>Date</b> 2016-05-01
Author Sebastian Calonico <scalonico@bus.miami.edu>, Matias D. Catta- neo <cattaneo@umich.edu>, Max H. Farrell, <max.farrell@chicagobooth.edu>, Ro- cio Titiunik <titiunik@umich.edu></titiunik@umich.edu></max.farrell@chicagobooth.edu></cattaneo@umich.edu></scalonico@bus.miami.edu>
Maintainer Sebastian Calonico <scalonico@bus.miami.edu></scalonico@bus.miami.edu>
Description Regression-discontinuity (RD) designs are quasi-experimental research designs popular in social, behavioral and natural sciences. The RD design is usually employed to study the (local) causal effect of a treatment, intervention or policy. This package provides tools for data-driven graphical and analytical statistical inference in RD designs: rdrobust to construct local-polynomial point estimators and robust confidence intervals for average treatment effects at the cutoff in Sharp, Fuzzy and Kink RD settings, rdbwselect to perform bandwidth selection for the different procedures implemented, and rdplot to conduct exploratory data analysis (RD plots).
<b>Depends</b> R (>= $3.1.1$ )
License GPL-2
NeedsCompilation no
Repository CRAN
<b>Date/Publication</b> 2016-05-08 09:52:01  R topics documented:
it topics documented.
rdrobust-package rdbwselect rdbwselect_2014 rdplot rdrobust 1 rdrobust-internal 1 rdrobust_RDsenate 1

rdrobust 11

Cattaneo, M. D., B. Frandsen, and R. Titiunik. 2015. Randomization Inference in the Regression Discontinuity Design: An Application to the Study of Party Advantages in the U.S. Senate. Journal of Causal Inference 3(1): 1-24. http://www-personal.umich.edu/~cattaneo/papers/Cattaneo-Frandsen-Titiunik\_2015\_JCI.pdf.

#### See Also

```
rdbwselect, rdrobust
```

### **Examples**

```
x<-runif(1000,-1,1)
y<-5+3*x+2*(x>=0)+rnorm(1000)
rdplot(y,x)
```

rdrobust

Local-Polynomial RD Estimation with Robust Confidence Intervals

### **Description**

rdrobust implements local polynomial Regression Discontinuity (RD) point estimators with robust bias-corrected confidence intervals and inference procedures developed in Calonico, Cattaneo and Titiunik (2014a), Calonico, Cattaneo and Farrell (2016a), and Calonico, Cattaneo, Farrell and Titiunik (2016). It also computes alternative estimation and inference procedures available in the literature.

Companion commands are: rdbwselect for data-driven bandwidth selection, and rdplot for data-driven RD plots (see Calonico, Cattaneo and Titiunik (2015a) for details).

A detailed introduction to this command is given in Calonico, Cattaneo and Titiunik (2015b), and Calonico, Cattaneo, Farrell and Titiunik (2016b). A companion Stata package is described in Calonico, Cattaneo and Titiunik (2014b).

For more details, and related Stata and R packages useful for analysis of RD designs, visit https://sites.google.com/site/rdpackages/

### Usage

### **Arguments**

- y is the dependent variable.
- x is the running variable (a.k.a. score or forcing variable).

covs specifies additional covariates to be used for estimation and inference.

1.3 Maimonides' Rule 7

# 1.2.6 Exiting a treatment group after treatment assignment

Randomized experiment: Once assigned to a treatment group, subjects do not exit. A subject who does not comply with the assigned treatment, or switches to another treatment, or is lost to follow-up, remains in the assigned treatment group with these characteristics noted. An analysis that compares the groups as randomly assigned, ignoring deviations between intended and actual treatment, is called an 'intention-to-treat' analysis, and it is one of the central analyses reported in a randomized trial. Randomization inference may partially address noncompliance with assigned treatment by viewing treatment assignment as an instrumental variable for treatment received; see §5.3 and [18].

Better observational study: Once assigned to a treatment group, subjects do not exit. A subject who does not comply with the assigned treatment, or switches to another treatment, or is lost to follow-up, remains in the assigned treatment group with these characteristics noted. Inference may partially address noncompliance by viewing treatment assignment as an instrumental variable for treatment received; see §5.3 and [22].

Poorer observational study: There is no clear distinction between assignment to treatment, acceptance of treatment, receipt of treatment, or switching treatments, so problems that arise in experiments seem to be avoided, when in fact they are simply ignored.

# 1.2.7 Study protocol

Randomized experiment: Before beginning the actual experiment, a written protocol describes the design, exclusion criteria, primary and secondary outcomes, and proposed analyses.

Better observational study: Before examining outcomes that will form the basis for the study's conclusions, a written protocol describes the design, exclusion criteria, primary and secondary outcomes, and proposed analyses; see Chapter 19.

Poorer observational study: If sufficiently many analyses are performed, something publishable will turn up sooner or later.

# 1.3 Maimonides' Rule

In 1999, Joshua Angrist and Victor Lavy [3] published an unusual and much admired study of the effects of class size on academic achievement. They wrote [3, pages 533-535]:

[C]ausal effects of class size on pupil achievement have proved very difficult to measure. Even though the level of educational inputs differs substantially both between and within schools, these differences are often associated with factors such as remedial training or students' socioeconomic background ... The great twelfth century Rabbinic scholar, Maimonides, interprets the Talmud's discussion of class size as follows: 'Twenty-five children may be put in charge of one teacher. If the number in the class exceeds twenty-five but is not more than forty, he should have an assistant to help with instruction. If there are more than forty, two teachers must be appointed.' ... The importance of Maimonides' rule for our purposes is that, since 1969, it has been used to determine the division of enrollment cohorts into classes in Israeli public schools.

In most places at most times, class size has been determined by the affluence or poverty of a community, its enthusiasm or skepticism about the value of education, the special needs of students for remedial or advanced instruction, the obscure, transitory, barely intelligible obsessions of bureaucracies, and each of these determinants of class size clouds its actual effect on academic performance. However, if adherence to Maimonides' rule were perfectly rigid, then what would separate a school with a single class of size 40 from the same school with two classes whose average size is 20.5 is the enrollment of a single student.

Maimonides' rule has the largest impact on a school with about 40 students in a grade cohort. With cohorts of size 40, 80, and 120 students, the steps down in average class size required by Maimonides' rule when an additional student enrolls are, respectively, from 40 to 20.5, from 40 to 27, and from 40 to 30.25. For this reason, we will look at schools with fifth grade cohorts in 1991 with between 31 and 50 students, where average class sizes might be cut in half by Maimonides' rule. There were 211 such schools, with 86 of these schools having between 31 and 40 students in fifth grade, and 125 schools having between 41 and 50 students in the fifth grade.

Adherence to Maimonides' rule is not perfectly rigid. In particular, Angrist and Lavy [3, page 538] note that the percentage of disadvantaged students in a school "is used by the Ministry of Education to allocate supplementary hours of instruction and other school resources." Among the 211 schools with between 31 and 50 students in fifth grade, the percentage disadvantaged has a slightly negative Kendall's correlation of -0.10 with average class size, which differs significantly from zero (P-value = 0.031), and it has more strongly negative correlations of -0.42 and -0.55, respectively, with performance on verbal and mathematics test scores. For this reason, 86 matched pairs of two schools were formed, matching to minimize to total absolute difference in percentage disadvantaged. Figure 1.1 shows the paired schools, 86 schools with 31 and 40 students in fifth grade, and 86 schools with between 41 and 50 students in the fifth grade. After matching, the upper left panel in Figure 1.1 shows that the percentage of disadvantaged students was balanced; indeed, the average absolute difference within a pair was less than 1%. The upper right panel in Figure 1.1 shows Maimonides' rule at work: with some exceptions, the slightly larger schools had substantially smaller class sizes. The bottom panels of Figure 1.1 show the average mathematics and verbal test performance of these fifth graders, with somewhat higher scores in the schools with between 41 and 50 fifth graders, where class sizes tended to be smaller.

The paper by Angrist and Lavy (1999) on 'Maimonides rule' pushed economists interest/revival (originally Don Campbell in psychology) of regression discontinuity designs. The goal is to estimate the 'causal' effect of class size on reading achievement of elementary school children in Israel. Maimonides rule is from the talmud and says class sizes above 40 must be broken into 2 smaller class. So a class of size 40 would be left intact, but a class of size 41 would be divided into two classes: sizes 20 and 21. So assignment into a class size reduction mechanism is a function of original class size, with cutoff at 40. [note in the real life data there are discrepancies and deviations from the rule, but that didn't phase the economists much, and for our purposes we will treat these data as following the design intent].

I obtained the Angrist dataset from the UCLA repository for the "Methods Matter" book: http://www.ats.ucla.edu/stat/examples/methods\_matter/ (but beware most of the links are broken, but with some modifications I could get the data). I formed a dataset named 'ang2' which contained classes of original size 36 through 45 (classes above size 40 are broken up). Dataset ang2 contains 180 classrooms (rows).

The 'read' variable is the class mean reading score (analysis is done at the classroom level).

- I formed the variable 'treat' in ang2 by
- > ang2\$treat = ang2\$size > 40

-----

# Analysis

> attach(ang2)

```
> angreg = lm(read ~ treat + size)
> summary(angreg)
Call: lm(formula = read ~ treat + size)
Coefficients:
          Estimate Std. Error t value Pr(>|t|)
(Intercept) 61.6978 16.6373
                                3.708 0.000279 ***
treatTRUE 3.8472
                       2.8112
                               1.368 0.172894
size
             0.1707
                        0.4360
                                0.391 0.695923
Residual standard error: 9.674 on 177 degrees of freedom
Multiple R-squared: 0.04579, Adjusted R-squared: 0.035
F-statistic: 4.246 on 2 and 177 DF, p-value: 0.0158
```

```
> library(rdd)
```

> summary(RDestimate(read ~ size, cutpoint = 40, data = ang2))

Call: RDestimate(formula = read ~ size, data = ang2, cutpoint = 40)

Type: sharp

## Estimates:

Bandwidth	Observations	Estimate	Std. Error	z value	Pr(> z )
3.175	115	1.0893	7.105	0.1533	0.8782
1.587	47	-0.9415	4.532	-0.2078	0.8354
6.350	180	1.0718	4.926	0.2176	0.8278
	3.175 1.587	3.175 115 1.587 47	3.175 115 1.0893 (1.587) 47 -0.9415	3.175     115     1.0893     7.105       1.587     47     -0.9415     4.532	1.587     47       -0.9415     4.532       -0.2078

> plot(RDestimate(read ~ size, cutpoint = 40, data = ang2)) #plot attached

