

# Pset3

Becky Cardinali, Yuen Ho, Sara Johns, and Jacob Lefler

12/4/2020

## Question 1: OLS Regressions

## Question 2: Regression Discontinuity Design

- (a) Consider the HRS score as the running variable for an RD research design. What assumptions are needed on the HRS score? How do each of the below “facts” impact the appropriateness of these assumptions?

In order for regression discontinuity to be a valid research design, we need to assume that the potential outcomes  $Y_i(0)$  and  $Y_i(1)$  (housing prices) are smooth functions of the running variable  $X_i$  (the HRS score) as it crosses the threshold  $c$  (28.5). In other words,  $E[Y_i(0)|X_i = x]$  and  $E[Y_i(1)|X_i = x]$  are continuous in  $x$ . If there is imperfect compliance, that is if the probability of treatment increases, but by less than 100 pp, when the running variable crosses the threshold, then we need to use a fuzzy RD design. In this case, we need to make an additional monotonicity assumption that  $D_i(x^*)$  is non-increasing in  $x^*$  at  $x^* = c$ , that is we need to assume there are no “defiers.”

Importantly, our first assumption is violated if there is manipulation based on the HRS score. In other words, if individuals understand the assignment mechanism and can manipulate the HRS score to place a census block just above (or below) the threshold, then there is selection into treatment so census tracts just above and below the threshold are no longer comparable. Thus we need to assume that individuals cannot game the assignment mechanism in order for this to be a valid research design. Relatedly, we also need to assume that covariates are smooth at the threshold, that is that covariates are balanced above and below the threshold. If this is not true, then we have selection into treatment and observations just above and below the threshold are again not comparable.

- (i) The EPA assertion that “the 28.5” cutoff was selected because it produced a manageable number of sites.”

This fact makes it more likely that our assumptions hold, because the threshold was not selected based on specific site characteristics, which would have potentially made covariates imbalanced across the threshold. For example, if instead the “28.5” cutoff was selected because a HRS rating of 28.5 or higher is especially (disproportionately) dangerous for human health, then our first assumption will no longer hold because houses close to sites above this threshold may benefit disproportionately from treatment.

- (ii) None of the individuals involved in identifying the site, testing the level of pollution, or running the 1982 HRS test knew the cutoff threshold score.

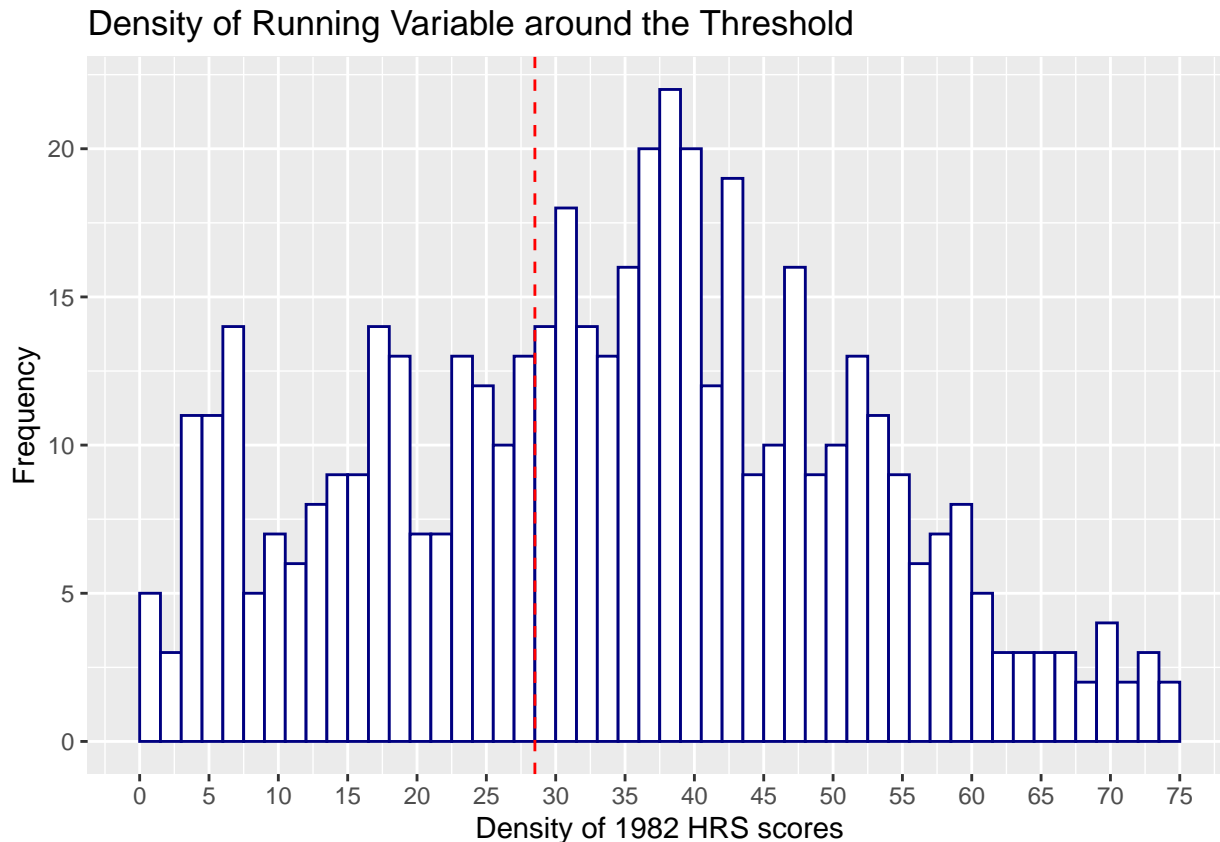
This fact makes it more likely that our assumptions hold. In particular, if none of the individuals involved knew the cutoff threshold score, it is less likely they were able to manipulate the test results to make certain census tracts be above (or below) the threshold. Even if individuals had an incentive to cheat, without knowing the assignment mechanism they would not have been able to game the system effectively.

- (iii) EPA documentation emphasizes that the HRS test is an imperfect scoring measure

Whether this fact violates our assumptions depends on the type of error associated with the HRS test. If this is classical measurement error then it should not affect our assumptions. However, if the error is correlated with our covariates or with our outcome variable (housing prices) then this would violate our first assumption.

- (b) Create a histogram of the distribution of the 1982 HRS scores by dividing the HRS score into non-overlapping bins. Include a vertical line at 28.5. Next run local linear regressions on either side of 28.5 using the midpoints of the bins as the data. What do you conclude?

```
## histogram of the density of 1982 HRS scores
ggplot(data, aes(x = hrs_82)) +
  geom_histogram(binwidth = 1.5, boundary = 0, closed = "left", col = "navy", fill = "white") +
  geom_vline(xintercept = 28.5, linetype = "dashed", color = "red") +
  theme_gray() +
  scale_x_continuous(breaks = seq(0,75,5)) +
  xlab("Density of 1982 HRS scores") +
  ylab("Frequency") +
  ggtitle("Density of Running Variable around the Threshold")
```



```
## Run local linear regressions on either side of threshold, using the midpoints of the bins as the data
range(data$hrs_82) # between 0 and 74.16
```

```
## [1] 0.00 74.16
```

```
h = 1.5 #set bandwidth
bins = seq(from = 0, to = 75, by = h) # set cutoffs for bins
length(bins)
```

```
## [1] 51
```

```
# returns the bin index for each observation
data$hrs_82_bin <- cut(data$hrs_82, breaks = bins, right = FALSE)

# calculate the midpoint of each bin
bins.midpoint = (bins[-1] + bins[-(length(bins))])/2

# assign a bin midpoint to each observation
data$hrs_82_binmid = bins.midpoint[ data$hrs_82_bin ]

# generate average of the treatment variable (NPL assignment) for each bin
npl2000_bin =tapply(data$npl2000, data$hrs_82_bin, mean)

# generate average outcome in each bin
lnmdvalhs0_nbr_bin = tapply(data$lnmdvalhs0_nbr, data$hrs_82_bin, mean)

# regression fitted on data below the cutoff
below_lm <- lm(lnmdvalhs0_nbr_bin ~ bins.midpoint, data.frame(lnmdvalhs0_nbr_bin, bins.midpoint)[1:19,])
summary(below_lm)
```

```
##
## Call:
## lm(formula = lnmdvalhs0_nbr_bin ~ bins.midpoint, data = data.frame(lnmdvalhs0_nbr_bin,
##     bins.midpoint)[1:19, ])
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -0.32629 -0.08308  0.03012  0.08214  0.30266
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept)   11.486342    0.080292 143.057  <2e-16 ***
## bins.midpoint  0.007585    0.004881   1.554   0.139
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 0.1748 on 17 degrees of freedom
## Multiple R-squared:  0.1244, Adjusted R-squared:  0.07284
## F-statistic: 2.414 on 1 and 17 DF,  p-value: 0.1387
```

```
# regression fitted on data above the cutoff
above_lm <- lm(lnmdvalhs0_nbr_bin ~ bins.midpoint, data.frame(lnmdvalhs0_nbr_bin, bins.midpoint)[20:50,])
summary(above_lm)
```

```
##
## Call:
## lm(formula = lnmdvalhs0_nbr_bin ~ bins.midpoint, data = data.frame(lnmdvalhs0_nbr_bin,
##     bins.midpoint)[20:50, ])
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -0.44550 -0.10299 -0.02503  0.11681  0.31197
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept)   11.657923    0.124356  93.746  <2e-16 ***
## bins.midpoint  0.001049    0.002326   0.451   0.655
```

```
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 0.1738 on 29 degrees of freedom
## Multiple R-squared:  0.006959,    Adjusted R-squared:  -0.02728
## F-statistic: 0.2032 on 1 and 29 DF,  p-value: 0.6555
```

We estimate the treatment effect  $\hat{\tau} = \hat{\alpha}_r - \hat{\alpha}_l = 11.658 - 11.486 = 0.172$ , the difference in the estimated intercepts from our local linear regressions above and below the threshold. This is suggestive evidence that being above the threshold, and therefore being more likely to be placed on the NPL, is associated with higher mean housing prices in 2000. Of course, a drawback of fitting separate local linear regressions on either side of the threshold is that we cannot conduct statistical inference on our estimated treatment effect.

### Question 3: First Stage of RD Design

- (a) Use a 2SLS (IV) econometric setup that uses whether or not a census tract has a site scoring above/below 28.5 as the instrument. Write down the 1st stage equation. Run the 1st stage regression experimenting with the same set of covariates used in question (1). In addition, run a second specification in which you limit the sample to only those census tracts with sites between 16.5 and 40.5 and run the specification using all of the control variables (we will use this as the size of the bandwidth for the “regression discontinuity” regression). Interpret the results.

We can write the first stage as

$$NPL_t = \delta_0 + \delta_1 1(HRS_t \geq 28.5) + \gamma X_t + \nu_t$$

where the instrument for NPL status is whether HRS is above 28.5 and we control for other covariates. Now we estimate this first stage regression, including controls for population density, education (college educated), children, poverty rate, and home characteristics (no full kitchen, 3 or more bedrooms, mobile).

```
first_stage <- lm(npl2000 ~ above_28pt5 + pop_den8_nbr + ba_or_better8_nbr + child8_nbr + povrat8_nbr +
  nofullkitchen80_nbr + bedrms3_80occ_nbr + mobile80occ_nbr,
  data = data)
summary(first_stage)
```

```
##
## Call:
## lm(formula = npl2000 ~ above_28pt5 + pop_den8_nbr + ba_or_better8_nbr +
##   child8_nbr + povrat8_nbr + nofullkitchen80_nbr + bedrms3_80occ_nbr +
##   mobile80occ_nbr, data = data)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -0.99305 -0.13976 -0.00303  0.01608  0.89830
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept)    2.613e-01  1.125e-01   2.323  0.0206 *
## above_28pt5     8.116e-01  2.311e-02  35.120 <2e-16 ***
## pop_den8_nbr   -4.754e-06  2.985e-06  -1.592  0.1120
## ba_or_better8_nbr  1.370e-01  1.721e-01   0.796  0.4264
## child8_nbr     -1.839e-01  2.699e-01  -0.681  0.4960
## povrat8_nbr    -1.979e-02  2.250e-01  -0.088  0.9300
## nofullkitchen80_nbr -8.361e-01  6.089e-01  -1.373  0.1704
## bedrms3_80occ_nbr -4.647e-02  1.462e-01  -0.318  0.7507
## mobile80occ_nbr  1.730e-02  1.547e-01   0.112  0.9110
## ---
```

```
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 0.2372 on 474 degrees of freedom
## Multiple R-squared:  0.7431, Adjusted R-squared:  0.7388
## F-statistic: 171.4 on 8 and 474 DF,  p-value: < 2.2e-16
```

As expected, having HRS above 28.5 is strongly predictive of NPL status. Now we rerun our IV analysis but focusing on census tracts with HRS between 16.5 and 40.5

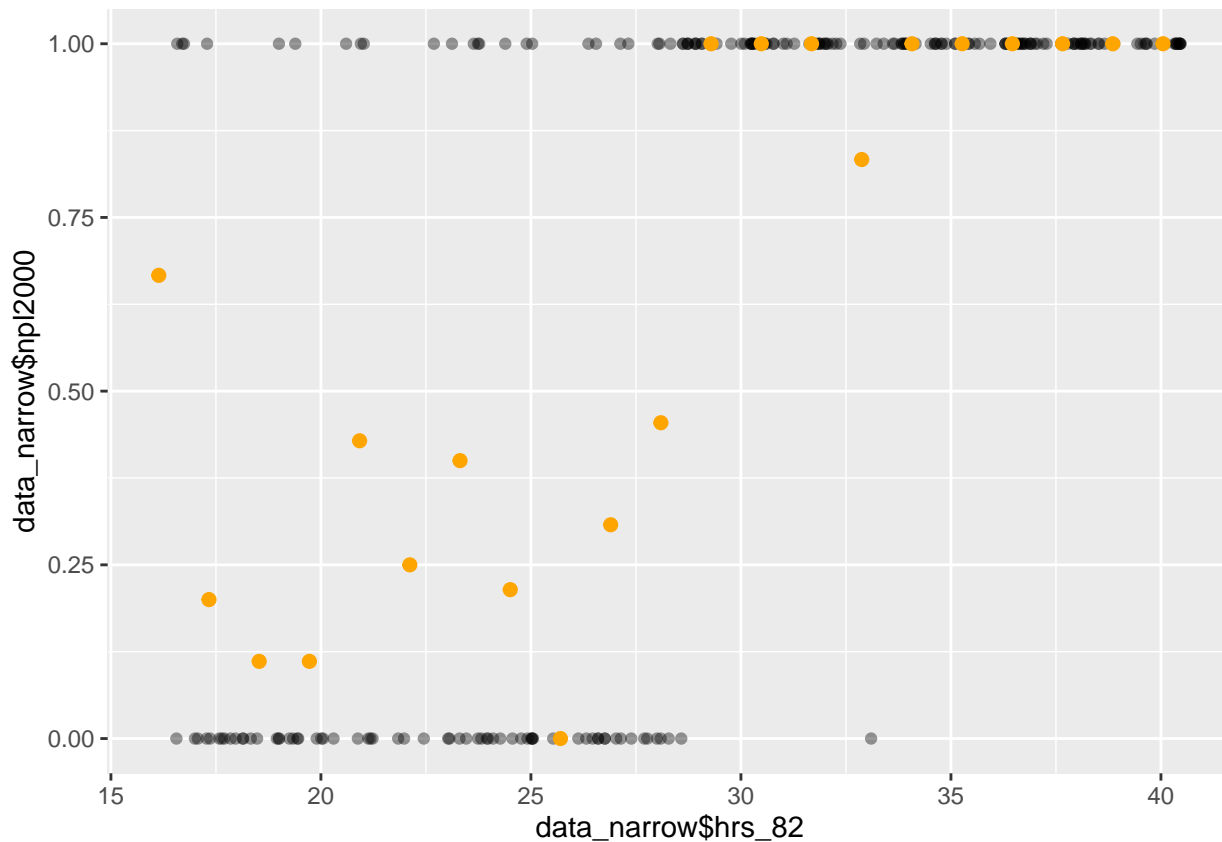
```
data_narrow <- data[data$hrs_82 >= 16.5 & data$hrs_82 <= 40.5,]
first_stage <- lm(npl2000 ~ above_28pt5 + pop_den8_nbr + ba_or_better8_nbr + child8_nbr + povrat8_nbr +
  nofullkitchen80_nbr + bedrms3_80occ_nbr + mobile80occ_nbr,
  data = data_narrow)
summary(first_stage)
```

```
##
## Call:
## lm(formula = npl2000 ~ above_28pt5 + pop_den8_nbr + ba_or_better8_nbr +
##   child8_nbr + povrat8_nbr + nofullkitchen80_nbr + bedrms3_80occ_nbr +
##   mobile80occ_nbr, data = data_narrow)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -0.98494 -0.22997 -0.00238  0.02202  0.86897
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept)    2.898e-01  2.017e-01   1.436  0.1523
## above_28pt5     7.080e-01  4.070e-02  17.396 <2e-16 ***
## pop_den8_nbr    -6.086e-06  5.464e-06  -1.114  0.2666
## ba_or_better8_nbr  1.010e-01  2.941e-01   0.344  0.7315
## child8_nbr      -3.404e-01  4.785e-01  -0.711  0.4776
## povrat8_nbr      3.785e-01  3.884e-01   0.975  0.3309
## nofullkitchen80_nbr -1.768e+00  9.395e-01  -1.882  0.0612 .
## bedrms3_80occ_nbr   1.770e-01  2.701e-01   0.655  0.5129
## mobile80occ_nbr    -9.817e-02  3.098e-01  -0.317  0.7516
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 0.296 on 217 degrees of freedom
## Multiple R-squared:  0.5968, Adjusted R-squared:  0.5819
## F-statistic: 40.14 on 8 and 217 DF,  p-value: < 2.2e-16
```

Here again the threshold is strongly predictive of NPL status.

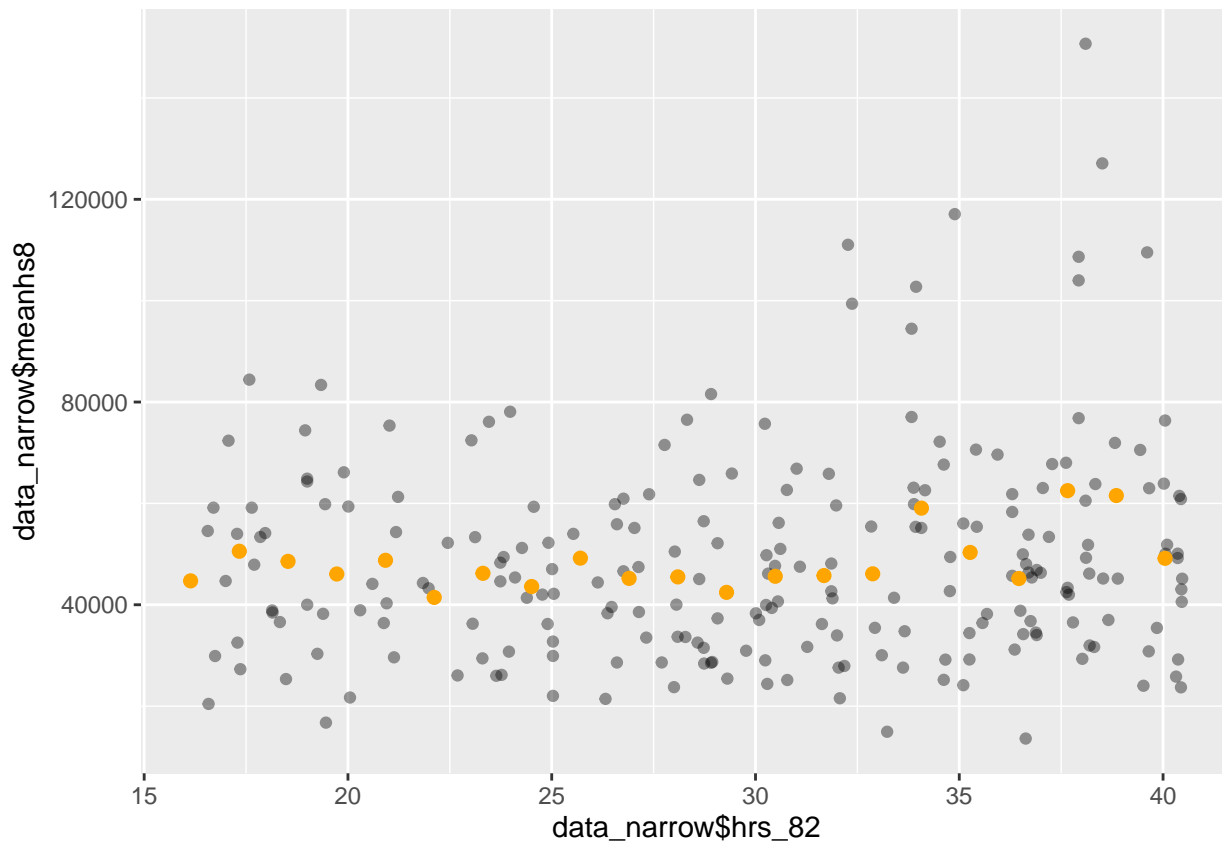
- (b) Create a graph plotting the the 1982 HRS score against whether a site is listed on the NPL by year 2000 (NPL on the y-axis, HRS on the x-axis). Briefly explain and interpret this graph.

```
(ggplot(data_narrow, aes(x=data_narrow$hrs_82,y=data_narrow$npl2000)) +
  geom_point(alpha = 0.4) +
  stat_summary_bin(fun='mean', bins=20,
    color='orange', size=2, geom='point'))
```



Here the yellow dots are the binned means and the black dots are the observed values. With an HRS score below 28.5, there is still a reasonable change (around 25%) that the site will be added to the NPL. With HRS above 28.5, it is almost guaranteed (the graph shows one exception).

```
(ggplot(data_narrow, aes(x=data_narrow$hrs_82,y=data_narrow$meanhs8)) +
  geom_point(alpha = 0.4) +
  stat_summary_bin(fun='mean', bins=20,
    color='orange', size=2, geom='point'))
```



There aren't any obvious differences in this range of HRS values. If anything, a higher HRS appears to be correlated with slightly higher housing values, which would cause our estimates to be downward biased, if anything. All in all, the values are largely comparable across this range.

#### Question 4: Second Stage of RD Design

#### Question 5: Conclusion