

# Honest inference in Sharp Regression Discontinuity

Michal Kolesár

2016-06-23

The package `RDHonest` implements confidence intervals for the regression discontinuity parameter considered in Armstrong and Kolesár (2016a) and Armstrong and Kolesár (2016b). In this vignette, we demonstrate the implementation of these confidence intervals using a dataset from Lee (2008), which is included in the package as a data frame `lee08`.

## Sharp RD model

In the sharp regression discontinuity model, we observe units  $i = 1, \dots, n$ , with the outcome  $y_i$  for the  $i$ th unit given by

$$y_i = f(x_i) + u_i,$$

where  $f(x_i)$  is the expectation of  $y_i$  conditional on the running variable  $x_i$  and  $u_i$  is the regression error. A unit is treated if and only if the running variable  $x_i$  lies above a known cutoff  $c_0$ . The parameter of interest is given by the jump of  $f$  at the cutoff,

$$\beta = \lim_{x \downarrow c_0} f(x) - \lim_{x \uparrow c_0} f(x).$$

Let  $\sigma^2(x_i)$  denote the conditional variance of  $u_i$ .

In the Lee dataset, the running variable corresponds to the margin of victory of a Democratic candidate in a US House election, and the treatment corresponds to winning the election. Therefore, the cutoff is zero. The outcome of interest is the Democratic vote share in the following election.

The package provides a function `plot_RDscatter` to plot the raw data. To remove some noise, the function plots averages over `avg` number of observations.

```
library("RDHonest")
## transform data to an RDdata object
dt <- RDData(lee08, cutoff = 0)
## plot 25-bin averages in a window equal to 50 around
## the cutoff, see Figure 1
plot_RDscatter(dt, avg = 25, window = 50, xlab = "Margin of victory",
  ylab = "Vote share in next election")
```

The function `RDHonest` constructs one- and two-sided confidence intervals (CIs) around local linear and local quadratic estimators using either a user-supplied bandwidth (which is allowed to differ on either side of the cutoff), or bandwidth that is optimized for a given performance criterion. The sense of honesty is that, if the regression errors are normally distributed with known variance, the CIs are guaranteed to achieve correct coverage *in finite samples*, and achieve correct coverage asymptotically uniformly over the parameter space otherwise. Furthermore, because the CIs explicitly take into account the possible bias of the estimators, the asymptotic approximation doesn't rely on the bandwidth to shrink to zero at a particular rate.

## Honest CIs

To describe the form of the CIs, let  $\hat{\beta}_{h_+, h_-}$  denote a local polynomial estimator with bandwidth equal to  $h_+$  above the cutoff and equal to  $h_-$  below the cutoff. Let  $\beta_{h_+, h_-}(f)$  denote its expectation conditional on the

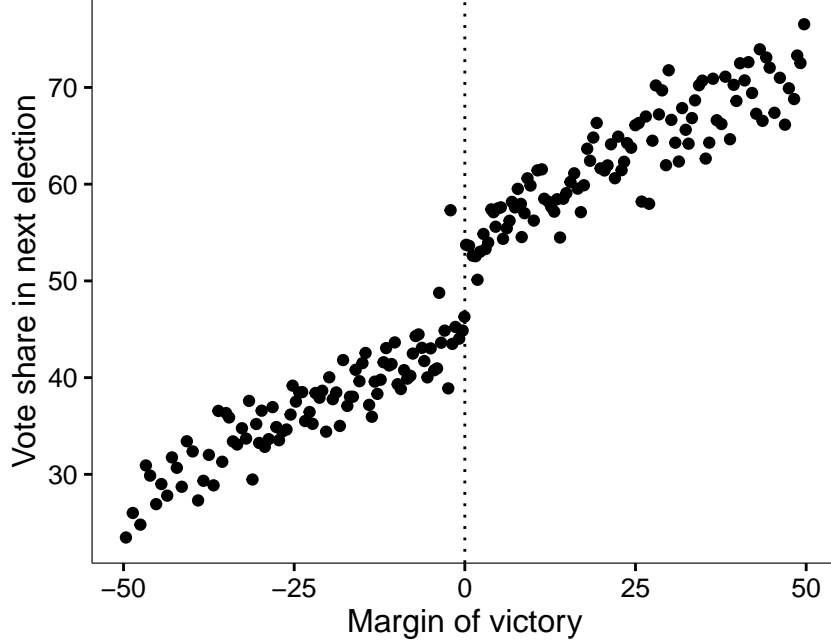


Figure 1: Lee (2008) data

covariates when the regression function equals  $f$ . Then the bias of the estimator is given by  $\beta_{h_+,h_-}(f) - \beta$ . Let

$$B(\hat{\beta}_{h_+,h_-}) = \sup_{f \in \mathcal{F}} |\beta_{h_+,h_-}(f) - \beta|$$

denote the worst-case bias over the parameter space  $\mathcal{F}$ . Then the lower limit of a one-sided CI is given by

$$\hat{\beta}_{h_+,h_-} - B(\hat{\beta}_{h_+,h_-}) - z_{1-\alpha} \widehat{se}(\hat{\beta}_{h_+,h_-}),$$

where  $z_{1-\alpha}$  is the  $1 - \alpha$  quantile of a standard normal distribution, and  $\widehat{se}(\hat{\beta}_{h_+,h_-})$  is the standard error (an estimate of the standard deviation of the estimator). Subtracting the worst-case bias in addition to the usual critical value times standard error ensures correct coverage at all points in the parameter space.

A two-sided CI is given by

$$\hat{\beta}_{h_+,h_-} \pm cv_{1-\alpha}(B(\hat{\beta}_{h_+,h_-})/\widehat{se}(\hat{\beta}_{h_+,h_-})) \times \widehat{se}(\hat{\beta}_{h_+,h_-}),$$

where the critical value function  $cv_{1-\alpha}(b)$  corresponds to the  $1 - \alpha$  quantile of the  $|N(b, 1)|$  distribution. To see why using this critical value ensures honesty, decompose the  $t$ -statistic as

$$\frac{\hat{\beta}_{h_+,h_-} - \beta}{\widehat{se}(\hat{\beta}_{h_+,h_-})} = \frac{\hat{\beta}_{h_+,h_-} - \beta_{h_+,h_-}(f)}{\widehat{se}(\hat{\beta}_{h_+,h_-})} + \frac{\beta_{h_+,h_-}(f) - \beta}{\widehat{se}(\hat{\beta}_{h_+,h_-})}$$

By a central limit theorem, the first term on the right-hand side will be distributed standard normal, irrespective of the bias. The second term is bounded in absolute value by  $B(\hat{\beta}_{h_+,h_-})/\widehat{se}(\hat{\beta}_{h_+,h_-})$ , so that, in large samples, the  $1 - \alpha$  quantile of the absolute value of the  $t$ -statistic will be bounded by  $cv_{1-\alpha}(B(\hat{\beta}_{h_+,h_-})/\widehat{se}(\hat{\beta}_{h_+,h_-}))$ . The function `CVb` gives these critical values:

```
## Usual critical value
CVb(0, alpha = 0.05)
#>   bias alpha      cv TeXDescription
#> 1    0  0.05 1.95996 $\\alpha=0.05$
```

```
## Tabulate critical values for different significance
## levels when bias-sd ratio equals 1/4
knitr::kable(CVb(1/4, alpha = c(0.01, 0.05, 0.1)))
```

bias	alpha	cv	TeXDescription
0.25	0.01	2.65224	$\alpha = 0.01$
0.25	0.05	2.01971	$\alpha = 0.05$
0.25	0.10	1.69558	$\alpha = 0.1$

The field `TeXDescription` is useful for plotting, or for exporting to  $\text{\LaTeX}$ , as in the table above.

## Parameter space

The package computes honest CIs when the parameter space  $\mathcal{F}$  corresponds to a second-order Taylor or second-order Hölder smoothness class, which capture two different types of smoothness restrictions. The second-order Taylor class assumes that  $f$  lies in the the class of functions

$$\mathcal{F}_{\text{Taylor}}(M) = \{f_+ - f_- : f_+ \in \mathcal{F}_T(M; [c_0, \infty)), f_- \in \mathcal{F}_T(M; (-\infty, c_0))\},$$

where  $\mathcal{F}_T(M; \mathcal{X})$  consists of functions  $f$  such that the approximation error from second-order Taylor expansion of  $f(x)$  about  $c_0$  is bounded by  $M|x|^2/2$ , uniformly over  $\mathcal{X}$ :

$$\mathcal{F}_T(M; \mathcal{X}) = \{f : |f(x) - f(c_0) - f'(c_0)x| \leq M|x|^2/2 \text{ all } x \in \mathcal{X}\}.$$

The class  $\mathcal{F}_T(C; \mathcal{X})$  formalizes the idea that the second derivative of  $f$  at zero should be bounded by  $M$ . See Section 5 in Armstrong and Kolesár (2016a). A disadvantage of this class is that it doesn't impose smoothness away from boundary, which may be undesirable in many empirical applications. The Hölder class addresses this problem directly by bounding the second derivative globally. In particular, it assumes that  $f$  lies in the class of functions

$$\mathcal{F}_{\text{Hölder}}(M) = \{f_+ - f_- : f_+ \in \mathcal{F}_H(M; [c_0, \infty)), f_- \in \mathcal{F}_H(M; (-\infty, c_0))\},$$

where

$$\mathcal{F}_H(M; \mathcal{X}) = \{f : |f'(x) - f'(y)| \leq M|x - y| \text{ } x, y \in \mathcal{X}\}.$$

## Honest CIs in Lee Dataset

CIs around a local linear estimator with bandwidth that equals to 10 on either side of the cutoff when the parameter space is given by a Taylor and Hölder smoothness class, respectively, with  $M = 0.1$ :

```
RDHonest(voteshare ~ margin, data = lee08, kern = "uniform",
          M = 0.1, hp = 10, sclass = "T")
```

```
#> Call:
```

```
#> RDHonest(formula = voteshare ~ margin, data = lee08, M = 0.1,
```

```
#>
```

```
kern = "uniform", hp = 10, sclass = "T")
```

```

#> Inference by se.method:
#>   Estimate Maximum Bias Std. Error
#> nn  6.05677      3.78224    1.19053
#>
#> Confidence intervals:
#> nn  (0.316293, 11.7973), (0.316293, Inf), (-Inf, 11.7973)
#>
#> Bandwidth below cutoff: 10
#> Bandwidth above cutoff: 10 (Bandwidths are the same)
#> Number of effective observations: 292.325
RDHonest(votesshare ~ margin, data = lee08, kern = "uniform",
  M = 0.1, hp = 10, sclass = "H")
#> Call:
#> RDHonest(formula = votesshare ~ margin, data = lee08, M = 0.1,      kern = "uniform", hp = 10, sclass = "H")
#>
#> Inference by se.method:
#>   Estimate Maximum Bias Std. Error
#> nn  6.05677      1.72377    1.19053
#>
#> Confidence intervals:
#> nn  (2.37473, 9.73882), (2.37476, Inf), (-Inf, 9.73878)
#>
#> Bandwidth below cutoff: 10
#> Bandwidth above cutoff: 10 (Bandwidths are the same)
#> Number of effective observations: 292.325

```

The confidence intervals use the nearest-neighbor method to estimate the standard error by default. The package reports two-sided as well one-sided CIs (with lower as well as upper limit) by default.

CIs around MSE-optimal bandwidth:

```

RDHonest(votesshare ~ margin, data = lee08, kern = "uniform",
  M = 0.1, opt.criterion = "MSE", sclass = "T")
#> Call:
#> RDHonest(formula = votesshare ~ margin, data = lee08, M = 0.1,      kern = "uniform", opt.criterion = "MSE", sclass = "T")
#>
#> Inference by se.method:
#>   Estimate Maximum Bias Std. Error
#> nn  4.94066      0.971353    1.53665
#>
#> Confidence intervals:
#> nn  (1.41581, 8.46551), (1.44175, Inf), (-Inf, 8.43958)
#>
#> Bandwidth below cutoff: 4.93675
#> Bandwidth above cutoff: 4.93675 (Bandwidths are the same)
#> Number of effective observations: 138.802
RDHonest(votesshare ~ margin, data = lee08, kern = "uniform",
  M = 0.1, opt.criterion = "MSE", sclass = "H")
#> Call:
#> RDHonest(formula = votesshare ~ margin, data = lee08, M = 0.1,      kern = "uniform", opt.criterion = "MSE", sclass = "H")
#>
#> Inference by se.method:
#>   Estimate Maximum Bias Std. Error
#> nn  5.61533      0.796411    1.4757

```

```
#>
#> Confidence intervals:
#> nn      ( 2.3484, 8.88226), (2.39161, Inf), (-Inf, 8.83905)
#>
#> Bandwidth below cutoff: 6.73296
#> Bandwidth above cutoff: 6.73296 (Bandwidths are the same)
#> Number of effective observations: 192.403
```

It is possible to compute the MSE-optimal bandwidth directly using the function `RDOpt.BW`

```
RDOptBW(votesshare ~ margin, data = lee08, kern = "uniform",
        M = 0.1, opt.criterion = "MSE", sclass = "T")
#> Call:
#> RDOptBW(formula = votesshare ~ margin, data = lee08, M = 0.1,      kern = "uniform", opt.criterion = "MSE", sclass = "T")
#>
#> Bandwidth below cutoff: 4.93675
#> Bandwidth above cutoff: 4.93675 (Bandwidths are the same)
RDOptBW(votesshare ~ margin, data = lee08, kern = "uniform",
        M = 0.1, opt.criterion = "MSE", sclass = "H")
#> Call:
#> RDOptBW(formula = votesshare ~ margin, data = lee08, M = 0.1,      kern = "uniform", opt.criterion = "MSE", sclass = "H")
#>
#> Bandwidth below cutoff: 6.73296
#> Bandwidth above cutoff: 6.73296 (Bandwidths are the same)
```

## References

- Armstrong, Timothy B., and Michal Kolesár. 2016a. “Optimal Inference in a Class of Regression Models.”
- . 2016b. “Simple and Honest Confidence Intervals in Nonparametric Regression.”
- Lee, David S. 2008. “Randomized Experiments from Non-Random Selection in U.S. House Elections.” *Journal of Econometrics* 142 (2): 675–97.