

# Lab #17 - Regression Discontinuity Basics

*Econ 224*

*November 6th, 2018*

## Introduction

1. Start with a varying slopes, varying intercepts model from Econ 103. Make sure they know how to simulate, fit, plot, and interpret
2. Then move on to a regression with a *break*, in other words there is one linear regression to the left of a cutoff and another to the right of that cutoff. Have them figure out how to write this as a *single* linear regression. Then simulate, fit, plot, and interpret.
3. Talk about the basic regression discontinuity idea: the only thing that jumps at the cutoff is treatment allocation.
4. Create an example where non-linearity masquerades as a discontinuity.
5. Have them repeat 2 but with a *quadratic* regression. Maybe talk about the non-parametric version?

## “Sharp” Regression Discontinuity

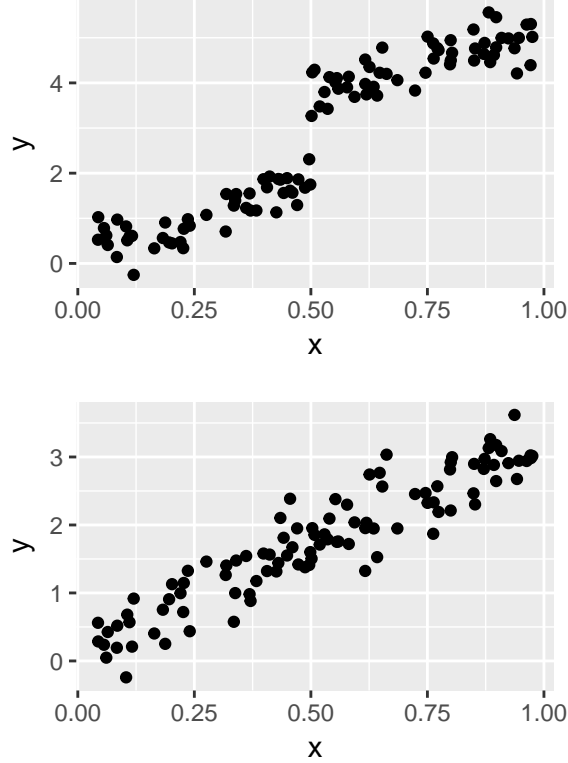
Suppose we are interested in learning the causal effect of a binary treatment  $D$  on an outcome  $Y$ . In some special settings, whether or not a person is treated is solely determined by a special covariate  $x$ , called the *running variable*

$$D_i = \begin{cases} 0 & \text{if } x_i < x_0 \\ 1 & \text{if } x_i \geq x_0 \end{cases}$$

The preceding expression says that  $D$  is a *deterministic function* of  $x$ : everyone who has  $x \geq x_0$  is treated, and no one who has  $x < x_0$  is treated. This setting is called a *sharp regression discontinuity design* and it provides us with a powerful tool for causal inference. We’ll distinguish this from another kind of regression discontinuity setup called a *fuzzy regression discontinuity design* below. I’ll use the shorthand RD to refer to regression discontinuity in these notes.

When we previously used regression to carry out causal inference, the idea was to compare two groups of people who had been *matched* using a set of covariates  $\mathbf{x}$ . One group was treated and the other was not, but both groups had exactly the same values of  $\mathbf{x}$ . Under the assumption that treatment is “as good as randomly assigned” after conditioning on  $\mathbf{x}$ , we could learn the causal effect of treatment by comparing the mean outcomes of the two groups.

Sharp RD is very different since there is no way to carry out matching using the running variable  $x$ . This is because everyone who has  $x < x_0$  is untreated while everyone who has  $x \geq x_0$  is treated. Instead of matching people who have the same covariate values, sharp RD *extrapolates* by comparing people with *different* covariate values. The basic idea is very simple: we compare people whose  $x$  is close to but slightly *below* the cutoff  $x_0$  to people whose  $x$  is close to but slightly *above* the cutoff. Both  $D$  and  $x$  could affect  $Y$ , but since  $D$  abruptly switches from 0 to 1 at  $x_0$ , a causal relationship between  $D$  and  $Y$  should show up as a “jump” in the relationship between  $x$  and  $Y$  at  $x_0$ . For example, in the left panel there is clearly a jump at  $x_0 = 0.5$  and in the right panel there is not:



If the threshold  $x_0$  equals 0.5, then the left figure suggests that  $D$  has a substantial causal effect on  $Y$ : when  $D$  switches from zero to one as  $x$  crosses the threshold, the mean of  $Y$  jumps from around 2 to around 4. In contrast, the right panel doesn't show evidence of a causal effect of  $D$  on  $Y$ : when  $D$  switches from zero to one, there is no discernible change in the mean of  $Y$ . Now we'll be a little more precise about this intuition by thinking about exactly how  $x$  and  $Y$  are related. The key will be to create a link to potential outcomes, since this is our main tool for thinking about causality. Let's stick with the same  $x$ -axis as in the preceding figures: imagine that the running variable  $x$  is between zero and one, and that the cutoff  $x_0$  equals 0.5. This means that anyone with  $x < 0.5$  is untreated while anyone with  $x \geq 0.5$  is treated.

To begin, suppose we had a data for a large random sample of people who all had  $x = 0.3$ . If we took the mean  $Y$  for these people we would get an unbiased estimate of  $E[Y_i|x_i = 0.3]$ . The crucial point about sharp RD is that treatment is *completely determined* by  $x$ . This means that there is *no selection bias* since individuals are not free to choose their treatment status. Since everyone with  $x_i = 0.3$  is untreated, we have  $E[Y_i|x_i = 0.3] = E[Y_{0i}|x_i = 0.3]$ . Note that this is *not* the same thing as  $E[Y_{0i}]$  since the running variable  $x$  could have a direct effect on  $Y$ . In words, a person's potential outcome when untreated could depend on her value of  $x$ . For example,  $E[Y_{0i}|x_i = 0.3]$  may not equal  $E[Y_{0i}|x_i = 0.4]$  even though neither someone with  $x$  equal to 0.3 nor someone with  $x$  equal to 0.4 is treated. But this is fine, since we know how to use *predictive* modeling tools to estimate a conditional mean function. Here is the key point: since there is no selection into treatment for people with  $x < x_0$ , we can use *predictive regression* to estimate  $E[Y_i|x_i]$  and this will give us an estimate of  $E[Y_{0i}|x]$  for any  $x$  below the threshold.

What about when  $x$  is above the threshold? Consider for example, a large group of people with  $x = 0.6$  and suppose as above that  $x_0 = 0.5$ . All of these people are treated since their  $x$  exceeds the threshold. If we take the average  $Y$  for this group of people, we will obtain an unbiased estimator of  $E[Y_i|x_i = 0.6]$ . But since this group of people could not possibly select *out* of treatment, there is once again no selection bias and hence  $E[Y_i|x_i = 0.6] = E[Y_{1i}|x_i = 0.6]$ . This is not the same thing as  $E[Y_{1i}]$  since  $x$  could affect  $Y$  directly. But, again, this doesn't present a problem: since there is no selection out of treatment for people with  $x \geq x_0$ , we can use *predictive regression* to estimate  $E[Y_i|x_i]$  and this will give us an estimate of  $E[Y_{1i}|x_i]$ . To summarize

the reasoning from this and the preceding paragraph,

$$E[Y_i|x_i] = \begin{cases} E[Y_{0i}|x_i], & \text{if } x_i < x_0 \\ E[Y_{1i}|x_i], & \text{if } x_i \geq x_0 \end{cases}$$

Again, this relationship holds because individuals are *not* free to choose their treatment: everyone with  $x \geq x_0$  is treated and no one with  $x < x_0$  is treated.

There is a key distinction you need to bear in mind:  $E[Y_i|x_i]$  includes the effect of *both*  $D$  and  $x$  while  $E[Y_{0i}|x_i]$  and  $E[Y_{1i}|x_i]$  hold  $D$  *fixed*. The function  $E[Y_i|x_i]$  answers the question “what value should we predict for  $Y$  for someone who has a covariate value of  $x_i$ ?” In contrast, the function  $E[Y_{0i}|x_i]$  answers the question “what would be the average outcome for a person with covariate value  $x_i$  if I randomly assigned her  $D = 0$ ?” Similarly,  $E[Y_{1i}|x_i]$  answers the question “what would be the average outcome for a person with covariate value  $x_i$  if I randomly assigned her  $D = 1$ ?” If we knew  $E[Y_{0i}|x_i]$  and  $E[Y_{1i}|x_i]$  for all values of  $x$ , then by taking the difference, we could learn how the ATE *varies* across people with different values of  $x$ :

$$ATE(x) = E[Y_{1i}|x_i = x] - E[Y_{0i}|x_i = x] = E[Y_{1i} - Y_{0i}|x_i = x]$$

using the linearity of expectation. The idea of estimating an ATE as a function of some covariate  $x$  is called “heterogeneous treatment effects.” For example, a treatment may be more effective for younger people than older people. If we had experimental data, we could estimate  $ATE(x)$ . But in the RD setting we only have *observational data*. Crucially, we never observe  $E[Y_{0i}|x_i]$  for anyone with  $x_i \geq x_0$ , and we never observe  $E[Y_{1i}|x_i]$  for anyone with  $x_i < x_0$ .

So how can we proceed? In RD, the key assumption is that  $E[Y_{0i}|x_i]$  and  $E[Y_{1i}|x_i]$  are *continuous functions* of  $x$ . In other words, while we allow for the possibility that people with different values of  $x$  will have different potential outcomes, we assume that people with values of  $x$  that are *very similar* have will have potential outcomes that are *nearly equal*. In particular, we assume that  $\lim_{\Delta \rightarrow 0} E[Y_{1i}|x_i = x_0 + \Delta] = E[Y_{1i}|x_i = x_0]$  and similarly that  $\lim_{\Delta \rightarrow 0} E[Y_{0i}|x_i = x_0 - \Delta] = E[Y_{0i}|x_i = x_0]$ . But since  $E[Y_i|x_i]$  equals  $E[Y_{0i}|x_i]$  for  $x$  above the threshold and  $E[Y_{1i}|x_i]$  for  $x$  below the threshold, this implies that

$$\lim_{\Delta \rightarrow 0} E[Y_i|x_i = x_0 + \Delta] - E[Y_i|x_i = x_0 - \Delta] = E[Y_{1i} - Y_{0i}|x_i = x_0] = ATE(x_0)$$

So by using predictive regression to estimate  $E[Y_i|x_i]$  for  $x_i$  *just below*  $x_0$  and comparing it to an estimate of  $E[Y_i|x_i]$  for  $x_i$  *just above*  $x_0$ , RD allows us to learn the average treatment effect for individuals with  $x = x_0$ .