

A Bayesian Nonparametric Approach to Geographic Regression Discontinuity Designs

Maxime Rischard¹, Zach Branson¹, Luke Miratrix², and Luke Bornn³

¹ Department of Statistics, Harvard University

² Graduate School of Education, Harvard University

³ Simon Fraser University

What is a Geographical Regression Discontinuity Design?

A GeoRDD is a natural experiment where units on one side of a geographical border are given a treatment while units on the other side are not. It's the spatial analog of univariate RDDs, which are an increasingly popular tool in econometrics and other social sciences.

What is your estimand?

The difference in expected potential outcomes from the control side of the border to the treatment side. It is defined along the border. Interestingly, the estimand is *functional*.

$\tau: \mathcal{B} \rightarrow \mathbb{R}$ defined as $\tau(\mathbf{b}) = \mathbb{E}[Y_{iT} - Y_{iC} \mid \mathbf{s}_i = \mathbf{b}]$

Annotations:
 - τ : treatment effect
 - \mathcal{B} : border
 - \mathbf{b} : point on the border
 - \mathbf{s}_i : spatial location of unit i
 - Y_{iT} : potential outcome for unit i under treatment...and under control

How do you estimate the treatment effect?

Our GeoRDD framework proceeds by analogy to the univariate RDD:

Univariate RDD

- ① fit a smooth **function** to the outcomes against the forcing variable on each side of the discontinuity,
- ② extrapolate the functions to the **discontinuity point**, and
- ③ take the difference between the two extrapolations to estimate the treatment effect at the threshold point.

GeoRDD

- ① fit a smooth **surface** to the outcomes against the geographical covariates in each region,
- ② extrapolate the surfaces to the **border curve**, and
- ③ take the pointwise difference between the two extrapolations to estimate the treatment effect along the border.

To implement this strategy, we use Gaussian process regression:

$Y_{iT} = m_T + f_T(\mathbf{s}_i) + \epsilon_i$ and $Y_{iC} = m_C + f_C(\mathbf{s}_i) + \epsilon_i$; $\epsilon_i \sim \text{iid normal noise}$

$f_T, f_C \stackrel{\text{Gaussian process}}{\sim} \mathcal{GP}(0, k(\mathbf{s}, \mathbf{s}'))$ with $k(\mathbf{s}, \mathbf{s}') = \sigma_{\text{GP}}^2 \exp\left(-\frac{(\mathbf{s} - \mathbf{s}')^\top (\mathbf{s} - \mathbf{s}')}{2\ell^2}\right)$.

Annotations:
 - f_T, f_C : smooth surface
 - \mathbf{s}_i : spatial location of unit i
 - $k(\mathbf{s}, \mathbf{s}')$: covariance between two locations
 - σ_{GP}^2 : mean shift
 - ℓ^2 : e.g. Squared Exponential kernel

Can we estimate the average treatment effect (ATE)?

We have to define which average. The geometry of the GeoRDD makes the choice non-trivial. We consider the class of weighted means over the border:

$$\tau^w = \frac{\int_{\mathcal{B}} w_{\mathcal{B}}(\mathbf{b}) \tau(\mathbf{b}) d\mathbf{b}}{\int_{\mathcal{B}} w_{\mathcal{B}}(\mathbf{b}) d\mathbf{b}} \approx \frac{\sum_{r=1}^R w_{\mathcal{B}}(\mathbf{b}_r) \tau(\mathbf{b}_r)}{\sum_{r=1}^R w_{\mathcal{B}}(\mathbf{b}_r)}$$

Annotations:
 - τ^w : ATE
 - $w_{\mathcal{B}}(\mathbf{b})$: weight function
 - \mathbf{b}_r : discretized at a set of R "sentinel" points

e.g. τ^{UNIF} : $w_{\mathcal{B}}(\mathbf{b}_r) = 1$ (population density)

τ^{ρ} : $w_{\mathcal{B}}(\mathbf{b}_r) = \rho(\mathbf{b}_r)$ (posterior covariance of treatment effect)

τ^{INV} : $w_{\mathcal{B}}(\mathbf{b}_{1:R}) = \Sigma_{\mathbf{b}_{1:R}|Y}^{-1} \mathbf{1}_R$

We prefer the inverse-variance weighted average treatment effect. Ask me why!

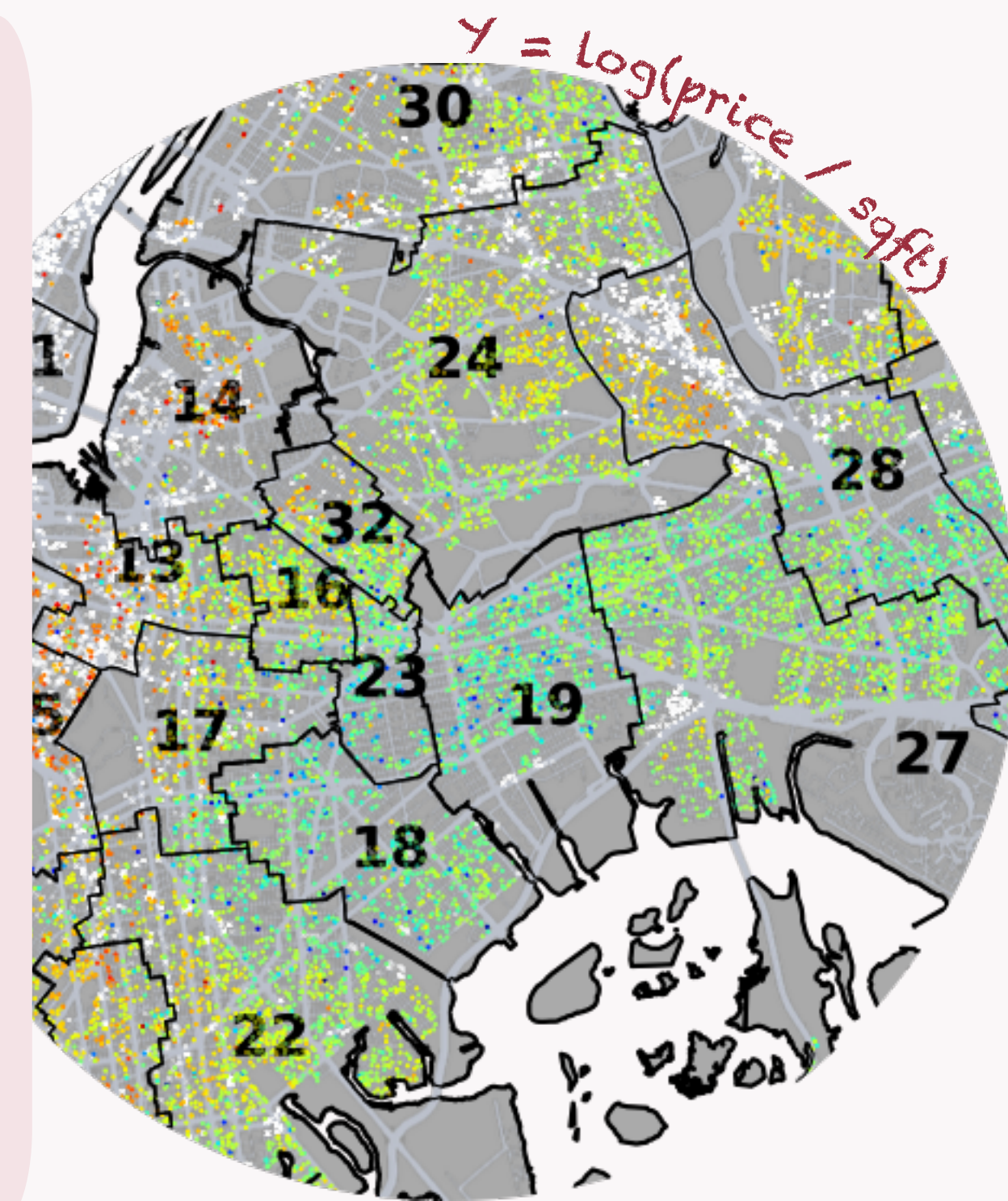
Illustration: Estimating the Cost of a NYC School District

The Data. A year of house sales in NYC, shown on the adjacent map. Each circle is a sale, coloured according to the price per square foot. Also shown are the outlines of the NYC school districts, numbered from 1 to 32. We also have the building class at time of sale.

The Question. It is a common belief that school districts impact real estate prices, as parents are willing to pay more to live in better districts.

Can we measure a discontinuous jump in house prices across the borders separating school districts?

Exploratory Analysis. Look closely at the border between districts 19 and 27. It looks like the dots in district 19 are bluer than those in district 27. But can we detect and quantify this difference at the border, and is it significant?

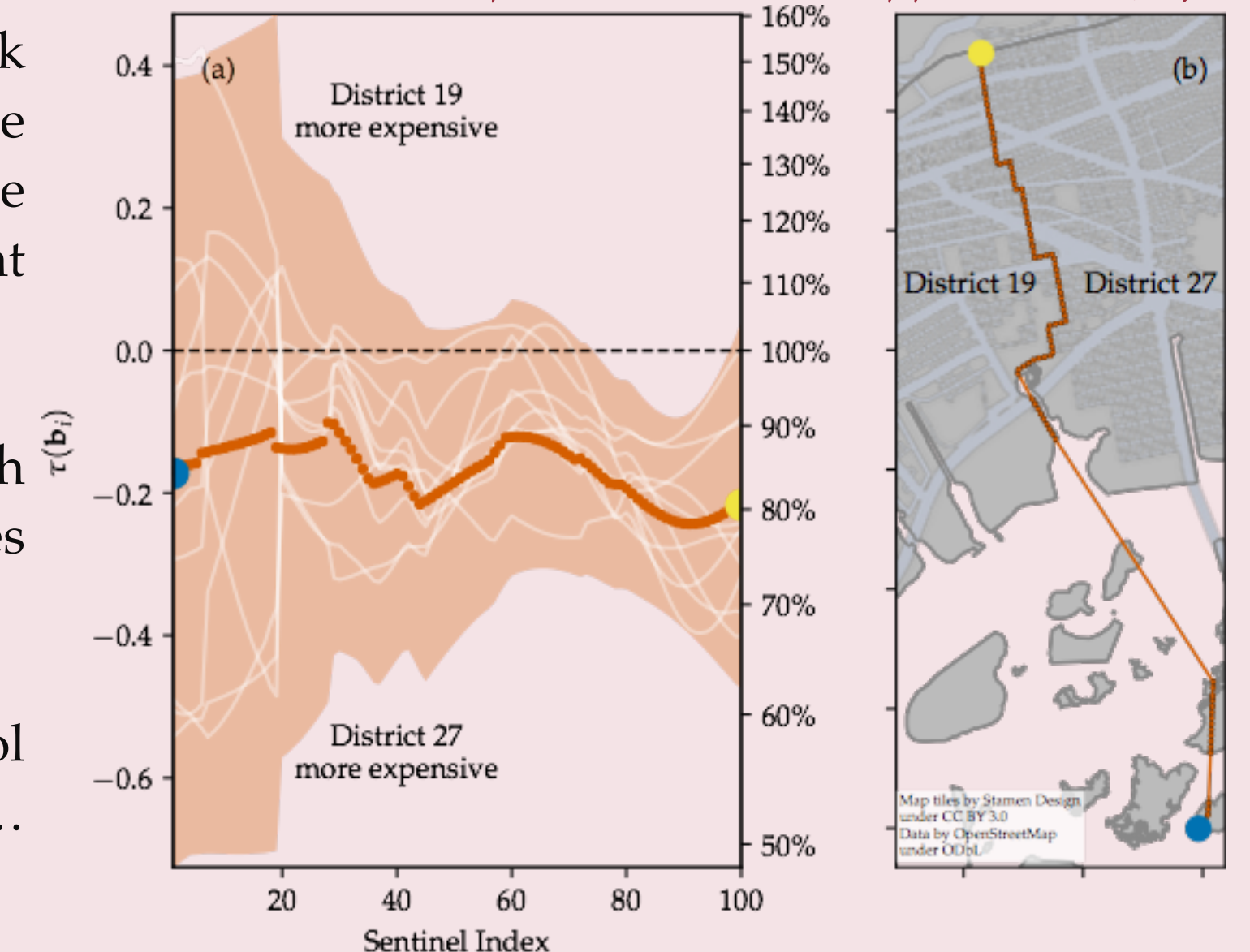


After fitting the hyperparameters by maximizing the marginal likelihood, and accounting for covariates (ask me how!), we extract the posterior mean and covariance of the treatment effect at 100 sentinel locations on the border between districts 19 and 27, shown on the right and in 3D below.

The inverse-variance weighted ATE is -0.19, which corresponds to an almost 20% increase in houses prices going from district 19 to district 27.

The difference is significant with $p=0.002$. So school districts cause differences in house prices? Not quite... the causal story is murkier. Ask me for details.

Plot of treatment effect $\tau(\mathbf{b}) \mid Y$



The posterior treatment effect is analytically tractable and easily computed.

$$\tau(\mathbf{b}_{1:R}) \mid Y, \theta \sim \mathcal{N}(\mu_{\mathbf{b}_{1:R}|Y}, \Sigma_{\mathbf{b}_{1:R}|Y})$$

Annotations:
 - $\tau(\mathbf{b}_{1:R})$: discretized at a set of R "sentinel" points
 - Y : data
 - θ : GP hyperparams
 - $\mu_{\mathbf{b}_{1:R}|Y}$: mean vector and covariance matrix

Is the treatment effect significant?

How do we know we're not just staring at noise? We can look at the tail probability of the ATE posterior distribution and derive a pseudo- p -value

$$\tilde{p}^w = 2\Phi\left(-\left|\mu_{\tau^w|Y}\right| / \sqrt{\Sigma_{\tau^w|Y}}\right)$$

Annotations:
 - \tilde{p}^w : normal CDF
 - $\mu_{\tau^w|Y}$: posterior mean and variance of avg. treatment effect

But it's not a valid p -value in the frequentist sense. We do better by treating the posterior mean ATE as a test statistic and deriving (analytically or using a bootstrap) its distribution under a parametric null distribution. This gives us a valid frequentist test derived from a Bayesian posterior!

