

# Regression Discontinuity Designs in Economics

DAVID S. LEE AND THOMAS LEMIEUX\*

*This paper provides an introduction and “user guide” to Regression Discontinuity (RD) designs for empirical researchers. It presents the basic theory behind the research design, details when RD is likely to be valid or invalid given economic incentives, explains why it is considered a “quasi-experimental” design, and summarizes different ways (with their advantages and disadvantages) of estimating RD designs and the limitations of interpreting these estimates. Concepts are discussed using examples drawn from the growing body of empirical research using RD. (JEL C21, C31)*

## 1. Introduction

Regression Discontinuity (RD) designs were first introduced by Donald L. Thistlethwaite and Donald T. Campbell (1960) as a way of estimating treatment effects in a nonexperimental setting where treatment is determined by whether an observed “assignment” variable (also referred to in the literature as the “forcing” variable or the “running” variable) exceeds a known cutoff point. In their initial application of RD designs, Thistlethwaite and Campbell

(1960) analyzed the impact of merit awards on future academic outcomes, using the fact that the allocation of these awards was based on an observed test score. The main idea behind the research design was that individuals with scores just below the cutoff (who did not receive the award) were good comparisons to those just above the cutoff (who did receive the award). Although this evaluation strategy has been around for almost fifty years, it did not attract much attention in economics until relatively recently.

Since the late 1990s, a growing number of studies have relied on RD designs to estimate program effects in a wide variety of economic contexts. Like Thistlethwaite and Campbell (1960), early studies by Wilbert van der Klaauw (2002) and Joshua D. Angrist and Victor Lavy (1999) exploited threshold rules often used by educational institutions to estimate the effect of financial aid and class size, respectively, on educational outcomes. Sandra E. Black (1999) exploited the presence of discontinuities at the geographical level (school district

\*Lee: Princeton University and NBER. Lemieux: University of British Columbia and NBER. We thank David Autor, David Card, John DiNardo, Guido Imbens, and Justin McCrary for suggestions for this article, as well as for numerous illuminating discussions on the various topics we cover in this review. We also thank two anonymous referees for their helpful suggestions and comments, and Damon Clark, Mike Geruso, Andrew Marder, and Zhuan Pei for their careful reading of earlier drafts. Diane Alexander, Emily Buchsbaum, Elizabeth Debraggio, Enkeleda Gjenci, Ashley Hodgson, Yan Lau, Pauline Leung, and Xiaotong Niu provided excellent research assistance.

boundaries) to estimate the willingness to pay for good schools. Following these early papers in the area of education, the past five years have seen a rapidly growing literature using RD designs to examine a range of questions. Examples include the labor supply effect of welfare, unemployment insurance, and disability programs; the effects of Medicaid on health outcomes; the effect of remedial education programs on educational achievement; the empirical relevance of median voter models; and the effects of unionization on wages and employment.

One important impetus behind this recent flurry of research is a recognition, formalized by Jinyong Hahn, Petra Todd, and van der Klaauw (2001), that RD designs require seemingly mild assumptions compared to those needed for other nonexperimental approaches. Another reason for the recent wave of research is the belief that the RD design is not “just another” evaluation strategy, and that causal inferences from RD designs are potentially more credible than those from typical “natural experiment” strategies (e.g., difference-in-differences or instrumental variables), which have been heavily employed in applied research in recent decades. This notion has a theoretical justification: David S. Lee (2008) formally shows that one need not *assume* the RD design isolates treatment variation that is “as good as randomized”; instead, such randomized variation is a *consequence* of agents’ inability to precisely control the assignment variable near the known cutoff.

So while the RD approach was initially thought to be “just another” program evaluation method with relatively little general applicability outside of a few specific problems, recent work in economics has shown quite the opposite.<sup>1</sup> In addition to providing

a highly credible and transparent way of estimating program effects, RD designs can be used in a wide variety of contexts covering a large number of important economic questions. These two facts likely explain why the RD approach is rapidly becoming a major element in the toolkit of empirical economists.

Despite the growing importance of RD designs in economics, there is no single comprehensive summary of what is understood about RD designs—when they succeed, when they fail, and their strengths and weaknesses.<sup>2</sup> Furthermore, the “nuts and bolts” of implementing RD designs in practice are not (yet) covered in standard econometrics texts, making it difficult for researchers interested in applying the approach to do so. Broadly speaking, the main goal of this paper is to fill these gaps by providing an up-to-date overview of RD designs in economics and creating a guide for researchers interested in applying the method.

A reading of the most recent research reveals a certain body of “folk wisdom” regarding the applicability, interpretation, and recommendations of practically implementing RD designs. This article represents our attempt at summarizing what we believe to be the most important pieces of this wisdom, while also dispelling misconceptions that could potentially (and understandably) arise for those new to the RD approach.

We will now briefly summarize the most important points about RD designs to set the stage for the rest of the paper where we systematically discuss identification, interpretation, and estimation issues. Here, and throughout the paper, we refer to the assignment variable as  $X$ . Treatment is, thus,

---

the RD design in economics is unique as it is still rarely used in other disciplines.

<sup>2</sup> See, however, two recent overview papers by van der Klaauw (2008b) and Guido W. Imbens and Thomas Lemieux (2008) that have begun bridging this gap.

<sup>1</sup> See Thomas D. Cook (2008) for an interesting history of the RD design in education research, psychology, statistics, and economics. Cook argues the resurgence of

assigned to individuals (or “units”) with a value of  $X$  greater than or equal to a cutoff value  $c$ .

- **RD designs can be invalid if individuals can precisely manipulate the “assignment variable.”**

When there is a payoff or benefit to receiving a treatment, it is natural for an economist to consider how an individual may behave to obtain such benefits. For example, if students could effectively “choose” their test score  $X$  through effort, those who chose a score  $c$  (and hence received the merit award) could be somewhat different from those who chose scores just below  $c$ . The important lesson here is that the existence of a treatment being a discontinuous function of an assignment variable is *not* sufficient to justify the validity of an RD design. Indeed, if anything, discontinuous rules may generate incentives, causing behavior that would *invalidate* the RD approach.

- **If individuals—even while having some influence—are unable to precisely manipulate the assignment variable, a consequence of this is that the variation in treatment near the threshold is randomized as though from a randomized experiment.**

This is a crucial feature of the RD design, since it is the reason RD designs are often so compelling. Intuitively, when individuals have imprecise control over the assignment variable, even if some are especially likely to have values of  $X$  near the cutoff, *every* individual will have approximately the same probability of having an  $X$  that is just above (receiving the treatment) or just below (being denied the treatment) the cutoff—similar to a coin-flip experiment. This result clearly differentiates the RD and

instrumental variables (IV) approaches. When using IV for causal inference, one must *assume* the instrument is exogenously generated as if by a coin-flip. Such an assumption is often difficult to justify (except when an actual lottery was run, as in Angrist (1990), or if there were some biological process, e.g., gender determination of a baby, mimicking a coin-flip). By contrast, the variation that RD designs isolate is randomized *as a consequence* of the assumption that individuals have imprecise control over the assignment variable.

- **RD designs can be analyzed—and tested—like randomized experiments.**

This is the key implication of the local randomization result. If variation in the treatment near the threshold is approximately randomized, then it follows that all “baseline characteristics”—all those variables determined prior to the realization of the assignment variable—should have the same distribution just above and just below the cutoff. If there is a discontinuity in these baseline covariates, then at a minimum, the underlying identifying assumption of individuals’ inability to precisely manipulate the assignment variable is unwarranted. Thus, the baseline covariates are used to *test* the validity of the RD design. By contrast, when employing an IV or a matching/regression-control strategy, assumptions typically need to be made about the relationship of these other covariates to the treatment and outcome variables.<sup>3</sup>

- **Graphical presentation of an RD design is helpful and informative, but the visual presentation should not be**

<sup>3</sup>Typically, one assumes that, *conditional on the covariates*, the treatment (or instrument) is essentially “as good as” randomly assigned.

**tilted toward either finding an effect or finding no effect.**

It has become standard to summarize RD analyses with a simple graph showing the relationship between the outcome and assignment variables. This has several advantages. The presentation of the “raw data” enhances the transparency of the research design. A graph can also give the reader a sense of whether the “jump” in the outcome variable at the cutoff is unusually large compared to the bumps in the regression curve away from the cutoff. Also, a graphical analysis can help identify why different functional forms give different answers, and can help identify outliers, which can be a problem in any empirical analysis. The problem with graphical presentations, however, is that there is some room for the researcher to construct graphs making it seem as though there are effects when there are none, or hiding effects that truly exist. We suggest later in the paper a number of methods to minimize such biases in presentation.

- **Nonparametric estimation does not represent a “solution” to functional form issues raised by RD designs. It is therefore helpful to view it as a complement to—rather than a substitute for—parametric estimation.**

When the analyst chooses a parametric functional form (say, a low-order polynomial) that is incorrect, the resulting estimator will, in general, be biased. When the analyst uses a nonparametric procedure such as local linear regression—essentially running a regression using only data points “close” to the cutoff—there will also be bias.<sup>4</sup> With a finite sample, it is impossible to know

which case has a smaller bias without knowing something about the true function. There will be some functions where a low-order polynomial is a very good approximation and produces little or no bias, and therefore it is efficient to use all data points—both “close to” and “far away” from the threshold. In other situations, a polynomial may be a bad approximation, and smaller biases will occur with a local linear regression. In practice, parametric and nonparametric approaches lead to the computation of the exact same statistic.<sup>5</sup> For example, the procedure of regressing the outcome  $Y$  on  $X$  and a treatment dummy  $D$  can be viewed as a parametric regression (as discussed above), or as a local linear regression with a very large bandwidth. Similarly, if one wanted to exclude the influence of data points in the tails of the  $X$  distribution, one could call the exact same procedure “parametric” after trimming the tails, or “nonparametric” by viewing the restriction in the range of  $X$  as a result of using a smaller bandwidth.<sup>6</sup> Our main suggestion in estimation is to not rely on one particular method or specification. In any empirical analysis, results that are stable across alternative

<sup>5</sup> See section 1.2 of James L. Powell (1994), where it is argued that is more helpful to view *models* rather than particular statistics as “parametric” or “nonparametric.” It is shown there how the same least squares estimator can simultaneously be viewed as a solution to parametric, semi-parametric, and nonparametric problems.

<sup>6</sup> The main difference, then, between a parametric and nonparametric approach is not in the actual estimation but rather in the discussion of the asymptotic behavior of the estimator as sample sizes tend to infinity. For example, standard nonparametric asymptotics considers what would happen if the bandwidth  $h$ —the width of the “window” of observations used for the regression—were allowed to shrink as the number of observations  $N$  tended to infinity. It turns out that if  $h \rightarrow 0$  and  $Nh \rightarrow \infty$  as  $N \rightarrow \infty$ , the bias will tend to zero. By contrast, with a parametric approach, when one is not allowed to make the model more flexible with more data points, the bias would generally remain—even with infinite samples.

<sup>4</sup> Unless the underlying function is exactly linear in the area being examined.

and equally plausible specifications are generally viewed as more reliable than those that are sensitive to minor changes in specification. RD is no exception in this regard.

- **Goodness-of-fit and other statistical tests can help rule out overly restrictive specifications.**

Often the consequence of trying many different specifications is that it may result in a wide range of estimates. Although there is no simple formula that works in all situations and contexts for weeding out inappropriate specifications, it seems reasonable, at a minimum, not to rely on an estimate resulting from a specification that can be rejected by the data when tested against a strictly more flexible specification. For example, it seems wise to place less confidence in results from a low-order polynomial model when it is rejected in favor of a less restrictive model (e.g., separate means for each discrete value of  $X$ ). Similarly, there seems little reason to prefer a specification that uses all the data if using the same specification, but restricting to observations closer to the threshold, gives a substantially (and statistically) different answer.

Although we (and the applied literature) sometimes refer to the RD “method” or “approach,” the RD design should perhaps be viewed as more of a *description* of a particular data generating process. All other things (topic, question, and population of interest) equal, we as researchers might prefer data from a randomized experiment or from an RD design. But in reality, like the randomized experiment—which is also more appropriately viewed as a particular data generating process rather than a “method” of analysis—an RD design will simply not exist to answer a great number of questions. That

said, as we show below, there has been an explosion of discoveries of RD designs that cover a wide range of interesting economic topics and questions.

The rest of the paper is organized as follows. In section 2, we discuss the origins of the RD design and show how it has recently been formalized in economics using the potential outcome framework. We also introduce an important theme that we stress throughout the paper, namely that RD designs are particularly compelling because they are close cousins of randomized experiments. This theme is more formally explored in section 3 where we discuss the conditions under which RD designs are “as good as a randomized experiment,” how RD estimates should be interpreted, and how they compare with other commonly used approaches in the program evaluation literature. Section 4 goes through the main “nuts and bolts” involved in implementing RD designs and provides a “guide to practice” for researchers interested in using the design. A summary “checklist” highlighting our key recommendations is provided at the end of this section. Implementation issues in several specific situations (discrete assignment variable, panel data, etc.) are covered in section 5. Based on a survey of the recent literature, section 6 shows that RD designs have turned out to be much more broadly applicable in economics than was originally thought. We conclude in section 7 by discussing recent progress and future prospects in using and interpreting RD designs in economics.

## 2. Origins and Background

In this section, we set the stage for the rest of the paper by discussing the origins and the basic structure of the RD design, beginning with the classic work of Thistlethwaite and Campbell (1960) and moving to the recent interpretation of the design using modern tools of program evaluation in economics (potential outcomes framework). One of

the main virtues of the RD approach is that it can be naturally presented using simple graphs, which greatly enhances its credibility and transparency. In light of this, the majority of concepts introduced in this section are represented in graphical terms to help capture the intuition behind the RD design.

### 2.1 Origins

The RD design was first introduced by Thistlethwaite and Campbell (1960) in their study of the impact of merit awards on the future academic outcomes (career aspirations, enrollment in postgraduate programs, etc.) of students. Their study exploited the fact that these awards were allocated on the basis of an observed test score. Students with test scores  $X$ , greater than or equal to a cutoff value  $c$ , received the award, while those with scores below the cutoff were denied the award. This generated a sharp discontinuity in the “treatment” (receiving the award) as a function of the test score. Let the receipt of treatment be denoted by the dummy variable  $D \in \{0, 1\}$ , so that we have  $D = 1$  if  $X \geq c$  and  $D = 0$  if  $X < c$ .

At the same time, there appears to be no reason, other than the merit award, for future academic outcomes,  $Y$ , to be a discontinuous function of the test score. This simple reasoning suggests attributing the discontinuous jump in  $Y$  at  $c$  to the causal effect of the merit award. Assuming that the relationship between  $Y$  and  $X$  is otherwise linear, a simple way of estimating the treatment effect  $\tau$  is by fitting the linear regression

$$(1) \quad Y = \alpha + D\tau + X\beta + \varepsilon,$$

where  $\varepsilon$  is the usual error term that can be viewed as a purely random error generating variation in the value of  $Y$  around the regression line  $\alpha + D\tau + X\beta$ . This case is depicted in figure 1, which shows both the true underlying function and numerous realizations of  $\varepsilon$ .

Thistlethwaite and Campbell (1960) provide some graphical intuition for why the coefficient  $\tau$  could be viewed as an estimate of the causal effect of the award. We illustrate their basic argument in figure 1. Consider an individual whose score  $X$  is exactly  $c$ . To get the causal effect for a person scoring  $c$ , we need guesses for what her  $Y$  would be with and without receiving the treatment.

If it is “reasonable” to assume that all factors (other than the award) are evolving “smoothly” with respect to  $X$ , then  $B'$  would be a reasonable guess for the value of  $Y$  of an individual scoring  $c$  (and hence receiving the treatment). Similarly,  $A''$  would be a reasonable guess for that same individual in the counterfactual state of not having received the treatment. It follows that  $B' - A''$  would be the causal estimate. This illustrates the intuition that the RD estimates should use observations “close” to the cutoff (e.g., in this case at points  $c'$  and  $c''$ ).

There is, however, a limitation to the intuition that “the closer to  $c$  you examine, the better.” In practice, one *cannot* “only” use data close to the cutoff. The narrower the area that is examined, the less data there are. In this example, examining data any closer than  $c'$  and  $c''$  will yield no observations at all! Thus, in order to produce a reasonable guess for the treated and untreated states at  $X = c$  with finite data, one has no choice but to use data *away* from the discontinuity.<sup>7</sup> Indeed, if the underlying function is truly linear, we know that the best linear unbiased estimator of  $\tau$  is the coefficient on  $D$  from OLS estimation (using all of the observations) of equation (1).

This simple heuristic presentation illustrates two important features of the RD

<sup>7</sup> Interestingly, the very first application of the RD design by Thistlethwaite and Campbell (1960) was based on discrete data (interval data for test scores). As a result, their paper clearly points out that the RD design is fundamentally based on an extrapolation approach.

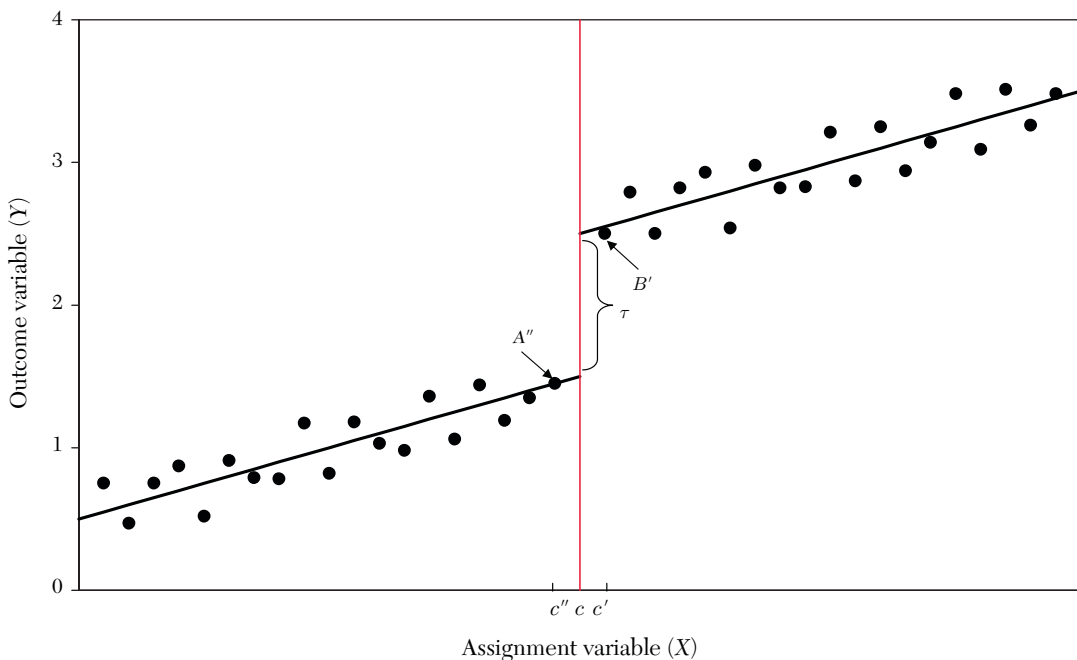


Figure 1. Simple Linear RD Setup

design. First, in order for this approach to work, “all other factors” determining  $Y$  must be evolving “smoothly” with respect to  $X$ . If the other variables also jump at  $c$ , then the gap  $\tau$  will potentially be biased for the treatment effect of interest. Second, since an RD estimate requires data away from the cut-off, the estimate will be dependent on the chosen functional form. In this example, if the slope  $\beta$  were (erroneously) restricted to equal zero, it is clear the resulting OLS coefficient on  $D$  would be a biased estimate of the true discontinuity gap.

## 2.2 RD Designs and the Potential Outcomes Framework

While the RD design was being imported into applied economic research by studies such as van der Klaauw (2002), Black (1999), and Angrist and Lavy (1999), the identification issues discussed above were formalized

in the theoretical work of Hahn, Todd, and van der Klaauw (2001), who described the RD evaluation strategy using the language of the treatment effects literature. Hahn, Todd, and van der Klaauw (2001) noted the key assumption of a valid RD design was that “all other factors” were “continuous” with respect to  $X$ , and suggested a nonparametric procedure for estimating  $\tau$  that did not assume underlying linearity, as we have in the simple example above.

The necessity of the continuity assumption is seen more formally using the “potential outcomes framework” of the treatment effects literature with the aid of a graph. It is typically imagined that, for each individual  $i$ , there exists a pair of “potential” outcomes:  $Y_i(1)$  for what would occur if the unit were exposed to the treatment and  $Y_i(0)$  if not exposed. The causal effect of the treatment is represented by the difference  $Y_i(1) - Y_i(0)$ .

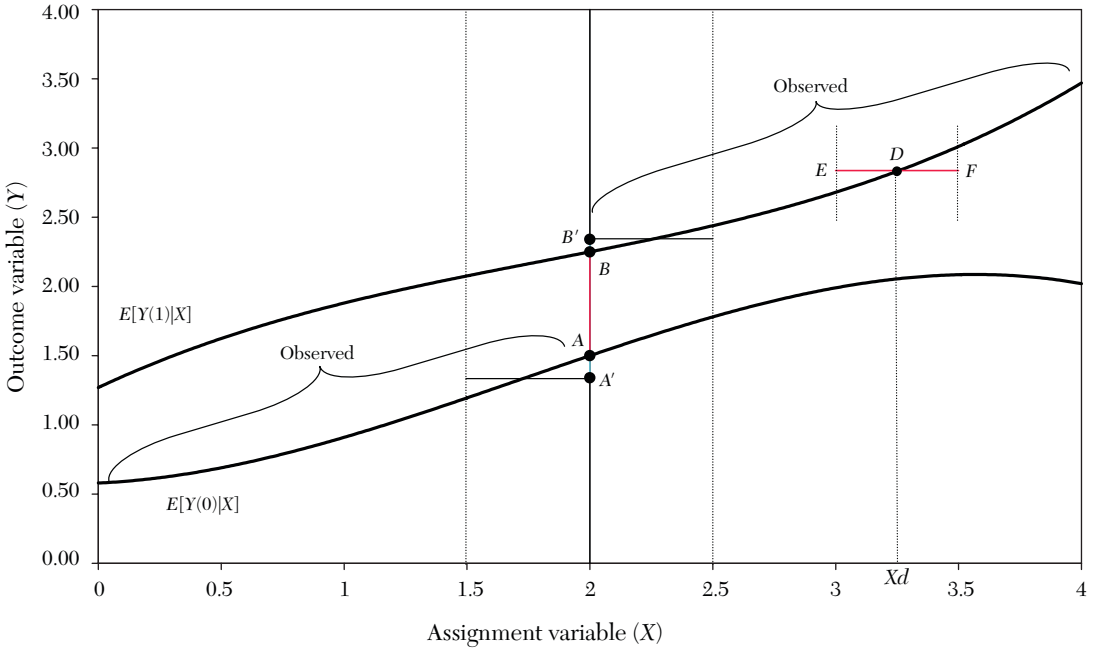


Figure 2. Nonlinear RD

The fundamental problem of causal inference is that we cannot observe the pair  $Y_i(0)$  and  $Y_i(1)$  simultaneously. 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.

In the RD setting, we can imagine there are two underlying relationships between average outcomes and  $X$ , represented by  $E[Y_i(1)|X]$  and  $E[Y_i(0)|X]$ , as in figure 2. But by definition of the RD design, all individuals to the right of the cutoff ( $c = 2$  in this example) are exposed to treatment and all those to the left are denied treatment. Therefore, we only observe  $E[Y_i(1)|X]$  to the right of the cutoff and  $E[Y_i(0)|X]$  to the left of the cutoff as indicated in the figure.

It is easy to see that with what is observable, we could try to estimate the quantity

$$B - A = \lim_{\varepsilon \downarrow 0} E[Y_i | X_i = c + \varepsilon] - \lim_{\varepsilon \uparrow 0} E[Y_i | X_i = c + \varepsilon],$$

which would equal

$$E[Y_i(1) - Y_i(0) | X = c].$$

This is the “average treatment effect” at the cutoff  $c$ .

This inference is possible because of the continuity of the underlying functions  $E[Y_i(1)|X]$  and  $E[Y_i(0)|X]$ .<sup>8</sup> In essence,

<sup>8</sup>The continuity of both functions is not the minimum that is required, as pointed out in Hahn, Todd, and van der Klaauw (2001). For example, identification is still possible even if only  $E[Y_i(0)|X]$  is continuous, and only continuous at  $c$ . Nevertheless, it may seem more natural to assume that the conditional expectations are continuous for all values of  $X$ , since cases where continuity holds at the cutoff point but not at other values of  $X$  seem peculiar.

this continuity condition enables us to use the average outcome of those right below the cutoff (who are denied the treatment) as a valid counterfactual for those right above the cutoff (who received the treatment).

Although the potential outcome framework is very useful for understanding how RD designs work in a framework applied economists are used to dealing with, it also introduces some difficulties in terms of interpretation. First, while the continuity assumption sounds generally plausible, it is not completely clear what it means from an economic point of view. The problem is that since continuity is not required in the more traditional applications used in economics (e.g., matching on observables), it is not obvious what assumptions about the behavior of economic agents are required to get continuity.

Second, RD designs are a fairly peculiar application of a “selection on observables” model. Indeed, the view in James J. Heckman, Robert J. Lalonde, and Jeffrey A. Smith (1999) was that “[r]egression discontinuity estimators constitute a special case of selection on observables,” and that the RD estimator is “a limit form of matching at one point.” In general, we need two crucial conditions for a matching/selection on observables approach to work. First, treatment must be randomly assigned conditional on observables (the *ignorability* or *unconfoundedness* assumption). In practice, this is typically viewed as a strong, and not particularly credible, assumption. For instance, in a standard regression framework this amounts to assuming that all relevant factors are controlled for, and that no omitted variables are correlated with the treatment dummy. In an RD design, however, this crucial assumption is trivially satisfied. When  $X \geq c$ , the treatment dummy  $D$  is always equal to 1. When  $X < c$ ,  $D$  is always equal to 0. Conditional on  $X$ , there is no variation left in  $D$ , so it

cannot, therefore, be correlated with any other factor.<sup>9</sup>

At the same time, the other standard assumption of *overlap* is violated since, strictly speaking, it is not possible to observe units with either  $D = 0$  or  $D = 1$  for a given value of the assignment variable  $X$ . This is the reason the continuity assumption is required—to compensate for the failure of the overlap condition. So while we cannot observe treatment and non-treatment for the same value of  $X$ , we can observe the two outcomes for values of  $X$  around the cutoff point that are arbitrarily close to each other.

### 2.3 RD Design as a Local Randomized Experiment

When looking at RD designs in this way, one could get the impression that they require some assumptions to be satisfied, while other methods such as matching on observables and IV methods simply require other assumptions.<sup>10</sup> From this point of view, it would seem that the assumptions for the RD design are just as arbitrary as those used for other methods. As we discuss throughout the paper, however, we do not believe this way of looking at RD designs does justice to their important advantages over most other existing methods. This point becomes much clearer once we compare the RD design to the “gold standard” of program evaluation methods, randomized experiments. We will show that the RD design is a much closer cousin of randomized experiments than other competing methods.

<sup>9</sup> In technical terms, the treatment dummy  $D$  follows a degenerate (concentrated at  $D = 0$  or  $D = 1$ ), but nonetheless random distribution conditional on  $X$ . Ignorability is thus trivially satisfied.

<sup>10</sup> For instance, in the survey of Angrist and Alan B. Krueger (1999), RD is viewed as an IV estimator, thus having essentially the same potential drawbacks and pitfalls.

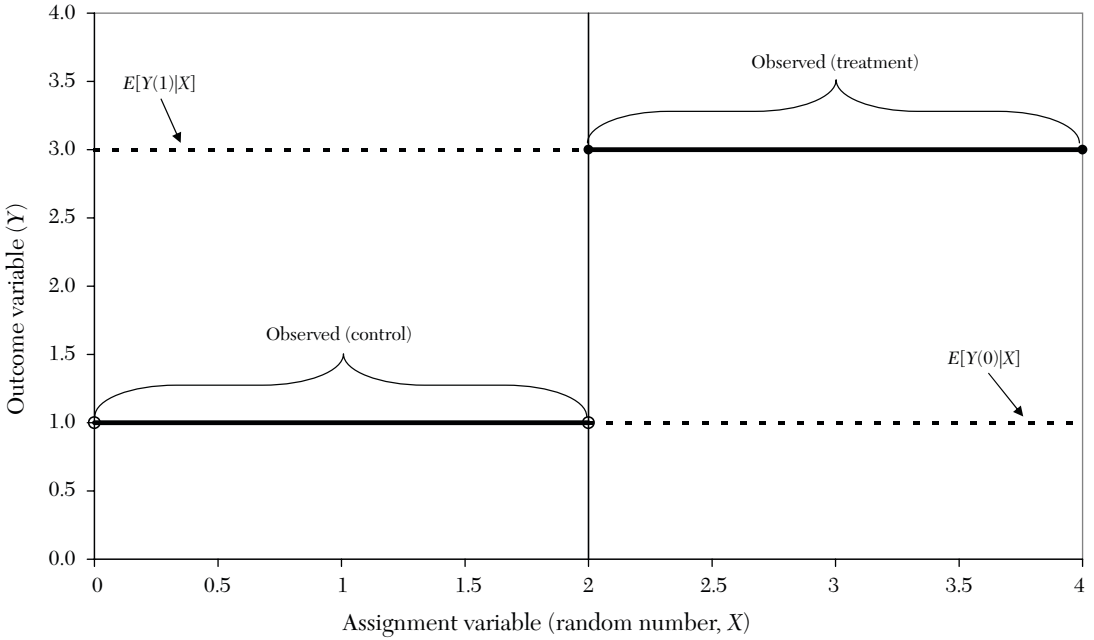


Figure 3. Randomized Experiment as a RD Design

In a randomized experiment, units are typically divided into treatment and control groups on the basis of a randomly generated number,  $\nu$ . For example, if  $\nu$  follows a uniform distribution over the range  $[0, 4]$ , units with  $\nu \geq 2$  are given the treatment while units with  $\nu < 2$  are denied treatment. So the randomized experiment can be thought of as an RD design where the assignment variable is  $X = \nu$  and the cutoff is  $c = 2$ . Figure 3 shows this special case in the potential outcomes framework, just as in the more general RD design case of figure 2. The difference is that because the assignment variable  $X$  is now completely random, it is independent of the potential outcomes  $Y_i(0)$  and  $Y_i(1)$ , and the curves  $E[Y_i(1)|X]$  and  $E[Y_i(0)|X]$  are flat. Since the curves are flat, it trivially follows that they are also continuous at the cutoff point  $X = c$ . In other

words, continuity is a direct consequence of randomization.

The fact that the curves  $E[Y_i(1)|X]$  and  $E[Y_i(0)|X]$  are flat in a randomized experiment implies that, as is well known, the average treatment effect can be computed as the difference in the mean value of  $Y$  on the right and left hand side of the cutoff. One could also use an RD approach by running regressions of  $Y$  on  $X$ , but this would be less efficient since we know that if randomization were successful, then  $X$  is an irrelevant variable in this regression.

But now imagine that, for ethical reasons, people are compensated for having received a “bad draw” by getting a monetary compensation inversely proportional to the random number  $X$ . For example, the treatment could be job search assistance for the unemployed, and the outcome whether one found a job

within a month of receiving the treatment. If people with a larger monetary compensation can afford to take more time looking for a job, the potential outcome curves will no longer be flat and will slope upward. The reason is that having a higher random number, i.e., a lower monetary compensation, increases the probability of finding a job. So in this “smoothly contaminated” randomized experiment, the potential outcome curves will instead look like the classical RD design case depicted in figure 2.

Unlike a classical randomized experiment, in this contaminated experiment a simple comparison of means no longer yields a consistent estimate of the treatment effect. By focusing right around the threshold, however, an RD approach would still yield a consistent estimate of the treatment effect associated with job search assistance. The reason is that since people just above or below the cutoff receive (essentially) the same monetary compensation, we still have locally a randomized experiment around the cutoff point. Furthermore, as in a randomized experiment, it is possible to test whether randomization “worked” by comparing the local values of baseline covariates on the two sides of the cutoff value.

Of course, this particular example is highly artificial. Since we know the monetary compensation is a continuous function of  $X$ , we also know the continuity assumption required for the RD estimates of the treatment effect to be consistent is also satisfied. The important result, due to Lee (2008), that we will show in the next section is that the conditions under which we locally have a randomized experiment (and continuity) right around the cutoff point are remarkably weak. Furthermore, in addition to being weak, the conditions for local randomization are testable in the same way global randomization is testable in a randomized experiment by looking at whether baseline covariates are balanced. It is in this sense

that the RD design is more closely related to randomized experiments than to other popular program evaluation methods such as matching on observables, difference-in-differences, and IV.

### 3. Identification and Interpretation

This section discusses a number of issues of identification and interpretation that arise when considering an RD design. Specifically, the applied researcher may be interested in knowing the answers to the following questions:

1. How do I know whether an RD design is appropriate for my context? When are the identification assumptions plausible or implausible?
2. Is there any way I can test those assumptions?
3. To what extent are results from RD designs generalizable?

On the surface, the answers to these questions seem straightforward: (1) “An RD design will be appropriate if it is plausible that all other unobservable factors are “continuously” related to the assignment variable,” (2) “No, the continuity assumption is necessary, so there are no tests for the validity of the design,” and (3) “The RD estimate of the treatment effect is only applicable to the subpopulation of individuals at the discontinuity threshold, and uninformative about the effect anywhere else.” These answers suggest that the RD design is no more compelling than, say, an instrumental variables approach, for which the analogous answers would be (1) “The instrument must be uncorrelated with the error in the outcome equation,” (2) “The identification assumption is ultimately untestable,” and (3) “The estimated treatment effect is applicable

to the subpopulation whose treatment was affected by the instrument.” After all, who’s to say whether one untestable design is more “compelling” or “credible” than another untestable design? And it would seem that having a treatment effect for a vanishingly small subpopulation (those at the threshold, in the limit) is hardly more (and probably much less) useful than that for a population “affected by the instrument.”

As we describe below, however, a closer examination of the RD design reveals quite different answers to the above three questions:

1. “When there is a continuously distributed stochastic error component to the assignment variable—which can occur when optimizing agents do not have *precise* control over the assignment variable—then the variation in the treatment will be as good as randomized in a neighborhood around the discontinuity threshold.”
2. “Yes. As in a randomized experiment, the distribution of observed baseline covariates should not change discontinuously at the threshold.”
3. “The RD estimand can be interpreted as a weighted average treatment effect, where the weights are the relative *ex ante* probability that the value of an individual’s assignment variable will be in the neighborhood of the threshold.”

Thus, in many contexts, the RD design may have more in common with randomized experiments (or circumstances when an instrument is truly randomized)—in terms of their “internal validity” and how to implement them in practice—than with regression control or matching methods, instrumental variables, or panel data approaches. We will return to this point after first discussing the above three issues in greater detail.

### 3.1 *Valid or Invalid RD?*

Are individuals able to influence the assignment variable, and if so, what is the nature of this control? This is probably the most important question to ask when assessing whether a particular application should be analyzed as an RD design. If individuals have a great deal of control over the assignment variable and if there is a perceived benefit to a treatment, one would certainly expect individuals on one side of the threshold to be systematically different from those on the other side.

Consider the test-taking RD example. Suppose there are two types of students: *A* and *B*. Suppose type *A* students are more able than *B* types, and that *A* types are also keenly aware that passing the relevant threshold (50 percent) will give them a scholarship benefit, while *B* types are completely ignorant of the scholarship and the rule. Now suppose that 50 percent of the questions are trivial to answer correctly but, due to random chance, students will sometimes make careless errors when they initially answer the test questions, but would certainly correct the errors if they checked their work. In this scenario, only type *A* students will make sure to check their answers before turning in the exam, thereby assuring themselves of a passing score. Thus, while we would expect those who barely passed the exam to be a mixture of type *A* and type *B* students, those who barely failed would exclusively be type *B* students. In this example, it is clear that the marginal failing students do *not* represent a valid counterfactual for the marginal passing students. Analyzing this scenario within an RD framework would be inappropriate.

On the other hand, consider the same scenario, except assume that questions on the exam are *not* trivial; there are no guaranteed passes, no matter how many times the students check their answers before turning in the exam. In this case, it seems more

plausible that, among those scoring near the threshold, it is a matter of “luck” as to which side of the threshold they land. Type A students can exert more effort—because they know a scholarship is at stake—but they do not know the exact score they will obtain. In this scenario, it would be reasonable to argue that those who marginally failed and passed would be otherwise comparable, and that an RD analysis *would* be appropriate and would yield credible estimates of the impact of the scholarship.

These two examples make it clear that one must have some knowledge about the mechanism generating the assignment variable beyond knowing that, if it crosses the threshold, the treatment is “turned on.” It is “folk wisdom” in the literature to judge whether the RD is appropriate based on whether individuals could manipulate the assignment variable and *precisely* “sort” around the discontinuity threshold. The key word here is “precise” rather than “manipulate.” After all, in both examples above, individuals do exert some control over the test score. And indeed, in virtually every known application of the RD design, it is easy to tell a plausible story that the assignment variable is to some degree influenced by *someone*. But individuals will not always be able to have *precise* control over the assignment variable. It should perhaps seem obvious that it is necessary to rule out precise sorting to justify the use of an RD design. After all, individual self-selection into treatment or control regimes is exactly why simple comparison of means is unlikely to yield valid causal inferences. Precise sorting around the threshold is self-selection.

What is not obvious, however, is that, when one formalizes the notion of having imprecise control over the assignment variable, there is a striking consequence: the variation in the treatment in a neighborhood of the threshold is “as good as randomized.” We explain this below.

### 3.1.1 Randomized Experiments from Nonrandom Selection

To see how the inability to precisely control the assignment variable leads to a source of randomized variation in the treatment, consider a simplified formulation of the RD design:<sup>11</sup>

$$(2) \quad Y = D\tau + W\delta_1 + U$$

$$D = 1[X \geq c]$$

$$X = W\delta_2 + V,$$

where  $Y$  is the outcome of interest,  $D$  is the binary treatment indicator, and  $W$  is the vector of all predetermined and observable characteristics of the individual that might impact the outcome and/or the assignment variable  $X$ .

This model looks like a standard endogenous dummy variable set-up, except that we observe the assignment variable,  $X$ . This allows us to relax most of the other assumptions usually made in this type of model. First, we allow  $W$  to be endogenously determined as long as it is determined prior to  $V$ . Second, we take no stance as to whether some elements of  $\delta_1$  or  $\delta_2$  are zero (exclusion restrictions). Third, we make no assumptions about the correlations between  $W$ ,  $U$ , and  $V$ .<sup>12</sup>

In this model, individual heterogeneity in the outcome is completely described by the pair of random variables  $(W, U)$ ; anyone with the same values of  $(W, U)$  will have one of two values for the outcome, depending on whether they receive treatment. Note that,

<sup>11</sup> We use a simple linear endogenous dummy variable setup to describe the results in this section, but all of the results could be stated within the standard potential outcomes framework, as in Lee (2008).

<sup>12</sup> This is much less restrictive than textbook descriptions of endogenous dummy variable systems. It is typically assumed that  $(U, V)$  is independent of  $W$ .

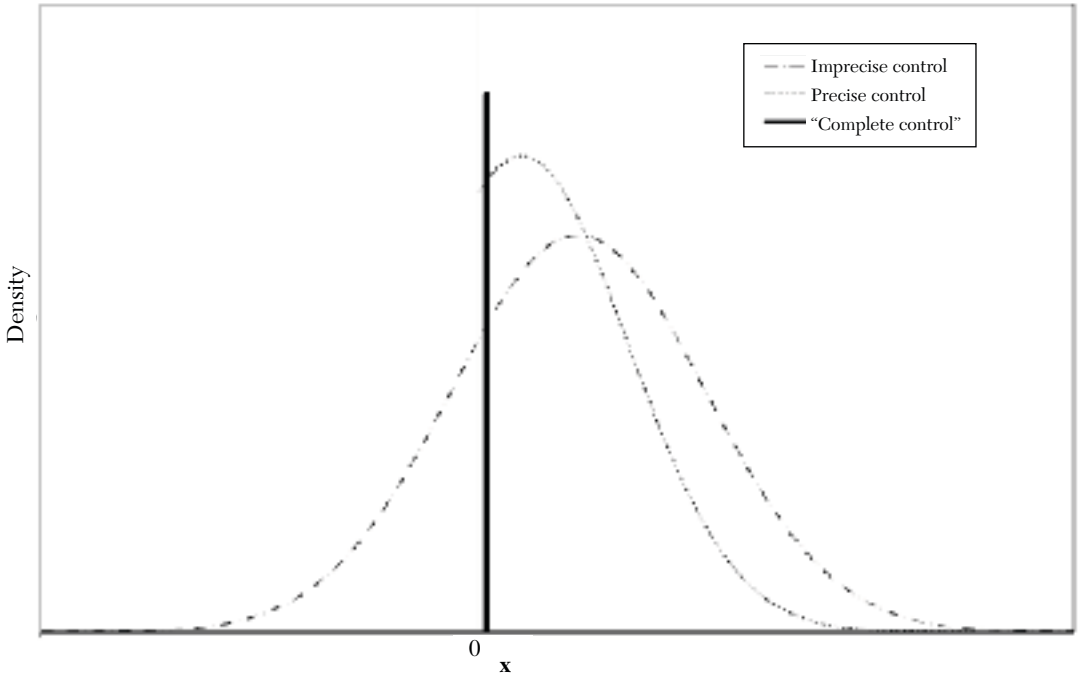


Figure 4. Density of Assignment Variable Conditional on  $W = w, U = u$

since RD designs are implemented by running regressions of  $Y$  on  $X$ , equation (2) looks peculiar since  $X$  is not included with  $W$  and  $U$  on the right hand side of the equation. We could add a function of  $X$  to the outcome equation, but this would not make a difference since we have not made any assumptions about the joint distribution of  $W, U$ , and  $V$ . For example, our setup allows for the case where  $U = X\delta_3 + U'$ , which yields the outcome equation  $Y = D\tau + W\delta_1 + X\delta_3 + U'$ . For the sake of simplicity, we work with the simple case where  $X$  is not included on the right hand side of the equation.<sup>13</sup>

Now consider the distribution of  $X$ , conditional on a particular pair of values  $W = w, U = u$ . It is equivalent (up to a translational shift) to the distribution of  $V$  conditional on  $W = w, U = u$ . If an individual has complete and exact control over  $X$ , we would model it as having a degenerate distribution, conditional on  $W = w, U = u$ . That is, in repeated trials, this individual would choose the same score. This is depicted in figure 4 as the thick line.

If there is some room for error but individuals can nevertheless have precise control about whether they will fail to receive the

<sup>13</sup> When RD designs are implemented in practice, the estimated effect of  $X$  on  $Y$  can either reflect a true causal effect of  $X$  on  $Y$  or a spurious correlation between  $X$  and the

unobservable term  $U$ . Since it is not possible to distinguish between these two effects in practice, we simplify the setup by implicitly assuming that  $X$  only comes into equation (2) indirectly through its (spurious) correlation with  $U$ .

treatment, then we would expect the density of  $X$  to be zero just below the threshold, but positive just above the threshold, as depicted in figure 4 as the truncated distribution. This density would be one way to model the first example described above for the type  $A$  students. Since type  $A$  students know about the scholarship, they will double-check their answers and make sure they answer the easy questions, which comprise 50 percent of the test. How high they score above the passing threshold will be determined by some randomness.

Finally, if there is stochastic error in the assignment variable and individuals do *not* have precise control over the assignment variable, we would expect the density of  $X$  (and hence  $V$ ), conditional on  $W = w, U = u$  to be continuous at the discontinuity threshold, as shown in figure 4 as the untruncated distribution.<sup>14</sup> It is important to emphasize that, in this final scenario, the individual still has control over  $X$ : through her efforts, she can choose to shift the distribution to the right. This is the density for someone with  $W = w, U = u$ , but may well be different—with a different mean, variance, or shape of the density—for other individuals, with different levels of ability, who make different choices. We are assuming, however, that all individuals are unable to precisely control the score just around the threshold.

**Definition:** We say individuals have imprecise control over  $X$  when conditional on  $W = w$  and  $U = u$ , the density of  $V$  (and hence  $X$ ) is continuous.

When individuals have imprecise control over  $X$  this leads to the striking implication that variation in treatment status will be

<sup>14</sup> For example, this would be plausible when  $X$  is a test score modeled as a sum of Bernoulli random variables, which is approximately normal by the central limit theorem.

randomized in a neighborhood of the threshold. To see this, note that by Bayes' Rule, we have

$$(3) \quad \Pr[W = w, U = u | X = x] \\ = f(x | W = w, U = u) \frac{\Pr[W = w, U = u]}{f(x)},$$

where  $f(\cdot)$  and  $f(\cdot | \cdot)$  are marginal and conditional densities for  $X$ . So when  $f(x | W = w, U = u)$  is continuous in  $x$ , the right hand side will be continuous in  $x$ , which therefore means that the distribution of  $W, U$  conditional on  $X$  will be continuous in  $x$ .<sup>15</sup> That is, *all observed and unobserved predetermined characteristics will have identical distributions on either side of  $x = c$ , in the limit, as we examine smaller and smaller neighborhoods of the threshold.*

In sum,

**Local Randomization:** If individuals have imprecise control over  $X$  as defined above, then  $\Pr[W = w, U = u | X = x]$  is continuous in  $x$ : the treatment is “as good as” randomly assigned around the cutoff.

In other words, the behavioral assumption that individuals do not precisely manipulate  $X$  around the threshold has the *prediction* that treatment is locally randomized.

This is perhaps why RD designs can be so compelling. A deeper investigation into the real-world details of how  $X$  (and hence  $D$ ) is determined can help assess whether it is plausible that individuals have precise or imprecise control over  $X$ . By contrast, with

<sup>15</sup> Since the potential outcomes  $Y(0)$  and  $Y(1)$  are functions of  $W$  and  $U$ , it follows that the distribution of  $Y(0)$  and  $Y(1)$  conditional on  $X$  is also continuous in  $x$  when individuals have imprecise control over  $X$ . This implies that the conditions usually invoked for consistently estimating the treatment effect (the conditional means  $E[Y(0) | X = x]$  and  $E[Y(1) | X = x]$  being continuous in  $x$ ) are also satisfied. See Lee (2008) for more detail.

most nonexperimental evaluation contexts, learning about how the treatment variable is determined will rarely lead one to conclude that it is “as good as” randomly assigned.

### 3.2 Consequences of Local Random Assignment

There are three practical implications of the above local random assignment result.

#### 3.2.1 Identification of the Treatment Effect

First and foremost, it means that the discontinuity gap at the cutoff identifies the treatment effect of interest. Specifically, we have

$$\begin{aligned} & \lim_{\varepsilon \downarrow 0} E[Y|X = c + \varepsilon] \\ & - \lim_{\varepsilon \uparrow 0} E[Y|X = c + \varepsilon] \\ & = \tau + \lim_{\varepsilon \downarrow 0} \sum_{w,u} (w\delta_1 + u) \\ & \quad \times \Pr[W = w, U = u | X = c + \varepsilon] \\ & - \lim_{\varepsilon \uparrow 0} \sum_{w,u} (w\delta_1 + u) \\ & \quad \times \Pr[W = w, U = u | X = c + \varepsilon] \\ & = \tau, \end{aligned}$$

where the last line follows from the continuity of  $\Pr[W = w, U = u | X = x]$ .

As we mentioned earlier, nothing changes if we augment the model by adding a direct impact of  $X$  itself in the outcome equation, as long as the effect of  $X$  on  $Y$  does not jump at the cutoff. For example, in the example of Thistlethwaite and Campbell (1960), we can allow higher test scores to improve future academic outcomes (perhaps by raising the probability of admission to higher quality schools) as long as that probability does not jump at precisely the same cutoff used to award scholarships.

#### 3.2.2 Testing the Validity of the RD Design

An almost equally important implication of the above local random assignment result is that it makes it possible to empirically assess the prediction that  $\Pr[W = w, U = u | X = x]$  is continuous in  $x$ . Although it is impossible to test this directly—since  $U$  is unobserved—it is nevertheless possible to assess whether  $\Pr[W = w | X = x]$  is continuous in  $x$  at the threshold. A discontinuity would indicate a failure of the identifying assumption.

This is akin to the tests performed to empirically assess whether the randomization was carried out properly in randomized experiments. It is standard in these analyses to demonstrate that treatment and control groups are similar in their observed baseline covariates. It is similarly impossible to test whether unobserved characteristics are balanced in the experimental context, so the most favorable statement that can be made about the experiment is that the data “failed to reject” the assumption of randomization.

Performing this kind of test is arguably more important in the RD design than in the experimental context. After all, the true nature of individuals’ control over the assignment variable—and whether it is precise or imprecise—may well be somewhat debatable even after a great deal of investigation into the exact treatment-assignment mechanism (which itself is always advisable to do). Imprecision of control will often be nothing more than a conjecture, but thankfully it has testable predictions.

There is a complementary, and arguably more direct and intuitive test of the imprecision of control over the assignment variable: examination of the density of  $X$  itself, as suggested in Justin McCrary (2008). If the density of  $X$  for each individual is continuous, then the marginal density of  $X$  over the population should be continuous as well. A jump in the density at the threshold is probably the most direct evidence of some degree

of sorting around the threshold, and should provoke serious skepticism about the appropriateness of the RD design.<sup>16</sup> Furthermore, one advantage of the test is that it can always be performed in a RD setting, while testing whether the covariates  $W$  are balanced at the threshold depends on the availability of data on these covariates.

This test is also a partial one. Whether each individual's ex ante density of  $X$  is continuous is fundamentally untestable since, for each individual, we only observe one realization of  $X$ . Thus, in principle, at the threshold some individuals' densities may jump up while others may sharply fall, so that in the aggregate, positives and negatives offset each other making the density appear continuous. In recent applications of RD such occurrences seem far-fetched. Even if this were the case, one would certainly expect to see, after stratifying by different values of the observable characteristics, some discontinuities in the density of  $X$ . These discontinuities could be detected by performing the local randomization test described above.

### 3.2.3 Irrelevance of Including Baseline Covariates

A consequence of a randomized experiment is that the assignment to treatment is, by construction, independent of the baseline covariates. As such, it is not necessary to include them to obtain consistent estimates of the treatment effect. In practice, however,

<sup>16</sup>Another possible source of discontinuity in the density of the assignment variable  $X$  is selective attrition. For example, John DiNardo and Lee (2004) look at the effect of unionization on wages several years after a union representation vote was taken. In principle, if firms that were unionized because of a majority vote are more likely to close down, then conditional on firm survival at a later date, there will be a discontinuity in  $X$  (the vote share) that could threaten the validity of the RD design for estimating the effect of unionization on wages (conditional on survival). In that setting, testing for a discontinuity in the density (conditional on survival) is similar to testing for selective attrition (linked to treatment status) in a standard randomized experiment.

researchers will include them in regressions, because doing so can reduce the sampling variability in the estimator. Arguably the greatest potential for this occurs when one of the baseline covariates is a pre-random-assignment observation on the dependent variable, which may likely be highly correlated with the post-assignment outcome variable of interest.

The local random assignment result allows us to apply these ideas to the RD context. For example, if the lagged value of the dependent variable was determined prior to the realization of  $X$ , then the local randomization result will imply that that lagged dependent variable will have a continuous relationship with  $X$ . Thus, performing an RD analysis on  $Y$  minus its lagged value should also yield the treatment effect of interest. The hope, however, is that the differenced outcome measure will have a sufficiently lower variance than the level of the outcome, so as to lower the variance in the RD estimator.

More formally, we have

$$\begin{aligned} & \lim_{\varepsilon \downarrow 0} E[Y - W\pi | X = c + \varepsilon] \\ & - \lim_{\varepsilon \uparrow 0} E[Y - W\pi | X = c + \varepsilon] \\ & = \tau + \lim_{\varepsilon \downarrow 0} \sum_{w,u} (w(\delta_1 - \pi) + u) \\ & \quad \times \Pr[W = w, U = u | X = c + \varepsilon] \\ & - \lim_{\varepsilon \uparrow 0} \sum_{w,u} (w(\delta_1 - \pi) + u) \\ & \quad \times \Pr[W = w, U = u | X = c + \varepsilon] \\ & = \tau, \end{aligned}$$

where  $W\pi$  is *any* linear function, and  $W$  can include a lagged dependent variable, for example. We return to how to implement this in practice in section 4.4.

### 3.3 Generalizability: The RD Gap as a Weighted Average Treatment Effect

In the presence of heterogeneous treatment effects, the discontinuity gap in an RD design can be interpreted as a *weighted* average treatment effect across *all* individuals. This is somewhat contrary to the temptation to conclude that the RD design only delivers a credible treatment effect for the subpopulation of individuals at the threshold and says nothing about the treatment effect “away from the threshold.” Depending on the context, this may be an overly simplistic and pessimistic assessment.

Consider the scholarship test example again, and define the “treatment” as “receiving a scholarship by scoring 50 percent or greater on the scholarship exam.” Recall that the pair  $W, U$  characterizes individual heterogeneity. We now let  $\tau(w, u)$  denote the treatment effect for an individual with  $W = w$  and  $U = u$ , so that the outcome equation in (2) is instead given by

$$Y = D\tau(W, U) + W\delta_1 + U.$$

This is essentially a model of completely unrestricted heterogeneity in the treatment effect. Following the same line of argument as above, we obtain

$$\begin{aligned} (5) \quad & \lim_{\varepsilon \downarrow 0} E[Y|X = c + \varepsilon] \\ & - \lim_{\varepsilon \uparrow 0} E[Y|X = c + \varepsilon] \\ & = \sum_{w,u} \tau(w, u) \Pr[W = w, U = u | X = c] \\ & = \sum_{w,u} \tau(w, u) \frac{f(c|W = w, U = u)}{f(c)} \\ & \quad \times \Pr[W = w, U = u], \end{aligned}$$

where the second line follows from equation (3).

The discontinuity gap then, is a particular kind of average treatment effect *across all individuals*. If not for the term  $f(c|W = w, U = u)/f(c)$ , it would be the average treatment effect for the entire population. The presence of the ratio  $f(c|W = w, U = u)/f(c)$  implies the discontinuity is instead a *weighted* average treatment effect where the weights are directly proportional to the ex ante likelihood that an individual’s realization of  $X$  will be close to the threshold. All individuals could get some weight, and the similarity of the weights across individuals is ultimately untestable, since again we only observe one realization of  $X$  per person and do not know anything about the ex ante probability distribution of  $X$  for any one individual. The weights may be relatively similar across individuals, in which case the RD gap would be closer to the overall average treatment effect; but, if the weights are highly varied and also related to the magnitude of the treatment effect, then the RD gap would be very different from the overall average treatment effect. While it is not possible to know how close the RD gap is from the overall average treatment effect, it remains the case that the treatment effect estimated using a RD design is averaged over a larger population than one would have anticipated from a purely “cut-off” interpretation.

Of course, we do not observe the density of the assignment variable at the individual level so we therefore do not know the weight for each individual. Indeed, if the signal to noise ratio of the test is extremely high, someone who scores a 90 percent may have almost a zero chance of scoring near the threshold, implying that the RD gap is almost entirely dominated by those who score near 50 percent. But if the reliability is lower, then the RD gap applies to a relatively broader subpopulation. It remains to be seen whether or not and how information on the reliability, or a second test measurement, or other

covariates that can predict the assignment could be used in conjunction with the RD gap to learn about average treatment effects for the overall population. The understanding of the RD gap as a weighted average treatment effect serves to highlight that RD causal evidence is not somehow fundamentally disconnected from the average treatment effect that is often of interest to researchers.

It is important to emphasize that the RD gap is not informative about the treatment if it were defined as “receipt of a scholarship that is awarded by scoring 90 percent or higher on the scholarship exam.” This is not so much a “drawback” of the RD design as a limitation shared with even a carefully controlled randomized experiment. For example, if we randomly assigned financial aid awards to low-achieving students, whatever treatment effect we estimate may not be informative about the effect of financial aid for high-achieving students.

In some contexts, the treatment effect “away from the discontinuity threshold” may not make much practical sense. Consider the RD analysis of incumbency in congressional elections of Lee (2008). When the treatment is “being the incumbent party,” it is implicitly understood that incumbency entails winning the previous election by obtaining at least 50 percent of the vote.<sup>17</sup> In the election context, the treatment “being the incumbent party by virtue of winning an election, whereby 90 percent of the vote is required to win” simply does not apply to any real-life situation. Thus, in this context, it is awkward to interpret the RD gap as “the effect of incumbency that exists at 50 percent vote-share threshold” (as if there is an effect at a 90 percent threshold). Instead it is more natural to interpret the RD gap as estimating a weighted average treatment effect of incumbency across all districts, where more

weight is given to those districts in which a close election race was expected.

### 3.4 Variations on the Regression Discontinuity Design

To this point, we have focused exclusively on the “classic” RD design introduced by Thistlethwaite and Campbell (1960), whereby there is a single binary treatment and the assignment variable perfectly predicts treatment receipt. We now discuss two variants of this base case: (1) when there is so-called “imperfect compliance” of the rule and (2) when the treatment of interest is a continuous variable.

In both cases, the notion that the RD design generates local variation in treatment that is “as good as randomly assigned” is helpful because we can apply known results for randomized instruments to the RD design, as we do below. The notion is also helpful for addressing other data problems, such as differential attrition or sample selection, whereby the treatment affects whether or not you observe the outcome of interest. The local random assignment result means that, in principle, one could extend the ideas of Joel L. Horowitz and Charles F. Manski (2000) or Lee (2009), for example, to provide bounds on the treatment effect, accounting for possible sample selection bias.

#### 3.4.1. Imperfect Compliance: The “Fuzzy” RD

In many settings of economic interest, treatment is determined partly by whether the assignment variable crosses a cutoff point. This situation is very important in practice for a variety of reasons, including cases of imperfect take-up by program participants or when factors other than the threshold rule affect the probability of program participation. Starting with William M. K. Trochim (1984), this setting has been referred to as a “fuzzy” RD design. In the case we have discussed so far—the “sharp” RD design—the

<sup>17</sup> For this example, consider the simplified case of a two-party system.

probability of treatment jumps from 0 to 1 when  $X$  crosses the threshold  $c$ . The fuzzy RD design allows for a smaller jump in the probability of assignment to the treatment at the threshold and only requires

$$\lim_{\varepsilon \downarrow 0} \Pr(D = 1 | X = c + \varepsilon) \\ \neq \lim_{\varepsilon \uparrow 0} \Pr(D = 1 | X = c + \varepsilon).$$

Since the probability of treatment jumps by less than one at the threshold, the jump in the relationship between  $Y$  and  $X$  can no longer be interpreted as an average treatment effect. As in an instrumental variable setting however, the treatment effect can be recovered by dividing the jump in the relationship between  $Y$  and  $X$  at  $c$  by the fraction induced to take-up the treatment at the threshold—in other words, the discontinuity jump in the relation between  $D$  and  $X$ . In this setting, the treatment effect can be written as

$$\tau_F = \frac{\lim_{\varepsilon \downarrow 0} E[Y | X = c + \varepsilon] - \lim_{\varepsilon \uparrow 0} E[Y | X = c + \varepsilon]}{\lim_{\varepsilon \downarrow 0} E[D | X = c + \varepsilon] - \lim_{\varepsilon \uparrow 0} E[D | X = c + \varepsilon]},$$

where the subscript “F” refers to the fuzzy RD design.

There is a close analogy between how the treatment effect is defined in the fuzzy RD design and in the well-known “Wald” formulation of the treatment effect in an instrumental variables setting. Hahn, Todd and van der Klaauw (2001) were the first to show this important connection and to suggest estimating the treatment effect using two-stage least-squares (TSLS) in this setting. We discuss estimation of fuzzy RD designs in greater detail in section 4.3.3.

Hahn, Todd and van der Klaauw (2001) furthermore pointed out that the interpretation of this ratio as a causal effect requires the same assumptions as in Imbens and Angrist (1994). That is, one must assume “monotonicity” (i.e.,  $X$  crossing the cutoff cannot simultaneously *cause* some units to take up and others to reject the treatment)

and “excludability” (i.e.,  $X$  crossing the cutoff cannot impact  $Y$  except through impacting receipt of treatment). When these assumptions are made, it follows that<sup>18</sup>

$$\tau_F = E[Y(1) - Y(0) | \text{unit is complier}, X = c],$$

where “compliers” are units that receive the treatment when they satisfy the cutoff rule ( $X_i \geq c$ ), but would not otherwise receive it.

In summary, if there is local random assignment (e.g., due to the plausibility of individuals’ imprecise control over  $X$ ), then we can simply apply all of what is known about the assumptions and interpretability of instrumental variables. The difference between the “sharp” and “fuzzy” RD design is exactly parallel to the difference between the randomized experiment with perfect compliance and the case of imperfect compliance, when only the “intent to treat” is randomized.

For example, in the case of imperfect compliance, even if a proposed binary instrument  $Z$  is randomized, it is necessary to rule out the possibility that  $Z$  affects the outcome, outside of its influence through treatment receipt,  $D$ . Only then will the instrumental variables estimand—the ratio of the reduced form effects of  $Z$  on  $Y$  and of  $Z$  on  $D$ —be properly interpreted as a causal effect of  $D$  on  $Y$ . Similarly, supposing that individuals do not have precise control over  $X$ , it is necessary to assume that whether  $X$  crosses the threshold  $c$  (the instrument) has no impact on  $y$  except by influencing  $D$ . Only then will the ratio of the two RD gaps in  $Y$  and  $D$  be properly interpreted as a causal effect of  $D$  on  $Y$ .

In the same way that it is important to verify a strong first-stage relationship in an IV design, it is equally important to verify

<sup>18</sup> See Imbens and Lemieux (2008) for a more formal exposition.

that a discontinuity exists in the relationship between  $D$  and  $X$  in a fuzzy RD design.

Furthermore, in this binary-treatment–binary-instrument context with unrestricted heterogeneity in treatment effects, the IV estimand is interpreted as the average treatment effect “for the subpopulation affected by the instrument,” (or LATE). Analogously, the ratio of the RD gaps in  $Y$  and  $D$  (the “fuzzy design” estimand) can be interpreted as a *weighted* LATE, where the weights reflect the ex ante likelihood the individual’s  $X$  is near the threshold. In both cases, the exclusion restriction and monotonicity condition must hold.

### 3.4.2 Continuous Endogenous Regressor

In a context where the “treatment” is a continuous variable—call it  $T$ —and there is a randomized binary instrument (that can additionally be excluded from the outcome equation), an IV approach is an obvious way of obtaining an estimate of the impact of  $T$  on  $Y$ . The IV estimand is the reduced-form impact of  $Z$  on  $Y$  divided by the first-stage impact of  $Z$  on  $T$ .

The same is true for an RD design when the regressor of interest is continuous. Again, the causal impact of interest will still be the ratio of the two RD gaps (i.e., the discontinuities in  $Y$  and  $T$ ).

To see this more formally, consider the model

$$\begin{aligned} (6) \quad Y &= T\gamma + W\delta_1 + U_1 \\ T &= D\phi + W\gamma + U_2 \\ D &= 1[X \geq c] \\ X &= W\delta_2 + V, \end{aligned}$$

which is the same set-up as before, except with the added second equation, allowing for imperfect compliance or other factors (observables  $W$  or unobservables  $U_2$ ) to

impact the continuous regressor of interest  $T$ . If  $\gamma = 0$  and  $U_2 = 0$ , then the model collapses to a “sharp” RD design (with a continuous regressor).

Note that we make no additional assumptions about  $U_2$  (in terms of its correlation with  $W$  or  $V$ ). We do continue to assume imprecise control over  $X$  (conditional on  $W$  and  $U_1$ , the density of  $X$  is continuous).<sup>19</sup>

Given the discussion so far, it is easy to show that

$$\begin{aligned} (7) \quad & \lim_{\varepsilon \downarrow 0} E[Y|X = c + \varepsilon] \\ & - \lim_{\varepsilon \uparrow 0} E[Y|X = c + \varepsilon] \\ & = \left\{ \lim_{\varepsilon \downarrow 0} E[T|X = c + \varepsilon] \right. \\ & \quad \left. - \lim_{\varepsilon \uparrow 0} E[T|X = c + \varepsilon] \right\} \gamma. \end{aligned}$$

The left hand side is simply the “reduced form” discontinuity in the relation between  $y$  and  $X$ . The term preceding  $\gamma$  on the right hand side is the “first-stage” discontinuity in the relation between  $T$  and  $X$ , which is also estimable from the data. Thus, analogous to the exactly identified instrumental variable case, the ratio of the two discontinuities yields the parameter  $\gamma$ : the effect of  $T$  on  $Y$ . Again, because of the added notion of imperfect compliance, it is important to assume that  $D$  ( $X$  crossing the threshold) does not directly enter the outcome equation.

In some situations, more might be known about the rule determining  $T$ . For example, in Angrist and Lavy (1999) and Miguel Urquiola and Eric A. Verhoogen (2009), class size is an increasing function of total school enrollment, except for discontinuities at various enrollment thresholds. But

<sup>19</sup> Although it would be unnecessary to do so for the identification of  $\gamma$ , it would probably be more accurate to describe the situation of imprecise control with the continuity of the density of  $X$  conditional on the three variables ( $W, U_1, U_2$ ). This is because  $U_2$  is now another variable characterizing heterogeneity in individuals.

additional information about characteristics such as the slope and intercept of the underlying function (apart from the magnitude of the discontinuity) generally adds nothing to the identification strategy.

To see this, change the second equation in (6) to  $T = D\phi + g(X)$  where  $g(\cdot)$  is any continuous function in the assignment variable. Equation (7) will remain the same and, thus, knowledge of the function  $g(\cdot)$  is irrelevant for identification.<sup>20</sup>

There is also no need for additional theoretical results in the case when there is individual-level heterogeneity in the causal effect of the continuous regressor  $T$ . The local random assignment result allows us to borrow from the existing IV literature and interpret the ratio of the RD gaps as in Angrist and Krueger (1999), except that we need to add the note that all averages are weighted by the ex ante relative likelihood that the individual's  $X$  will land near the threshold.

### 3.5 Summary: A Comparison of RD and Other Evaluation Strategies

We conclude this section by comparing the RD design with other evaluation approaches. We believe it is helpful to view the RD design as a distinct approach rather than as a special case of either IV or matching/regression-control. Indeed, in important ways the RD design is more similar to a randomized experiment, which we illustrate below.

Consider a randomized experiment where subjects are assigned a random number  $X$  and are given the treatment if  $X \geq c$ . By construction,  $X$  is independent and not systematically related to any observable or unobservable characteristic determined prior to the randomization. This situation is illustrated in panel A of figure 5. The first column shows

the relationship between the treatment variable  $D$  and  $X$ , a step function, going from 0 to 1 at the  $X = c$  threshold. The second column shows the relationship between the observables  $W$  and  $X$ . This is flat because  $X$  is completely randomized. The same is true for the unobservable variable  $U$ , depicted in the third column. These three graphs capture the appeal of the randomized experiment: treatment varies while all other factors are kept constant (on average). And even though we cannot directly test whether there are no treatment-control differences in  $U$ , we can test whether there are such differences in the observable  $W$ .

Now consider an RD design (panel B of figure 5) where individuals have imprecise control over  $X$ . Both  $W$  and  $U$  may be systematically related to  $X$ , perhaps due to the actions taken by units to increase their probability of receiving treatment. Whatever the shape of the relation, as long as individuals have imprecise control over  $X$ , the relationship will be continuous. And therefore, as we examine  $Y$  near the  $X = c$  cutoff, we can be assured that like an experiment, treatment varies (the first column) while other factors are kept constant (the second and third columns). And, like an experiment, we can test this prediction by assessing whether observables truly are continuous with respect to  $X$  (the second column).<sup>21</sup>

We now consider two other commonly used nonexperimental approaches, referring to the model (2):

$$Y = D\tau + W\delta_1 + U$$

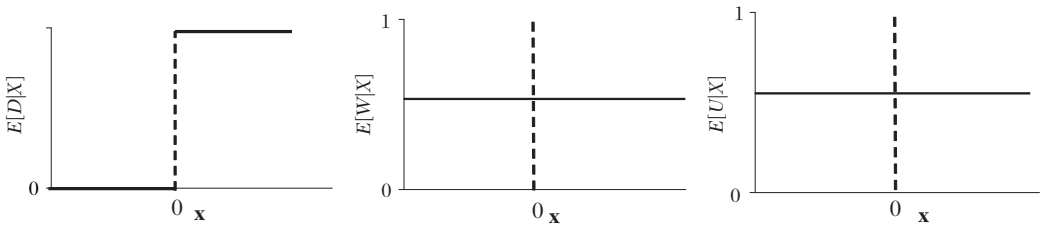
$$D = 1[X \geq c]$$

$$X = W\delta_2 + V.$$

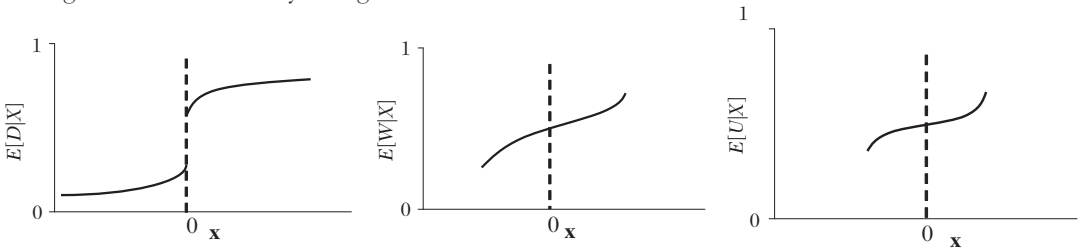
<sup>20</sup> As discussed in 3.2.1, the inclusion of a direct effect of  $X$  in the outcome equation will not change identification of  $\tau$ .

<sup>21</sup> We thank an anonymous referee for suggesting these illustrative graphs.

A. Randomized Experiment



B. Regression Discontinuity Design



C. Matching on Observables



D. Instrumental Variables

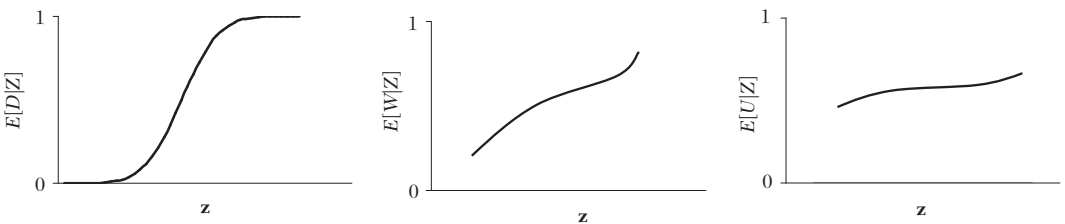


Figure 5. Treatment, Observables, and Unobservables in Four Research Designs

### 3.5.1 Selection on Observables: Matching/ Regression Control

The basic idea of the “selection on observables” approach is to adjust for differences in the  $W$ 's between treated and control individuals. It is usually motivated by the fact that it seems “implausible” that the unconditional mean  $Y$  for the control group represents a valid counterfactual for the treatment group. So it is argued that, *conditional on  $W$* , treatment-control contrasts may identify the ( $W$ -specific) treatment effect.

The underlying assumption is that conditional on  $W$ ,  $U$  and  $V$  are independent. From this it is clear that

$$\begin{aligned} E[Y|D = 1, W = w] \\ & - E[Y|D = 0, W = w] \\ & = \tau + E[U|W = w, V \geq c - w\delta_2] \\ & \quad - E[U|W = w, V < c - w\delta_2] \\ & = \tau. \end{aligned}$$

Two issues arise when implementing this approach. The first is one of functional form: how exactly to control for the  $W$ 's? When the  $W$ 's take on discrete values, one possibility is to compute treatment effects for each distinct value of  $W$ , and then average these effects across the constructed “cells.” This will not work, however, when  $W$  has continuous elements, in which case it is necessary to implement multivariate matching, propensity score, reweighting procedures, or nonparametric regressions.<sup>22</sup>

Regardless of the functional form issue, there is arguably a more fundamental question of which  $W$ 's to use in the analysis. While it is tempting to answer “all of them” and

hope that more  $W$ 's will lead to less biased estimates, this is obviously not necessarily the case. For example, consider estimating the economic returns to graduating high school (versus dropping out). It seems natural to include variables like parents' socioeconomic status, family income, year, and place of birth in the regression. Including more and more family-level  $W$ 's will ultimately lead to a “within-family” sibling analysis; extending it even further by including date of birth leads to a “within-twin-pair” analysis. And researchers have been critical—justifiably so—of this source of variation in education. The same reasons causing discomfort about the twin analyses should also cause skepticism about “kitchen sink” multivariate matching/propensity score/regression control analyses.<sup>23</sup>

It is also tempting to believe that, if the  $W$ 's do a “good job” in predicting  $D$ , the selection on observables approach will “work better.” But the opposite is true: in the extreme case when the  $W$ 's perfectly predict  $X$  (and hence  $D$ ), it is *impossible* to construct a treatment-control contrast for virtually all observations. For each value of  $W$ , the individuals will either all be treated or all control. In other words, there will be literally no overlap in the support of the propensity score for the treated and control observations. The propensity score would take the values of either 1 or 0.

The “selection on observables” approach is illustrated in panel C of figure 5. Observables  $W$  can help predict the probability of treatment (first column), but ultimately one must assume that unobservable factors  $U$  must be the same for treated and control units for

<sup>22</sup> See Hahn (1998) on including covariates directly with nonparametric regression.

<sup>23</sup> Researchers question the twin analyses on the grounds that it is not clear why one twin ends up having more education than the other, and that the assumption that education differences among twins is purely random (as ignorability would imply) is viewed as far-fetched. We thank David Card for pointing out this connection between twin analyses and matching approaches.

every value of  $W$ . That is, the crucial assumption is that the two lines in the third column be on top of each other. Importantly, there is no comparable graph in the second column because there is no way to test the design since all the  $W$ 's are used for estimation.

### 3.5.2 Selection on Unobservables:

#### *Instrumental Variables and "Heckit"*

A less restrictive modeling assumption is to allow  $U$  and  $V$  to be correlated, conditional on  $W$ . But because of the arguably "more realistic"/flexible data generating process, another assumption is needed to identify  $\tau$ . One such assumption is that some elements of  $W$  (call them  $Z$ ) enter the selection equation, but not the outcome equation and are also uncorrelated with  $U$ . An instrumental variables approach utilizes the fact that

$$\begin{aligned} E[Y|W^* = w^*, Z = z] & \\ &= E[D|W^* = w^*, Z = z]\tau + w^* \gamma \\ &\quad + E[U|W^* = w^*, Z = z] \\ &= E[D|W^* = w^*, Z = z]\tau + w^* \gamma \\ &\quad + E[U|W^* = w^*], \end{aligned}$$

where  $W$  has been split up into  $W^*$  and  $Z$  and  $\gamma$  is the corresponding coefficient for  $w^*$ . Conditional on  $W^* = w^*$ ,  $Y$  only varies with  $Z$  because of how  $D$  varies with  $Z$ . Thus, one identifies  $\tau$  by "dividing" the reduced form quantity  $E[D|W^* = w^*, Z = z]\tau$  (which can be obtained by examining the expectation of  $Y$  conditional on  $Z$  for a particular value  $w^*$  of  $W^*$ ) by  $E[D|W^* = w^*, Z = z]$ , which is also provided by the observed data. It is common to model the latter quantity as a linear function in  $Z$ , in which case the IV estimator is (conditional on  $W^*$ ) the ratio of coefficients from regressions of  $Y$  on  $Z$  and  $D$  on  $Z$ . When  $Z$  is binary, this appears to be the only way to identify  $\tau$  without imposing further assumptions.

When  $Z$  is continuous, there is an additional approach to identifying  $\tau$ . The "Heckit" approach uses the fact that

$$\begin{aligned} E[Y|W^* = w^*, Z = z, D = 1] & \\ &= \tau + w^* \gamma \\ &\quad + E[U|W^* = w^*, Z = z, V \geq c - w\delta_2] \\ E[Y|W^* = w^*, Z = z, D = 0] & \\ &= w^* \gamma \\ &\quad + E[U|W^* = w^*, Z = z, V < c - w\delta_2]. \end{aligned}$$

If we further assume a functional form for the joint distribution of  $U, V$ , conditional on  $W^*$  and  $Z$ , then the "control function" terms  $E[U|W = w, V \geq c - w\delta_2]$  and  $E[U|W = w, V < c - w\delta_2]$  are functions of observed variables, with the parameters then estimable from the data. It is then possible, for any value of  $W = w$ , to identify  $\tau$  as

$$\begin{aligned} (8) \quad &E[Y|W^* = w^*, Z = z, D = 1] \\ &- E[Y|W^* = w^*, Z = z, D = 0]) \\ &- (E[U|W^* = w^*, Z = z, V \geq c - w\delta_2] \\ &- E[U|W^* = w^*, Z = z, V < c - w\delta_2]). \end{aligned}$$

Even if the joint distribution of  $U, V$  is unknown, in principle it is still possible to identify  $\tau$ , if it were possible to choose two different values of  $Z$  such that  $c - w\delta_2$  approaches  $-\infty$  and  $\infty$ . If so, the last two terms in (8) approach  $E[U|W^* = w^*]$  and, hence, cancel one another. This is known as "identification at infinity."

Perhaps the most important assumption that any of these approaches require is the existence of a variable  $Z$  that is (conditional

on  $W^*$ ) independent of  $U$ .<sup>24</sup> There does not seem to be any way of testing the validity of this assumption. Different, but equally “plausible”  $Z$ ’s may lead to different answers, in the same way that including different sets of  $W$ ’s may lead to different answers in the selection on observables approach.

Even when there is a mechanism that justifies an instrument  $Z$  as “plausible,” it is often unclear which covariates  $W^*$  to include in the analysis. Again, when different sets of  $W^*$  lead to different answers, the question becomes which is more plausible:  $Z$  is independent of  $U$  conditional on  $W^*$  or  $Z$  is independent of  $U$  conditional on a *subset* of the variables in  $W^*$ ? While there may be some situations where knowledge of the mechanism dictates which variables to include, in other contexts, it may not be obvious.

The situation is illustrated in panel D of figure 5. It is necessary that the instrument  $Z$  is related to the treatment (as in the first column). The crucial assumption is regarding the relation between  $Z$  and the unobservables  $U$  (the third column). In order for an IV or a “Heckit” approach to work, the function in the third column needs to be flat. Of course, we cannot observe whether this is true. Furthermore, in most cases, it is unclear how to interpret the relation between  $W$  and  $Z$  (second column). Some might argue the observed relation between  $W$  and  $Z$  should be flat if  $Z$  is truly exogenous, and that if  $Z$  is highly correlated with  $W$ , then it casts doubt on  $Z$  being uncorrelated with  $U$ . Others will argue that using the second graph as a test is only appropriate when  $Z$  is truly randomized,

and that the assumption invoked is that  $Z$  is uncorrelated with  $U$ , *conditional on  $W$* . In this latter case, the design seems fundamentally untestable, since all the remaining observable variables (the  $W$ ’s) are being “used up” for identifying the treatment effect.

### 3.5.3 RD as “Design” not “Method”

RD designs can be valid under the more general “selection on unobservables” environment, allowing an arbitrary correlation among  $U$ ,  $V$ , and  $W$ , but at the same time not requiring an instrument. As discussed above, all that is needed is that conditional on  $W$ ,  $U$ , the density of  $V$  is continuous, and the local randomization result follows.

How is an RD design able to achieve this, given these weaker assumptions? The answer lies in what is absolutely necessary in an RD design: observability of the latent index  $X$ . Intuitively, given that both the “selection on observables” and “selection on unobservables” approaches rely heavily on modeling  $X$  and its components (e.g., which  $W$ ’s to include, and the properties of the unobservable error  $V$  and its relation to other variables, such as an instrument  $Z$ ), actually *knowing* the value of  $X$  ought to help.

In contrast to the “selection on observables” and “selection on unobservables” modeling approaches, with the RD design the researcher can avoid taking any strong stance about what  $W$ ’s to include in the analysis, since the design *predicts* that the  $W$ ’s are irrelevant and unnecessary for identification. Having data on  $W$ ’s is, of course, of some use, as they allow testing of the underlying assumption (described in section 4.4).

For this reason, it may be more helpful to consider RD designs as a description of a particular *data generating process*, rather than a “method” or even an “approach.” In virtually any context with an outcome variable  $Y$ , treatment status  $D$ , and other observable variables  $W$ , in principle a researcher can construct a regression-control or instrumental variables

<sup>24</sup> For IV, violation of this assumption essentially means that  $Z$  varies with  $Y$  for reasons other than its influence on  $D$ . For the textbook “Heckit” approach, it is typically assumed that  $U, V$  have the same distribution for any value of  $Z$ . It is also clear that the “identification at infinity” approach will only work if  $Z$  is uncorrelated with  $U$ , otherwise the last two terms in equation (8) would not cancel. See also the framework of Heckman and Edward Vytlacil (2005), which maintains the assumption of the independence of the error terms and  $Z$ , conditional on  $W^*$ .

(after designating one of the  $W$  variables a valid instrument) estimator, and state that the identification assumptions needed are satisfied.

This is not so with an RD design. Either the situation is such that  $X$  is observed, or it is not. If not, then the RD design simply does not apply.<sup>25</sup> If  $X$  is observed, then one has little choice but to attempt to estimate the expectation of  $Y$  conditional on  $X$  on either side of the cutoff. In this sense, the RD design *forces* the researcher to analyze it in a particular way, and there is little room for researcher discretion—at least from an identification standpoint. The design also predicts that the inclusion of  $W$ 's in the analysis should be irrelevant. Thus it naturally leads the researcher to examine the density of  $X$  or the distribution of  $W$ 's, conditional on  $X$ , for discontinuities as a test for validity.

The analogy of the truly randomized experiment is again helpful. Once the researcher is faced with what she thinks is a properly carried out randomized controlled trial, the analysis is quite straightforward. Even before running the experiment, most researchers agree it would be helpful to display the treatment-control contrasts in the  $W$ 's to test whether the randomization was carried out properly, then to show the simple mean comparisons, and finally to verify the inclusion of the  $W$ 's make little difference in the analysis, even if they might reduce sampling variability in the estimates.

#### 4. *Presentation, Estimation, and Inference*

In this section, we systematically discuss the nuts and bolts of implementing RD designs in practice. An important virtue of RD designs is that they provide a very

transparent way of graphically showing how the treatment effect is identified. We thus begin the section by discussing how to graph the data in an informative way. We then move to arguably the most important issue in implementing an RD design: the choice of the regression model. We address this by presenting the various possible specifications, discussing how to choose among them, and showing how to compute the standard errors.

Next, we discuss a number of other practical issues that often arise in RD designs. Examples of questions discussed include whether we should control for other covariates and what to do when the assignment variable is discrete. We discuss a number of tests to assess the validity of the RD designs, which examine whether covariates are “balanced” on the two sides of the threshold, and whether the density of the assignment variable is continuous at the threshold. Finally, we summarize our recommendations for implementing the RD design.

Throughout this section, we illustrate the various concepts using an empirical example from Lee (2008) who uses an RD design to estimate the causal effect of incumbency in U.S. House elections. We use a sample of 6,558 elections over the 1946–98 period (see Lee 2008 for more detail). The assignment variable in this setting is the fraction of votes awarded to Democrats in the previous election. When the fraction exceeds 50 percent, a Democrat is elected and the party becomes the incumbent party in the next election. Both the share of votes and the probability of winning the next election are considered as outcome variables.

##### 4.1 *Graphical Presentation*

A major advantage of the RD design over competing methods is its transparency, which can be illustrated using graphical methods. A standard way of graphing the data is to divide the assignment variable into a number of bins, making sure there are two separate

<sup>25</sup> Of course, sometimes it may seem at first that an RD design does not apply, but a closer inspection may reveal that it does. For example, see Per Pettersson (2000), which eventually became the RD analysis in Pettersson-Lidbom (2008b).

bins on each side of the cutoff point (to avoid having treated and untreated observations mixed together in the same bin). Then, the average value of the outcome variable can be computed for each bin and graphed against the mid-points of the bins.

More formally, for some bandwidth  $h$ , and for some number of bins  $K_0$  and  $K_1$  to the left and right of the cutoff value, respectively, the idea is to construct bins  $(b_k, b_{k+1}]$ , for  $k = 1, \dots, K = K_0 + K_1$ , where

$$b_k = c - (K_0 - k + 1)h.$$

The average value of the outcome variable in the bin is

$$\bar{Y}_k = \frac{1}{N_k} \sum_{i=1}^N Y_i 1\{b_k < X_i \leq b_{k+1}\}.$$

It is also useful to calculate the number of observations in each bin

$$N_k = \sum_{i=1}^N 1\{b_k < X_i \leq b_{k+1}\}$$

to detect a possible discontinuity in the assignment variable at the threshold, which would suggest manipulation.

There are several important advantages in graphing the data this way before starting to run regressions to estimate the treatment effect. First, the graph provides a simple way of visualizing what the functional form of the regression function looks like on either side of the cutoff point. Since the mean of  $Y$  in a bin is, for nonparametric kernel regression estimators, evaluated at the bin mid-point using a rectangular kernel, the set of bin means literally represent nonparametric estimates of the regression function. Seeing what the nonparametric regression looks like can then provide useful guidance in choosing the functional form of the regression models.

A second advantage is that comparing the mean outcomes just to the left and right of the cutoff point provides an indication of the magnitude of the jump in the regression function

at this point, i.e., of the treatment effect. Since an RD design is “as good as a randomized experiment” right around the cutoff point, the treatment effect could be computed by comparing the average outcomes in “small” bins just to the left and right of the cutoff point. If there is no visual evidence of a discontinuity in a simple graph, it is unlikely the formal regression methods discussed below will yield a significant treatment effect.

A third advantage is that the graph also shows whether there are unexpected comparable jumps at other points. If such evidence is clearly visible in the graph and cannot be explained on substantive grounds, this calls into question the interpretation of the jump at the cutoff point as the causal effect of the treatment. We discuss below several ways of testing explicitly for the existence of jumps at points other than the cutoff.

Note that the visual impact of the graph is typically enhanced by also plotting a relatively flexible regression model, such as a polynomial model, which is a simple way of smoothing the graph. The advantage of showing both the flexible regression line and the unrestricted bin means is that the regression line better illustrates the shape of the regression function and the size of the jump at the cutoff point, and laying this over the unrestricted means gives a sense of the underlying noise in the data.

Of course, if bins are too narrow the estimates will be highly imprecise. If they are too wide, the estimates may be biased as they fail to account for the slope in the regression line (negligible for very narrow bins). More importantly, wide bins make the comparisons on both sides of the cutoff less credible, as we are no longer comparing observations just to the left and right of the cutoff point.

This raises the question of how to choose the bandwidth (the width of the bin). In practice, this is typically done informally by trying to pick a bandwidth that makes the graphs look informative in the sense that bins

are wide enough to reduce the amount of noise, but narrow enough to compare observations “close enough” on both sides of the cutoff point. While it is certainly advisable to experiment with different bandwidths and see how the corresponding graphs look, it is also useful to have some formal guidance in the selection process.

One approach to bandwidth choice is based on the fact that, as discussed above, the mean outcomes by bin correspond to kernel regression estimates with a rectangular kernel. Since the standard kernel regression is a special case of a local linear regression where the slope term is equal to zero, the cross-validation procedure described in more detail in section 4.3.1 can also be used here by constraining the slope term to equal zero.<sup>26</sup> For reasons we discuss below, however, one should not solely rely on this approach to select the bandwidth since other reasonable subjective goals should be considered when choosing how to plot the data.

Furthermore, a range of bandwidths often yield similar values of the cross-validation function in practical applications (see below). A researcher may, therefore, want to use some discretion in choosing a bandwidth that provides a particularly compelling illustration of the RD design. An alternative approach is to choose a bandwidth based on a more heuristic visual inspection of the data, and then perform some tests to make sure this informal choice is not clearly rejected.

We suggest two such tests. Consider the case where one has decided to use  $K'$  bins based on a visual inspection of the data. The

first test is a standard  $F$ -test comparing the fit of a regression model with  $K'$  bin dummies to one where we further divide each bin into two equal sized smaller bins, i.e., increase the number of bins to  $2K'$  (reduce the bandwidth from  $h'$  to  $h'/2$ ). Since the model with  $K'$  bins is nested in the one with  $2K'$  bins, a standard  $F$ -test with  $K'$  degrees of freedom can be used. If the null hypothesis is not rejected, this provides some evidence that we are not oversmoothing the data by using only  $K'$  bins.

Another test is based on the idea that if the bins are “narrow enough,” then there should not be a systematic relationship between  $Y$  and  $X$ , that we capture using a simple regression of  $Y$  on  $X$ , within each bin. Otherwise, this suggests the bin is too wide and that the mean value of  $Y$  over the whole bin is not representative of the mean value of  $Y$  at the boundaries of the bin. In particular, when this happens in the two bins next to the cutoff point, a simple comparison of the two bin means yields a biased estimate of the treatment effect. A simple test for this consists of adding a set of interactions between the bin dummies and  $X$  to a base regression of  $Y$  on the set of bin dummies, and testing whether the interactions are jointly significant. The test statistic once again follows a  $F$  distribution with  $K'$  degrees of freedom.

Figures 6–11 show the graphs for the share of Democrat vote in the next election and the probability of Democrats winning the next election, respectively. Three sets of graphs with different bandwidths are reported using a bandwidth of 0.02 in figures 6 and 9, 0.01 in figures 7 and 10, and 0.005

<sup>26</sup> In section 4.3.1, we consider the cross-validation function  $CV_Y(h) = (1/N) \sum_{i=1}^N (Y_i - \hat{Y}(X_i))^2$  where  $\hat{Y}(X_i)$  is the predicted value of  $Y_i$  based on a regression using observations with a bin of width  $h$  on either the left (for observations on left of the cutoff) or the right (for observations on the right of the cutoff) of observation  $i$ , but not including observation  $i$  itself. In the context of the graph discussed here, the only modification to the cross-validation function is that the predicted value  $\hat{Y}(X_i)$  is based only

on a regression with a constant term, which means  $\hat{Y}(X_i)$  is the average value of  $Y$  among all observations in the bin (excluding observation  $i$ ). Note that this is slightly different from the standard cross-validation procedure in kernel regressions where the left-out observation is in the middle instead of the edge of the bin (see, for example, Richard Blundell and Alan Duncan 1998). Our suggested procedure is arguably better suited to the RD context since estimation of the treatment effect takes place at boundary points.

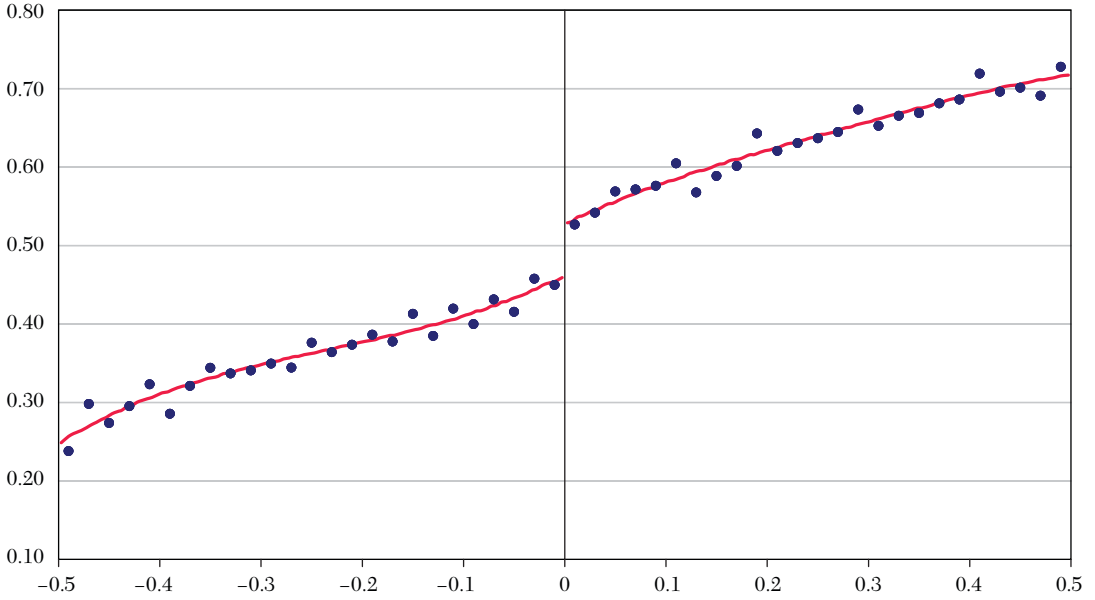


Figure 6. Share of Vote in Next Election, Bandwidth of 0.02 (50 bins)

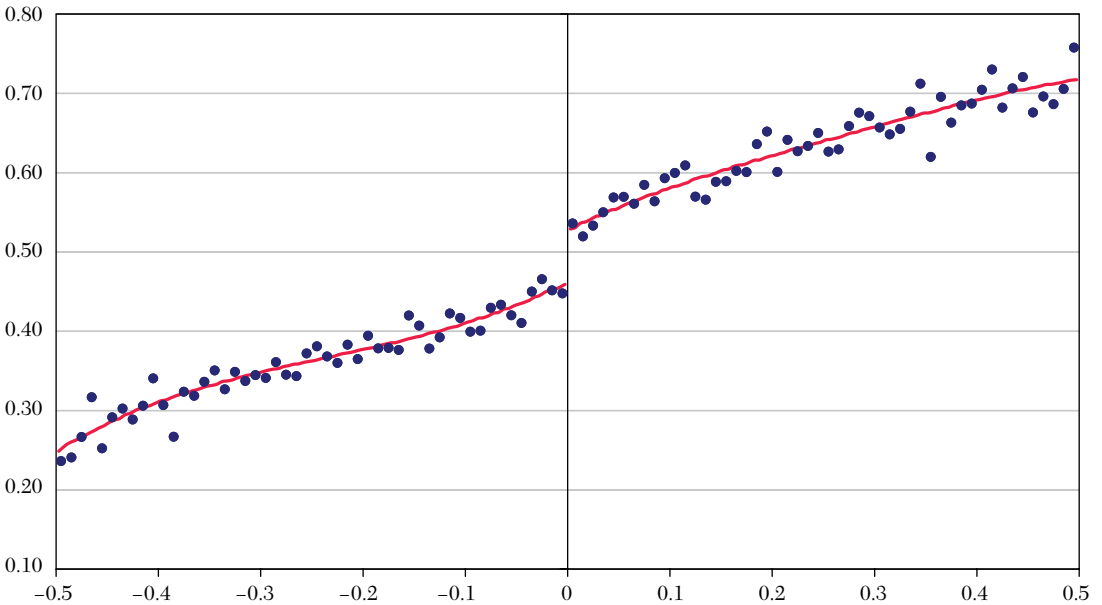


Figure 7. Share of Vote in Next Election, Bandwidth of 0.01 (100 bins)

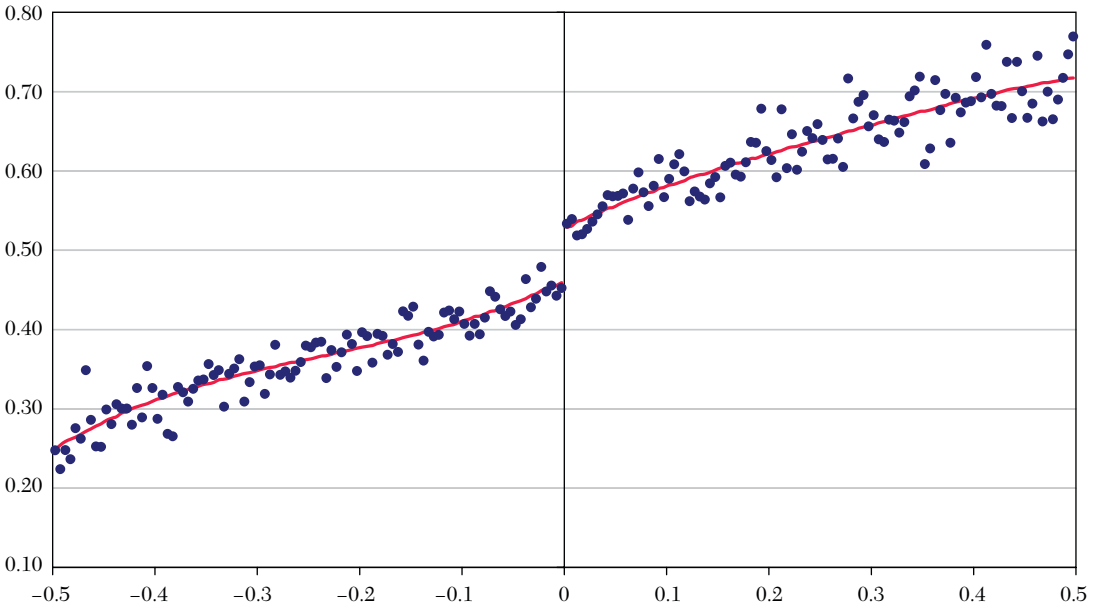


Figure 8. Share of Vote in Next Election, Bandwidth of 0.005 (200 bins)

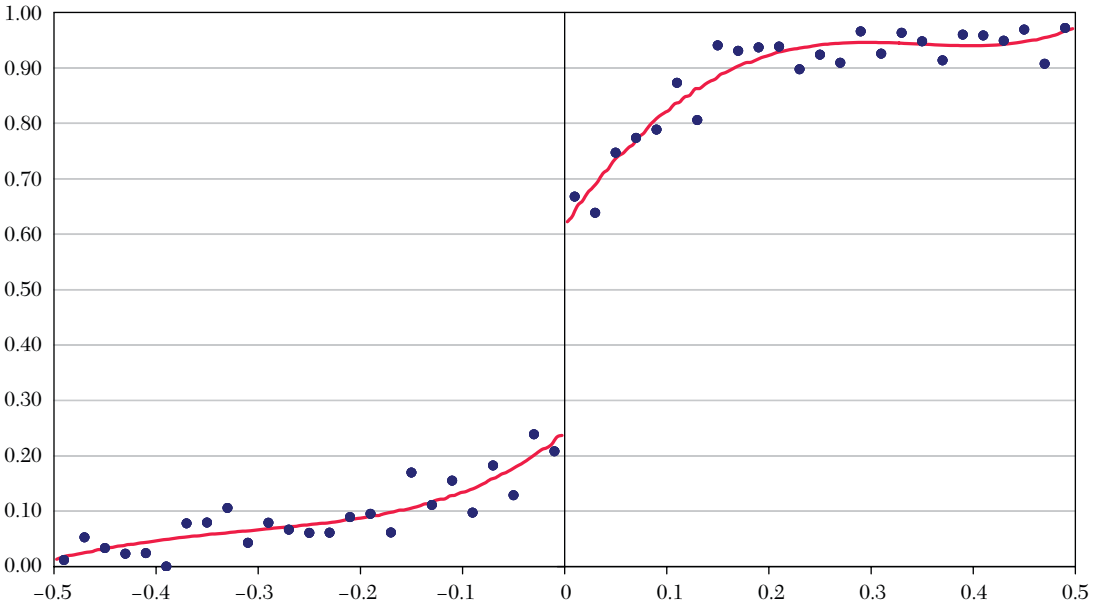


Figure 9. Winning the Next Election, Bandwidth of 0.02 (50 bins)

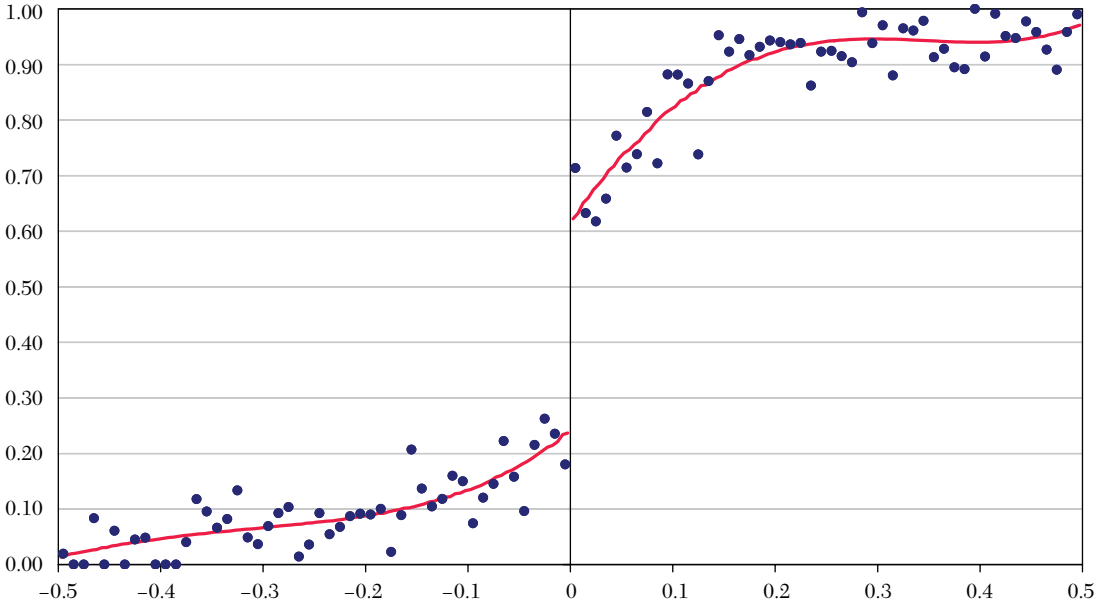


Figure 10. Winning the Next Election, Bandwidth of 0.01 (100 bins)

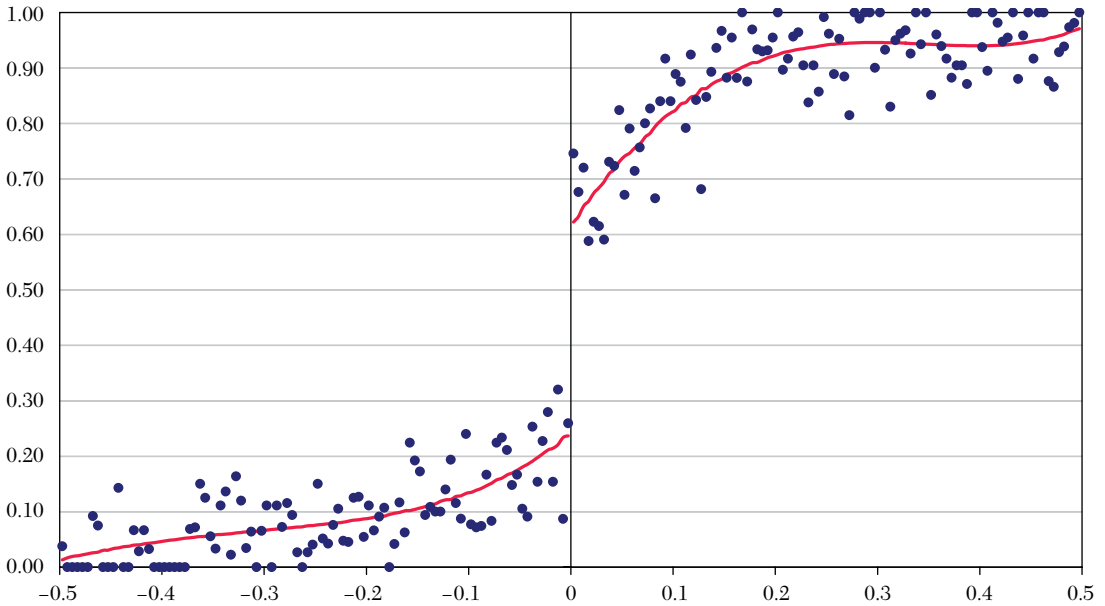


Figure 11. Winning the Next Election, Bandwidth of 0.005 (200 bins)

in figures 8 and 11. In all cases, we also show the fitted values from a quartic regression model estimated separately on each side of the cutoff point. Note that the assignment variable is normalized as the difference between the share of vote to Democrats and Republicans in the previous election. This means that a Democrat is the incumbent when the assignment variable exceeds zero. We also limit the range of the graphs to winning margins of 50 percent or less (in absolute terms) as data become relatively sparse for larger winning (or losing) margins.

All graphs show clear evidence of a discontinuity at the cutoff point. While the graphs are all quite informative, the ones with the smallest bandwidth (0.005, figures 8 and 11) are more noisy and likely provide too many data points (200) for optimal visual impact.

The results of the bandwidth selection procedures are presented in table 1. Panel A shows the cross-validation procedure always suggests using a bandwidth of 0.02 or more, which corresponds to similar or wider bins than those used in figures 6 and 9 (those with the largest bins). This is true irrespective of whether we pick a separate bandwidth on each side of the cutoff (first two rows of the panel), or pick the bandwidth that minimizes the cross-validation function for the entire date range on both the left and right sides of the cutoff. In the case where the outcome variable is winning the next election, the cross-validation procedure for the data to the right of the cutoff point and for the entire range suggests using a very wide bin (0.049) that would only yield about ten bins on each side of the cutoff.

As it turns out, the cross-validation function for the entire data range has two local minima at 0.021 and 0.049 that correspond to the optimal bandwidths on the left and right hand side of the cutoff. This is illustrated in figure 12, which plots the cross-validation function as a function of the bandwidth. By contrast, the cross-validation function is better behaved and shows a global minimum around 0.020

when the outcome variable is the vote share (figure 13). For both outcome variables, the value of the cross-validation function grows quickly for bandwidths smaller than 0.02, suggesting that the graphs with narrower bins (figures 7, 8, 10, and 11) are too noisy.

Panel B of table 1 shows the results of our two suggested specification tests. The tests based on doubling the number of bins and running regressions within each bin yield remarkably similar results. Generally speaking, the results indicate that only fairly wide bins are rejected. Looking at both outcome variables, the tests systematically reject models with bandwidths of 0.05 or more (twenty bins over the  $-0.5$  to  $0.5$  range). The models are never rejected for either outcome variable once we hit bandwidths of 0.02 (fifty bins) or less. In practice, the testing procedure rules out bins that are larger than those reported in figures 6–11.

At first glance, the results in the two panels of table 1 appear to be contradictory. The cross-validation procedure suggests bandwidths ranging from 0.02 to 0.05, while the bin and regression tests suggest that almost all bandwidths of less than 0.05 are acceptable. The reason for this discrepancy is that while the cross-validation procedure tries to balance precision and bias, the bin and regression tests only deal with the “bias” part of the equation by checking whether the value of  $Y$  is more or less constant within a given bin. Models with small bins easily pass this kind of test, although they may yield a very noisy graph. One alternative approach is to choose the largest possible bandwidth that passes the bin and the regression test, which turns out to be 0.033 in table 1, a bandwidth that is within the range of those suggested by the cross-validation procedure.

From a practical point of view, it seems to be the case that formal procedures, and in particular cross-validation, suggest bandwidths that are wider than those one would likely choose based on a simple visual examination

TABLE 1  
CHOICE OF BANDWIDTH IN GRAPH FOR VOTING EXAMPLE

**A. Optimal bandwidth selected by cross-validation**

Side of cutoff	Share of vote	Win next election
Left	0.021	0.049
Right	0.026	0.021
Both	0.021	0.049

**B. P-values of tests for the numbers of bins in RD graph**

No. of bins	Bandwidth	Share of vote		Win next election	
		Bin test	Regr. test	Bin test	Regr. test
10	0.100	0.000	0.000	0.001	0.000
20	0.050	0.000	0.000	0.026	0.049
30	0.033	0.163	0.390	0.670	0.129
40	0.025	0.157	0.296	0.024	0.020
50	0.020	0.957	0.721	0.477	0.552
60	0.017	0.159	0.367	0.247	0.131
70	0.014	0.596	0.130	0.630	0.743
80	0.013	0.526	0.740	0.516	0.222
90	0.011	0.815	0.503	0.806	0.803
100	0.010	0.787	0.976	0.752	0.883

*Notes:* Estimated over the range of the forcing variable (Democrat to Republican difference in the share of vote in the previous election) ranging between  $-0.5$  and  $0.5$ . The “bin test” is computed by comparing the fit of a model with the number of bins indicated in the table to an alternative where each bin is split in 2. The “regression test” is a joint test of significance of bin-specific regression estimates of the outcome variable on the share of vote in the previous election.

of the data. In particular, both figures 7 and 10 (bandwidth of 0.01) look visually acceptable but are clearly not recommended on the basis of the cross-validation procedure. This likely reflects the fact that one important goal of the graph is to show how the raw data look, and too much smoothing would defy the purpose of such a data illustration exercise. Furthermore, the regression estimates of the treatment effect accompanying the graphical results are a formal way of smoothing the data to get precise estimates. This suggests that there is probably little harm in under-

smoothing (relative to what formal bandwidth selection procedures would suggest) to better illustrate the variation in the raw data when graphically illustrating an RD design.

## 4.2 Regression Methods

### 4.2.1 Parametric or Nonparametric Regressions?

When we introduced the RD design in section 2, we followed Thistlethwaite and Campbell (1960) in assuming that the

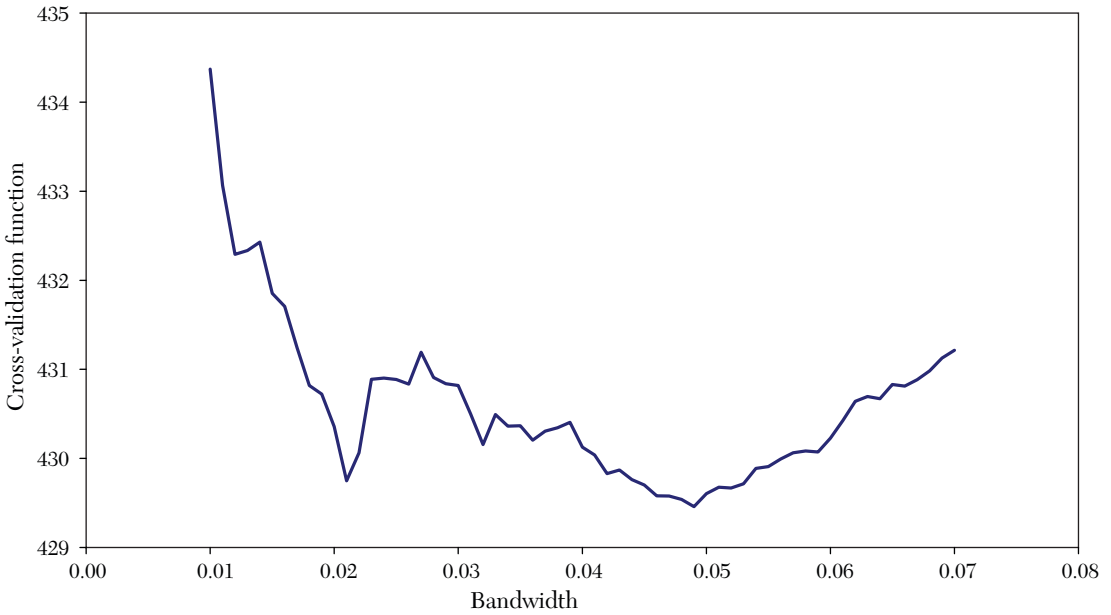


Figure 12. Cross-Validation Function for Choosing the Bandwidth in a RD Graph: Winning the Next Election

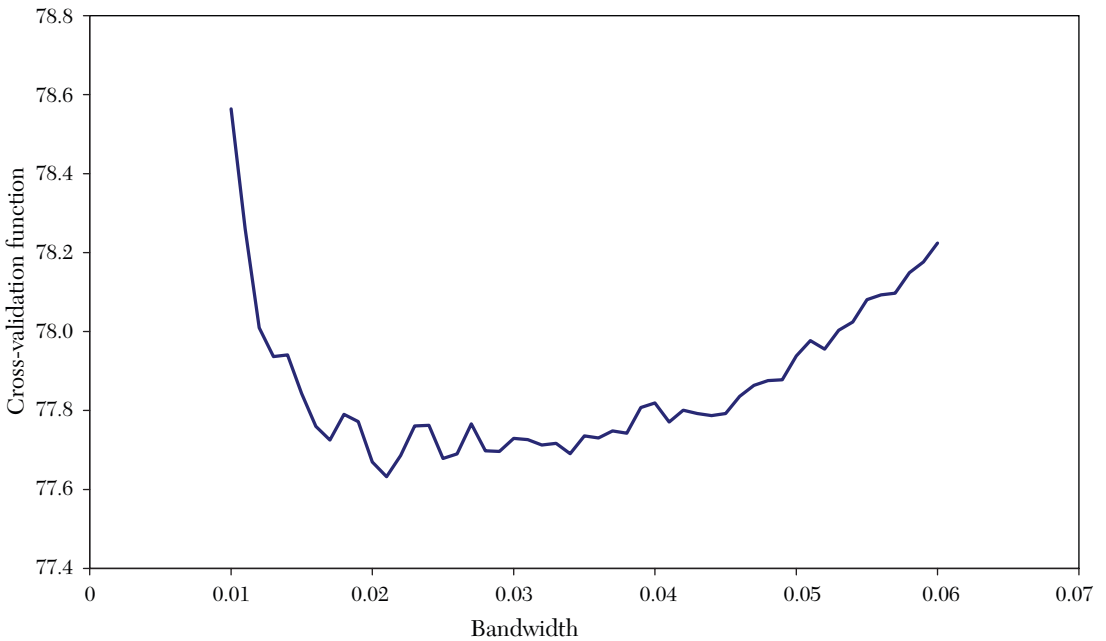


Figure 13. Cross-Validation Function for Choosing Bandwidth in a RD Graph: Share of Vote at Next Election

underlying regression model was linear in the assignment variable  $X$ :

$$Y = \alpha + D\tau + X\beta + \varepsilon.$$

In general, as in any other setting, there is no particular reason to believe that the true model is linear. The consequences of using an incorrect functional form are more serious in the case of RD designs however, since misspecification of the functional form typically generates a bias in the treatment effect,  $\tau$ .<sup>27</sup> This explains why, starting with Hahn, Todd, and van der Klaauw (2001), the estimation of RD designs have generally been viewed as a nonparametric estimation problem.

This being said, applied papers using the RD design often just report estimates from parametric models. Does this mean that these estimates are incorrect? Should all studies use nonparametric methods instead? As we pointed out in the introduction, we think that the distinction between parametric and nonparametric methods has sometimes been a source of confusion to practitioners. Before covering in detail the practical issues involved in the estimation of RD designs, we thus provide some background to help clarify the insights provided by nonparametric analysis, while also explaining why, in practice, RD designs can still be implemented using “parametric” methods.

Going beyond simple parametric linear regressions when the true functional form is unknown is a well-studied problem in econometrics and statistics. A number of nonparametric methods have been suggested to provide flexible estimates of the regression

function. As it turns out, however, the RD setting poses a particular problem because we need to estimate regressions at the cutoff point. This results in a “boundary problem” that causes some complications for nonparametric methods.

From an applied perspective, a simple way of relaxing the linearity assumption is to include polynomial functions of  $X$  in the regression model. This corresponds to the series estimation approach often used in nonparametric analysis. A possible disadvantage of the approach, however, is that it provides global estimates of the regression function over all values of  $X$ , while the RD design depends instead on local estimates of the regression function at the cutoff point. The fact that polynomial regression models use data far away from the cutoff point to predict the value of  $Y$  at the cutoff point is not intuitively appealing. That said, trying more flexible specification by adding polynomials in  $X$  as regressors is an important and useful way of assessing the robustness of the RD estimates of the treatment effect.

The other leading nonparametric approach is kernel regressions. Unlike series (polynomial) estimators, the kernel regression is fundamentally a local method well suited for estimating the regression function at a particular point. Unfortunately, this property does not help very much in the RD setting because the cutoff represents a boundary point where kernel regressions perform poorly.

These issues are illustrated in figure 2, which shows a situation where the relationship between  $Y$  and  $X$  (under treatment or control) is nonlinear. First, consider the point  $D$  located away from the cutoff point. The kernel estimate of the regression of  $Y$  on  $X$  at  $X = X_d$  is simply a local mean of  $Y$  for values of  $X$  close to  $X_d$ . The kernel function provides a way of computing this local average by putting more weight on observations with values of  $X$  close to  $X_d$  than on observations with values of  $X$  far away from  $X_d$ . Following Imbens

<sup>27</sup> By contrast, when one runs a linear regression in a model where the true functional form is nonlinear, the estimated model can still be interpreted as a linear predictor that minimizes specification errors. But since specification errors are only minimized globally, we can still have large specification errors at specific points including the cutoff point and, therefore, a large bias in RD estimates of the treatment effect.

and Lemieux (2008), we focus on the convenient case of the rectangular kernel. In this setting, computing kernel regressions simply amounts to computing the average value of  $Y$  in the bin illustrated in figure 2. The resulting local average is depicted as the horizontal line  $EF$ , which is very close to true value of  $Y$  evaluated at  $X = X_d$  on the regression line.

Applying this local averaging approach is problematic, however, for the RD design. Consider estimating the value of the regression function just on the right of the cutoff point. Clearly, only observations on the right of the cutoff point that receive the treatment should be used to compute mean outcomes on the right hand side. Similarly, only observations on the left of the cutoff point that do not receive the treatment should be used to compute mean outcomes on the left hand side. Otherwise, regression estimates would mix observations with and without the treatment, which would invalidate the RD approach.

In this setting, the best thing is to compute the average value of  $Y$  in the bin just to the right and just to the left of the cutoff point. These two bins are shown in figure 2. The RD estimate based on kernel regressions is then equal to  $B' - A'$ . In this example where the regression lines are upward sloping, it is clear, however, that the estimate  $B' - A'$  overstates the true treatment effect represented as the difference  $B - A$  at the cutoff point. In other words, there is a systematic bias in kernel regression estimates of the treatment effect. Hahn, Todd, and van der Klaauw (2001) provide a more formal derivation of the bias (see also Imbens and Lemieux 2008 for a simpler exposition when the kernel is rectangular). In practical terms, the problem is that in finite samples the bandwidth has to be large enough to encompass enough observations to get a reasonable amount of precision in the estimated average values of  $Y$ . Otherwise, attempts to reduce the bias by shrinking the bandwidth will result in

extremely noisy estimates of the treatment effect.<sup>28</sup>

As a solution to this problem, Hahn, Todd, and van der Klaauw (2001) suggests running local linear regressions to reduce the importance of the bias. In our setup with a rectangular kernel, this suggestion simply amounts to running standard linear regressions within the bins on both sides of the cutoff point to better predict the value of the regression function right at the cutoff point. In this example, the regression lines within the bins around the cutoff point are close to linear. It follows that the predicted values of the local linear regressions at the cutoff point are very close to the true values of  $A$  and  $B$ . Intuitively, this means that running local linear regressions instead of just computing averages within the bins reduces the bias by an order of magnitude. Indeed, Hahn, Todd, and van der Klaauw (2001) show that the remaining bias is of an order of magnitude lower, and is comparable to the usual bias in kernel estimation at interior points like  $D$  (the small difference between the horizontal line  $EF$  and the true value of the regression line evaluated at  $D$ ).

In the literature on nonparametric estimation at boundary points, local linear regressions have been introduced as a means of reducing the bias in standard kernel regression methods.<sup>29</sup> One of the several contributions of Hahn, Todd, and van der Klaauw (2001) is to show how the same bias-reducing

<sup>28</sup> The trade-off between bias and precision is a fundamental feature of kernel regressions. A larger bandwidth yields more precise, but potentially biased, estimates of the regression. In an interior point like  $D$ , however, we see that the bias is of an order of magnitude lower than at the cutoff (boundary) point. In more technical terms, it can be shown (see Hahn, Todd, and van der Klaauw 2001 or Imbens and Lemieux 2008) that the usual bias is of order  $h^2$  at interior points, but of order  $h$  at boundary points, where  $h$  is the bandwidth. In other words, the bias dies off much more quickly when  $h$  goes to zero when we are at interior, as opposed to boundary, points.

<sup>29</sup> See Jianqing Fan and Irene Gijbels (1996).

procedure should also be applied to the RD design. We have shown here that, in practice, this simply amounts to applying the original insight of Thistlethwaite and Campbell (1960) to a narrower window of observations around the cutoff point. When one is concerned that the regression function is not linear over the whole range of  $X$ , a highly sensible procedure is, thus, to restrict the estimation range to values closer to the cutoff point where the linear approximation of the regression line is less likely to result in large biases in the RD estimates. In practice, many applied papers present RD estimates with varying window widths to illustrate the robustness (or lack thereof) of the RD estimates to specification issues. It is comforting to know that this common empirical practice can be justified on more formal econometric grounds like those presented by Hahn, Todd, and van der Klaauw (2001). The main conclusion we draw from this discussion of nonparametric methods is that it is essential to explore how RD estimates are robust to the inclusion of higher order polynomial terms (the series or polynomial estimation approach) and to changes in the window width around the cutoff point (the local linear regression approach).

### 4.3 Estimating the Regression

A simple way of implementing RD designs in practice is to estimate two separate regressions on each side of the cutoff point. In terms of computations, it is convenient to subtract the cutoff value from the covariate, i.e., transform  $X$  to  $X - c$ , so the intercepts of the two regressions yield the value of the regression functions at the cutoff point.

The regression model on the left hand side of the cutoff point ( $X < c$ ) is

$$Y = \alpha_l + f_l(X - c) + \varepsilon,$$

while the regression model on the right hand side of the cutoff point ( $X \geq c$ ) is

$$Y = \alpha_r + f_r(X - c) + \varepsilon,$$

where  $f_l(\cdot)$  and  $f_r(\cdot)$  are functional forms that we discuss later. The treatment effect can then be computed as the difference between the two regressions intercepts,  $\alpha_r$  and  $\alpha_l$ , on the two sides of the cutoff point. A more direct way of estimating the treatment effect is to run a pooled regression on both sides of the cutoff point:

$$Y = \alpha_l + \tau D + f(X - c) + \varepsilon,$$

where  $\tau = \alpha_r - \alpha_l$  and  $f(X - c) = f_l(X - c) + D [f_r(X - c) - f_l(X - c)]$ . One advantage of the pooled approach is that it directly yields estimates and standard errors of the treatment effect  $\tau$ . Note, however, that it is recommended to let the regression function differ on both sides of the cutoff point by including interaction terms between  $D$  and  $X$ . For example, in the linear case where  $f_l(X - c) = \beta_l(X - c)$  and  $f_r(X - c) = \beta_r(X - c)$ , the pooled regression would be

$$Y = \alpha_l + \tau D + \beta_l(X - c) + (\beta_r - \beta_l) D (X - c) + \varepsilon.$$

The problem with constraining the slope of the regression lines to be the same on both sides of the cutoff ( $\beta_r = \beta_l$ ) is best illustrated by going back to the separate regressions above. If we were to constrain the slope to be identical on both sides of the cutoff, this would amount to using data on the right hand side of the cutoff to estimate  $\alpha_l$ , and vice versa. Remember from section 2 that in an RD design, the treatment effect is obtained by comparing conditional expectations of  $Y$  when approaching from the left ( $\alpha_l = \lim_{x \uparrow c} E[Y_i | X_i = x]$ ) and from the right ( $\alpha_r = \lim_{x \downarrow c} E[Y_i | X_i = x]$ ) of the cutoff. Constraining the slope to be the same would thus be inconsistent with the spirit of

the RD design, as data from the right of the cutoff would be used to estimate  $\alpha_l$ , which is defined as a limit when approaching from the left of the cutoff, and vice versa.

In practice, however, estimates where the regression slope or, more generally, the regression function  $f(X - c)$  are constrained to be the same on both sides of the cutoff point are often reported. One possible justification for doing so is that if the functional form is indeed the same on both sides of the cutoff, then more efficient estimates of the treatment effect  $\tau$  are obtained by imposing that constraint. Such a constrained specification should only be viewed, however, as an additional estimate to be reported for the sake of completeness. It should not form the core basis of the empirical approach.

#### 4.3.1 Local Linear Regressions and Bandwidth Choice

As discussed above, local linear regressions provide a nonparametric way of consistently estimating the treatment effect in an RD design (Hahn, Todd, and van der Klaauw (2001), Jack Porter (2003)). Following Imbens and Lemieux (2008), we focus on the case of a rectangular kernel, which amounts to estimating a standard regression over a window of width  $h$  on both sides of the cutoff point. While other kernels (triangular, Epanechnikov, etc.) could also be used, the choice of kernel typically has little impact in practice. As a result, the convenience of working with a rectangular kernel compensates for efficiency gains that could be achieved using more sophisticated kernels.<sup>30</sup>

The regression model on the left hand side of the cutoff point is

$$Y = \alpha_l + \beta_l(X - c) + \varepsilon,$$

$$\text{where } c - h \leq X < c,$$

while the regression model on the right hand side of the cutoff point is

$$Y = \alpha_r + \beta_r(X - c) + \varepsilon,$$

$$\text{where } c \leq X \leq c + h.$$

As before, it is also convenient to estimate the pooled regression

$$Y = \alpha_l + \tau D + \beta_l(X - c) + (\beta_r - \beta_l) D(X - c) + \varepsilon,$$

$$\text{where } c - h \leq X \leq c + h,$$

since the standard error of the estimated treatment effect can be directly obtained from the regression.

While it is straightforward to estimate the linear regressions within a given window of width  $h$  around the cutoff point, a more difficult question is how to choose this bandwidth. In general, choosing a bandwidth in nonparametric estimation involves finding an optimal balance between precision and bias. On the one hand, using a larger bandwidth yields more precise estimates as more observations are available to estimate the regression. On the other hand, the linear specification is less likely to be accurate

<sup>30</sup> It has been shown in the statistics literature (Fan and Gijbels 1996) that a triangular kernel is optimal for estimating local linear regressions at the boundary. As it turns out, the only difference between regressions using a rectangular or a triangular kernel is that the latter puts more weight (in a linear way) on observations closer to the cutoff point. It thus involves estimating a weighted, as opposed to an unweighted, regression within a bin of width  $h$ . An

arguably more transparent way of putting more weight on observations close to the cutoff is simply to reestimate a model with a rectangular kernel using a smaller bandwidth. In practice, it is therefore simpler and more transparent to just estimate standard linear regressions (rectangular kernel) with a variety of bandwidths, instead of trying out different kernels corresponding to particular weighted regressions that are more difficult to interpret.

when a larger bandwidth is used, which can bias the estimate of the treatment effect. If the underlying conditional expectation is not linear, the linear specification will provide a close approximation over a limited range of values of  $X$  (small bandwidth), but an increasingly bad approximation over a larger range of values of  $X$  (larger bandwidth).

As the number of observations available increases, it becomes possible to use an increasingly small bandwidth since linear regressions can be estimated relatively precisely over even a small range of values of  $X$ . As it turns out, Hahn, Todd, and van der Klaauw (2001) show the optimal bandwidth is proportional to  $N^{-1/5}$ , which corresponds to a fairly slow rate of convergence to zero. For example, this suggests that the bandwidth should only be cut in half when the sample size increases by a factor of 32 ( $2^5$ ). For technical reasons, however, it would be preferable to undersmooth by shrinking the bandwidth at a faster rate requiring that  $h \propto N^{-\delta}$  with  $1/5 < \delta < 2/5$ , in order to eliminate an asymptotic bias that would remain when  $\delta = 1/5$ . In the presence of this bias, the usual formula for the variance of a standard least square estimator would be invalid.<sup>31</sup>

In practice however, knowing at what rate the bandwidth should shrink in the limit does not really help since only one actual sample with a given number of observations is

available. The importance of undersmoothing only has to do with a thought experiment of how much the bandwidth should shrink if the sample size were larger so that one obtains asymptotically correct standard errors, and does not help one choose a particular bandwidth in a particular sample.<sup>32</sup>

In the econometrics and statistics literature, two procedures are generally considered for choosing bandwidths. The first procedure consists of characterizing the optimal bandwidth in terms of the unknown joint distribution of all variables. The relevant components of this distribution can then be estimated and plugged into the optimal bandwidth function.<sup>33</sup> In the context of local linear regressions, Fan and Gijbels (1996) show this involves estimating a number of parameters including the curvature of the regression function. In practice, this can be done in two steps. In step one, a rule-of-thumb (ROT) bandwidth is estimated over the whole relevant data range. In step two, the ROT bandwidth is used to estimate the optimal bandwidth right at the cutoff point. For the rectangular kernel, the ROT bandwidth is given by:

$$h_{ROT} = 2.702 \left[ \frac{\hat{\sigma}^2 R}{\sum_{i=1}^N \left\{ \tilde{m}''(x_i) \right\}^2} \right]^{1/5},$$

<sup>31</sup> See Hahn, Todd, and van der Klaauw (2001) and Imbens and Lemieux (2008) for more details.

<sup>32</sup> The main purpose of asymptotic theory is to use the large sample properties of estimators to approximate the distribution of an estimator in the real sample being considered. The issue is a little more delicate in a nonparametric setting where one also has to think about how fast the bandwidth should shrink when the sample size approaches infinity. The point about undersmoothing is simply that one unpleasant property of the optimal bandwidth is that it does not yield the convenient least squares variance formula. But this can be fixed by shrinking the bandwidth a little faster as the sample size goes to infinity. Strictly speaking, this is only a technical issue with how to perform the thought experiment (what happens when the sample size goes to infinity?) required for using asymptotics to

approximate the variance of the RD estimator in the actual sample. This does not say anything about what bandwidth should be chosen in the actual sample available for implementing the RD design.

<sup>33</sup> A well known example of this procedure is the "rule-of-thumb" bandwidth selection formula in kernel density estimation where an estimate of the dispersion in the variable (standard deviation or the interquartile range),  $\hat{\sigma}$ , is plugged into the formula  $0.9 \cdot \hat{\sigma} \cdot N^{-1/5}$ . Bernard W. Silverman (1986) shows that this formula is the closed form solution for the optimal bandwidth choice problem when both the actual density and the kernel are Gaussian. See also Imbens and Karthik Kalyanaraman (2009), who derive an optimal bandwidth for this RD setting, and propose a data-dependent method for choosing the bandwidth.

where  $\tilde{m}''(\cdot)$  is the second derivative (curvature) of an estimated regression of  $Y$  on  $X$ ,  $\tilde{\sigma}$  is the estimated standard error of the regression,  $R$  is the range of the assignment variable over which the regression is estimated, and the constant 2.702 is a number specific to the rectangular kernel. A similar formula can be used for the optimal bandwidth, except both the regression standard error and the average curvature of the regression function are estimated locally around the cutoff point. For the sake of simplicity, we only compute the ROT bandwidth in our empirical example. Following the common practice in studies using these bandwidth selection methods, we also use a quartic specification for the regression function.<sup>34</sup>

The second approach is based on a cross-validation procedure. In the case considered here, Jens Ludwig and Douglas Miller (2007) and Imbens and Lemieux (2008) have proposed a “leave one out” procedure aimed specifically at estimating the regression function at the boundary. The basic idea behind this procedure is the following. Consider an observation  $i$ . To see how well a linear regression with a bandwidth  $h$  fits the data, we run a regression with observation  $i$  left out and use the estimates to predict the value of  $Y$  at  $X = X_i$ . In order to mimic the fact that RD estimates are based on regression estimates at the boundary, the regression is estimated using only observations with values of  $X$  on the left of  $X_i$  ( $X_i - h \leq X < X_i$ ) for observations on the left of the cutoff point ( $X_i < c$ ). For observations on the right of the cutoff point ( $X_i \geq c$ ), the regression is estimated

using only observations with values of  $X$  on the right of  $X_i$  ( $X_i < X \leq X_i + h$ ).

Repeating the exercise for each and every observation, we get a whole set of predicted values of  $Y$  that can be compared to the actual values of  $Y$ . The optimal bandwidth can be picked by choosing the value of  $h$  that minimizes the mean square of the difference between the predicted and actual value of  $Y$ .

More formally, let  $\hat{Y}(X_i)$  represent the predicted value of  $Y$  obtained using the regressions described above. The cross-validation criterion is defined as

$$(9) \quad \text{CV}_Y(h) = \frac{1}{N} \sum_{i=1}^N (Y_i - \hat{Y}(X_i))^2$$

with the corresponding cross-validation choice for the bandwidth

$$h_{\text{CV}}^{\text{opt}} = \arg \min_h \text{CV}_Y(h).$$

Imbens and Lemieux (2008) discuss this procedure in more detail and point out that since we are primarily interested in what happens around the cutoff, it may be advisable to only compute  $\text{CV}_Y(h)$  for a subset of observations with values of  $X$  close enough to the cutoff point. For instance, only observations with values of  $X$  between the median value of  $X$  to the left and right of the cutoff could be used to perform the cross-validation.

The second rows of tables 2 and 3 show the local linear regression estimates of the treatment effect for the two outcome variables (share of vote and winning the next election). We show the estimates for a wide range of bandwidths going from the entire data range (bandwidth of 1 on each side of the cutoff) to a very small bandwidth of 0.01 (winning margins of one percent or less). As expected, the precision of the estimates declines quickly as we approach smaller and smaller bandwidths. Notice also that estimates based

<sup>34</sup> See McCrary and Heather Royer (2003) for an example where the bandwidth is selected using the ROT procedure (with a triangular kernel), and Stephen L. Desjardins and Brian P. McCall (2008) for an example where the second step optimal bandwidth is computed (for the Epanechnikov kernel). Both papers use a quartic regression function  $m(x) = \beta_0 + \beta_1 x + \dots + \beta_4 x^4$ , which means that  $m''(x) = 2\beta_2 + 6\beta_3 x + 12\beta_4 x^2$ . Note that the quartic regressions are estimated separately on both sides of the cutoff.

TABLE 2  
RD ESTIMATES OF THE EFFECT OF WINNING THE PREVIOUS ELECTION ON THE  
SHARE OF VOTES IN THE NEXT ELECTION

Bandwidth:	1.00	0.50	0.25	0.15	0.10	0.05	0.04	0.03	0.02	0.01
Polynomial of order:										
Zero	0.347 (0.003) [0.000]	0.257 (0.004) [0.000]	0.179 (0.004) [0.000]	0.143 (0.005) [0.000]	0.125 (0.006) [0.003]	0.096 (0.009) [0.047]	0.080 (0.011) [0.778]	0.073 (0.012) [0.821]	0.077 (0.014) [0.687]	0.088 (0.015)
One	0.118 (0.006) [0.000]	0.090 (0.007) [0.332]	0.082 (0.008) [0.423]	0.077 (0.011) [0.216]	0.061 (0.013) [0.543]	0.049 (0.019) [0.168]	0.067 (0.022) [0.436]	0.079 (0.026) [0.254]	0.098 (0.029) [0.935]	0.096 (0.028)
Two	0.052 (0.008) [0.000]	0.082 (0.010) [0.335]	0.069 (0.013) [0.371]	0.050 (0.016) [0.385]	0.057 (0.020) [0.458]	0.100 (0.029) [0.650]	0.101 (0.033) [0.682]	0.119 (0.038) [0.272]	0.088 (0.044) [0.943]	0.098 (0.045)
Three	0.111 (0.011) [0.001]	0.068 (0.013) [0.335]	0.057 (0.017) [0.524]	0.061 (0.022) [0.421]	0.072 (0.028) [0.354]	0.112 (0.037) [0.603]	0.119 (0.043) [0.453]	0.092 (0.052) [0.324]	0.108 (0.062) [0.915]	0.082 (0.063)
Four	0.077 (0.013) [0.014]	0.066 (0.017) [0.325]	0.048 (0.022) [0.385]	0.074 (0.027) [0.425]	0.103 (0.033) [0.327]	0.106 (0.048) [0.560]	0.088 (0.056) [0.497]	0.049 (0.067) [0.044]	0.055 (0.079) [0.947]	0.077 (0.063)
Optimal order of the polynomial	6	3	1	2	1	2	0	0	0	0
Observations	6,558	4,900	2,763	1,765	1,209	610	483	355	231	106

*Notes:* Standard errors in parentheses. *P*-values from the goodness-of-fit test in square brackets. The goodness-of-fit test is obtained by jointly testing the significance of a set of bin dummies included as additional regressors in the model. The bin width used to construct the bin dummies is 0.01. The optimal order of the polynomial is chosen using Akaike's criterion (penalized cross-validation).

on very wide bandwidths (0.5 or 1) are systematically larger than those for the smaller bandwidths (in the 0.05 to 0.25 range) that are still large enough for the estimates to be reasonably precise. A closer examination of figures 6–11 also suggests that the estimates for very wide bandwidths are larger than what the graphical evidence would suggest.<sup>35</sup> This is consistent with a substantial bias for

these estimates linked to the fact that the linear approximation does not hold over a wide data range. This is particularly clear in the case of winning the next election where figures 9–11 show some clear curvature in the regression function.

Table 4 shows the optimal bandwidth obtained using the ROT and cross-validation procedure. Consistent with the above

<sup>35</sup> In the case of the vote share, the quartic regression shown in figures 6–8 implies a treatment effect of 0.066, which is substantially smaller than the local linear regression estimates with a bandwidth of 0.5 (0.090) or 1 (0.118).

Similarly, the quartic regression shown in figures 9–11 for winning the next election implies a treatment effect of 0.375, which is again smaller than the local linear regression estimates with a bandwidth of 0.5 (0.566) or 1 (0.689).

TABLE 3  
RD ESTIMATES OF THE EFFECT OF WINNING THE PREVIOUS ELECTION ON  
PROBABILITY OF WINNING THE NEXT ELECTION

Bandwidth:	1.00	0.50	0.25	0.15	0.10	0.05	0.04	0.03	0.02	0.01
Polynomial of order:										
Zero	0.814 (0.007) [0.000]	0.777 (0.009) [0.000]	0.687 (0.013) [0.000]	0.604 (0.018) [0.000]	0.550 (0.023) [0.011]	0.479 (0.035) [0.201]	0.428 (0.040) [0.852]	0.423 (0.047) [0.640]	0.459 (0.058) [0.479]	0.533 (0.082)
One	0.689 (0.011) [0.000]	0.566 (0.016) [0.000]	0.457 (0.026) [0.126]	0.409 (0.036) [0.269]	0.378 (0.047) [0.336]	0.378 (0.073) [0.155]	0.472 (0.083) [0.400]	0.524 (0.099) [0.243]	0.567 (0.116) [0.125]	0.453 (0.157)
Two	0.526 (0.016) [0.075]	0.440 (0.023) [0.145]	0.375 (0.039) [0.253]	0.391 (0.055) [0.192]	0.450 (0.072) [0.245]	0.607 (0.110) [0.485]	0.586 (0.124) [0.367]	0.589 (0.144) [0.191]	0.440 (0.177) [0.134]	0.225 (0.246)
Three	0.452 (0.021) [0.818]	0.370 (0.031) [0.277]	0.408 (0.052) [0.295]	0.435 (0.075) [0.115]	0.472 (0.096) [0.138]	0.566 (0.143) [0.536]	0.547 (0.166) [0.401]	0.412 (0.198) [0.234]	0.266 (0.247) [0.304]	0.172 (0.349)
Four	0.385 (0.026) [0.965]	0.375 (0.039) [0.200]	0.424 (0.066) [0.200]	0.529 (0.093) [0.173]	0.604 (0.119) [0.292]	0.453 (0.183) [0.593]	0.331 (0.214) [0.507]	0.134 (0.254) [0.150]	0.050 (0.316) [0.244]	0.168 (0.351)
Optimal order of the polynomial	4	3	2	1	1	2	0	0	0	1
Observations	6,558	4,900	2,763	1,765	1,209	610	483	355	231	106

*Notes:* Standard errors in parentheses. *P*-values from the goodness-of-fit test in square brackets. The goodness-of-fit test is obtained by jointly testing the significance of a set of bin dummies included as additional regressors in the model. The bin width used to construct the bin dummies is 0.01. The optimal order of the polynomial is chosen using Akaike's criterion (penalized cross-validation).

discussion, the suggested bandwidth ranges from 0.14 to 0.28, which is large enough to get precise estimates, but narrow enough to minimize the bias. Two interesting patterns can be observed in table 4. First, the bandwidth chosen by cross-validation tends to be a bit larger than the one based on the rule-of-thumb. Second, the bandwidth is generally smaller for winning the next election (second column) than for the vote share (first column). This is particularly clear when the optimal bandwidth is constrained to be the same on both sides of the cutoff point. This is consistent with the graphical evidence

showing more curvature for winning the next election than the vote share, which calls for a smaller bandwidth to reduce the estimation bias linked to the linear approximation.

Figures 14 and 15 plot the value of the cross-validation function over a wide range of bandwidths. In the case of the vote share where the linearity assumption appears more accurate (figures 6–8), the cross-validation function is fairly flat over a sizable range of values for the bandwidth (from about 0.16 to 0.29). This range includes the optimal bandwidth suggested by cross-validation (0.282) at the upper end, and the ROT

TABLE 4  
OPTIMAL BANDWIDTH FOR LOCAL LINEAR REGRESSIONS,  
VOTING EXAMPLE

	Share of vote	Win next election
<b>A. Rule-of-thumb bandwidth</b>		
Left	0.162	0.164
Right	0.208	0.130
Both	0.180	0.141
<b>B. Optimal bandwidth selected by cross-validation</b>		
Left	0.192	0.247
Right	0.282	0.141
Both	0.282	0.172

*Notes:* Estimated over the range of the forcing variable (Democrat to Republican difference in the share of vote in the previous election) ranging between  $-0.5$  and  $0.5$ . See the text for a description of the rule-of-thumb and cross-validation procedures for choosing the optimal bandwidth.

bandwidth (0.180) at the lower end. In the case of winning the next election (figure 15), the cross-validation procedure yields a sharper suggestion of optimal bandwidth around 0.15, which is quite close to both the optimal cross-validation bandwidth (0.172) and the ROT bandwidth (0.141).

The main difference between the two outcome variables is that larger bandwidths start getting penalized more quickly in the case of winning the election (figure 15) than in the case of the vote share (figure 14). This is consistent with the graphical evidence in figures 6–11. Since the regression function looks fairly linear for the vote share, using larger bandwidths does not get penalized as much since they improve efficiency without generating much of a bias. But in the case of winning the election where the regression function exhibits quite a bit of curvature, larger bandwidths are quickly penalized for introducing an estimation bias. Since there is a real tradeoff between precision and bias, the cross-validation procedure is quite informative. By contrast, there is not much of a tradeoff when the regression function is more or less linear, which explains why the

optimal bandwidth is larger in the case of the vote share.

This example also illustrates the importance of first graphing the data before running regressions and trying to choose the optimal bandwidth. When the graph shows a more or less linear relationship, it is natural to expect different bandwidths to yield similar results and the bandwidth selection procedure not to be terribly informative. But when the graph shows substantial curvature, it is natural to expect the results to be more sensitive to the choice of bandwidth and that bandwidth selection procedures will play a more important role in selecting an appropriate empirical specification.

#### 4.3.2 *Order of Polynomial in Local Polynomial Modeling*

In the case of polynomial regressions, the equivalent to bandwidth choice is the choice of the order of the polynomial regressions. As in the case of local linear regressions, it is advisable to try and report a number of specifications to see to what extent the results are sensitive to the order of the polynomial. For the same reason

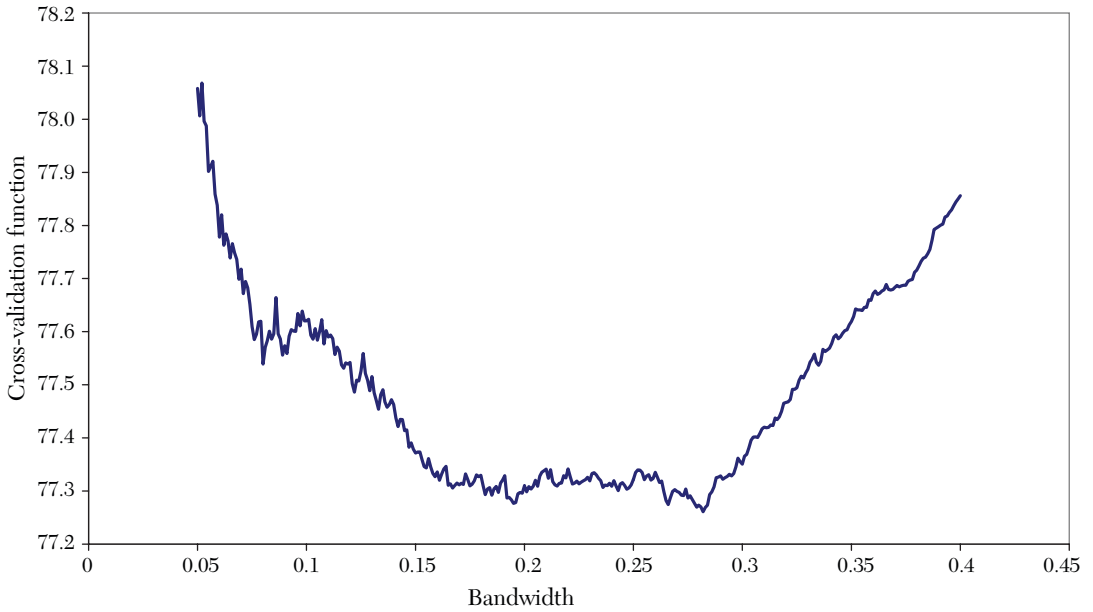


Figure 14. Cross-Validation Function for Local Linear Regression: Share of Vote at Next Election

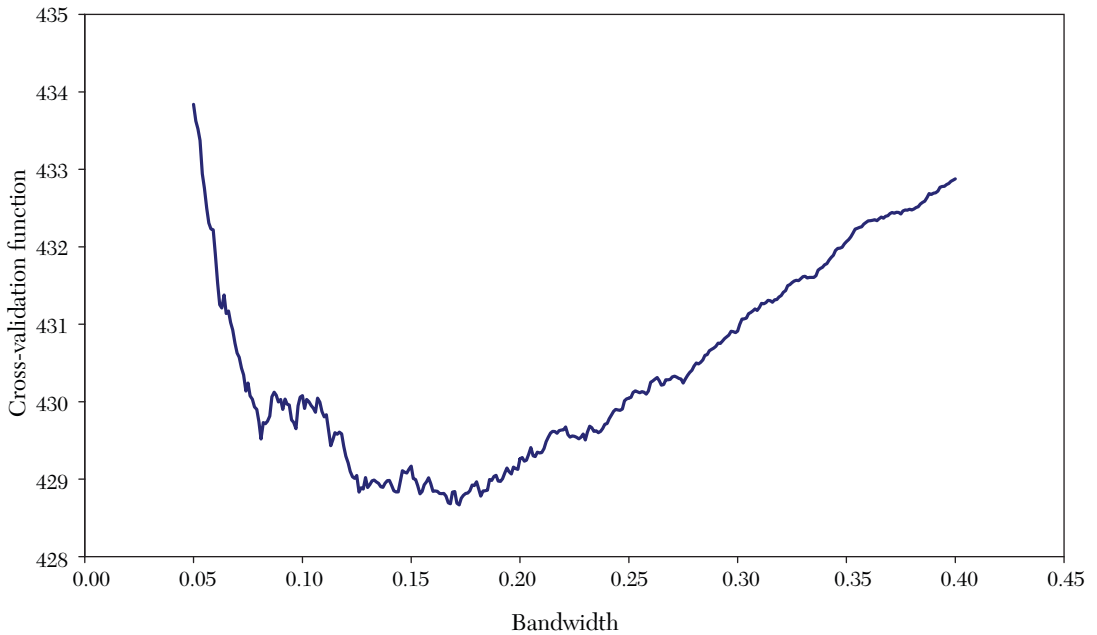


Figure 15. Cross-Validation Function for Local Linear Regression: Winning the Next Election

mentioned earlier, it is also preferable to estimate separate regressions on the two sides of the cutoff point.

The simplest way of implementing polynomial regressions and computing standard errors is to run a pooled regression. For example, in the case of a third order polynomial regression, we would have

$$\begin{aligned} Y = & \alpha_l + \tau D + \beta_{l1}(X - c) \\ & + \beta_{l2}(X - c)^2 + \beta_{l3}(X - c)^3 \\ & + (\beta_{r1} - \beta_{l1})D(X - c) \\ & + (\beta_{r2} - \beta_{l2})D(X - c)^2 \\ & + (\beta_{r3} - \beta_{l3})D(X - c)^3 + \varepsilon. \end{aligned}$$

While it is important to report a number of specifications to illustrate the robustness of the results, it is often useful to have some more formal guidance on the choice of the order of the polynomial. Starting with van der Klaauw (2002), one approach has been to use a generalized cross-validation procedure suggested in the literature on nonparametric series estimators.<sup>36</sup> One special case of generalized cross-validation (used by Dan A. Black, Jose Galdo, and Smith (2007a), for example), which we also use in our empirical example, is the well known Akaike information criterion (AIC) of model selection. In a regression context, the AIC is given by

$$AIC = N \ln(\hat{\sigma}^2) + 2p,$$

where  $\hat{\sigma}^2$  is the mean squared error of the regression, and  $p$  is the number of parameters in the regression model (order of the polynomial plus one for the intercept).

One drawback of this approach is that it does not provide a very good sense of how

a particular parametric model (say a cubic model) compares relative to a more general nonparametric alternative. In the context of the RD design, a natural nonparametric alternative is the set of unrestricted means of the outcome variable by bin used to graphically depict the data in section 4.1. Since one virtue of polynomial regressions is that they provide a smoothed version of the graph, it is natural to ask how well the polynomial model fits the unrestricted graph. A simple way of implementing the test is to add the set of bin dummies to the polynomial regression and jointly test the significance of the bin dummies. For example, in a first order polynomial model (linear regression), the test can be computed by including  $K - 2$  bin dummies  $B_k$ , for  $k = 2$  to  $K - 1$ , in the model

$$\begin{aligned} Y = & \alpha_l + \tau D + \beta_{l1}(X - c) \\ & + (\beta_{r1} - \beta_{l1})D(X - c) \\ & + \sum_{k=2}^{K-1} \phi_k B_k + \varepsilon \end{aligned}$$

and testing the null hypothesis that  $\phi_2 = \phi_3 = \dots = \phi_{K-1} = 0$ . Note that two of the dummies are excluded because of collinearity with the constant and the treatment dummy,  $D$ .<sup>37</sup> In terms of specification choice procedure, the idea is to add a higher order term to the polynomial until the bin dummies are no longer jointly significant.

Another major advantage of this procedure is that testing whether the bin dummies are significant turns out to be a test for

<sup>36</sup> See Blundell and Duncan (1998) for a more general discussion of series estimators.

<sup>37</sup> While excluding dummies for the two bins next to the cutoff point yields more interpretable results ( $\tau$  remains the treatment effect), the test is invariant to the excluded bin dummies, provided that one excluded dummy is on the left of the cutoff point and the other one is on the right (something standard regression packages will automatically do if all  $K$  dummies are included in the regression).

the presence of discontinuities in the regression function at points other than the cutoff point. In that sense, it provides a falsification test of the RD design by examining whether there are other unexpected discontinuities in the regression function at randomly chosen points (the bin thresholds). To see this, rewrite  $\sum_{k=1}^K \phi_k B_k$  as

$$\sum_{k=1}^K \phi_k B_k = \phi_1 + \sum_{k=2}^K (\phi_k - \phi_{k-1}) B_k^+,$$

where  $B_k^+ = \sum_{j=k}^K B_j$  is a dummy variable indicating that the observation is in bin  $k$  or above, i.e., that the assignment variable  $X$  is above the bin cutoff  $b_k$ . Testing whether all the  $\phi_k - \phi_{k-1}$  are equal to zero is equivalent to testing that all the  $\phi_k$  are the same (the above test), which amounts to testing that the regression line does not jump at the bin thresholds  $b_k$ .

Tables 2 and 3 show the estimates of the treatment effect for the voting example. For the sake of completeness, a wide range of bandwidths and specifications are presented, along with the corresponding  $p$ -values for the goodness-of-fit test discussed above (a bandwidth of 0.01 is used for the bins used to construct the test). We also indicate at the bottom of the tables the order of the polynomial selected for each bandwidth using the AIC. Note that the estimates of the treatment effect for the “order zero” polynomials are just comparisons of means on the two sides of the cutoff point, while the estimates for the “order one” polynomials are based on (local) linear regressions.

Broadly speaking, the goodness-of-fit tests do a very good job ruling out clearly misspecified models, like the zero order polynomials with large bandwidths that yield upward biased estimates of the treatment effect. Estimates from models that pass the goodness-of-fit test mostly fall in the 0.05–0.10 range for the vote share (table 2)

and 0.37–0.57 for the probability of winning (table 3). One set of models the goodness-of-fit test does not rule out, however, is higher order polynomial models with small bandwidths that tend to be imprecisely estimated as they “overfit” the data.

Looking informally at both the fit of the model (goodness-of-fit test) and the precision of the estimates (standard errors) suggests the following strategy: use higher order polynomials for large bandwidths of 0.50 and more, lower order polynomials for bandwidths between 0.05 and 0.50, and zero order polynomials (comparisons of means) for bandwidths of less than 0.05, since the latter specification passes the goodness-of-fit test for these very small bandwidths. Interestingly, this informal approach more or less corresponds to what is suggested by the AIC. In this specific example, it seems that given a specific bandwidth, the AIC provides reasonable suggestions on which order of the polynomial to use.

#### 4.3.3 Estimation in the Fuzzy RD Design

As discussed earlier, in both the “sharp” and the “fuzzy” RD designs, the probability of treatment jumps discontinuously at the cutoff point. Unlike the case of the sharp RD where the probability of treatment jumps from 0 to 1 at the cutoff, in the fuzzy RD case, the probability jumps by less than one. In other words, treatment is not solely determined by the strict cutoff rule in the fuzzy RD design. For example, even if eligibility for a treatment solely depends on a cutoff rule, not all the eligibles may get the treatment because of imperfect compliance. Similarly, program eligibility may be extended in some cases even when the cutoff rule is not satisfied. For example, while Medicare eligibility is mostly determined by a cutoff rule (age 65 or older), some disabled individuals under the age of 65 are also eligible.

Since we have already discussed the interpretation of estimates of the treatment effect

in a fuzzy RD design in section 3.4.1, here we just focus on estimation and implementation issues. The key message to remember from the earlier discussion is that, as in a standard IV framework, the estimated treatment effect can be interpreted as a local average treatment effect, provided monotonicity holds.

In the fuzzy RD design, we can write the probability of treatment as

$$\Pr(D = 1 | X = x) = \gamma + \delta T + g(x - c),$$

where  $T = 1[X \geq c]$  indicates whether the assignment variable exceeds the eligibility threshold  $c$ .<sup>38</sup> Note that the sharp RD is a special case where  $\gamma = 0$ ,  $g(\cdot) = 0$ , and  $\delta = 1$ . It is advisable to draw a graph for the treatment dummy  $D$  as a function of the assignment variable  $X$  using the same procedure discussed in section 4.1. This provides an informal way of seeing how large the jump in the treatment probability  $\delta$  is at the cutoff point, and what the functional form  $g(\cdot)$  looks like.

Since  $D = \Pr(D = 1 | X = x) + \nu$ , where  $\nu$  is an error term independent of  $X$ , the fuzzy RD design can be described by the two equation system:

$$(10) \quad Y = \alpha + \tau D + f(X - c) + \varepsilon,$$

$$(11) \quad D = \gamma + \delta T + g(X - c) + \nu.$$

Looking at these equations suggests estimating the treatment effect  $\tau$  by instrumenting the treatment dummy  $D$  with  $T$ . Note also that substituting the treatment determining equation into the outcome equation yields the reduced form

$$(12) \quad Y = \alpha_r + \tau_r T + f_r(X - c) + \varepsilon_r,$$

where  $\tau_r = \tau \delta$ . In this setting,  $\tau_r$  can be interpreted as an “intent-to-treat” effect.

Estimation in the fuzzy RD design can be performed using either the local linear regression approach or polynomial regressions. Since the model is exactly identified, 2SLS estimates are numerically identical to the ratio of reduced form coefficients  $\tau_r/\delta$ , provided that the same bandwidth is used for equations (11) and (12) in the local linear regression case, and that the same order of polynomial is used for  $g(\cdot)$  and  $f(\cdot)$  in the polynomial regression case.

In the case of the local linear regression, Imbens and Lemieux (2008) recommend using the same bandwidth in the treatment and outcome regression. When we are close to a sharp RD design, the function  $g(\cdot)$  is expected to be very flat and the optimal bandwidth to be very wide. In contrast, there is no particular reason to expect the function  $f(\cdot)$  in the outcome equation to be flat or linear, which suggests the optimal bandwidth would likely be less than the one for the treatment equation. As a result, Imbens and Lemieux (2008) suggest focusing on the outcome equation for selecting bandwidth, and then using the same bandwidth for the treatment equation.

While using a wider bandwidth for the treatment equation may be advisable on efficiency grounds, there are two practical reasons that suggest not doing so. First, using different bandwidths complicates the computation of standard errors since the outcome and treatment samples used for the estimation are no longer the same, meaning the usual 2SLS standard errors are no longer valid. Second, since it is advisable to explore the sensitivity of results to changes in the bandwidth, “trying out” separate bandwidths for each of the two equations would lead to a large and difficult-to-interpret number of specifications.

<sup>38</sup> Although the probability of treatment is modeled as a linear probability model here, this does not impose any restrictions on the probability model since  $g(x - c)$  is unrestricted on both sides of the cutoff  $c$ , while  $T$  is a dummy variable. So there is no need to write the model using a probit or logit formulation.

The same broad arguments can be used in the case of local polynomial regressions. In principle, a lower order of polynomial could be used for the treatment equation (11) than for the outcome equation (12). In practice, however, it is simpler to use the same order of polynomial and just run 2SLS (and use 2SLS standard errors).

#### 4.3.4 How to Compute Standard Errors?<sup>9</sup>

As discussed above, for inference in the sharp RD case we can use standard least squares methods. As usual, it is recommended to use heteroskedasticity-robust standard errors (Halbert White 1980) instead of standard least squares standard errors. One additional reason for doing so in the RD case is to ensure the standard error of the treatment effect is the same when either a pooled regression or two separate regressions on each side of the cutoff are used to compute the standard errors. As we just discussed, it is also straightforward to compute standard errors in the fuzzy RD case using 2SLS methods, although robust standard errors should also be used in this case. Imbens and Lemieux (2008) propose an alternative way of computing standard errors in the fuzzy RD case, but nonetheless suggest using 2SLS standard errors readily available in econometric software packages.

One small complication that arises in the nonparametric case of local linear regressions is that the usual (robust) standard errors from least squares are only valid provided that  $h \propto N^{-\delta}$  with  $1/5 < \delta < 2/5$ . As we mentioned earlier, this is not a very important point in practice, and the usual standard errors can be used with local linear regressions.

#### 4.4 Implementing Empirical Tests of RD Validity and Using Covariates

In this part of the section, we describe how to implement tests of the validity of the RD design and how to incorporate covariates in the analysis.

##### 4.4.1 Inspection of the Histogram of the Assignment Variable

Recall that the underlying assumption that generates the local random assignment result is that each individual has imprecise control over the assignment variable, as defined in section 3.1.1. We cannot test this directly (since we will only observe one observation on the assignment variable per individual at a given point in time), but an intuitive test of this assumption is whether the *aggregate* distribution of the assignment variable is discontinuous, since a mixture of individual-level continuous densities is itself a continuous density.

McCrary (2008) proposes a simple two-step procedure for testing whether there is a discontinuity in the density of the assignment variable. In the first step, the assignment variable is partitioned into equally spaced bins and frequencies are computed within those bins. The second step treats the frequency counts as a dependent variable in a local linear regression. See McCrary (2008), who adopts the nonparametric framework for asymptotics, for details on this procedure for inference.

As McCrary (2008) points out, this test can fail to detect a violation of the RD identification condition if for some individuals there is a “jump” up in the density, offset by jumps “down” for others, making the aggregate density continuous at the threshold. McCrary (2008) also notes it is possible the RD estimate could remain unbiased, even when there is important manipulation of the assignment variable causing a jump in the density. It should be noted, however, that in order to rely upon the RD estimate as unbiased, one needs to invoke other identifying assumptions and cannot rely upon the mild conditions we focus on in this article.<sup>39</sup>

<sup>39</sup> McCrary (2008) discusses an example where students who barely fail a test are given extra points so that they barely pass. The RD estimator can remain unbiased if one assumes that those who are given extra points were chosen randomly from those who barely failed.

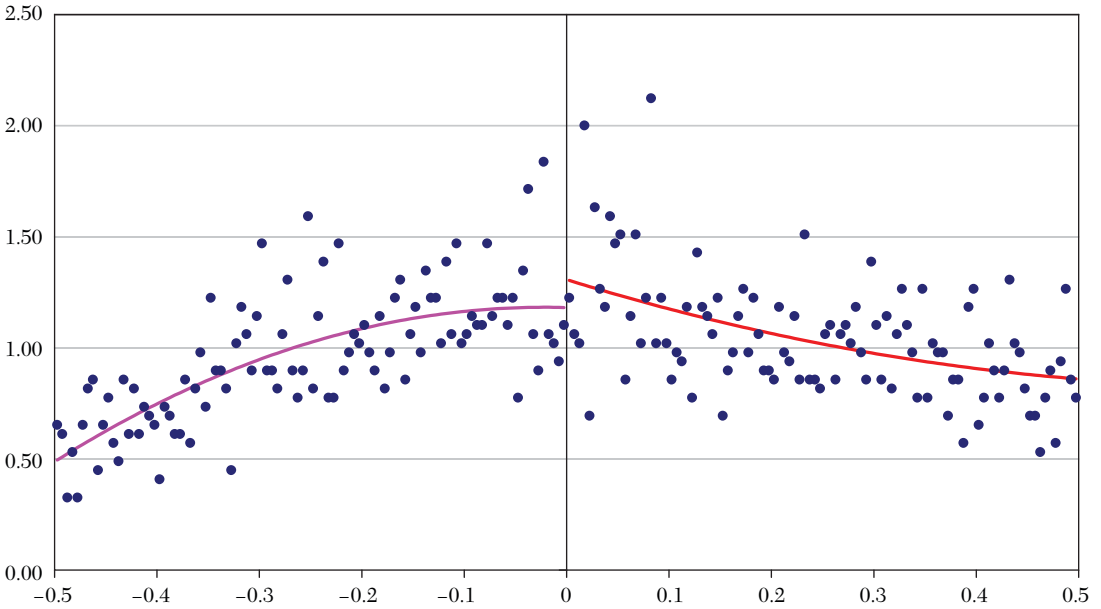


Figure 16. Density of the Forcing Variable (Vote Share in Previous Election)

One of the examples McCrary uses for his test is the voting model of Lee (2008) that we used in the earlier empirical examples. Figure 16 shows a graph of the raw densities computed over bins with a bandwidth of 0.005 (200 bins in the graph), along with a smooth second order polynomial model. Consistent with McCrary (2008), the graph shows no evidence of discontinuity at the cutoff. McCrary also shows that a formal test fails to reject the null hypothesis of no discontinuity in the density at the cutoff.

#### 4.4.2 Inspecting Baseline Covariates

An alternative approach for testing the validity of the RD design is to examine whether the observed baseline covariates are “locally” balanced on either side of the threshold, which should be the case if the treatment indicator is locally randomized.

A natural thing to do is conduct both a graphical RD analysis and a formal

estimation, replacing the dependent variable with each of the observed baseline covariates in  $W$ . A discontinuity would indicate a violation in the underlying assumption that predicts local random assignment. Intuitively, if the RD design is valid, we *know* that the treatment variable cannot influence variables determined prior to the realization of the assignment variable and treatment assignment; if we observe it does, something is wrong in the design.

If there are many covariates in  $W$ , even abstracting from the possibility of misspecification of the functional form, some discontinuities will be statistically significant by random chance. It is thus useful to combine the multiple tests into a single test statistic to see if the data are consistent with no discontinuities for any of the observed covariates. A simple way to do this is with a Seemingly Unrelated Regression (SUR) where each equation represents a different baseline

covariate, and then perform a  $\chi^2$  test for the discontinuity gaps in all questions being zero. For example, supposing the underlying functional form is linear, one would estimate the system

$$\begin{aligned} w_1 &= \alpha_1 + D\beta_1 + X\gamma_1 + \varepsilon_1 \\ &\vdots \\ w_K &= \alpha_K + D\beta_K + X\gamma_K + \varepsilon_K \end{aligned}$$

and test the hypothesis that  $\beta_1, \dots, \beta_K$  are jointly equal to zero, where we allow the  $\varepsilon$ 's to be correlated across the  $K$  equations. Alternatively, one can simply use the OLS estimates of  $\beta_1, \dots, \beta_K$  obtained from a “stacked” regression where all the equations for each covariate are pooled together, while  $D$  and  $X$  are fully interacted with a set of  $K$  dummy variables (one for each covariate  $w_k$ ). Correlation in the error terms can then be captured by clustering the standard errors on individual observations (which appear in the stacked dataset  $K$  times). Under the null hypothesis of no discontinuities, the Wald test statistic  $N\hat{\beta}'\hat{V}^{-1}\hat{\beta}$  (where  $\hat{\beta}$  is the vector of estimates of  $\beta_1, \dots, \beta_K$ , and  $\hat{V}$  is the cluster-and-heteroskedasticity consistent estimate of the asymptotic variance of  $\hat{\beta}$ ) converges in distribution to a  $\chi^2$  with  $K$  degrees of freedom.

Of course, the importance of functional form for RD analysis means a rejection of the null hypothesis tells us either that the underlying assumptions for the RD design are invalid, or that at least some of the equations are sufficiently misspecified and too restrictive, so that nonzero discontinuities are being estimated, even though they do not exist in the population. One could use the parametric specification tests discussed earlier for each of the individual equations to see if misspecification of the functional form is an important problem. Alternatively, the test could be performed only for observations

within a narrower window around the cut-off point, such as the one suggested by the bandwidth selection procedures discussed in section 4.3.1.

Figure 17 shows the RD graph for a baseline covariate, the Democratic vote share in the election prior to the one used for the assignment variable (four years prior to the current election). Consistent with Lee (2008), there is no indication of a discontinuity at the cutoff. The actual RD estimate using a quartic model is  $-0.004$  with a standard error of 0.014. Very similar results are obtained using winning the election as the outcome variable instead (RD estimate of  $-0.003$  with a standard error of 0.017).

#### 4.5 Incorporating Covariates in Estimation

If the RD design is valid, the other use for the baseline covariates is to reduce the sampling variability in the RD estimates. We discuss two simple ways to do this. First, one can “residualize” the dependent variable—subtract from  $Y$  a prediction of  $Y$  based on the baseline covariates  $W$ —and then conduct an RD analysis on the residuals. Intuitively, this procedure nets out the portion of the variation in  $Y$  we could have predicted using the predetermined characteristics, making the question whether the treatment variable can explain the remaining residual variation in  $Y$ . The important thing to keep in mind is that if the RD design is valid, this procedure provides a consistent estimate of the same RD parameter of interest. Indeed, any combination of covariates can be used, and abstracting from functional form issues, the estimator will be consistent for the same parameter, as discussed above in equation (4). Importantly, this two-step approach also allows one to perform a graphical analysis of the residual.

To see this more formally in the parametric case, suppose one is willing to assume that the expectation of  $Y$  as a function of  $X$  is a polynomial, and the expectation of each

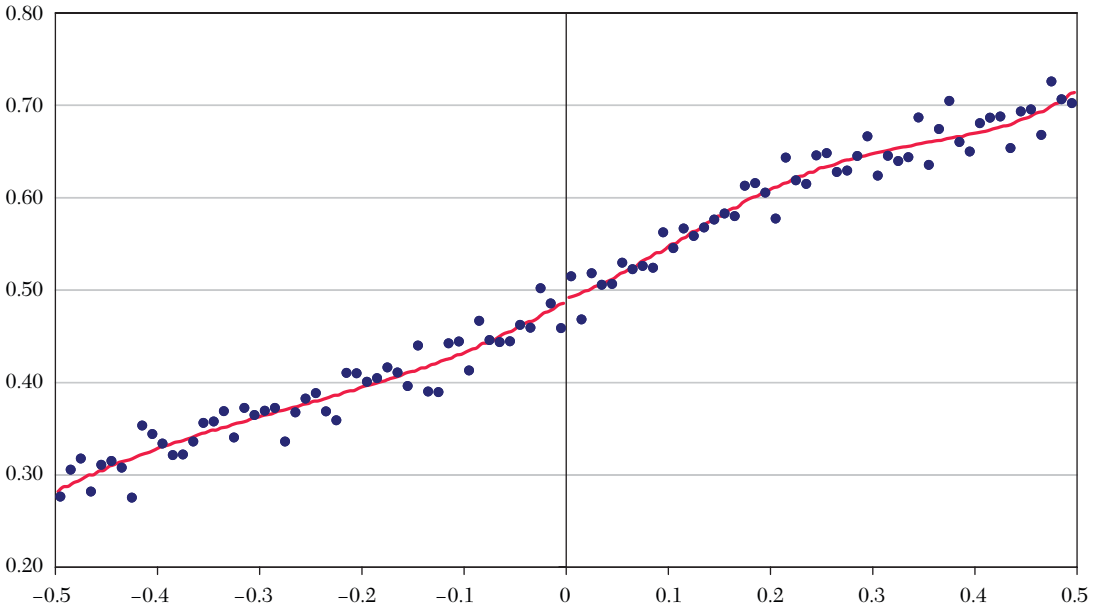


Figure 17. Discontinuity in Baseline Covariate (Share of Vote in Prior Election)

element of  $W$  is also a polynomial function of  $X$ . This implies

$$(13) \quad Y = D\tau + \tilde{X}\tilde{\gamma} + \varepsilon$$

$$W = \tilde{X}\delta + u,$$

where  $\tilde{X}$  is a vector of polynomial terms in  $X$ ,  $\delta$  and  $u$  are of conformable dimension, and  $\varepsilon$  and  $u$  are by construction orthogonal to  $D$  and  $\tilde{X}$ . It follows that

$$(14) \quad \begin{aligned} Y - W\pi &= D\tau + \tilde{X}\tilde{\gamma} - W\pi + \varepsilon \\ &= D\tau + \tilde{X}(\tilde{\gamma} - \delta\pi) - u\pi + \varepsilon \\ &= D\tau + \tilde{X}\tilde{\gamma} - u\pi + \varepsilon. \end{aligned}$$

This makes clear that a regression of  $Y - W\pi$  on  $D$  and  $\tilde{X}$  will give consistent estimates of  $\tau$  and  $\tilde{\gamma}$ . This is true no matter the value of  $\pi$ . Furthermore, as long as the specification in

equation (13) is correct, in computing estimated standard errors in the second step, one can ignore the fact that the first step was estimated.<sup>40</sup>

The second approach—which uses the same assumptions implicit in equation (13)—is to simply add  $W$  to the regression. While this may seem to impose linearity in how  $W$

<sup>40</sup>The two-step procedure solves the sample analogue to the following set of moment equations:

$$\begin{aligned} E\left[\begin{pmatrix} D \\ \tilde{X} \end{pmatrix}(Y - W\pi_0 - D\tau - \tilde{X}\tilde{\gamma})\right] &= 0 \\ E[W(Y - W\pi_0)] &= 0. \end{aligned}$$

As noted above, the second-step estimator for  $\tau$  is consistent for any value of  $\pi$ . Letting  $\theta \equiv \begin{pmatrix} \tau \\ \tilde{\gamma} \end{pmatrix}$ , and using the notation of Whitney K. Newey and Daniel L. McFadden (1994), this means that the first row of  $\nabla_{\pi}\theta(\pi_0) = -G_{\theta}^{-1}G_{\pi}$  is a row of zeros. It follows from their theorem 6.1, with the 1,1 element of  $V$  being the asymptotic variance of the estimator for  $\tau$ , that the 1,1 element of  $V$  is equal to the 1,1 element of  $G_{\theta}^{-1}E[g(z)g(z)']G_{\theta}^{-1}$ , which is the asymptotic covariance matrix of the second stage estimator ignoring estimation in the first step.

affects  $Y$ , it can be shown that the inclusion of these regressors will not affect the consistency of the estimator for  $\tau$ .<sup>41</sup> The advantage of this second approach is that under these functional form assumptions and with homoskedasticity, the estimator for  $\tau$  is guaranteed to have a lower asymptotic variance.<sup>42</sup> By contrast, the “residualizing” approach can in some cases *raise* standard errors.<sup>43</sup>

The disadvantage of solely relying upon this second approach, however, is that it does not help distinguish between an inappropriate functional form and discontinuities in  $W$ , as both could potentially cause the estimates of  $\tau$  to change significantly when  $W$  is included.<sup>44</sup> On the other hand, the “residualizing” approach allows one to examine how well the residuals fit the assumed order of polynomial (using, for example, the methods described in subsection 4.3.2). If it does not fit well, then it suggests that the use of that order of polynomial with the second approach is not justified. Overall, one sensible approach is to directly enter the covariates, but then to use the “residualizing” approach as an additional diagnostic check on whether the assumed order of the polynomial is justified.

As discussed earlier, an alternative approach to estimating the discontinuity

involves limiting the estimation to a window of data around the threshold and using a linear specification within that window.<sup>45</sup> We note that as the neighborhood shrinks, the true expectation of  $W$  conditional on  $X$  will become closer to being linear, and so equation (13) (with  $\tilde{X}$  containing only the linear term) will become a better approximation.

For the voting example used throughout this paper, Lee (2008) shows that adding a set of covariates essentially has no impact on the RD estimates in the model where the outcome variable is winning the next election. Doing so does not have a large impact on the standard errors either, at least up to the third decimal. Using the procedure based on residuals instead actually slightly increases the second step standard errors—a possibility mentioned above. Therefore in this particular example, the main advantage of using baseline covariates is to help establish the validity of the RD design, as opposed to improving the efficiency of the estimators.

#### 4.6 A Recommended “Checklist” for Implementation

Below is a brief summary of our recommendations for the analysis, presentation, and estimation of RD designs.

<sup>41</sup> To see this, rewrite equation (13) as  $Y = D\tau + \tilde{X}\gamma + Da + \tilde{X}b + Wc + \mu$ , where  $a, b, c$ , and  $\mu$  are linear projection coefficients and the residual from a population regression  $\varepsilon$  on  $D, \tilde{X}$ , and  $W$ . If  $a = 0$ , then adding  $W$  will not affect the coefficient on  $D$ . This will be true—applying the Frisch–Waugh theorem—when the covariance between  $\varepsilon$  and  $D - \tilde{X}d - We$  (where  $d$  and  $e$  are coefficients from projecting  $D$  on  $\tilde{X}$  and  $W$ ) is zero. This will be true when  $e = 0$ , because  $\varepsilon$  is by assumption orthogonal to both  $D$  and  $\tilde{X}$ . Applying the Frisch–Waugh theorem again,  $e$  is the coefficient obtained by regressing  $D$  on  $W - \tilde{X}\delta \equiv u$ ; by assumption  $u$  and  $D$  are uncorrelated, so  $e = 0$ .

<sup>42</sup> The asymptotic variance for the least squares estimator (without including  $W$ ) of  $\tau$  is given by the ratio  $V(\varepsilon)/V(\tilde{D})$  where  $\tilde{D}$  is the residual from the population regression of  $D$  on  $\tilde{X}$ . If  $W$  is included, then the least squares estimator has asymptotic variance of  $\sigma^2/V(D - \tilde{X}d - We)$ , where  $\sigma^2$  is the variance of the error when  $W$  is included, and  $d$  and  $e$  are coefficients from projecting  $D$  on  $\tilde{X}$  and  $W$ .  $\sigma^2$  cannot exceed  $V(\varepsilon)$ , and as shown in

the footnote above,  $e = 0$ , and thus  $D - \tilde{X}d = \tilde{D}$ , implying that the denominator in the ratio does not change when  $W$  is included.

<sup>43</sup> From equation (14), the regression error variance will increase if  $V(\varepsilon - u\pi) > V(\varepsilon) \Leftrightarrow V(u\pi) - 2C(\varepsilon, u\pi) > 0$ , which will hold when, for example,  $\varepsilon$  is orthogonal to  $u$  and  $\pi$  is nonzero.

<sup>44</sup> If the true equation for  $W$  contains more polynomial terms than  $\tilde{X}$ , then  $e$ , as defined in the preceding footnotes (the coefficient obtained by regressing  $D$  on the residual from projecting  $W$  on  $\tilde{X}$ ), will not be zero. This implies that including  $W$  will generally lead to inconsistent estimates of  $\tau$ , and may cause the asymptotic variance to increase (since  $V(D - \tilde{X}d - We) \leq V(\tilde{D})$ ).

<sup>45</sup> And we have noted that one can justify this by assuming that in that specified neighborhood, the underlying function is in fact linear, and make standard parametric inferences. Or one can conduct a nonparametric inference approach by making assumptions about the rate at which the bandwidth shrinks as the sample size grows.

1. **To assess the possibility of manipulation of the assignment variable, show its distribution.** The most straightforward thing to do is to present a histogram of the assignment variable, using a fixed number of bins. The bin widths should be as small as possible, without compromising the ability to visually see the overall shape of the distribution. For an example, see figure 16. The bin-to-bin jumps in the frequencies can provide a sense in which any jump at the threshold is “unusual.” For this reason, we recommend *against* plotting a smooth function comprised of kernel density estimates. A more formal test of a discontinuity in the density can be found in McCrary (2008).
2. **Present the main RD graph using binned local averages.** As with the histogram, we recommend using a fixed number of nonoverlapping bins, as described in subsection 4.1. For examples, see figures 6–11. The nonoverlapping nature of the bins for the local averages is important; we recommend *against* simply presenting a continuum of nonparametric estimates (with a single break at the threshold), as this will naturally tend to give the impression of a discontinuity even if there does not exist one in the population. We recommend reporting bandwidths implied by cross-validation, as well as the range of widths that are not statistically rejected in favor of strictly less restrictive alternatives (for an example, see table 1). We recommend generally “undersmoothing,” while at the same time avoiding “too narrow” bins that produce a scatter of data points, from which it is difficult to see the shape of the underlying function. Indeed, we recommend *against* simply plotting the raw data without a minimal amount of local averaging.
3. **Graph a benchmark polynomial specification.** Superimpose onto the graph the predicted values from a low-order polynomial specification (see figures 6–11). One can often informally assess by comparing the two functions whether a simple polynomial specification is an adequate summary of the data. If the local averages represent the most flexible “nonparametric” representation of the function, the polynomial represents a “best case” scenario in terms of the variance of the RD estimate, since if the polynomial specification is correct, under certain conditions, the least squares estimator is efficient.
4. **Explore the sensitivity of the results to a range of bandwidths, and a range of orders to the polynomial.** For an example, see tables 2 and 3. The tables should be supplemented with information on the implied rule-of-thumb bandwidth and cross-validation bandwidths for local linear regression (as in table 4), as well as the AIC-implied optimal order of the polynomial. The specification tests that involve adding bin dummies to the polynomial specifications can help rule out overly restrictive specifications. Among all the specifications that are not rejected by the bin-dummy tests, and among the polynomial orders recommended by the AIC, and the estimates given by both rule of thumb and CV bandwidths, report a “typical” point estimate and a range of point estimates. A useful graphical device for illustrating the sensitivity of the results to bandwidths is to plot the local linear discontinuity estimate against a continuum of bandwidths (within a range of bandwidths that are not ruled out by the above specification tests). For an example

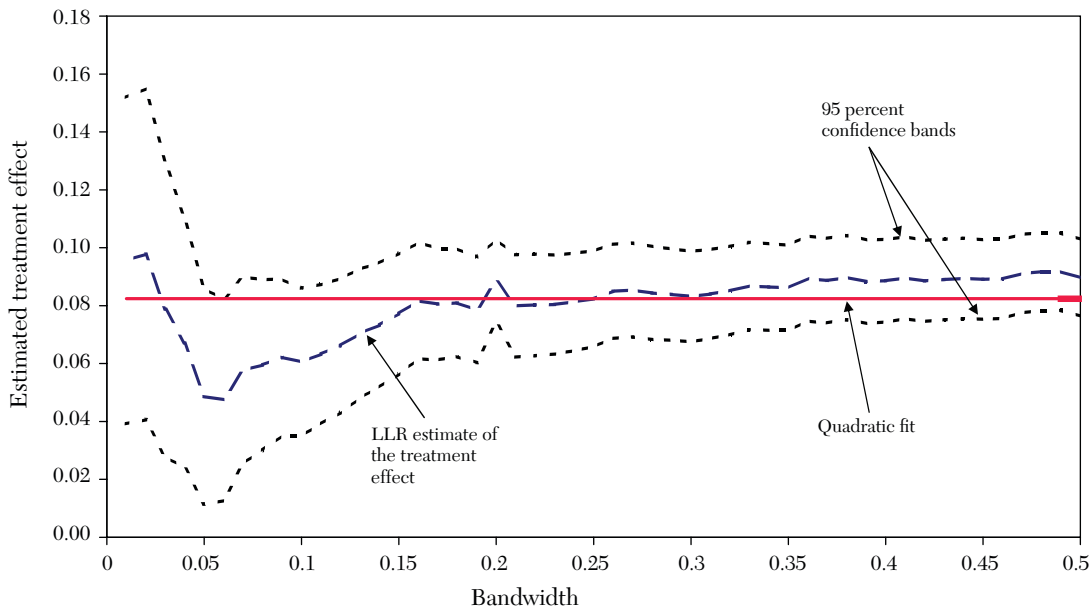


Figure 18. Local Linear Regression with Varying Bandwidth: Share of Vote at Next Election

of such a presentation, see the online appendix to Card, Carlos Dobkin, and Nicole Maestas (2009), and figure 18.

5. **Conduct a parallel RD analysis on the baseline covariates.** As discussed earlier, if the assumption that there is no precise manipulation or sorting of the assignment variable is valid, then there should be no discontinuities in variables that are determined prior to the assignment. See figure 17, for example.
6. **Explore the sensitivity of the results to the inclusion of baseline covariates.** As discussed above, the inclusion of baseline covariates—no matter how highly correlated they are with the outcome—should not affect the estimated discontinuity, if the no-manipulation

assumption holds. If the estimates do change in an important way, it may indicate a potential sorting of the assignment variable that may be reflected in a discontinuity in one or more of the baseline covariates. In terms of implementation, in subsection 4.5, we suggest simply including the covariates directly, after choosing a suitable order of polynomial. Significant changes in the estimated effect or increases in the standard errors may be an indication of a misspecified functional form. Another check is to perform the “residualizing” procedure suggested there, to see if that same order of polynomial provides a good fit for the residuals, using the specification tests from point 4.

We recognize that, due to space limitations, researchers may be unable to present every

permutation of presentation (e.g., points 2–4 for every one of 20 baseline covariates) within a published article. Nevertheless, we do believe that documenting the sensitivity of the results to these array of tests and alternative specifications—even if they only appear in unpublished, online appendices—is an important component of a thorough RD analysis.

### 5. *Special Cases*

In this section, we discuss how the RD design can be implemented in a number of specific cases beyond the one considered up to this point (that of a single cross-section with a continuous assignment variable).

#### 5.1 *Discrete Assignment Variable and Specification Errors*

Up until now, we have assumed the assignment variable was continuous. In practice, however,  $X$  is often discrete. For example, age or date of birth are often only available at a monthly, quarterly, or annual frequency level. Studies relying on an age-based cutoff thus typically rely on discrete values of the age variable when implementing an RD design.

Lee and Card (2008) study this case in detail and make a number of important points. First, with a discrete assignment variable, it is not possible to compare outcomes in very narrow bins just to the right and left of the cutoff point. Consequently, one must use regressions to estimate the conditional expectation of the outcome variable at the cutoff point by extrapolation. As discussed in section 4, however, in practice we always extrapolate to some extent, even in the case of a continuous assignment variable. So the fact we must do so in the case of a discrete assignment variable does not introduce particular complications from an econometric point of view, provided the discrete variable is not too coarsely distributed.

Additionally, the various estimation and graphing techniques discussed in section 4 can readily be used in the case of a discrete assignment variable. For instance, as with a continuous assignment variable, either local linear regressions or polynomial regressions can be used to estimate the jump in the regression function at the cutoff point. Furthermore, the discreteness of the assignment variable simplifies the problem of bandwidth choice when graphing the data since, in most cases, one can simply compute and graph the mean of the outcome variable for each value of the discrete assignment variable. The fact that the variable is discrete also provides a natural way of testing whether the regression model is well specified by comparing the fitted model to the raw dispersion in mean outcomes at each value of the assignment variable. Lee and Card (2008) show that, when errors are homoskedastic, the model specification can be tested using the standard goodness-of-fit statistic

$$G \equiv \frac{(ESS_R - ESS_{UR})/(J - K)}{ESS_{UR}/(N - J)},$$

where  $ESS_R$  is the estimated sum of squares of the restricted model (e.g., low order polynomial), while  $ESS_{UR}$  is the estimated sum of squares of the unrestricted model where a full set of dummies (for each value of the assignment variable) are included. In this unrestricted model, the fitted regression corresponds to the mean outcome in each cell.  $G$  follows a  $F(J - K, N - J)$  distribution where  $J$  is the number of values taken by the assignment variables and  $K$  is the number of parameters of the restricted model.

This test is similar to the test in section 4 where we suggested including a full set of bin dummies in the regression model and testing whether the bin dummies were jointly significant. The procedure is even

simpler here, as bin dummies are replaced by dummies for each value of the discrete assignment variable. In the presence of heteroskedasticity, the goodness-of-fit test can be computed by estimating the model and testing whether a set of dummies for each value of the discrete assignment variable are jointly significant. In that setting, the test statistic follows a chi-square distribution with  $J - K$  degrees of freedom.

In Lee and Card (2008), the difference between the true conditional expectation  $E[Y|X = x]$  and the estimated regression function forming the basis of the goodness-of-fit test is interpreted as a random specification error that introduces a group structure in the standard errors. One way of correcting the standard errors for group structure is to run the model on cell means.<sup>46</sup> Another way is to “cluster” the standard errors. Note that in this setting, the goodness-of-fit test can also be interpreted as a test of whether standard errors should be adjusted for the group structure. In practice, it is nonetheless advisable to either group the data or cluster the standard errors in micro-data models irrespective of the results of the goodness-of-fit test. The main purpose of the test should be to help choose a reasonably accurate regression model.

Lee and Card (2008) also discuss a number of issues including what to do when specification errors under treatment and control are correlated, and how to possibly adjust the RD estimates in the presence of specification errors. Since these issues are beyond the scope of this paper, interested readers should consult Lee and Card (2008) for more detail.

<sup>46</sup>When the discrete assignment variable—and the “treatment” dummy solely dependent on this variable—is the only variable used in the regression model, standard OLS estimates will be numerically equivalent to those obtained by running a weighted regression on the cell means, where the weights are the number of observations (or the sum of individual weights) in each cell.

## 5.2 Panel Data and Fixed Effects

In some situations, the RD design will be embedded in a panel context, whereby period by period, the treatment variable is determined according to the realization of the assignment variable  $X$ . Again, it seems natural to propose the model

$$Y_{it} = D_{it}\tau + f(X_{it}; \gamma) + a_i + \varepsilon_{it}$$

(where  $i$  and  $t$  denote the individuals and time, respectively), and simply estimate a fixed effects regression by including individual dummy variables to capture the unit-specific error component,  $a_i$ . It is important to note, however, that including fixed effects is unnecessary for identification in an RD design. This sharply contrasts with a more traditional panel data setting where the error component  $a_i$  is allowed to be correlated with the observed covariates, including the treatment variable  $D_{it}$ , in which case including fixed effects is essential for consistently estimating the treatment effect  $\tau$ .

An alternative is to simply conduct the RD analysis for the entire pooled-cross-section dataset, taking care to account for within-individual correlation of the errors over time using clustered standard errors. The source of identification is a comparison between those just below and above the threshold, and can be carried out with a single cross-section. Therefore, imposing a specific dynamic structure introduces more restrictions without any gain in identification.

Time dummies can also be treated like any other baseline covariate. This is apparent by applying the main RD identification result: conditional on what period it is, we are assuming the density of  $X$  is continuous at the threshold and, hence, conditional on  $X$ , the probability of an individual observation coming from a particular period is also continuous.

We note that it becomes a little bit more awkward to use the justification proposed in subsection 4.5 for directly including dummies for individuals and time periods on the right hand side of the regression. This is because the assumption would have to be that the probability that an observation belonged to each individual (or the probability that an observation belonged to each time period) is a polynomial function in  $X$  and, strictly speaking, nontrivial polynomials are not bounded between 0 and 1.

A more practical concern is that inclusion of individual dummy variables may lead to an *increase* in the variance of the RD estimator for another reason. If there is little “within-unit” variability in treatment status, then the variation in the main variable of interest (treatment after partialling out the individual heterogeneity) may be quite small. Indeed, seeing standard errors rise when including fixed effects may be an indication of a misspecified functional form.<sup>47</sup>

Overall, since the RD design is still valid ignoring individual or time effects, then the only rationale for including them is to reduce sampling variance. But there are other ways to reduce sampling variance by exploiting the structure of panel data. For instance, we can treat the lagged dependent variable  $Y_{it-1}$  as simply another baseline covariate in period  $t$ . In cases where  $Y_{it}$  is highly persistent over time,  $Y_{it-1}$  may well be a very good predictor and has a very good chance of reducing the sampling error. As we have also discussed earlier, looking at possible discontinuities in baseline covariates is an important test of the validity of the RD design. In this particular case, since  $Y_{it}$  can be highly correlated with  $Y_{it-1}$ , finding a discontinuity in  $Y_{it}$  but not in  $Y_{it-1}$  would be a strong piece of evidence supporting the validity of the RD design.

In summary, one can utilize the panel nature of the data by conducting an RD analysis on the entire dataset, using lagged variables as baseline covariates for inclusion as described in subsection 4.5. The primary caution in doing this is to ensure that for each period, the included covariates are the variables determined *prior* to the present period’s realization of  $X_{it}$ .

## 6. *Applications of RD Designs in Economics*

In what areas has the RD design been applied in economic research? Where do discontinuous rules come from and where might we expect to find them? In this section, we provide some answers to these questions by providing a survey of the areas of applied economic research that have employed the RD design. Furthermore, we highlight some examples from the literature that illustrate what we believe to be the most important elements of a compelling, “state-of-the-art” implementation of RD.

### 6.1 *Areas of Research Using RD*

As we suggested in the introduction, the notion that the RD design has limited applicability to a few specific topics is inconsistent with our reading of existing applied research in economics. Table 5 summarizes our survey of empirical studies on economic topics that have utilized the RD design. In compiling this list, we searched economics journals as well as listings of working papers from economists, and chose any study that recognized the potential use of an RD design in their given setting. We also included some papers from non-economists when the research was closely related to economic work.

Even with our undoubtedly incomplete compilation of over sixty studies, table 5 illustrates that RD designs have been applied in many different contexts. Table 5 summarizes the context of the study, the outcome

<sup>47</sup> See the discussion in section 4.5.

TABLE 5  
REGRESSION DISCONTINUITY APPLICATIONS IN ECONOMICS

Study	Context	Outcome(s)	Treatment(s)	Assignment variable(s)
<b>Education</b>				
Angrist and Lavy (1999)	Public Schools (Grades 3–5), Israel	Test scores	Class size	Student enrollment
Asadullah (2005)	Secondary schools, Bangladesh	Examination pass rate	Class size	Student enrollment
Bayer, Ferreira, and McMillan (2007)	Valuation of schools and neighborhoods, Northern California	Housing prices, school test scores, demographic characteristics	Inclusion in school attendance region	Geographic location
Black (1999)	Valuation of school quality, Massachusetts	Housing prices	Inclusion in school attendance region	Geographic location
Canton and Blom (2004)	Higher education, Mexico	University enrollment, GPA, Part-time employment, Career choice	Student loan receipt	Economic need index
Cascio and Lewis (2006)	Teenagers, United States	AFQT test scores	Age at school entry	Birthdate
Chay, McEwan, and Urquiola (2005)	Elementary schools, Chile	Test scores	Improved infrastructure, more resources	School averages of test scores
Chiang (2009)	School accountability, Florida	Test scores, education quality	Threat of sanctions	School's assessment score
Clark (2009)	High schools, U.K.	Examination pass rates	“Grant maintained” school status	Vote share
Ding and Lehrer (2007)	Secondary school students, China	Academic achievement (Test scores)	School assignment	Entrance examination scores
Figlio and Kenny (2009)	Elementary and middle schools, Florida	Private donations to school	D or F grade in school performance measure	Grading points
Goodman (2008)	College enrollment, Massachusetts	School choice	Scholarship offer	Test scores
Goolsbee and Guryan (2006)	Public schools, California	Internet access in classrooms, test scores	E-Rate subsidy amount	Proportion of students eligible for lunch program
Guryan (2001)	State-level equalization: elementary, middle schools, Massachusetts	Spending on schools, test scores	State education aid	Relative average property values
Hoxby (2000)	Elementary schools, Connecticut	Test scores	Class size	Student enrollment
Kane (2003)	Higher education, California	College attendance	Financial aid receipt	Income, assets, GPA
Lavy (2002)	Secondary schools, Israel	Test scores, drop out rates	Performance based incentives for teachers	Frequency of school type in community
Lavy (2004)	Secondary schools, Israel	Test scores	Pay-for-performance incentives	School matriculation rates
Lavy (2006)	Secondary schools, Tel Aviv	Dropout rates, test scores	School choice	Geographic location
Jacob and Lefgren (2004a)	Elementary schools, Chicago	Test scores	Teacher training	School averages on test scores

TABLE 5 (continued)  
REGRESSION DISCONTINUITY APPLICATIONS IN ECONOMICS

Study	Context	Outcome(s)	Treatment(s)	Assignment variable(s)
Jacob and Lefgren (2004b)	Elementary schools, Chicago	Test scores	Summer school attendance, grade retention	Standardized test scores
Leuven, Lindahl, Oosterbeek, and Webbink (2007)	Primary schools, Netherlands	Test scores	Extra funding	Percent disadvantaged minority pupils
Matsudaira (2008)	Elementary schools, Northeastern United States	Test scores	Summer school, grade promotion	Test scores
Urquiola (2006)	Elementary schools, Bolivia	Test scores	Class size	Student enrollment
Urquiola and Verhoogen (2009)	Class size sorting- RD violations, Chile	Test scores	Class size	Student enrollment
Van der Klaauw (2002, 1997)	College enrollment, East Coast College	Enrollment	Financial Aid offer	SAT scores, GPA
Van der Klaauw (2008a)	Elementary/middle schools, New York City	Test scores, student attendance	Title I federal funding	Poverty rates
<b>Labor Market</b>				
Battistin and Rettore (2002)	Job training, Italy	Employment rates	Training program (computer skills)	Attitudinal test score
Behaghel, Crepon, and Sedillot (2008)	Labor laws, France	Hiring among age groups	Tax exemption for hiring firm	Age of worker
Black, Smith, Berger, and Noel (2003); Black, Galdo, and Smith (2007b)	UI claimants, Kentucky	Earnings, benefit receipt/duration	Mandatory reemployment services (job search assistance)	Profiling score (expected benefit duration)
Card, Chetty, and Weber (2007)	Unemployment benefits, Austria	Unemployment duration	Lump-sum severance pay, extended UI benefits	Months employed, job tenure
Chen and van der Klaauw (2008)	Disability insurance beneficiaries, United States	Labor force participation	Disability insurance benefits	Age at disability decision
De Giorgi (2005)	Welfare-to-work program, United Kingdom	Re-employment probability	Job search assistance, training, education	Age at end of unemployment spell
DiNardo and Lee (2004)	Unionization, United States	Wages, employment, output	Union victory in NLRB election	Vote share
Dobkin and Ferreira (2009)	Individuals, California and Texas	Educational attainment, wages	Age at school entry	Birthdate
Edmonds (2004)	Child labor supply and school attendance, South Africa	Child labor supply, school attendance	Pension receipt of oldest family member	Age
Hahn, Todd, and van der Klaauw (1999)	Discrimination, United States	Minority employment	Coverage of federal antidiscrimination law	Number of employees at firm
Lalive (2008)	Unemployment Benefits, Austria	Unemployment duration	Maximum benefit duration	Age at start of unemployment spell, geographic location

TABLE 5 (continued)  
REGRESSION DISCONTINUITY APPLICATIONS IN ECONOMICS

Study	Context	Outcome(s)	Treatment(s)	Assignment variable(s)
Lalive (2007)	Unemployment, Austria	Unemployment duration, duration of job search, quality of post-unemployment jobs	Benefits duration	Age at start of unemployment spell
Lalive, Van Ours, and Zweimüller (2006)	Unemployment, Austria	Unemployment duration	Benefit replacement rate, potential benefit duration	Pre-unemployment income, age
Leuven and Oosterbeek (2004)	Employers, Netherlands	Training, wages	Business tax deduction, training	Age of employee
Lemieux and Milligan (2008)	Welfare, Canada	Employment, marital status, living arrangements	Cash benefit	Age
Oreopoulos (2006)	Returns to education, U.K.	Earnings	Coverage of compulsory schooling law	Birth year
<b>Political Economy</b>				
Albouy (2009)	Congress, United States	Federal expenditures	Party control of seat	Vote share in election
Albouy (2008)	Senate, United States	Roll call votes	Incumbency	Initial vote share
Ferreira and Gyourko (2009)	Mayoral Elections, United States	Local expenditures	Incumbency	Initial vote share
Lee (2008, 2001)	Congressional elections, United States	Vote share in next election	Incumbency	Initial vote share
Lee, Moretti, and Butler (2004)	House of Representatives, United States	Roll call votes	Incumbency	Initial vote share
McCrary (2008)	House of Representatives, United States	N/A	Passing of resolution	Share of roll call vote "Yeay"
Pettersson-Lidbom (2006)	Local Governments, Sweden and Finland	Expenditures, tax revenues	Number of council seats	Population
Pettersson-Lidbom (2008)	Local Governments, Sweden	Expenditures, tax revenues	Left-, right-wing bloc	Left-wing parties' share
<b>Health</b>				
Card and Shore-Sheppard (2004)	Medicaid, United States	Overall insurance coverage	Medicaid eligibility	Birthdate
Card, Dobkin, and Maestas (2008)	Medicare, United States	Health care utilization	Coverage under Medicare	Age
Card, Dobkin, and Maestas (2009)	Medicare, California	Insurance coverage, Health services, Mortality	Medicare coverage	Age
Carpenter and Dobkin (2009)	Alcohol and mortality, United States	Mortality	Attaining minimum legal drinking age	Age
Ludwig and Miller (2007)	Head Start, United States	Child mortality, educational attainment	Head Start funding	County poverty rates
McCrary and Royer (2003)	Maternal education, United States, California and Texas	Infant health, fertility timing	Age of school entry	Birthdate
Snyder and Evans (2006)	Social Security recipients, United States	Mortality	Social security payments (\$)	Birthdate

TABLE 5 (continued)  
REGRESSION DISCONTINUITY APPLICATIONS IN ECONOMICS

Study	Context	Outcome(s)	Treatment(s)	Assignment variable(s)
<b>Crime</b>				
Berk and DeLeeuw (1999)	Prisoner behavior in California	Inmate misconduct	Prison security levels	Classification score
Berk and Rauma (1983)	Ex-prisoners recidivism, California	Arrest, parole violation	Unemployment insurance benefit	Reported hours of work
Chen and Shapiro (2004)	Ex-prisoners recidivism, United States	Arrest rates	Prison security levels	Classification score
Lee and McCrary (2005)	Criminal offenders, Florida	Arrest rates	Severity of sanctions	Age at arrest
Hjalmarsson (2009)	Juvenile offenders, Washington State	Recidivism	Sentence length	Criminal history score
<b>Environment</b>				
Chay and Greenstone (2003)	Health effects of pollution, United States	Infant mortality	Regulatory status	Pollution levels
Chay and Greenstone (2005)	Valuation of air quality, United States	Housing prices	Regulatory status	Pollution levels
Davis (2008)	Restricted driving policy, Mexico	Hourly air pollutant measures	Restricted automobile use	Time
Greenstone and Gallagher (2008)	Hazardous waste, United States	Housing prices	Superfund clean-up status	Ranking of level of hazard
<b>Other</b>				
Battistin and Rettore (2008)	Mexican anti-poverty program (PROGRESA)	School attendance	Cash grants	Pre-assigned probability of being poor
Baum-Snow and Marion (2009)	Housing subsidies, United States	Residents' characteristics, new housing construction	Increased subsidies	Percentage of eligible households in area
Buddelmeyer and Skoufias (2004)	Mexican anti-poverty program (PROGRESA)	Child labor and school attendance	Cash grants	Pre-assigned probability of being poor
Buettner (2006)	Fiscal equalization across municipalities, Germany	Business tax rate	Implicit marginal tax rate on grants to localities	Tax base
Card, Mas, and Rothstein (2008)	Racial segregation, United States	Changes in census tract racial composition	Minority share exceeding "tipping" point	Initial minority share
Cole (2009)	Bank nationalization, India	Share of credit granted by public banks	Nationalization of private banks	Size of bank
Edmonds, Mammen, and Miller (2005)	Household structure, South Africa	Household composition	Pension receipt of oldest family member	Age
Ferreira (2007)	Residential mobility, California	Household mobility	Coverage of tax benefit	Age
Pence (2006)	Mortgage credit, United States	Size of loan	State mortgage credit laws	Geographic location
Pitt and Khandker (1998)	Poor households, Bangladesh	Labor supply, children school enrollment	Group-based credit program	Acreage of land
Pitt, Khandker, McKernan, and Latif (1999)	Poor households, Bangladesh	Contraceptive use, Childbirth	Group-based credit program	Acreage of land

variable, the treatment of interest, and the assignment variable employed.

While the categorization of the various studies into broad areas is rough and somewhat arbitrary, it does appear that a large share come from the area of education, where the outcome of interest is often an achievement test score and the assignment variable is also a test score, either at the individual or group (school) level. The second clearly identifiable group are studies that deal with labor market issues and outcomes. This probably reflects that, within economics, the RD design has so far primarily been used by labor economists, and that the use of quasi-experiments and program evaluation methods in documenting causal relationships is more prevalent in labor economics research.

There is, of course, nothing in the structure of the RD design tying it specifically to labor economics applications. Indeed, as the rest of the table shows, the remaining half of the studies are in the areas of political economy, health, crime, environment, and other areas.

## 6.2 Sources of Discontinuous Rules

Where do discontinuous rules come from, and in what situations would we expect to encounter them? As table 5 shows, there is a wide variety of contexts where discontinuous rules determine treatments of interest. There are, nevertheless, some patterns that emerge. We organize the various discontinuous rules below.

Before doing so, we emphasize that a good RD analysis—as with any other approach to program evaluation—is careful in clearly spelling out exactly what the treatment is, and whether it is of any real salience, independent of whatever effect it might have on the outcome. For example, when a pretest score is the assignment variable, we could always define a “treatment” as being “having passed the exam” (with a test score of 50 percent or higher), but this is not a very inter-

esting “treatment” to examine since it seems nothing more than an arbitrary label. On the other hand, if failing the exam meant not being able to advance to the next grade in school, the actual experience of treated and control individuals is observably different, no matter how large or small the impact on the outcome.

As another example, in the U.S. Congress, a Democrat obtaining the most votes in an election means something real—the Democratic candidate becomes a representative in Congress; otherwise, the Democrat has no official role in the government. But in a three-way electoral race, the treatment of the Democrat receiving the *second-most* number of votes (versus receiving the lowest number) is not likely a treatment of interest: only the first-place candidate is given any legislative authority. In principle, stories could be concocted about the psychological effect of placing second rather than third in an election, but this would be an example where the salience of the treatment is more speculative than when treatment is a concrete and observable event (e.g., a candidate becoming the sole representative of a constituency).

### 6.2.1 Necessary Discretization

Many discontinuous rules come about because resources cannot, for all practical purposes, be provided in a continuous manner. For example, a school can only have a whole number of classes per grade. For a fixed level of enrollment, the moment a school adds a single class, the average class size drops. As long as the number of classes is an increasing function of enrollment, there will be discontinuities at enrollments where a teacher is added. If there is a mandated maximum for the student to teacher ratio, this means that these discontinuities will be expected at enrollments that are exact multiples of the maximum. This is the

essence of the discontinuous rules used in the analyses of Angrist and Lavy (1999), M. Niaz Asadullah (2005), Caroline M. Hoxby (2000), Urquiola (2006), and Urquiola and Verhoogen (2009).

Another example of necessary discretization arises when children begin their schooling years. Although there are certainly exceptions, school districts typically follow a guideline that aims to group children together by age, leading to a grouping of children born in year-long intervals, determined by a single calendar date (e.g., September 1). This means children who are essentially of the same age (e.g., those born on August 31 and September 1), start school one year apart. This allocation of students to grade cohorts is used in Elizabeth U. Cascio and Ethan G. Lewis (2006), Dobkin and Fernando Ferreira (2009), and McCrary and Royer (2003).

Choosing a single representative by way of an election is yet another example. When the law or constitution calls for a single representative of some constituency and there are many competing candidates, the choice can be made via a “first-past-the-post” or “winner-take-all” election. This is the typical system for electing government officials at the local, state, and federal level in the United States. The resulting discontinuous relationship between win/loss status and the vote share is used in the context of the U.S. Congress in Lee (2001, 2008), Lee, Enrico Moretti and Matthew J. Butler (2004), David Albouy (2009), Albouy (2008), and in the context of mayoral elections in Ferreira and Joseph Gyourko (2009). The same idea is used in examining the impacts of union recognition, which is also decided by a secret ballot election (DiNardo and Lee 2004).

### 6.2.2 *Intentional Discretization*

Sometimes resources could potentially be allocated on a continuous scale but, in practice, are instead done in discrete levels. Among the studies we surveyed, we

identified three broad motivations behind the use of these discontinuous rules.

First, a number of rules seem driven by a compensatory or equalizing motive. For example, in Kenneth Y. Chay, Patrick J. McEwan and Urquiola (2005), Edwin Leuven et al. (2007), and van der Klaauw (2008a), extra resources for schools were allocated to the neediest communities, either on the basis of school-average test scores, disadvantaged minority proportions, or poverty rates. Similarly, Ludwig and Miller (2007), Erich Battistin and Enrico Rettore (2008), and Hielke Buddelmeyer and Emmanuel Skoufias (2004) study programs designed to help poor communities, where the eligibility of a community is based on poverty rates. In each of these cases, one could imagine providing the most resources to the neediest and gradually phasing them out as the need index declines, but in practice this is not done, perhaps because it was impractical to provide very small levels of the treatment, given the fixed costs in administering the program.

A second motivation for having a discontinuous rule is to allocate treatments on the basis of some measure of merit. This was the motivation behind the merit award from the analysis of Thistlethwaite and Campbell (1960), as well as recent studies of the effect of financial aid awards on college enrollment, where the assignment variable is some measure of student achievement or test score, as in Thomas J. Kane (2003) and van der Klaauw (2002).

Finally, we have observed that a number of discontinuous rules are motivated by the need to most effectively target the treatment. For example, environmental regulations or clean-up efforts naturally will focus on the most polluted areas, as in Chay and Michael Greenstone (2003), Chay and Greenstone (2005), and Greenstone and Justin Gallagher (2008). In the context of criminal behavior, prison security levels are often assigned based on an underlying score that quantifies

potential security risks, and such rules were used in Richard A. Berk and Jan de Leeuw (1999) and M. Keith Chen and Jesse M. Shapiro (2004).

### 6.3 Nonrandomized Discontinuity Designs

Throughout this article, we have focused on regression discontinuity designs that follow a certain structure and timing in the assignment of treatment. First, individuals or communities—potentially in anticipation of the assignment of treatment—make decisions and act, potentially altering their probability of receiving treatment. Second, there is a stochastic shock due to “nature,” reflecting that the units have incomplete control over the assignment variable. And finally, the treatment (or the intention to treat) is assigned on the basis of the assignment variable.

We have focused on this structure because in practice most RD analyses can be viewed along these lines, and also because of the similarity to the structure of a randomized experiment. That is, subjects of a randomized experiment may or may not make decisions in anticipation to participating in a randomized controlled trial (although their actions will ultimately have no influence on the probability of receiving treatment). Then the stochastic shock is realized (the randomization). Finally, the treatment is administered to one of the groups.

A number of the studies we surveyed though, did not seem to fit the spirit or essence of a randomized experiment. Since it is difficult to think of the treatment as being locally randomized in these cases, we will refer to the two research designs we identified in this category as “nonrandomized” discontinuity designs.

#### 6.3.1 Discontinuities in Age with Inevitable Treatment

Sometimes program status is turned on when an individual reaches a certain

age. Receipt of pension benefits is typically tied to reaching a particular age (see Eric V. Edmonds 2004; Edmonds, Kristin Mammen, and Miller 2005) and, in the United States, eligibility for the Medicare program begins at age 65 (see Card, Dobkin, and Maestas 2008) and young adults reach the legal drinking age at 21 (see Christopher Carpenter and Dobkin 2009). Similarly, one is subject to the less punitive juvenile justice system until the age of majority (typically, 18) (see Lee and McCrary 2005).

These cases stand apart from the typical RD designs discussed above because here assignment to treatment is essentially inevitable, as all subjects will eventually age into the program (or, conversely, age out of the program). One cannot, therefore, draw any parallels with a randomized experiment, which necessarily involves some *ex ante* uncertainty about whether a unit ultimately receives treatment (or the intent to treat).

Another important difference is that the tests of smoothness in baseline characteristics will generally be uninformative. Indeed, if one follows a single cohort over time, all characteristics determined prior to reaching the relevant age threshold are *by construction* identical just before and after the cutoff.<sup>48</sup> Note that in this case, *time* is the assignment variable, and therefore cannot be manipulated.

This design and the standard RD share the necessity of interpreting the discontinuity as the combined effect of *all* factors that switch on at the threshold. In the example of Thistlethwaite and Campbell (1960), if passing a scholarship exam provides the symbolic

<sup>48</sup> There are exceptions to this. There could be attrition over time, so that in principle, the number of observations could discontinuously drop at the threshold, changing the composition of the remaining observations. Alternatively, when examining a cross-section of different birth cohorts at a given point in time, it is possible to have sharp changes in the characteristics of individuals with respect to birthdate.

honor of passing the exam *as well as* a monetary award, the true treatment is a package of the two components, and one cannot attribute any effect to only one of the two. Similarly, when considering an age-activated treatment, one must consider the possibility that the age of interest is causing eligibility for potentially many other programs, which could affect the outcome.

There are at least two new issues that are irrelevant for the standard RD but are important for the analysis of age discontinuities. First, even if there is truly an effect on the outcome, if the effect is not immediate, it generally will not generate a discontinuity in the outcome. For example, suppose the receipt of Social Security benefits has no immediate impact but does have a long-run impact on labor force participation. Examining the labor force behavior as a function of age will not yield a discontinuity at age 67 (the full retirement age for those born after 1960), even though there may be a long-run effect. It is infeasible to estimate long-run effects because by the time we examine outcomes five years after receiving the treatment, for example, those individuals who were initially just below and just above age 67 will be exposed to essentially the same length of time of treatment (e.g., five years).<sup>49</sup>

The second important issue is that, because treatment is inevitable with the passage of time, individuals may fully anticipate the change in the regime and, therefore, may behave in certain ways prior to the time when treatment is turned on. Optimizing behavior in anticipation of a sharp regime change may either accentuate or mute observed effects. For example, simple life-cycle theories, assuming no liquidity constraints, suggest that the path of consumption will exhibit

no discontinuity at age 67, when Social Security benefits commence payment. On the other hand, some medical procedures are too expensive for an under-65-year-old but would be covered under Medicare upon turning 65. In this case, individuals' greater awareness of such a predicament will tend to *increase* the size of the discontinuity in utilization of medical procedures with respect to age (e.g., see Card, Dobkin, and Maestas 2008).

At this time we are unable to provide any more specific guidelines for analyzing these age/time discontinuities since it seems that how one models expectations, information, and behavior in anticipation of sharp changes in regimes will be highly context-dependent. But it does seem important to recognize these designs as being distinct from the standard RD design.

We conclude by emphasizing that when distinguishing between age-triggered treatments and a standard RD design, the involvement of age as an assignment variable is not as important as whether the receipt of treatment—or analogously, entering the control state—is inevitable. For example, on the surface, the analysis of the Medicaid expansions in Card and Lara D. Shore-Sheppard (2004) appears to be an age-based discontinuity since, effective July 1991, U.S. law requires states to cover children born after September 30, 1983, implying a discontinuous relationship between coverage and age, where the discontinuity in July 1991 was around 8 years of age. This design, however, actually fits quite easily into the standard RD framework we have discussed throughout this paper.

First, note that treatment receipt is *not* inevitable for those individuals born near the September 30, 1983, threshold. Those born strictly after that date were covered from July 1991 until their 18th birthday, while those born on or before the date received no such coverage. Second, the data generating process does follow the structure discussed

<sup>49</sup> By contrast, there is no such limitation with the standard RD design. One can examine outcomes defined at an arbitrarily long time period after the assignment to treatment.

above. Parents do have some influence regarding when their children are born, but with only imprecise control over the exact date (and at any rate, it seems implausible that parents would have anticipated that such a Medicaid expansion would have occurred eight years in the future, with the particular birthdate cutoff chosen). Thus the treatment is assigned based on the assignment variable, which is the birthdate in this context.

Examples of other age-based discontinuities where neither the treatment nor control state is guaranteed with the passage of time that can also be viewed within the standard RD framework include studies by Cascio and Lewis (2006), McCrary and Royer (2003), Dobkin and Ferreira (2009), and Phillip Oreopoulos (2006).

### 6.3.2 Discontinuities in Geography

Another “nonrandomized” RD design is one involving the location of residences, where the discontinuity threshold is a boundary that demarcates regions. Black (1999) and Patrick Bayer, Ferreira, and Robert McMillan (2007) examine housing prices on either side of school attendance boundaries to estimate the implicit valuation of different schools. Lavy (2006) examines adjacent neighborhoods in different cities, and therefore subject to different rules regarding student busing. Rafael Lalive (2008) compares unemployment duration in regions in Austria receiving extended benefits to adjacent control regions. Karen M. Pence (2006) examines census tracts along state borders to examine the impact of more borrower-friendly laws on mortgage loan sizes.

In each of these cases, it is awkward to view either houses or families as locally randomly assigned. Indeed this is a case where economic agents have quite precise control over where to place a house or where to live. The location of houses will be planned in response to geographic features (rivers, lakes, hills) and in conjunction with the planning of

streets, parks, commercial development, etc. In order for this to resemble a more standard RD design, one would have to imagine the relevant boundaries being set in a “random” way, so that it would be simply luck determining whether a house ended up on either side of the boundary. The concern over the endogeneity of boundaries is clearly recognized by Black (1999), who “. . . [b]ecause of concerns about neighborhood differences on opposite sides of an attendance district boundary, . . . was careful to omit boundaries from [her] sample if the two attendance districts were divided in ways that seemed to clearly divide neighborhoods; attendance districts divided by large rivers, parks, golf courses, or any large stretch of land were excluded.” As one could imagine, the selection of which boundaries to include could quickly turn into more of an art than a science.

We have no uniform advice on how to analyze geographic discontinuities because it seems that the best approach would be particularly context-specific. It does, however, seem prudent for the analyst, in assessing the internal validity of the research design, to carefully consider three sets of questions. First, what is the process that led to the location of the boundaries? Which came first: the houses or the boundaries? Were the boundaries a response to some preexisting geographical or political constraint? Second, how might sorting of families or the endogenous location of houses affect the analysis? And third, what are all the things differing between the two regions *other than the treatment of interest*? An exemplary analysis and discussion of these latter two issues in the context of school attendance zones is found in Bayer, Ferreira, and McMillan (2007).

## 7. Concluding Remarks on RD Designs in Economics: Progress and Prospects

Our reading of the existing and active literature is that—after being largely ignored

by economists for almost forty years—there have been significant inroads made in understanding the properties, limitations, interpretability, and perhaps most importantly, in the useful application of RD designs to a wide variety of empirical questions in economics. These developments have, for the most part, occurred within a short period of time, beginning in the late 1990s.

Here we highlight what we believe are the most significant recent contributions of the economics literature to the understanding and application of RD designs. We believe these are helpful developments in guiding applied researchers who seek to implement RD designs, and we also illustrate them with a few examples from the literature.

- **Sorting and Manipulation of the Assignment Variable:** Economists consider how self-interested individuals or optimizing organizations may behave in response to rules that allocate resources. It is therefore unsurprising that the discussion of how endogenous sorting around the discontinuity threshold *can invalidate* the RD design has been found (to our knowledge, exclusively) in the economics literature. By contrast, textbook treatments outside economics on RD do not discuss this sorting or manipulation, and give the impression that the knowledge of the assignment rule is sufficient for the validity of the RD.<sup>50</sup>

<sup>50</sup> For example, Trochim (1984) characterizes the three central assumptions of the RD design as: (1) perfect adherence to the cutoff rule, (2) having the correct functional form, and (3) no other factors (other than the program of interest) cause the discontinuity. More recently, William R. Shadish, Cook, and Campbell (2002) claim on page 243 that the proof of the unbiasedness of RD primarily follows from the fact that treatment is known perfectly once the assignment variable is known. They go on to argue that this deterministic rule implies omitted variables will not pose a problem. But Hahn, Todd, and van der Klaauw (2001)

We believe a “state-of-the-art” RD analysis today will consider carefully the possibility of endogenous sorting. A recent analysis that illustrates this standard is that of Urquiola and Verhoogen (2009), who examine the class size cap RD design pioneered by Angrist and Lavy (1999) in the context of Chile’s highly liberalized market for primary schools. In a certain segment of the private market, schools receive a fixed payment per student from the government. However, each school faces a very high marginal cost (hiring one extra teacher) for crossing a multiple of the class size cap. Perhaps unsurprisingly, they find striking discontinuities in the *histogram* of the assignment variable (total enrollment in the grade), with an undeniable “stacking” of schools at the relevant class size cap cutoffs. They also provide evidence that those families in schools just to the left and right of the thresholds are systematically different in family income, suggesting some degree of sorting. For this reason, they conclude that an RD analysis in this particular context is most likely inappropriate.<sup>51</sup> This study, as well as the analysis of Bayer, Ferreira, and McMillan (2007) reflects a heightened awareness of a sorting issue recognized since the beginning of the recent wave of RD applications in economics.<sup>52</sup> From a practitioner’s perspective, an important recent development

make it clear that the existence of a deterministic rule for the assignment of treatment is *not* sufficient for unbiasedness, and it is necessary to *assume* the influence of all other factors (omitted variables) are the same on either side of the discontinuity threshold (i.e., their continuity assumption).

<sup>51</sup> Urquiola and Verhoogen (2009) emphasize the sorting issues may well be specific to the liberalized nature of the Chilean primary school market, and that they may or may not be present in other countries.

<sup>52</sup> See, for example, footnote 23 in van der Klaauw (1997) and page 549 in Angrist and Lavy (1999).

is the notion that we can empirically examine the degree of sorting, and one way of doing so is suggested in McCrary (2008).

- **RD Designs as Locally Randomized Experiments:** Economists are hesitant to apply methods that have not been rigorously formalized within an econometric framework, and where crucial identifying assumptions have not been clearly specified. This is perhaps one of the reasons why RD designs were underutilized by economists for so long, since it is only relatively recently that the underlying assumptions needed for the RD were formalized.<sup>53</sup> In the recent literature, RD designs were initially viewed as a special case of matching (Heckman, Lalonde, and Smith 1999), or alternatively as a special case of IV (Angrist and Krueger 1999), and these perspectives may have provided empirical researchers a familiar econometric framework within which identifying assumptions could be more carefully discussed.

Today, RD is increasingly recognized in applied research as a distinct design that is a close relative to a randomized experiment. Formally shown in Lee (2008), even when individuals have

some control over the assignment variable, as long as this control is imprecise—that is, the *ex ante* density of the assignment variable is continuous—the consequence will be local randomization of the treatment. So in a number of nonexperimental contexts where resources are allocated based on a sharp cutoff rule, there may indeed be a hidden randomized experiment to utilize. And furthermore, as in a randomized experiment, this implies that all observable baseline covariates will locally have the same distribution on either side of the discontinuity threshold—an empirically testable proposition.

We view the testing of the continuity of the baseline covariates as an important part of assessing the validity of any RD design—particularly in light of the incentives that can potentially generate sorting—and as something that truly sets RD apart from other evaluation strategies. Examples of this kind of testing of the RD design include Jordan D. Matsudaira (2008), Card, Raj Chetty and Andrea Weber (2007), DiNardo and Lee (2004), Lee, Moretti and Butler (2004), McCrary and Royer (2003), Greenstone and Gallagher (2008), and Urquiola and Verhoogen (2009).

<sup>53</sup> An example of how economists/econometricians' notion of a proof differs from that in other disciplines is found in Cook (2008), who views the discussion in Arthur S. Goldberger (1972a) and Goldberger (1972b) as the first "proof of the basic design," quoting the following passage in Goldberger (1972a) (brackets from Cook 2008): "The explanation for this serendipitous result [no bias when selection is on an observed pretest score] is not hard to locate. Recall that  $z$  [a binary variable representing the treatment contrast at the cutoff] is completely determined by pretest score  $x$  [an obtained ability score]. It cannot contain any information about  $x^*$  [true ability] that is not contained within  $x$ . Consequently, when we control on  $x$  as in the multiple regression,  $z$  has no explanatory power with respect to  $y$  [the outcome measured with error]. More formally, the partial correlation of  $y$  and  $z$  controlling on  $x$  vanishes although the simple correlation of  $y$  and  $z$  is

nonzero" (p. 647). After reading the article, an econometrician will recognize the discussion above not as a proof of the validity of the RD, but rather as a restatement of the consequence of  $z$  being an indicator variable determined by an observed variable  $x$ , in a specific parameterized example. Today we know the existence of such a rule is *not sufficient* for a valid RD design, and a crucial necessary assumption is the continuity of the influence of all other factors, as shown in Hahn, Todd, and van der Klaauw (2001). In Goldberger (1972a), the role of the continuity of omitted factors was not mentioned (although it is implicitly assumed in the stylized model of test scores involving normally distributed and independent errors). Indeed, apparently Goldberger himself later clarified that he did not set out to propose the RD design, and was instead interested in the issues related to selection on observables and unobservables (Cook 2008).

- **Graphical Analysis and Presentation:**

The graphical presentation of an RD analysis is not a contribution of economists,<sup>54</sup> but it is safe to say that the body of work produced by economists has led to a kind of “industry standard” that the transparent identification strategy of the RD be accompanied by an equally transparent graph showing the empirical relation between the outcome and the assignment variable. Graphical presentations of RD are so prevalent in applied research, it is tempting to guess that studies not including the graphical evidence are ones where the graphs are not compelling or well-behaved.

In an RD analysis, the graph is indispensable because it can summarize a great deal of information in one picture. It can give a rough sense of the range of the both the assignment variable and the outcome variable as well as the overall shape of the relationship between the two, thus indicating what functional forms are likely to make sense. It can also alert the researcher to potential outliers in both the assignment and outcome variables. A graph of the raw means—in nonoverlapping intervals, as discussed in section 4.1—also gives a rough sense of the likely sampling variability of the RD gap estimate itself, since one can compare the size of the jump at the discontinuity to natural “bumpiness” in the graph away from the discontinuity. Our reading of the literature is that the most informative graphs are ones that simultaneously allow the raw data “to speak for themselves” in revealing a discontinuity if there is one, yet at the same time treat data near the threshold the same as data away from the

threshold.<sup>55</sup> There are many examples that follow this general principle; recent ones include Matsudaira (2008), Card, Chetty and Weber (2007), Card, Dobkin, and Maestas (2009), McCrary and Royer (2003), Lee (2008), and Ferreira and Gyourko (2009).

- **Applicability:** Soon after the introduction of RD, in a chapter in a book on research methods, Campbell and Julian C. Stanley (1963) wrote that the RD design was “very limited in range of possible applications.” The emerging body of research produced by economists in recent years has proven quite the opposite. Our survey of the literature suggests that there are many kinds of discontinuous rules that can help answer important questions in economics and related areas. Indeed, one may go so far as to guess that whenever a scarce resource is rationed for individual entities, if the political climate demands a transparent way of distributing that resource, it is a good bet there is an RD design lurking in the background. In addition, it seems that the approach of using changes in laws that disqualify older birth cohorts based on their date of birth (as in Card and Shore-Sheppard (2004) or Oreopoulos (2006)) may well have much wider applicability.

One way to understand both the applicability and limitations of the RD design is to recognize its relation to a standard econometric policy evaluation framework, where the main variable of interest is a potentially endogenous binary treatment variable (as considered in Heckman 1978 or more recently discussed in Heckman and Vytlačil

<sup>54</sup> Indeed the original article of Thistlethwaite and Campbell (1960) included a graphical analysis of the data.

<sup>55</sup> For example, graphing a smooth conditional expectation function everywhere *except* at the discontinuity threshold violates this principle.

2005). This selection model applies to a great deal of economic problems. As we pointed out in section 3, the RD design describes a situation where you are able to *observe* the latent variable that determines treatment. As long as the density of that variable is continuous for each individual, the benefit of observing the latent index is that one neither needs to make exclusion restrictions nor assume any variable (i.e., an instrument) is independent of errors in the outcome equation.

From this perspective, for the class of problems that fit into the standard treatment evaluation problem, RD designs can be seen as a subset since there is an institutional, index-based rule playing a role in determining treatment. Among this subset, the binding constraint of RD lies in obtaining the necessary data: readily available public-use household survey data, for example, will often only contain variables that are correlated with the true assignment variable (e.g., reported income in a survey, as opposed to the income used for allocation of benefits), or are measured too coarsely (e.g., years rather than months or weeks) to detect a discontinuity in the presence of a regression function with significant curvature. This is where there can be a significant payoff to investing in securing high quality data, which is evident in most of the studies listed in table 5.

### 7.1 Extensions

We conclude by discussing two natural directions in which the RD approach can be extended. First, we have discussed the “fuzzy” RD design as an important departure from the “classic” RD design where treatment is a deterministic function of the assignment variable, but there are other departures that could be practically relevant but not as well understood. For example, even if there

is perfect compliance of the discontinuous rule, it may be that the researcher does not directly observe the assignment variable, but instead possesses a slightly noisy measure of the variable. Understanding the effects of this kind of measurement error could further expand the applicability of RD. In addition, there may be situations where the researcher both suspects and statistically detects some degree of precise sorting around the threshold, but that the sorting may appear to be relatively minor, even if statistically significant (based on observing discontinuities in baseline characteristics). The challenge, then, is to specify under what conditions one can correct for small amounts of this kind of contamination.

Second, so far we have discussed the sorting or manipulation issue as a potential problem or nuisance to the general program evaluation problem. But there is another way of viewing this sorting issue. The observed sorting may well be evidence of economic agents responding to incentives, and may help identify economically interesting phenomena. That is, economic behavior may be what is driving discontinuities in the frequency distribution of grade enrollment (as in Urquiola and Verhoogen 2009), or in the distribution of roll call votes (as in McCrary 2008), or in the distribution of age at offense (as in Lee and McCrary 2005), and those behavioral responses may be of interest.

These cases, as well as the age/time and boundary discontinuities discussed above, do not fit into the “standard” RD framework, but nevertheless can tell us something important about behavior, and further expand the kinds of questions that can be addressed by exploiting discontinuous rules to identify meaningful economic parameters of interest.

### REFERENCES

- Albouy, David. 2008. “Do Voters Affect or Elect Policies? A New Perspective With Evidence from the U.S. Senate.” Unpublished.

- Albouy, David. 2009. "Partisan Representation in Congress and the Geographic Distribution of Federal Funds." National Bureau of Economic Research Working Paper 15224.
- Angrist, Joshua D. 1990. "Lifetime Earnings and the Vietnam Era Draft Lottery: Evidence from Social Security Administrative Records." *American Economic Review*, 80(3): 313–36.
- Angrist, Joshua D., and Alan B. Krueger. 1999. "Empirical Strategies in Labor Economics." In *Handbook of Labor Economics*, Volume 3A, ed. Orley Ashenfelter and David Card, 1277–1366. Amsterdam; New York and Oxford: Elsevier Science, North-Holland.
- Angrist, Joshua D., and Victor Lavy. 1999. "Using Maimonides' Rule to Estimate the Effect of Class Size on Scholastic Achievement." *Quarterly Journal of Economics*, 114(2): 533–75.
- Asadullah, M. Niaz. 2005. "The Effect of Class Size on Student Achievement: Evidence from Bangladesh." *Applied Economics Letters*, 12(4): 217–21.
- Battistin, Erich, and Enrico Rettore. 2002. "Testing for Programme Effects in a Regression Discontinuity Design with Imperfect Compliance." *Journal of the Royal Statistical Society: Series A (Statistics in Society)*, 165(1): 39–57.
- Battistin, Erich, and Enrico Rettore. 2008. "Ineligibles and Eligible Non-participants as a Double Comparison Group in Regression-Discontinuity Designs." *Journal of Econometrics*, 142(2): 715–30.
- Baum-Snow, Nathaniel, and Justin Marion. 2009. "The Effects of Low Income Housing Tax Credit Developments on Neighborhoods." *Journal of Public Economics*, 93(5–6): 654–66.
- Bayer, Patrick, Fernando Ferreira, and Robert McMillan. 2007. "A Unified Framework for Measuring Preferences for Schools and Neighborhoods." *Journal of Political Economy*, 115(4): 588–638.
- Behaghel, Luc, Bruno Crépon, and Béatrice Sédillot. 2008. "The Perverse Effects of Partial Employment Protection Reform: The Case of French Older Workers." *Journal of Public Economics*, 92(3–4): 696–721.
- Berk, Richard A., and Jan de Leeuw. 1999. "An Evaluation of California's Inmate Classification System Using a Generalized Regression Discontinuity Design." *Journal of the American Statistical Association*, 94(448): 1045–52.
- Berk, Richard A., and David Rauma. 1983. "Capitalizing on Nonrandom Assignment to Treatments: A Regression-Discontinuity Evaluation of a Crime-Control Program." *Journal of the American Statistical Association*, 78(381): 21–27.
- Black, Dan A., Jose Galdo, and Jeffrey A. Smith. 2007a. "Evaluating the Regression Discontinuity Design Using Experimental Data." Unpublished.
- Black, Dan A., Jose Galdo, and Jeffrey A. Smith. 2007b. "Evaluating the Worker Profiling and Reemployment Services System Using a Regression Discontinuity Approach." *American Economic Review*, 97(2): 104–07.
- Black, Dan A., Jeffrey A. Smith, Mark C. Berger, and Brett J. Noel. 2003. "Is the Threat of Reemployment Services More Effective Than the Services Themselves? Evidence from Random Assignment in the UI System." *American Economic Review*, 93(4): 1313–27.
- Black, Sandra E. 1999. "Do Better Schools Matter? Parental Valuation of Elementary Education." *Quarterly Journal of Economics*, 114(2): 577–99.
- Blundell, Richard, and Alan Duncan. 1998. "Kernel Regression in Empirical Microeconomics." *Journal of Human Resources*, 33(1): 62–87.
- Buddelmeyer, Hielke, and Emmanuel Skoufias. 2004. "An Evaluation of the Performance of Regression Discontinuity Design on PROGRESA." World Bank Policy Research Working Paper 3386.
- Buettner, Thiess. 2006. "The Incentive Effect of Fiscal Equalization Transfers on Tax Policy." *Journal of Public Economics*, 90(3): 477–97.
- Campbell, Donald T., and Julian C. Stanley. 1963. "Experimental and Quasi-experimental Designs for Research on Teaching." In *Handbook of Research on Teaching*, ed. N. L. Gage, 171–246. Chicago: Rand McNally.
- Canton, Erik, and Andreas Blom. 2004. "Can Student Loans Improve Accessibility to Higher Education and Student Performance? An Impact Study of the Case of SOFES, Mexico." World Bank Policy Research Working Paper 3425.
- Card, David, Raj Chetty, and Andrea Weber. 2007. "Cash-on-Hand and Competing Models of Intertemporal Behavior: New Evidence from the Labor Market." *Quarterly Journal of Economics*, 122(4): 1511–60.
- Card, David, Carlos Dobkin, and Nicole Maestas. 2008. "The Impact of Nearly Universal Insurance Coverage on Health Care Utilization: Evidence from Medicare." *American Economic Review*, 98(5): 2242–58.
- Card, David, Carlos Dobkin, and Nicole Maestas. 2009. "Does Medicare Save Lives?" *Quarterly Journal of Economics*, 124(2): 597–636.
- Card, David, Alexandre Mas, and Jesse Rothstein. 2008. "Tipping and the Dynamics of Segregation." *Quarterly Journal of Economics*, 123(1): 177–218.
- Card, David, and Lara D. Shore-Sheppard. 2004. "Using Discontinuous Eligibility Rules to Identify the Effects of the Federal Medicaid Expansions on Low-Income Children." *Review of Economics and Statistics*, 86(3): 752–66.
- Carpenter, Christopher, and Carlos Dobkin. 2009. "The Effect of Alcohol Consumption on Mortality: Regression Discontinuity Evidence from the Minimum Drinking Age." *American Economic Journal: Applied Economics*, 1(1): 164–82.
- Cascio, Elizabeth U., and Ethan G. Lewis. 2006. "Schooling and the Armed Forces Qualifying Test: Evidence from School-Entry Laws." *Journal of Human Resources*, 41(2): 294–318.
- Chay, Kenneth Y., and Michael Greenstone. 2003. "Air Quality, Infant Mortality, and the Clean Air Act of 1970." National Bureau of Economic Research Working Paper 10053.
- Chay, Kenneth Y., and Michael Greenstone. 2005.

- "Does Air Quality Matter? Evidence from the Housing Market." *Journal of Political Economy*, 113(2): 376–424.
- Chay, Kenneth Y., Patrick J. McEwan, and Miguel Urquiola. 2005. "The Central Role of Noise in Evaluating Interventions That Use Test Scores to Rank Schools." *American Economic Review*, 95(4): 1237–58.
- Chen, M. Keith, and Jesse M. Shapiro. 2004. "Does Prison Harden Inmates? A Discontinuity-Based Approach." Yale University Cowles Foundation Discussion Paper 1450.
- Chen, Susan, and Wilbert van der Klaauw. 2008. "The Work Disincentive Effects of the Disability Insurance Program in the 1990s." *Journal of Econometrics*, 142(2): 757–84.
- Chiang, Hanley. 2009. "How Accountability Pressure on Failing Schools Affects Student Achievement." *Journal of Public Economics*, 93(9–10): 1045–57.
- Clark, Damon. 2009. "The Performance and Competitive Effects of School Autonomy." *Journal of Political Economy*, 117(4): 745–83.
- Cole, Shawn. 2009. "Financial Development, Bank Ownership, and Growth: Or, Does Quantity Imply Quality?" *Review of Economics and Statistics*, 91(1): 33–51.
- Cook, Thomas D. 2008. "Waiting for Life to Arrive: A History of the Regression-Discontinuity Design in Psychology, Statistics and Economics." *Journal of Econometrics*, 142(2): 636–54.
- Davis, Lucas W. 2008. "The Effect of Driving Restrictions on Air Quality in Mexico City." *Journal of Political Economy*, 116(1): 38–81.
- De Giorgi, Giacomo. 2005. "Long-Term Effects of a Mandatory Multistage Program: The New Deal for Young People in the UK." Institute for Fiscal Studies Working Paper 05/08.
- DesJardins, Stephen L., and Brian P. McCall. 2008. "The Impact of the Gates Millennium Scholars Program on the Retention, College Finance- and Work-Related Choices, and Future Educational Aspirations of Low-Income Minority Students." Unpublished.
- DiNardo, John, and David S. Lee. 2004. "Economic Impacts of New Unionization on Private Sector Employers: 1984–2001." *Quarterly Journal of Economics*, 119(4): 1383–1441.
- Ding, Weili, and Steven F. Lehrer. 2007. "Do Peers Affect Student Achievement in China's Secondary Schools?" *Review of Economics and Statistics*, 89(2): 300–312.
- Dobkin, Carlos, and Fernando Ferreira. 2009. "Do School Entry Laws Affect Educational Attainment and Labor Market Outcomes?" National Bureau of Economic Research Working Paper 14945.
- Edmonds, Eric V. 2004. "Does Illiquidity Alter Child Labor and Schooling Decisions? Evidence from Household Responses to Anticipated Cash Transfers in South Africa." National Bureau of Economic Research Working Paper 10265.
- Edmonds, Eric V., Kristin Mammen, and Douglas L. Miller. 2005. "Rearranging the Family? Income Support and Elderly Living Arrangements in a Low-Income Country." *Journal of Human Resources*, 40(1): 186–207.
- Fan, Jianqing, and Irene Gijbels. 1996. *Local Polynomial Modelling and Its Applications*. London; New York and Melbourne: Chapman and Hall.
- Ferreira, Fernando. Forthcoming. "You Can Take It With You: Proposition 13 Tax Benefits, Residential Mobility, and Willingness to Pay for Housing Amenities." *Journal of Public Economics*.
- Ferreira, Fernando, and Joseph Gyourko. 2009. "Do Political Parties Matter? Evidence from U.S. Cities." *Quarterly Journal of Economics*, 124(1): 399–422.
- Figlio, David N., and Lawrence W. Kenny. 2009. "Public Sector Performance Measurement and Stakeholder Support." *Journal of Public Economics*, 93(9–10): 1069–77.
- Goldberger, Arthur S. 1972a. "Selection Bias in Evaluating Treatment Effects: Some Formal Illustrations." Unpublished.
- Goldberger, Arthur S. 1972b. *Selection Bias in Evaluating Treatment Effects: The Case of Interaction*. Unpublished.
- Goodman, Joshua. 2008. "Who Merits Financial Aid?: Massachusetts' Adams Scholarship." *Journal of Public Economics*, 92(10–11): 2121–31.
- Goolsbee, Austan, and Jonathan Guryan. 2006. "The Impact of Internet Subsidies in Public Schools." *Review of Economics and Statistics*, 88(2): 336–47.
- Greenstone, Michael, and Justin Gallagher. 2008. "Does Hazardous Waste Matter? Evidence from the Housing Market and the Superfund Program." *Quarterly Journal of Economics*, 123(3): 951–1003.
- Guryan, Jonathan. 2001. "Does Money Matter? Regression-Discontinuity Estimates from Education Finance Reform in Massachusetts." National Bureau of Economic Research Working Paper 8269.
- Hahn, Jinyong. 1998. "On the Role of the Propensity Score in Efficient Semiparametric Estimation of Average Treatment Effects." *Econometrica*, 66(2): 315–31.
- Hahn, Jinyong, Petra Todd, and Wilbert van der Klaauw. 1999. "Evaluating the Effect of an Antidiscrimination Law Using a Regression-Discontinuity Design." National Bureau of Economic Research Working Paper 7131.
- Hahn, Jinyong, Petra Todd, and Wilbert van der Klaauw. 2001. "Identification and Estimation of Treatment Effects with a Regression-Discontinuity Design." *Econometrica*, 69(1): 201–09.
- Heckman, James J. 1978. "Dummy Endogenous Variables in a Simultaneous Equation System." *Econometrica*, 46(4): 931–59.
- Heckman, James J., Robert J. Lalonde, and Jeffrey A. Smith. 1999. "The Economics and Econometrics of Active Labor Market Programs." In *Handbook of Labor Economics*, Volume 3A, ed. Orley Ashenfelter and David Card, 1865–2097. Amsterdam; New York and Oxford: Elsevier Science, North-Holland.
- Heckman, James J., and Edward Vytlacil. 2005. "Structural Equations, Treatment Effects, and Econometric

- Policy Evaluation." *Econometrica*, 73(3): 669–738.
- Hjalmarsson, Randi. 2009. "Juvenile Jails: A Path to the Straight and Narrow or to Hardened Criminality?" *Journal of Law and Economics*, 52(4): 779–809.
- Horowitz, Joel L., and Charles F. Manski. 2000. "Non-parametric Analysis of Randomized Experiments with Missing Covariate and Outcome Data." *Journal of the American Statistical Association*, 95(449): 77–84.
- Hoxby, Caroline M. 2000. "The Effects of Class Size on Student Achievement: New Evidence from Population Variation." *Quarterly Journal of Economics*, 115(4): 1239–85.
- Imbens, Guido W., and Joshua D. Angrist. 1994. "Identification and Estimation of Local Average Treatment Effects." *Econometrica*, 62(2): 467–75.
- Imbens, Guido W., and Karthik Kalyanaraman. 2009. "Optimal Bandwidth Choice for the Regression Discontinuity Estimator." National Bureau of Economic Research Working Paper 14726.
- Imbens, Guido W., and Thomas Lemieux. 2008. "Regression Discontinuity Designs: A Guide to Practice." *Journal of Econometrics*, 142(2): 615–35.
- Jacob, Brian A., and Lars Lefgren. 2004a. "The Impact of Teacher Training on Student Achievement: Quasi-experimental Evidence from School Reform Efforts in Chicago." *Journal of Human Resources*, 39(1): 50–79.
- Jacob, Brian A., and Lars Lefgren. 2004b. "Remedial Education and Student Achievement: A Regression-Discontinuity Analysis." *Review of Economics and Statistics*, 86(1): 226–44.
- Kane, Thomas J. 2003. "A Quasi-experimental Estimate of the Impact of Financial Aid on College-Going." National Bureau of Economic Research Working Paper 9703.
- Lalive, Rafael. 2007. "Unemployment Benefits, Unemployment Duration, and Post-unemployment Jobs: A Regression Discontinuity Approach." *American Economic Review*, 97(2): 108–12.
- Lalive, Rafael. 2008. "How Do Extended Benefits Affect Unemployment Duration? A Regression Discontinuity Approach." *Journal of Econometrics*, 142(2): 785–806.
- Lalive, Rafael, Jan C. van Ours, and Josef Zweimüller. 2006. "How Changes in Financial Incentives Affect the Duration of Unemployment." *Review of Economic Studies*, 73(4): 1009–38.
- Lavy, Victor. 2002. "Evaluating the Effect of Teachers' Group Performance Incentives on Pupil Achievement." *Journal of Political Economy*, 110(6): 1286–1317.
- Lavy, Victor. 2004. "Performance Pay and Teachers' Effort, Productivity and Grading Ethics." National Bureau of Economic Research Working Paper 10622.
- Lavy, Victor. 2006. "From Forced Busing to Free Choice in Public Schools: Quasi-Experimental Evidence of Individual and General Effects." National Bureau of Economic Research Working Paper 11969.
- Lee, David S. 2001. "The Electoral Advantage to Incumbency and Voters' Valuation of Politicians' Experience: A Regression Discontinuity Analysis of Close Elections." University of California Berkeley Center for Labor Economics Working Paper 31.
- Lee, David S. 2008. "Randomized Experiments from Non-random Selection in U.S. House Elections." *Journal of Econometrics*, 142(2): 675–97.
- Lee, David S. 2009. "Training, Wages, and Sample Selection: Estimating Sharp Bounds on Treatment Effects." *Review of Economic Studies*, 76(3): 1071–1102.
- Lee, David S., and David Card. 2008. "Regression Discontinuity Inference with Specification Error." *Journal of Econometrics*, 142(2): 655–74.
- Lee, David S., and Justin McCrary. 2005. "Crime, Punishment, and Myopia." National Bureau of Economic Research Working Paper 11491.
- Lee, David S., Enrico Moretti, and Matthew J. Butler. 2004. "Do Voters Affect or Elect Policies? Evidence from the U.S. House." *Quarterly Journal of Economics*, 119(3): 807–59.
- Lemieux, Thomas, and Kevin Milligan. 2008. "Incentive Effects of Social Assistance: A Regression Discontinuity Approach." *Journal of Econometrics*, 142(2): 807–28.
- Leuven, Edwin, Mikael Lindahl, Hessel Oosterbeek, and Dinand Webbink. 2007. "The Effect of Extra Funding for Disadvantaged Pupils on Achievement." *Review of Economics and Statistics*, 89(4): 721–36.
- Leuven, Edwin, and Hessel Oosterbeek. 2004. "Evaluating the Effect of Tax Deductions on Training." *Journal of Labor Economics*, 22(2): 461–88.
- Ludwig, Jens, and Douglas L. Miller. 2007. "Does Head Start Improve Children's Life Chances? Evidence from a Regression Discontinuity Design." *Quarterly Journal of Economics*, 122(1): 159–208.
- Matsudaira, Jordan D. 2008. "Mandatory Summer School and Student Achievement." *Journal of Econometrics*, 142(2): 829–50.
- McCrary, Justin. 2008. "Manipulation of the Running Variable in the Regression Discontinuity Design: A Density Test." *Journal of Econometrics*, 142(2): 698–714.
- McCrary, Justin, and Heather Royer. 2003. "Does Maternal Education Affect Infant Health? A Regression Discontinuity Approach Based on School Age Entry Laws." Unpublished.
- Newey, Whitney K., and Daniel L. McFadden. 1994. "Large Sample Estimation and Hypothesis Testing." In *Handbook of Econometrics, Volume 4*, ed. Robert F. Engle and Daniel L. McFadden, 2111–2245. Amsterdam; London and New York: Elsevier, North-Holland.
- Oreopoulos, Philip. 2006. "Estimating Average and Local Average Treatment Effects of Education When Compulsory Schooling Laws Really Matter." *American Economic Review*, 96(1): 152–75.
- Pence, Karen M. 2006. "Foreclosing on Opportunity: State Laws and Mortgage Credit." *Review of Economics and Statistics*, 88(1): 177–82.
- Pettersson, Per. 2000. "Do Parties Matter for Fiscal Policy Choices?" *Econometric Society World*

- Congress 2000 Contributed Paper 1373.
- Petterson-Lidbom, Per. 2008a. "Does the Size of the Legislature Affect the Size of Government? Evidence from Two Natural Experiments." Unpublished.
- Petterson-Lidbom, Per. 2008b. "Do Parties Matter for Economic Outcomes? A Regression-Discontinuity Approach." *Journal of the European Economic Association*, 6(5): 1037–56.
- Pitt, Mark M., and Shahidur R. Khandker. 1998. "The Impact of Group-Based Credit Programs on Poor Households in Bangladesh: Does the Gender of Participants Matter?" *Journal of Political Economy*, 106(5): 958–96.
- Pitt, Mark M., Shahidur R. Khandker, Signe-Mary McKernan, and M. Abdul Latif. 1999. "Credit Programs for the Poor and Reproductive Behavior in Low-Income Countries: Are the Reported Causal Relationships the Result of Heterogeneity Bias?" *Demography*, 36(1): 1–21.
- Porter, Jack. 2003. "Estimation in the Regression Discontinuity Model." Unpublished.
- Powell, James L. 1994. "Estimation of Semiparametric Models." In *Handbook of Econometrics, Volume 4*, ed. Robert F. Engle and Daniel L. McFadden, 2443–2521. Amsterdam; London and New York: Elsevier, North-Holland.
- Shadish, William R., Thomas D. Cook, and Donald T. Campbell. 2002. *Experimental and Quasi-Experimental Designs for Generalized Causal Inference*. Boston: Houghton Mifflin.
- Silverman, Bernard W. 1986. *Density Estimation for Statistics and Data Analysis*. London and New York: Chapman and Hall.
- Snyder, Stephen E., and William N. Evans. 2006. "The Effect of Income on Mortality: Evidence from the Social Security Notch." *Review of Economics and Statistics*, 88(3): 482–95.
- Thistlethwaite, Donald L., and Donald T. Campbell. 1960. "Regression-Discontinuity Analysis: An Alternative to the Ex Post Facto Experiment." *Journal of Educational Psychology*, 51(6): 309–17.
- Trochim, William M. K. 1984. *Research Design for Program Evaluation: The Regression-Discontinuity Approach*. Beverly Hills: Sage Publications.
- Urquiola, Miguel. 2006. "Identifying Class Size Effects in Developing Countries: Evidence from Rural Bolivia." *Review of Economics and Statistics*, 88(1): 171–77.
- Urquiola, Miguel, and Eric A. Verhoogen. 2009. "Class-Size Caps, Sorting, and the Regression-Discontinuity Design." *American Economic Review*, 99(1): 179–215.
- van der Klaauw, Wilbert. 1997. "A Regression-Discontinuity Evaluation of the Effect of Financial Aid Offers on College Enrollment." New York University C.V. Starr Center for Applied Economics Working Paper 10.
- van der Klaauw, Wilbert. 2002. "Estimating the Effect of Financial Aid Offers on College Enrollment: A Regression-Discontinuity Approach." *International Economic Review*, 43(4): 1249–87.
- van der Klaauw, Wilbert. 2008a. "Breaking the Link between Poverty and Low Student Achievement: An Evaluation of Title I." *Journal of Econometrics*, 142(2): 731–56.
- van der Klaauw, Wilbert. 2008b. "Regression-Discontinuity Analysis: A Survey of Recent Developments in Economics." *Labour*, 22(2): 219–45.
- White, Halbert. 1980. "A Heteroskedasticity-Consistent Covariance Matrix Estimator and a Direct Test for Heteroskedasticity." *Econometrica*, 48(4): 817–38.