

# A Bayesian Non-parametric Approach to Geographic Regression Discontinuity Designs: Do School Districts Affect NYC House Prices?

Maxime Rischard<sup>a</sup>, Zach Branson<sup>a</sup>, Luke Miratrix<sup>b</sup>, and Luke Bornn<sup>c</sup>

<sup>a</sup>Department of Statistics, Harvard University

<sup>b</sup>Graduate School of Education, Harvard University

<sup>c</sup>?

April 23, 2018

## Abstract

Regression discontinuity designs (RDDs) are natural experiments characterized by the treatment assignment being fully determined by covariates. Most research has focused on one-dimensional cases, where units with a “forcing” variable lying on one side of a threshold value receive a treatment that the rest do not receive. More recently, situations with multiple forcing variables have garnered interest. When these variables are spatial covariates—that is, when a treatment is applied to a region but not its neighbor—the resulting natural experiment is termed a geographic regression discontinuity design (GeoRDD). In this paper, we propose a general framework for analysing GeoRDDs, which we implement using Gaussian process regression. We address multiple nuances of having a functional estimand defined on a border with potentially intricate topology, particularly when defining and estimating causal estimands of the average treatment effect, and when testing for non-zero treatment effects. Finally, we use our methodology on a dataset of property sales in New York City, and show evidence of a discontinuity in house prices between some pairs of neighboring school districts.

## 1 Introduction

The original theory and methods for regression discontinuity designs (RDDs) date from the 1960s, starting with Thistlethwaite and Campbell (1960). Cook (2008) trace the history of how interest in RDDs subse-

quently waned until the late 1990s when the design saw renewed attention, theoretical progress, and popularity in the social sciences. In particular, beginning with Papay et al. (2011), methods have been recently developed to analyse RDDs with multiple forcing variables. Imbens and Zajonc (2011) extend the local linear regression methods (see Imbens and Lemieux, 2008) that are popular for analysing RDDs with a single forcing variable (1D RDDs) to settings with multiple forcing variables.

Geographical regression discontinuity designs (GeoRDDs) are such RDDs where the forcing variables are spatial coordinates (latitude and longitude), meaning that units within a certain region are assigned to treatment, while units in a neighboring region are assigned to control. For example, in MacDonald et al. (2015), a private police force patrols a neighborhood, but stays out of surrounding areas, and a causal effect on crime rates is sought. In Chen et al. (2013), a policy applies south of the Huai River in China but not in the north, and pollution levels and life expectancies are measured to infer environmental and health impacts of the policy. In our example, we seek to estimate the effect of school districts on house prices in New York City. Practitioners often wish to use the well-established methods and software developed for 1D RDDs for their spatial problem. It is therefore tempting to reduce a GeoRDD problem to a 1D RDD by using the signed distance from the boundary (positive for treatment and negative for control) as the forcing variable, a method that we will refer to as “projected 1D RDD,” and which is used by both examples cited above. However, this method can fail to fully capture the spatial variation in the outcomes, resulting in spatial confoundedness of the estimator. We demonstrate this in a simple example in Appendix A. See also Section 4.2 of Keele and Titiunik (2015) for a discussion of this issue.

A more principled treatment of the GeoRDD is offered by Keele and Titiunik (2015), who build theoretical foundations for the analysis of GeoRDDs. They extend the identification assumptions that were first formalized by Hahn et al. (2001) for 1D RDDs, and Imbens and Zajonc (2011) for multivariate RDDs, to GeoRDDs. The main requirement for identification of the treatment effect is two-dimensional continuity of the conditional regression functions near the border. Notably, this is violated if units can sort around the border, crossing the border to seek or avoid the treatment, a particular concern in geographical RDDs. Keele and Titiunik (2015) also discuss further pitfalls of GeoRDDs, such as the issue of compound treatments—when a geographical border determines the assignment of the treatment of interest, but also of other differences.

Several methods for estimating the treatment effect in GeoRDDs have been proposed. Keele and Titiunik (2015) and Keele et al. (2017) estimate the treatment effect using a modification of the projected 1D RDD method by applying it locally around each point along the border, thus alleviating the problem of spatial confounding. This is slightly different from the method of Imbens and Zajonc (2011) for multivariate regression discontinuities that are not necessarily geographical: they use a multivariate local linear regression

to estimate the treatment effect at each boundary point, thus avoiding the need to specify a distance metric in the space of covariates. Keele et al. (2015) propose an alternative methodology, deploying the matching methodology of Zubizarreta (2012) to match units on opposite sides of the border that are near each other geographically and in other covariates, and to obtain estimates and confidence intervals for the average treatment effect. Their method requires the selection of a “buffer” distance from the border inside of which units on opposite sides of the border are sufficiently similar, so that the observational study on the subset of matched units can be interpreted and analyzed as a randomized experiment.

In this paper, we propose a framework for analysing GeoRDDs that is analogous to their 1D counterpart. Broadly, 1D RDD methodologies are composed of three steps:

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.

Reusing the same methodological skeleton and applying it to geographical RDDs, our framework proceeds analogously:

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.

Recently, Branson et al. (2017) proposed a Gaussian process regression (GPR) methodology that exhibits promising coverage and MSE properties compared to local linear regression for 1D RDDs. We believe this approach to be particularly suitable to GeoRDDs, as GPR is a well-established tool in spatial statistics (where it is known as kriging) for fitting smoothly varying spatial processes. See Banerjee et al. (2014) for a textbook introduction to kriging for spatial data, and Rasmussen and Williams (2006) for a machine learning perspective on GPR.

In Section 2, we use GPR to estimate the treatment effect along the border by extending the model of Branson et al. (2017) to geographical settings. A peculiarity of GeoRDDs is that the estimand is a function defined everywhere along the border, which is a one-dimensional manifold embedded in two-dimensional space. Furthermore, geographical borders, whether they be political or natural, are rarely simple straight lines. The topology of borders complicates the definition and interpretation of estimands for the average treatment effect (ATE), which we address in Section 3. We obtain Bayesian estimators for multiple possible

ATE estimands and illustrate their properties through a simulation study. In Section 4 we turn to hypothesis testing, and propose methods to test against the null hypothesis of no treatment effect along the border. Here too, the functional estimand brings interesting conceptual and computational challenges.

In Section 5, we apply our methodology to a publicly available dataset of property sales in NYC to determine whether school districts affect property prices. First focusing on a single border between school districts 19 and 27, we estimate the treatment effect everywhere along the border, obtain estimates of the ATE, and perform and validate hypothesis tests. Extending the analysis to all pairs of adjacent districts in Brooklyn and Queens, we discuss the complications and pitfalls of interpreting the estimated average treatment effect as a causal effect.

## 2 GeoRDD Modeling with Gaussian processes

### 2.1 GeoRDD Setup and Notation

We largely adopt the setup and notation for GeoRDDs laid out in Keele and Titiunik (2015). The outcomes  $Y_i$  of  $n$  units with spatial coordinates  $s_i$  are observed within an area  $\mathcal{A}$  of 2-dimensional coordinate space. The units are separated into  $n_T$  treatment units in area  $\mathcal{A}^T \subset \mathcal{A}$  and  $n_C$  units in the control area  $\mathcal{A}^C$ . The defining characteristic of GeoRDDs is that the two areas are adjacent but non-overlapping, intersecting only at a border  $\mathcal{B}$ . Throughout this paper, points on the border are denoted by  $b$ . For computational reasons, we will often represent the border as a set  $b_{1:\mathcal{R}} = \{b_1, \dots, b_R\}$ ,  $b_r \in \mathcal{B}$  of  $R$  “sentinel points” along the border. Under the potential outcomes framework for causal inference, each unit  $i$  has a potential outcomes  $Y_{iT}$  and  $Y_{iC}$ , which represent the outcome under treatment and control respectively. Let  $Z_i$  denote the treatment indicator, which is equal to one if unit  $i$  is in the area group, and zero if it is in the control area. Unlike traditional randomized experiments, the treatment assignment is a deterministic function of the unit’s geographical coordinates  $s_i$ :  $Z_i = \mathbb{I}\{s_i \in \mathcal{A}^T\}$ . The observed outcome for unit  $i$  is  $Y_i = Z_i Y_{iT} + (1 - Z_i) Y_{iC}$ . We denote the vector of observed outcomes of the treatment units and control units respectively by  $\mathbf{Y}_T$  and  $\mathbf{Y}_C$ .

For 1D RDDs, because the treatment and control regions do not overlap, the treatment effect is typically only inferred at the threshold  $X = b$ . As was already recognized by Thistlethwaite and Campbell (1960), this choice requires the least extrapolation of the fitted regression functions, which makes the estimated treatment more credible. The estimand at the threshold can be obtained as the difference of the two limits

of the expectation of the conditional regression functions

$$\tau = \mathbb{E}[Y_{iT} | X_i = b] - \mathbb{E}[Y_{iC} | X_i = b] = \lim_{x \downarrow b} \mathbb{E}[Y | X = x] - \lim_{x \uparrow b} \mathbb{E}[Y | X = x], \quad (1)$$

where the second equality requires the assumption that the conditional regression functions  $\mathbb{E}[Y_{iT} | X_i = x]$  and  $\mathbb{E}[Y_{iC} | X_i = x]$  are continuous in  $x$  (see Assumption 2.1 in [Imbens and Lemieux \(2008\)](#) and the discussion that follows). Analogously, we focus on the treatment effect at the border  $\mathcal{B}$  between the treatment and control regions. Here,  $\tau : \mathcal{B} \rightarrow \mathbb{R}$  is a function with

$$\tau(\mathbf{b}) = \mathbb{E}[Y_{iT} - Y_{iC} | \mathbf{s}_i = \mathbf{b}] \quad (2)$$

This is the estimand defined in [Imbens and Zajonc \(2011\)](#) and [Keele and Titiunik \(2015\)](#). For any  $\mathbf{b} \in \mathcal{B}$ ,  $\tau(\mathbf{b})$  can be obtained as the difference of the two limits of the expected outcomes, approaching  $\mathbf{b}$  from the treatment or the control side of the border, given the assumption that the conditional regression functions  $\mathbb{E}[Y_{iT} | \mathbf{s}_i = \mathbf{s}]$  and  $\mathbb{E}[Y_{iC} | \mathbf{s}_i = \mathbf{s}]$  are continuous in  $\mathbf{s}$  within  $\mathcal{A}$ . This result is formalized under Assumption 2.2.2 by [Imbens and Zajonc \(2011\)](#) and Assumption 1 in [Keele and Titiunik \(2015\)](#).

## 2.2 Model Specification

Our GeoRDD framework allows any method to be used to fit the outcomes on either side of the border. In this paper we will use Gaussian process regression (GPR) for this purpose. GPR, known as kriging in the spatial statistics literature, is a Bayesian non-parametric method for fitting smooth functions, that was shown by [Branson et al. \(2017\)](#) to be a promising tool for fitting 1D RDDs. Further inspired by the popularity of GPR in spatial statistics, we extend the model of [Branson et al. \(2017\)](#) to geographical RDDs.

On each side of the border, we model the observed outcomes  $Y_i$  at location  $\mathbf{s}_i$  as the sum of an intercept  $m$ , a spatial Gaussian process  $f(\mathbf{s})$ , and iid normal noise  $\epsilon$ . The Gaussian process has zero mean, and its covariance function is a modeling choice. There is a rich literature of possible covariance functions, known as “kernels” in machine learning, but in this paper, we will use the squared exponential kernel, for its ease of understanding and its prevalence in applied spatial statistics (see for example [Banerjee et al. \(2014\)](#) and [Rasmussen and Williams \(2006\)](#) for other possible choices of covariance functions). This yields the

outcomes model:

$$\begin{aligned}
Y_{iT} &= \underbrace{m_T + f_T(s_i)}_{g_T(s_i)} + \epsilon_i \\
Y_{iC} &= \underbrace{m_C + f_C(s_i)}_{g_C(s_i)} + \epsilon_i \\
f_T, f_C &\stackrel{\perp}{\sim} \mathcal{GP}(0, k(s, s')) \\
k(s, s') &= \sigma_{GP}^2 \exp\left(-\frac{(s - s')^\top (s - s')}{2\ell^2}\right)
\end{aligned} \tag{3}$$

The treatment effect at a location  $\mathbf{b}$  on the border is derived as the difference between the two (noise-free) surfaces  $g_T$  and  $g_C$

$$\tau(\mathbf{b}) = [m_T + f_T(\mathbf{b})] - [m_C + f_C(\mathbf{b})]. \tag{4}$$

This can be visualized as the height of a cliff separating the two smooth plains of the treatment and control regions.

In this specification, the hyperparameters  $\ell$ ,  $\sigma_{GP}$ , and  $\sigma_\epsilon$  are the same in the treatment and control regions, so we assume that the spatial smoothness of the responses is not affected by the treatment. We expect that this assumption will be reasonable in many applications, but it can be easily relaxed, as discussed in [Branson et al. \(2017\)](#).

## 2.3 Inference

If  $m_T$  and  $m_C$  are given normal priors with variance  $\sigma_\mu$ , then the model specification (3) can be used to obtain covariances between the observations, the Gaussian processes, and the mean parameters. Given hyperparameters  $\theta \equiv (\ell, \sigma_{GP}, \sigma_\epsilon, \sigma_\mu)$ , any vector with entries consisting of observations, points on the potential outcomes surface  $f_T$  and  $f_C$ , and the mean parameters  $m_C, m_T$  is jointly multivariate normal. Therefore the distribution of any such vector conditioned on another is also multivariate normal, with mean and covariances analytically tractable, and easily computed.

In accordance with the framework laid out in Section 1, we proceed by extrapolating both Gaussian processes to the border, and then taking the difference of the predictions to obtain the posterior treatment effect  $\tau(\mathcal{B})$  along the border. Computationally, we need to represent this border as a set  $\mathbf{b}_{1:\mathcal{R}} = \{\mathbf{b}_1, \dots, \mathbf{b}_{\mathcal{R}}\}$  of  $\mathcal{R}$  “sentinel” units dotted along  $\mathcal{B}$ . The extrapolation step then follows mechanically through multivariate

normal theory. On the treatment side, for example, we have:

$$\begin{aligned} g_T(\mathbf{b}_{1:\mathcal{R}}) \mid \mathbf{Y}_T, S_T, \boldsymbol{\theta} &\sim \mathcal{N}\left(\mu_{\mathbf{b}_{1:\mathcal{R}}|T}, \Sigma_{\mathbf{b}_{1:\mathcal{R}}|T}\right) \\ \mu_{\mathbf{b}_{1:\mathcal{R}}|T} &\equiv \mathbf{K}_{BT}\Sigma_{TT}^{-1}\mathbf{Y}_T \\ \Sigma_{\mathbf{b}_{1:\mathcal{R}}|T} &\equiv \mathbf{K}_{BB} - \mathbf{K}_{BT}\Sigma_{TT}^{-1}\mathbf{K}_{BT}^T \end{aligned} \quad (5)$$

with all the covariance matrices derived from the model specification (see Appendix B). Analogously, predictions for  $g_C(\mathbf{b}_{1:\mathcal{R}})$  are obtained using the data in the control region. Denote their posterior mean and covariance respectively by  $\mu_{\mathbf{b}_{1:\mathcal{R}}|C}$  and  $\Sigma_{\mathbf{b}_{1:\mathcal{R}}|C}$ . Since the two surfaces are modeled as independent, the treatment effect  $\tau(\mathbf{b}_{1:\mathcal{R}}) = g_T(\mathbf{b}_{1:\mathcal{R}}) - g_C(\mathbf{b}_{1:\mathcal{R}})$  has posterior

$$\begin{aligned} \tau(\mathbf{b}_{1:\mathcal{R}}) \mid \mathbf{Y}, \boldsymbol{\theta} &\sim \mathcal{N}\left(\mu_{\mathbf{b}_{1:\mathcal{R}}|Y}, \Sigma_{\mathbf{b}_{1:\mathcal{R}}|Y}\right) \\ \mu_{\mathbf{b}_{1:\mathcal{R}}|Y} &= \mu_{\mathbf{b}_{1:\mathcal{R}}|T} - \mu_{\mathbf{b}_{1:\mathcal{R}}|C} \\ \Sigma_{\mathbf{b}_{1:\mathcal{R}}|Y} &= \Sigma_{\mathbf{b}_{1:\mathcal{R}}|T} + \Sigma_{\mathbf{b}_{1:\mathcal{R}}|C}. \end{aligned} \quad (6)$$

The posterior mean and covariance of the “cliff height”  $\tau(\mathbf{b}_{1:\mathcal{R}})$  are the primary output of our GeoRDD analysis, and we refer to (6) as the “cliff face” estimator.

This leaves the choice of the  $\boldsymbol{\theta}$  hyperparameters:  $\ell$ ,  $\sigma_{GP}$ ,  $\sigma_\epsilon$ , and  $\sigma_\mu$ . For  $\sigma_\mu$ , we arbitrarily pick a large number, so that the prior on the mean parameters is weak. The Gaussian process and noise hyperparameters are optimized by maximizing the marginal likelihood of the observations  $\mathbb{P}(\mathbf{Y}_T, \mathbf{Y}_C \mid \ell, \sigma_{GP}, \sigma_\epsilon)$ , which is available analytically and easily computed for GPR. This empirical Bayes approach is common in spatial and machine learning applications of Gaussian processes. An alternative would be to also specify a prior on the hyperparameters, which is preferable in order to fully account for the uncertainty in the model, but fully Bayesian inference of large Gaussian process models tends to be computationally expensive.

## 2.4 Handling Covariates

The Gaussian process specification also makes it easy to incorporate a linear model on non-spatial covariates, both mathematically and computationally. The models are modified by the addition of the linear regression term  $\mathbf{D}\beta$  on the  $n \times p$  matrix of covariates  $\mathbf{D}$ . We recommend placing a normal prior  $\mathcal{N}(0, \sigma_\beta^2)$  on the regression coefficients, as this preserves the multivariate normality of the model, with the simple addition of a term  $\sigma_\beta^2 \mathbf{D}\mathbf{D}^T$  to the covariance of  $\mathbf{Y}$ .

Our model becomes

$$\begin{aligned}
Y_{iT} &= \underbrace{m_T + f_T(s_i)}_{g_T(s_i)} + d_i^T \beta + \epsilon_i \\
Y_{iC} &= \underbrace{m_C + f_C(s_i)}_{g_C(s_i)} + d_i^T \beta + \epsilon_i \\
f_T, f_C &\stackrel{\perp}{\sim} \mathcal{GP}(0, k(s, s')) \\
k(s, s') &= \sigma_{GP}^2 \exp\left(-\frac{(s - s')^T (s - s')}{2\ell^2}\right) \\
\beta_j &\stackrel{\perp}{\sim} \mathcal{N}\left(0, \sigma_\beta^2\right) \text{ for } j = 1, 2, \dots, p
\end{aligned} \tag{7}$$

Unfortunately, the linear term induces a covariance between the treatment and control region, which quadruples the computational cost of the analysis. When the two regions are independent, fitting the Gaussian processes required the inversion of an  $n_T \times n_T$  covariance matrix, and of an  $n_C \times n_C$  matrix. But with the additional covariates, the covariance of  $Y$  is no longer block diagonal. Thus the inversion of an  $(n_T + n_C) \times (n_T + n_C)$  is now required. Matrix inversion algorithms generally have computational complexity  $O(n^3)$ . Therefore, if the units are evenly split between the two regions, the overall complexity of the model fitting increases fourfold.

The introduction of the linear term modifies the cliff face estimator (6) so that its posterior mean and covariance become:

$$\begin{aligned}
\mu_{b_{1:R}|Y,D} &= \begin{bmatrix} \mathbf{K}_{BT} & -\mathbf{K}_{BC} \end{bmatrix} \begin{bmatrix} \Sigma_{TT} + \sigma_\beta^2 D_T D_T^T & \sigma_\beta^2 D_T D_C^T \\ \sigma_\beta^2 D_C D_T^T & \Sigma_{CC} + \sigma_\beta^2 D_C D_C^T \end{bmatrix}^{-1} \begin{pmatrix} Y_T \\ Y_C \end{pmatrix}, \text{ and} \\
\Sigma_{b_{1:R}|Y,D} &= 2\mathbf{K}_{BB} - \begin{bmatrix} \mathbf{K}_{BT} & -\mathbf{K}_{BC} \end{bmatrix} \begin{bmatrix} \Sigma_{TT} + \sigma_\beta^2 D_T D_T^T & \sigma_\beta^2 D_T D_C^T \\ \sigma_\beta^2 D_C D_T^T & \Sigma_{CC} + \sigma_\beta^2 D_C D_C^T \end{bmatrix}^{-1} \begin{bmatrix} \mathbf{K}_{BT}^T \\ -\mathbf{K}_{BC}^T \end{bmatrix}.
\end{aligned} \tag{8}$$

To avoid the complexity caused by the correlation between  $Y_T$  and  $Y_C$  that the linear term induces, we suggest first obtaining an estimate  $\hat{\beta}$  of the coefficients. We show how to obtain the posterior mean of  $\beta$  in Appendix C. We can then proceed with the GeoRDD analysis on the residuals, which are decorrelated conditionally on  $\beta$ . This is an approximation, as it ignores the uncertainty in the estimate of  $\beta$ , but if the number of samples in the treatment and control areas is high (even away from the border), the approximation has negligible effect on the estimate of the treatment effect, and simplifies the subsequent GeoRDD analysis.

### 3 Characterizing and Estimating the Average Treatment Effect

Once we obtain the posterior on the treatment effect function  $\tau(\mathcal{B})$ , estimating the average treatment effect (ATE) along the border will often be of interest. We consider the class of weighted means of the functional treatment effect  $\tau(\mathbf{b})$ , with weight function  $w_{\mathcal{B}}(\mathbf{b})$  defined everywhere on the border  $\mathcal{B}$ . The weighted mean integral can be approximated as a weighted sum at the sentinels  $\mathbf{b}_{1:\mathbb{R}}$ :

$$\begin{aligned}\tau^w &= \frac{\oint_{\mathcal{B}} w_{\mathcal{B}}(\mathbf{b})\tau(\mathbf{b}) d\mathbf{s}}{\oint_{\mathcal{B}} w_{\mathcal{B}}(\mathbf{b}) d\mathbf{b}}, \\ &\approx \frac{\sum_{r=1}^R w_{\mathcal{B}}(\mathbf{b}_r)\tau(\mathbf{b}_r)}{\sum_{r=1}^R w_{\mathcal{B}}(\mathbf{b}_r)}.\end{aligned}\tag{9}$$

Note that the approximation assumes that the sentinels are evenly spaced, otherwise each term in the sum needs to be re-weighted by the length of the border that the sentinel occupies. We have shown the posterior distribution of  $\tau(\mathbf{b}_{1:\mathbb{R}})$  to be multivariate normal, with mean  $\mu_{\mathbf{b}_{1:\mathbb{R}}|Y}$  and covariance  $\Sigma_{\mathbf{b}_{1:\mathbb{R}}|Y}$  given in (6). Since  $\tau^w$  is a linear transformation of  $\tau(\mathbf{b}_{1:\mathbb{R}})$ , its posterior is also multivariate normal, with mean  $\mu_{\tau^w|Y}$  and covariance  $\Sigma_{\tau^w|Y}$  given by

$$\begin{aligned}\mu_{\tau^w|Y} &= \frac{w_{\mathcal{B}}(\mathbf{b}_{1:\mathbb{R}})^T \mu_{\mathbf{b}_{1:\mathbb{R}}|Y}}{w_{\mathcal{B}}(\mathbf{b}_{1:\mathbb{R}})^T \mathbf{1}_R} \\ \Sigma_{\tau^w|Y} &= \frac{w_{\mathcal{B}}(\mathbf{b}_{1:\mathbb{R}})^T \Sigma_{\mathbf{b}_{1:\mathbb{R}}|Y} w_{\mathcal{B}}(\mathbf{b}_{1:\mathbb{R}})}{(w_{\mathcal{B}}(\mathbf{b}_{1:\mathbb{R}})^T \mathbf{1}_R)^2}\end{aligned}\tag{10}$$

where  $w_{\mathcal{B}}(\mathbf{b}_{1:\mathbb{R}})$  is the R-vector of weights evaluated at the sentinels, and  $\mathbf{1}_R$  is an R-vector of ones. For each estimator obtained in (10) as a weighted mean of  $\mu_{\mathbf{b}_{1:\mathbb{R}}|Y}$ , we consider the “natural” estimand to be the same weighted mean applied to the truth  $\tau(\mathcal{B})$ , given by (9).

An alternative perspective on these estimators is given by the weights induced on the observations. Indeed, combining equations (5), (6), and (10), we obtain that the posterior mean of  $\tau^w$  is a linear combination

$$\mathbb{E}(\tau^w | Y) = \mathbf{w}_T^T \mathbf{Y}_T + \mathbf{w}_C^T \mathbf{Y}_C\tag{11}$$

of the observed data, with “unit weights” given by

$$\begin{aligned}\mathbf{w}_T &= \frac{1}{w_{\mathcal{B}}(\mathbf{b}_{1:\mathbb{R}})^T \mathbf{1}_R} \Sigma_{TT}^{-1} \mathbf{K}_{BT}^T w_{\mathcal{B}}(\mathbf{b}_{1:\mathbb{R}}), \text{ and} \\ \mathbf{w}_C &= -\frac{1}{w_{\mathcal{B}}(\mathbf{b}_{1:\mathbb{R}})^T \mathbf{1}_R} \Sigma_{CC}^{-1} \mathbf{K}_{BC}^T w_{\mathcal{B}}(\mathbf{b}_{1:\mathbb{R}}),\end{aligned}\tag{12}$$

for treatment and control units respectively.

The question remains: what is the most appropriate choice of weights? We next motivate and consider six possible choices of  $w_{\mathcal{B}}(\mathbf{b})$ , and explore interpretations, advantages, and disadvantages. A summary of their properties is provided in Table 2.

### 3.1 Uniform ATE

The simplest choice is uniform weights  $w_{\mathcal{B}}(\mathbf{b}) = 1$ , a seemingly reasonable and unopinionated decision. We estimate  $\tau^{\text{UNIF}}$ , the uniformly weighted mean of  $\tau(\mathcal{B})$ , by averaging the entries of the mean posterior at the sentinels. Following (9) and (10):

$$\begin{aligned}\tau^{\text{UNIF}} &\equiv \frac{\oint_{\mathcal{B}} \tau(x) d\mathbf{s}}{\oint_{\mathcal{B}} d\mathbf{x}} \\ \tau^{\text{UNIF}} | Y, \theta &\sim \mathcal{N}(\mu_{\tau^{\text{UNIF}}|Y}, \Sigma_{\tau^{\text{UNIF}}|Y}) \\ \mu_{\tau^{\text{UNIF}}|Y} &= (\mathbf{1}^\top \mu_{\mathbf{b}_{1:\mathbb{R}}|Y}) / R \\ \Sigma_{\tau^{\text{UNIF}}|Y} &= (\mathbf{1}^\top \Sigma_{\mathbf{b}_{1:\mathbb{R}}|Y} \mathbf{1}) / R^2\end{aligned}\tag{13}$$

The uniformly weighted estimand takes on a geometric interpretation: equal-length segments of the border are given equal weight. Unfortunately, uniform weights suffer from several issues that we describe and address in Sections 3.2 and 3.3.

### 3.2 Density Weighted ATE

With uniform border weights, parts of the border adjoining dense populations are given equal weights to those in sparsely populated areas. But if the border goes through an unpopulated area, like a lake or a public park, then the treatment effect there has little meaning and importance. Furthermore,  $\tau(\mathbf{b})$  in those empty areas will have large posterior variances, which will dominate the posterior variance of  $\tau^{\text{UNIF}}$ , potentially jeopardizing the successful detection of otherwise strong treatment effects.

We can address this issue by weighting the treatment effect at each sentinel location by the local density. That is we choose  $w_{\mathcal{B}}(\mathbf{b}) = \rho(\mathbf{b})$ , where  $\rho$  is the local population density. Attractively, the estimand is interpretable as the average treatment effect for the superpopulation of units that live on the border:

$$\tau^\rho = \mathbb{E}[Y_{iT} - Y_{iC} | \mathbf{s}_i \in \mathcal{B}].\tag{14}$$

It therefore better captures the “typical” treatment effect received by a unit than the uniformly weighted estimand. This is the estimand used by Keele and Titiunik (2015), who themselves follow in the footsteps

of Imbens and Zajonc (2011).

In practice, the local density needs to be estimated. A simple kernel density estimator can be used, though one could also deploy a more sophisticated spatial point process model. Strictly speaking, the uncertainty of the local density estimate should then be propagated to the estimate of  $\tau^0$ , which may therefore no longer have a normally distributed or analytically tractable posterior. These inconveniences certainly reduce the appeal of the density-weighted estimator, but there is a deeper issue affecting this choice of estimand: its susceptibility to the topology of the border.

### 3.3 Inverse-variance Weighted ATE

The unweighted and density-weighted mean treatment estimands are both affected by the shape of the border between the treatment and control regions, giving higher weight to wigglier sections of the border. We illustrate this with the border separating two American States: Louisiana and Mississippi. From North to South, the border follows the meandering Mississippi river, then takes a sharp turn to the East and becomes a straight line, until it meets the even more sinuous Pearl river, which it then follows until it reaches the Gulf of Mexico. Sentinels placed at equal distance intervals along this border will therefore be more densely packed along the rivers, and sparsest along the straight segment (see Figure 1). When averaging a function over the border, those sections will therefore be overrepresented. Troublingly, the sinuousness of the border therefore determines the estimand, even though the outcomes of interest will generally have nothing to do with river topologies. In population terms, the result is that units near wigglier segments receive more weight. Worse, the resolution of the map used in the analysis affects the estimated ATE.

This unwelcome dependence of the  $\tau^{\text{UNIF}}$  estimand on the border topology is a symptom of the geometry of the problem: the 1-dimensional treatment function  $\tau(\mathcal{B})$  is embedded in a Euclidean 2-dimensional space. The dependencies induced by this geometric fact are reflected in the covariance  $\Sigma_{b_{1:R}|Y}$ : sentinels in the straight segment of the border will be less strongly correlated than in the sinuous segments. The more correlated sentinels individually carry less information about the local treatment effect. Instead of averaging the posterior treatment effect along the border based on geometry or population, we consider averaging the information contained therein. This motivates the inverse-variance weighted mean  $\tau^{\text{INV}}$ :

$$\begin{aligned} \tau^{\text{INV}} | Y, \theta &\sim \mathcal{N}\left(\mu_{\tau^{\text{INV}}|Y}, \Sigma_{\tau^{\text{INV}}|Y}\right), \\ \mu_{\tau^{\text{INV}}|Y} &= \left(\mathbf{1}^\top \Sigma_{b_{1:R}|Y}^{-1} \mu_{b_{1:R}|Y}\right) / \left(\mathbf{1}^\top \Sigma_{b_{1:R}|Y}^{-1} \mathbf{1}\right), \\ \Sigma_{\tau^{\text{INV}}|Y} &= 1 / \left(\mathbf{1}^\top \Sigma_{b_{1:R}|Y}^{-1} \cdot \mathbf{1}\right) \end{aligned} \quad (15)$$

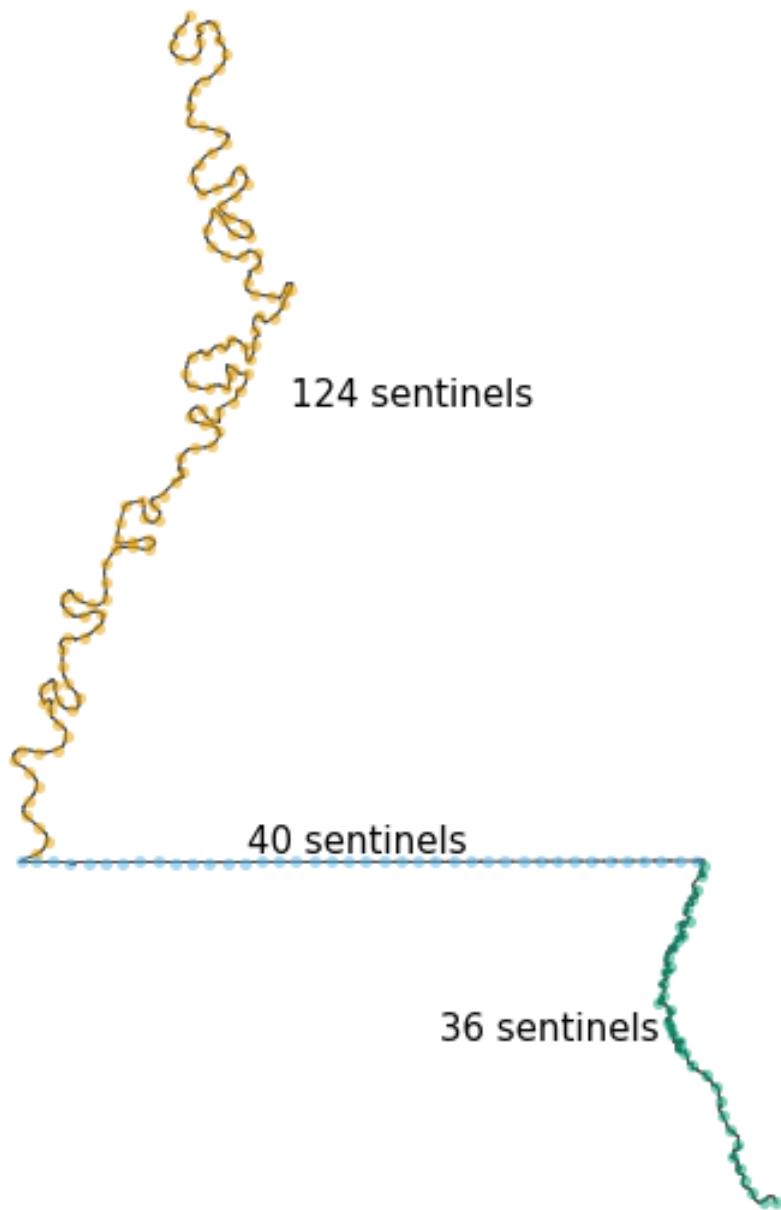


Figure 1: Evenly spaced sentinels along the border between Mississippi and Louisiana.

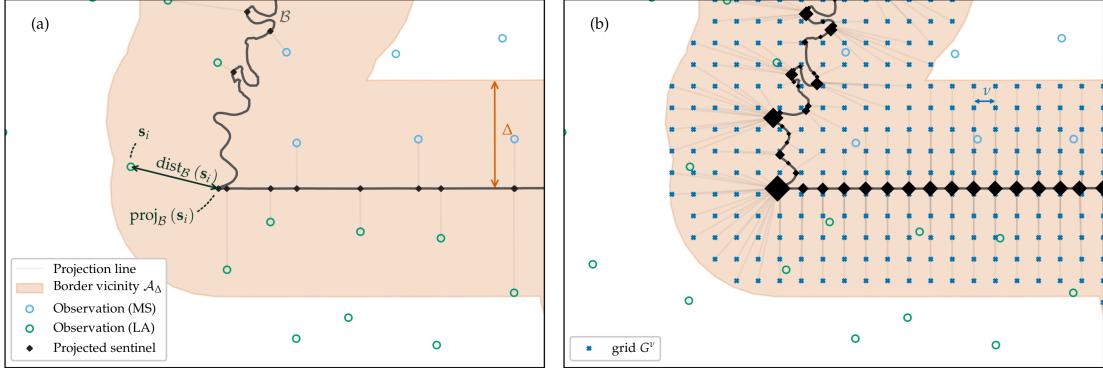


Figure 2: Illustration of (a) projected finite-population ATE  $\tau^{\text{PROJ}}$ , and (b) projected land ATE  $\tau^{\text{GEO}}$ , using the border separating Mississippi and Louisiana near Baton Rouge, with units at the centroid of each county. The border vicinity  $\mathcal{A}_\Delta$  is defined as all land within  $\Delta = 50$  km of the border. With both methods, every projected sentinel has equal weight in the ATE, but the tight grid in (b) causes sentinels to coincide or nearly coincide, which we depict by scaling up the size of the marker by the number of coinciding sentinels.

This estimator efficiently extracts the information from the posterior treatment effect, as it minimizes the posterior variance amongst weighted averages of the form (9). It automatically gives more weight to sentinels in dense areas (as the variance will be lower there), and to sentinels in straight sections of the border (as the correlations between sentinels will be lower).

The estimand is still a weighted mean, with weights for the sentinels given by  $w_{\mathcal{B}}(\mathbf{b}_{1:\mathbb{R}}) = \Sigma_{\mathbf{b}_{1:\mathbb{R}}|\mathcal{Y}}^{-1} \mathbf{1}$ . This can put negative weights on some sentinels, as seen in a simulated example in Figure 5(c), and generally this estimand doesn't lend itself to an intuitive interpretation. This estimand isn't chosen on scientific grounds, but rather it is dictated by the observed data. This is counter to the conventional wisdom in causal inference, that the estimand should be chosen based on substantive grounds, ideally before collecting any data. While perhaps unorthodox, analogous "estimands of convenience" have been proposed in other settings, for example matching methods that exclude some unmatched units from the analysis (discussed in Crump et al., 2009), or in the context of balancing treatment and control populations with little overlap in their covariate distributions (Li et al., 2016). The classical RDD could be said to provide another example, as the estimand (1) focuses on the treatment effect near the threshold not because those units are of particular substantive interest, but because the available data restricts estimation of the treatment effect elsewhere.

### 3.4 Projected Finite-Population ATE

All average treatment effect estimators considered so far presuppose evenly spaced sentinel points, which are then given weights. Alternatively, we can project the positions of treatment and control units that are within a distance  $\Delta$  of the border onto the border, and use those projected sentinel positions without

weights. This is illustrated in Figure 2(a). For any point  $\mathbf{s}$ , we use the notation  $\text{proj}_{\mathcal{B}}(\mathbf{s})$  to give the coordinates of the point on the border  $\mathcal{B}$  that is closest to  $\mathbf{s}$  (assuming uniqueness), and  $\text{dist}_{\mathcal{B}}(\mathbf{s})$  for the distance between the point and the border. Let  $\mathbb{I}^{\Delta}(\mathbf{s})$  indicate inclusion in the border vicinity by returning one if  $\text{dist}_{\mathcal{B}}(\mathbf{s}) < \Delta$  and zero otherwise. The projected finite-population  $\tau^{\text{PROJ}}$  is then the uniformly weighted mean applied with the projected sentinels instead of the evenly spaced sentinels. We can therefore modify (13), replacing the cliff face mean vector  $\mu_{b_{1:\mathbf{R}}|Y}$  and covariance matrix  $\Sigma_{b_{1:\mathbf{R}}|Y}$  with equivalent quantities obtained at the projected sentinels, to obtain the posterior mean and covariance of  $\tau^{\text{PROJ}}$ :

$$\begin{aligned}\tau^{\text{PROJ}} | Y_T, Y_C, \theta &\sim \mathcal{N}\left(\mu_{\tau^{\text{PROJ}}|Y}, \Sigma_{\tau^{\text{PROJ}}|Y}\right), \\ \mu_{\tau^{\text{PROJ}}|Y} &= \frac{1}{\sum_{i=1}^{n_T+n_C} \mathbb{I}^{\Delta}(\mathbf{s}_i)} \sum_{i=1}^{n_T+n_C} \mathbb{I}^{\Delta}(\mathbf{s}_i) \mathbb{E}[\tau(\text{proj}_{\mathcal{B}}(\mathbf{s}_i)) | Y, \theta], \\ \Sigma_{\tau^{\text{PROJ}}|Y} &= \frac{1}{\left(\sum_{i=1}^{n_T+n_C} \mathbb{I}^{\Delta}(\mathbf{s}_i)\right)^2} \sum_{i=1}^{n_T+n_C} \sum_{j=1}^{n_T+n_C} \mathbb{I}^{\Delta}(\mathbf{s}_i) \mathbb{I}^{\Delta}(\mathbf{s}_j) \text{Cov}[\tau(\text{proj}_{\mathcal{B}}(\mathbf{s}_i)), \tau(\text{proj}_{\mathcal{B}}(\mathbf{s}_j)) | Y, \theta].\end{aligned}\tag{16}$$

The posterior expectations and covariances in (16) can be obtained as in (6), but using the projected sentinels. Note that  $\tau^{\text{PROJ}}$  is in the class of weighted mean estimands (9), with weight function  $w_{\mathcal{B}}(\mathbf{b}) = \sum_{i=1}^{n_T+n_C} \mathbb{I}^{\Delta}(\mathbf{s}_i) \delta(\mathbf{b} - \text{proj}_{\mathcal{B}}(\mathbf{s}_i))$ , where  $\delta$  is the Dirac delta function.

The resulting estimator has desirable properties: densely populated regions receive proportionately more sentinels, but wigglier segments of the border do not. While it lacks the information efficiency of the inverse-variance estimator, the projected estimand is easier to understand and interpret, and may feel more familiar to practitioners used to finite-population inference. The averaging is over the observed units, although with their locations projected to the border.

In our experience, the choice of  $\Delta$  does not have a large effect on the estimate yielded by (16). A reasonable heuristic is to set  $\Delta$  to a small multiple of the Gaussian process lengthscale  $\ell$ . It should be noted that this choice only affects the location of sentinels on the border, the Gaussian process always gives low unit weights (11) to units far away from the border.

### 3.5 Projected Land ATE

In certain applications, population-based estimands can be undesirable, especially if the locations at which measurements are made are not representative of the population of interest. In such cases, geography-weighted estimands can be more natural. See [Antonelli et al. \(2016\)](#) for a discussion of this distinction in the context of preferential sampling. Remember that the “geometry-based” estimand  $\tau^{\text{UNIF}}$  places uniform

weights along the border. Instead, the “geography-based” projected land ATE estimand  $\tau^{\text{GEO}}$ , illustrated in Figure 2(b), begins by placing uniform weights on the treatment and control areas  $\mathcal{A}_T$  and  $\mathcal{A}_C$  that are within distance  $\Delta$  of the border  $\mathcal{B}$ , but then projects them onto the border to derive border weights. In other words, the projection method from  $\tau^{\text{PROJ}}$  is applied to an infinite population of uniform density on both sides of the border, instead of the finite population of observed units.

We denote the border vicinity area by  $\mathcal{A}_\Delta$ , defined as all points  $s$  such that  $s \in \mathcal{A}_T \cup \mathcal{A}_C$ , and  $\text{dist}_{\mathcal{B}}(s) < \Delta$ . To estimate  $\tau^{\text{GEO}}$ , a tight grid  $G^\nu$  of evenly spaced points separated by  $\nu$  is first generated covering  $\mathcal{A}_\Delta$ . Denote the number of grid points by  $L_\nu$ . Each point  $G_l^\nu$ ,  $l = 1, \dots, L_\nu$  in  $G^\nu$  is then projected onto the border to become a sentinel. The treatment effect at these positions is then estimated as before, yielding a mean vector and covariance matrix akin to (6). The mean of the mean vector then gives an estimate of  $\tau^{\text{GEO}}$ . In other words,  $\tau^{\text{GEO}}$  is estimated by applying the  $\tau^{\text{UNIF}}$  procedure with sentinels obtained by projecting the grid points, instead of equispaced sentinels.  $\tau^{\text{GEO}}$  remains in the category of weighted-mean estimands, with the weight function  $w_{\mathcal{B}}(\mathbf{b})$  in (9) proportional to the area of  $\mathcal{A}^T$  and  $\mathcal{A}^C$  that  $\mathbf{b}$  is nearest to, which can be written as the limit as the grid spacing goes to zero of point masses at the grid locations projected onto the border:

$$w_{\mathcal{B}}(\mathbf{b}) = \lim_{\nu \rightarrow 0} \frac{1}{L_\nu} \sum_{l=1}^{L_\nu} \mathbb{I}\{\mathbf{b} = \text{proj}_{\mathcal{B}}(G_l^\nu)\}, \quad (17)$$

where  $\mathbb{I}$  is the indicator function that returns one if its argument is true and zero otherwise.

For certain applications, it may be desirable to further restrict  $\mathcal{A}_\Delta$  to only certain types of land, for example residential areas in social studies, or farmland in agricultural studies. However, it is important to note that  $\tau^{\text{GEO}}$  is never interpretable as the average treatment effect in the vicinity of the border, that is  $\tau^{\text{GEO}} \neq \int_{\mathcal{A}_\Delta} \tau(s) d s$ . Estimating the latter estimand would require predicting the conditional regression function at grid locations within the treatment or control region using only observations on the *other* side of the border, which increases the extent of extrapolation required and thus makes the analysis more vulnerable to model misspecification.

### 3.6 Projected Super-Population ATE

Lastly, the purely geographical estimand  $\tau^{\text{GEO}}$  can be modified by weighting the grid points  $G_l^\nu$ ,  $l = 1, \dots, L_\nu$  by the population density  $\rho(G_l^\nu)$ . This gives the projected superpopulation ATE  $\tau^{\text{POP}}$ . Similarly to the density-weighted ATE  $\tau^\rho$ , estimating  $\tau^{\text{POP}}$  requires an estimate of the density  $\rho(G_l^\nu)$  at every grid point. As before, the uncertainty in the estimate of  $\rho$  should in principle be propagated to the estimate of  $\tau^{\text{POP}}$ , which generally will make the posterior distribution of  $\tau^{\text{POP}}$  neither normal nor analytically tractable.

The estimand  $\tau^{\text{POP}}$  can be interpreted as giving equal weight to each unit in the superpopulation of units

within the border vicinity  $\mathcal{A}_\Delta$ , but then moving each unit from its original location to its nearest location along the border, where the GeoRDD setting allows for the estimation of the treatment effect without undue extrapolation, and finally averaging the treatment effect of each unit in this displaced superpopulation.

### 3.7 Wiggly Border Simulation

We illustrate the above ATE estimators with a simulation. 200 units are placed in a square area delimited by spatial coordinates  $S_1 \in \{0, 2\}$  and  $S_2 \in \{-1, 1\}$ . A border at  $S_2 = 0$  divides units vertically into a control and treatment region, which are then further divided horizontally at  $S_1 = 0.5$  and  $S_1 = 1.5$  into three bands:

- The leftmost band  $S_1 < 0.5$  has a weak treatment effect.
- The middle band  $0.5 \geq S_1 < 1.5$  has a much lower population density, and a stronger treatment effect.
- The rightmost band  $S_1 \geq 1.5$ , has a much higher population density, and a very strong treatment effect.

Furthermore, the border in the leftmost band is a triangular wave, to create “wiggleness.” We increase the number of wiggles from 0 to 1000 to observe the effect on the estimates. The simulation setting is summarized in Table 1. We draw a single set of spatial coordinates, shown in Figure 3, then draw 10,000 simulations of the outcomes  $Y$  from a Gaussian process with squared exponential kernel ( $\ell = 0.4$ ,  $\sigma = 0.5$ ). To units above the border we add a treatment effect  $\tau(S_1, S_2) = S_1$ .

Table 1: Summary of wiggly border simulation setup.

	Left $S_1 < 0.5$	Middle $0.5 \geq S_1 < 1.5$	Right $1.5 \geq S_1$
Border	wiggly	straight	straight
Density	low $\rho = 1.0$	very low $\rho = 0.3$	high $\rho = 2.0$
$\tau$	weak	medium	strong

We fit the Gaussian process model (3), using the known hyperparameters of the covariance kernel and a weak prior on the mean parameter of each region, and estimate the average treatment effect using the six methods proposed above. In Figure 4(a) we show, for each estimator, the corresponding estimand and average posterior mean estimate evolving as the number of border wiggles increases. The behavior of the posterior standard deviation is shown in Figure 4(b).

As the border is a straight line and  $\mathcal{A}^T$  and  $\mathcal{A}^C$  are rectangles, and as the treatment effect does not depend on the vertical axis  $S_2$ , the density-weighted estimand  $\tau^\rho$  equals the projected superpopulation estimand

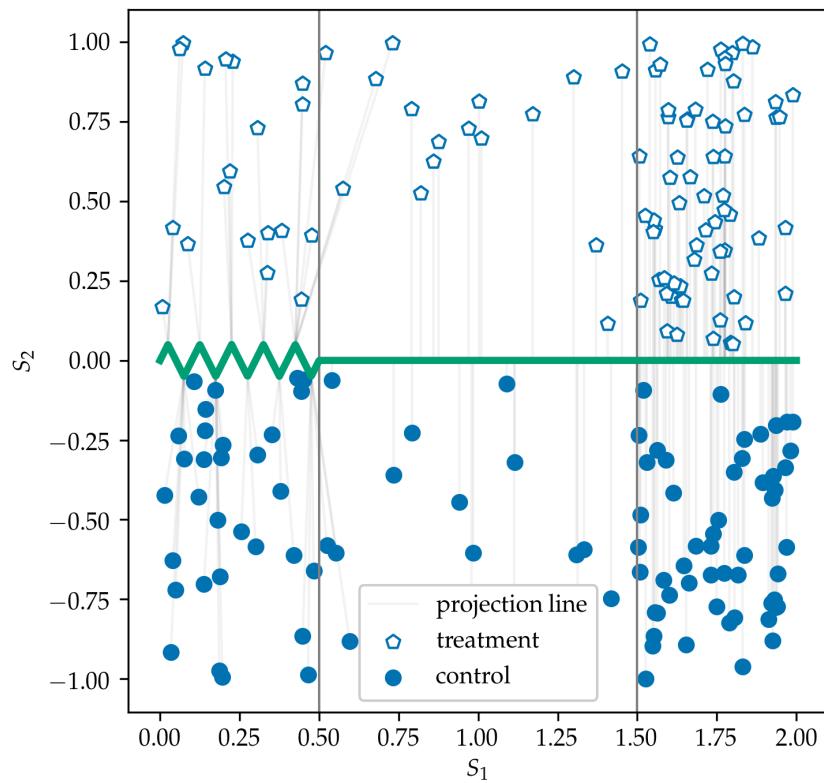


Figure 3: Spatial positions of units and border for the wiggly border simulation of Section 3.7. Projection lines for the projected finite population ATE are shown in light gray.

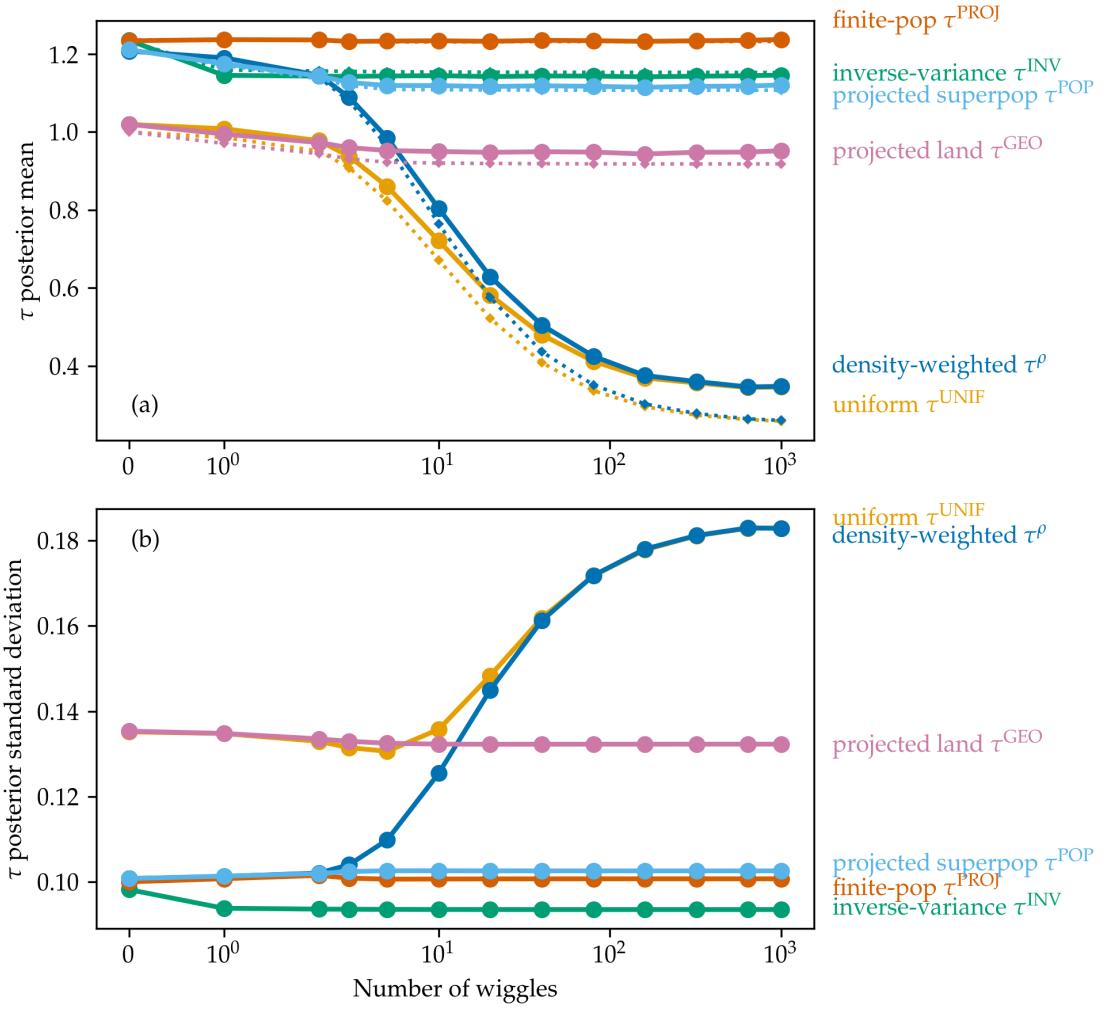


Figure 4: Results of the simulations of Section 3.7, showing for each ATE estimator as the leftmost section of the border gets wigglier (a) the estimate (posterior mean) averaged over 10,000 simulations with the corresponding estimand shown as a dotted line of the same color, and (b) the posterior standard deviation.

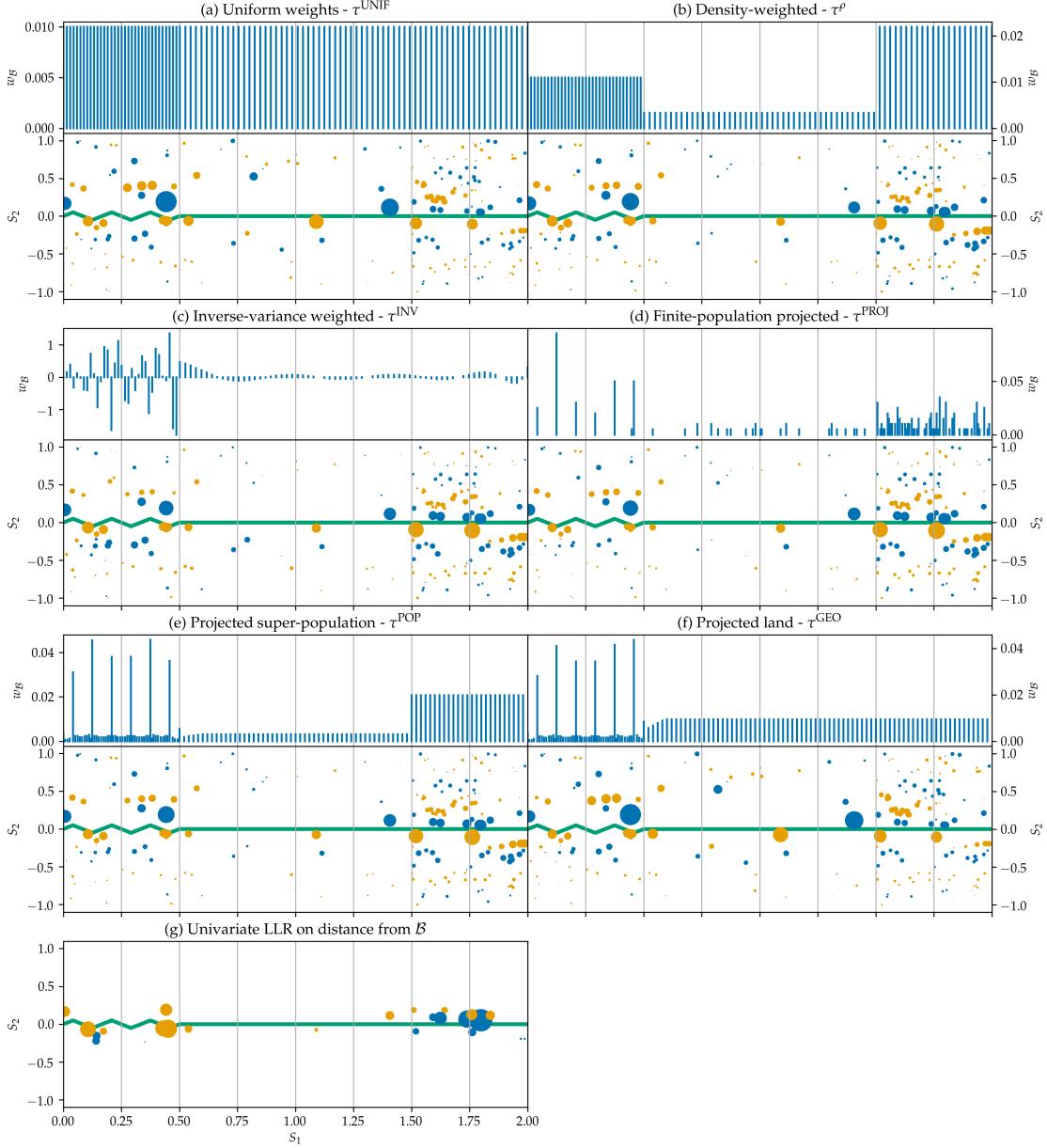


Figure 5: Weight functions and induced weights on the observations for the six weight functions proposed in this paper. The weight function plots show the weight  $w_B(b)$  against each sentinel's  $S_1$  coordinate. Sentinels with coinciding or nearly coinciding (within 0.005 of each other) coordinate  $S_1$  were merged and their weights summed. The induced weight plots show a circle for each unit, with the area of the circle proportional to its weight ( $w_T$  and  $w_C$  for treatment and control units respectively), and colored in blue for positive weights and orange for negative weights.

$\tau^{\text{POP}}$ , and they are in fact both equal to the infinite-population average treatment effect. Correspondingly, the posteriors of  $\tau^\rho$  and  $\tau^{\text{POP}}$  are identical. With 200 units,  $\tau^{\text{POP}}$  and the finite-population projected ATE  $\tau^{\text{PROJ}}$  are also similar, but the latter has the advantage of not require estimating the population density.

The geometry- and geography-based ATE  $\tau^{\text{UNIF}}$  and  $\tau^{\text{GEO}}$  are also equivalent when the border is a straight line. They give equal weight to the sparsely populated middle band, which produces a lower estimate with higher variance than the posteriors of  $\tau^\rho$  and  $\tau^{\text{POP}}$ .

Lastly, the information-based inverse-variance estimand  $\tau^{\text{INV}}$  does not coincide with any others. The estimand and mean estimate change slightly from 0 to 1 wiggles, but remains stable thereafter, demonstrating the robustness of this estimator to border topology. Weighting by the inverse variance gives the lowest posterior variance within the class of ATEs under consideration, which can indeed be seen in Figure 4(b).

As we introduce wiggles into the leftmost band,  $\tau^\rho$  and  $\tau^{\text{UNIF}}$  show their susceptibility to the border topology. Proportionally more sentinels are packed into the leftmost section of the border, upweighting the lower treatment effect of that band, and resulting in a drop of the two estimates and estimands. Meanwhile,  $\tau^{\text{INV}}$  remains stable despite the wiggles, because the additional sentinels in the leftmost band get automatically downweighted as their correlation rises. The estimators that rely on projection  $\tau^{\text{PROJ}}$ ,  $\tau^{\text{GEO}}$ , and  $\tau^{\text{POP}}$  also remain stable, because the projected sentinels hardly move. These robust estimands show only a slight displacement when the first wiggles are introduced, caused by the presence of some sentinels nearer to the observed units.

In Figure 5(a-f), we illustrate the behavior of border weights  $w_B(b)$  and unit weights ( $w_T$  and  $w_C$ ) in this simulation setting with 3 wiggles. Note how evenly spaced sentinels (for  $\tau^{\text{UNIF}}$ ,  $\tau^\rho$ , and  $\tau^{\text{INV}}$ ) are more densely packed along  $S_1$  in the leftmost area because of the zig-zagging border. The inverse-variance weighted estimator border weights can be seen to respond to this change in the border topology, though it is difficult to interpret their oscillating behavior. While these border-weights look unreasonable and unstable, the induced unit weights for  $\tau^{\text{INV}}$  are well-behaved, and in fact quite similar to those of the projected finite- and infinite-population estimators. Furthermore, note that all estimators can give some small negative weights  $w_T$  to treatment units, and small positive weights  $w_C$  to control units. For Gaussian processes, this can be understood in terms of the negative side-lobes of the equivalent kernel (see [Rasmussen and Williams \(2006\)](#) Section 2.6). The high variance of  $\tau^{\text{UNIF}}$  and  $\tau^{\text{GEO}}$  manifests itself as large weights given to a small number of units. All other estimators spread the weights more evenly amongst the units near the border, which reduces their variance.

For comparison, the weights placed on units by the projected 1D RDD are shown in Figure 5(g). A triangular kernel in  $S_2$  was used with bandwidth selected using the MSE-minimizing method proposed by [Imbens and Kalyanaraman \(2012\)](#). The Projected 1D RDD estimator can also be written as a linear

Table 2: Summary of average treatment effect estimator and estimand properties.

Symbol	Description	$\mathcal{B}$ Topology	Sentinels	Principle	Variance
$\tau^{\text{UNIF}}$	Uniform ATE	Sensitive	Equispaced	Geometry	High
$\tau^{\rho}$	Density-weighted ATE	Sensitive	Equispaced	Population	Low
$\tau^{\text{INV}}$	Inverse-var. weighted ATE	Robust	Equispaced	Information	Lowest
$\tau^{\text{PROJ}}$	Projected finite pop. ATE	Robust	Projected Units	Finite-pop.	Low
$\tau^{\text{GEO}}$	Projected land ATE	Robust	Projected Grid	Geography	High
$\tau^{\text{POP}}$	Projected superpop. ATE	Robust	Projected Grid	Population	Low

combination of the observed outcomes (11), and the unit weight vectors can be derived as:

$$\begin{aligned} \mathbf{w}_T &= \mathbf{X}_b (\mathbf{X}_T^\top \mathbf{W}_T \mathbf{X}_T)^{-1} \mathbf{X}_T^\top \mathbf{W}_T, \text{ and} \\ \mathbf{w}_C &= -\mathbf{X}_b (\mathbf{X}_C^\top \mathbf{W}_C \mathbf{X}_C)^{-1} \mathbf{X}_C^\top \mathbf{W}_C, \end{aligned} \quad (18)$$

where  $\mathbf{X}_b \equiv (1 \ 0)$ ,  $\mathbf{X}_T$  is the  $n_T \times 2$  design matrix with the first column filled with ones and the second column containing the distance from the border of each treatment unit, and  $\mathbf{W}_T$  is an  $n_T \times n_T$  diagonal matrix where the  $i^{\text{th}}$  diagonal element is the triangular kernel evaluated on the  $i^{\text{th}}$  unit's distance from the border. The  $\mathbf{X}_C$  and  $\mathbf{W}_C$  matrices are analogously defined for control units. By construction, the unit weights drop to zero outside of the support of the kernel. Within the support, Projected 1D RDD can also give negative weights to treatment units, and positive weights to control units. This results from the negative influence on the prediction  $\widehat{y}^*$  at  $x^*$  that univariate linear regression can give to an observation  $Y_i$  at  $X_i$  sufficiently far away on the opposite side of the mean  $\bar{X}$  of all observations. Strikingly, almost all of the positive weights are given to units in the rightmost treatment area that are closest to the border, and almost all the negative weights are given to units in the leftmost control area. Consequently, any trend in the outcomes across  $S_1$  would confound the estimated treatment effect.

In most applications, we recommend the use of the finite population or inverse-variance-weighted estimators, to prevent the undesirable influence of border topology. The projected finite population method is simplest to understand and interpret in the tradition of finite population estimators, and unlike the projected superpopulation estimator  $\tau^{\text{POP}}$  it does not require estimating population density. Meanwhile, the inverse-variance estimator is the most efficient (lowest posterior variance) weighted mean estimator, and avoids the potential complication of the choice of a distance cutoff for projected units.

## 4 Testing for Non-Zero Effect

Once we have obtained the “cliff face” estimate (6) and estimated an average treatment effect, we might also naturally wonder whether we can claim to have detected a significant treatment effect at the border. In the

hypothesis testing framework, we have two possible choices of null hypotheses. The **sharp null** specifies that the treatment effect is zero everywhere along the border:  $\tau(\mathcal{B}) = 0$ , while the **weak null** only requires the average treatment effect to be zero.

## 4.1 Marginal Likelihood Test

To target the sharp null hypothesis, we first define a parametric null model  $\mathcal{M}_0$ , specified as a single Gaussian process spanning the control and treatment regions, with the same kernel and hyperparameters obtained in the 2GP procedure.  $\mathcal{M}_0$  is smooth and continuous at the border, and therefore accords with the sharp null hypothesis. Intuitively, if there is a treatment effect, the likelihood of the observations should be lower under  $\mathcal{M}_0$  than under  $\mathcal{M}_1$ , the 2GP model as specified in equation (3). We therefore choose the difference in log-likelihoods as our test statistic

$$t = \log \mathbb{P}(Y_T, Y_C | \mathcal{M}_1) - \log \mathbb{P}(Y_T, Y_C | \mathcal{M}_0) \quad (19)$$

and wish to reject the sharp null hypothesis when its observed value  $t_{\text{obs}}$  is high.

A parametric bootstrap approach is used to quantify what “high” means. We draw  $Y_T^*, Y_C^*$  from  $\mathcal{M}_0$ , using the same spatial locations as in the original data, and then fit the two competing models to the simulated data in order to obtain the bootstrapped test statistic

$$t^* = \log \mathbb{P}(Y_T^*, Y_C^* | \mathcal{M}_1) - \log \mathbb{P}(Y_T^*, Y_C^* | \mathcal{M}_0) \quad (20)$$

Repeating this procedure, we obtain a distribution of  $t$  under  $\mathcal{M}_0$ , which we can then compare to the observed  $t$ . More precisely, we can interpret the proportion of  $t^*$  drawn above  $t_{\text{obs}}$  as a p-value.

$$p^{\text{lik}} = \mathbb{P}(t^* > t_{\text{obs}} | \mathcal{M}_0) \quad (21)$$

Computationally, because the hyperparameters and locations of the units are held constant during the bootstrap, we can reuse the Cholesky decomposition of the covariance matrix, allowing the test to be performed in seconds even with hundreds of units and thousands of bootstrap samples.

## 4.2 “Chi-squared” Test

The likelihood-based sharp null above is valid and easy to understand. But it may seem odd that the test aims to detect a non-zero treatment effect at the border, without any explicit reference to the border  $\mathcal{B}$ . The

test statistic and p-values can be computed without access to the sentinel positions, using only the treatment and control indicators. If the test is significant, there is no guarantee that this is due to a discontinuity at the border.

To address this oddity, we can derive a test statistic directly from the cliff face estimator (6). We will use  $\mu$  and  $\Sigma$  as shorthand for the posterior mean  $\mu_{b_{1:R}|Y}$  and covariance matrix  $\Sigma_{b_{1:R}|Y}$  throughout this section. If a  $k$ -vector  $y$  is distributed  $\mathcal{N}(\mu, \Sigma)$ , with mean vector  $\mu$  unknown and covariance  $\Sigma$  known, then under the null hypothesis that  $\mu = 0$ , the test statistic  $y^\top \Sigma^{-1} y$  has distribution  $\chi_k^2$ . See for example Rencher (2003) Section 5.2.2 for a classical derivation of this test. This suggests that we could use  $S^2 = \mu^\top \Sigma^{-1} \mu$  as a test statistic, and obtain a p-value from a  $\chi_R^2$  distribution function evaluated at  $S^2$ , where  $R$  is the number of sentinels. However, we face two problems. Firstly, this test, obtained heuristically from a Bayesian posterior by analogy with the classical multivariate normal result, is not a valid frequentist test. Secondly, while  $\Sigma$  is mathematically full-rank, it is typically numerically rank-deficient. Therefore,  $R$  overestimates the true degrees of freedom of the null distribution.

Benavoli and Mangili (2015), developing a test for function equality, address the second problem by trimming the  $\Sigma$  eigenvalues  $\lambda_i$  lower than  $\epsilon \sum_{j=1}^k \lambda_j$ , with  $\epsilon$  a pre-specified small number (they use 0.01). They address the first problem by showing that the resulting p-value is always conservative in their simulations. However, in our work, we found the resulting p-value to be sensitive to the arbitrarily chosen  $\epsilon$  tolerance parameter, which makes it difficult to trust its validity.

We therefore again take the parametric bootstrap approach, this time using  $S^2$  as the test statistic instead of the likelihood ratio. With  $B$  bootstrap samples, the p-value is obtained as

$$p = \frac{1}{B} \sum_{t=1}^T \mathbb{I}\left\{S_{(b)}^2 < S^2\right\}, \quad (22)$$

$$S_{(b)}^2 = (\mu_{(b)})^\top \Sigma^{-1} \mu_{(b)}$$

where  $\mu_{(b)}$  is the result of applying (6) to  $Y_T^{(b)}$  and  $Y_C^{(b)}$ , themselves drawn from  $\mathcal{M}_0$  at the same locations as the observations  $Y_T$  and  $Y_C$ .

Because calculating  $S^2$  involves inverting a matrix  $\Sigma$  that is mathematically of full rank, but numerically of low rank, we may worry about the numerical stability of computing  $S$ . We verified in simulated examples that regularizing  $\Sigma$  by adding a small constant to its diagonal does not greatly affect the computed  $S^2$ . The parametric bootstrap ensures the frequentist validity of the test regardless of the regularization.

### 4.3 Inverse-Variance Weighted Test

We now turn to the weak null hypothesis that the average treatment effect along the border is zero, but otherwise allowing the treatment effect to be positive along some parts of the border and negative in others. As we saw in Section 3, the “average” treatment effect can be defined in multiple ways. If we choose the inverse-variance weighted mean, then  $\tau^{\text{INV}}$  has posterior given by (15). While the posterior is a Bayesian object, we can use it heuristically to derive a pseudo-p-value

$$\begin{aligned} Z_0 &\sim \mathcal{N}(0, \Sigma_{\tau^{\text{INV}}|Y}) \\ p^{\text{INV}} &= \mathbb{P}\left(|Z_0| > |\mu_{\tau^{\text{INV}}|Y}|\right) \\ &= 2\Phi\left(-\frac{|\mu_{\tau^{\text{INV}}|Y}|}{\sqrt{\Sigma_{\tau^{\text{INV}}|Y}}}\right) \end{aligned} \tag{23}$$

This “p-value” obtained from the Bayesian posterior may not have good frequentist properties. In particular, there is no guarantee that under the null hypothesis,  $p^{\text{INV}}$  is below 0.05 less than 5% of the time. To turn it into a valid frequentist test, it can be calibrated using the same parametric bootstrap approach that was used for the likelihood and  $\chi^2$  tests, now using  $\mu_{\tau^{\text{INV}}|Y}$  as the test statistic. The calibration can also be achieved analytically, since  $\mu_{\tau^{\text{INV}}|Y}$  is normally distributed under the null hypothesis. We derive the analytical calibration of the inverse-variance test in Appendix D. Note that the p-value for all three tests defined in this section is derived under the same parametric null model  $\mathcal{M}_0$ , which accords with both the sharp null and weak null hypotheses. The calibrated inverse-variance test targets the weak null hypothesis in the sense that the test statistic is an estimate of the ATE (and thus the test is sensitive to deviations of the ATE from zero) rather than its p-value being derived under the weak null (like the classical t-test for example).

### 4.4 Power in Simulated Example

The three tests we developed leverage different aspects of the problem, and target two different null hypotheses. One may wonder how their power compares in the presence of a treatment effect. Considering once more the border between Louisiana and Mississippi, we imagine an experiment where the unit of analysis is the county, located at its centroid, as shown in Figure 6. We will simulate outcomes from a single Gaussian process covering both states. For simplicity, we fix the hyperparameters to arbitrary values:  $\sigma_\epsilon = \sigma_{\text{GP}} = 1.0$  and  $\ell = 100 \text{ km}$ . We then add a constant treatment effect  $\tau$  to all the outcomes in Louisiana. The results of the three tests proposed so far are shown in the first three rows of Table 3 for  $\tau = 0$  (null

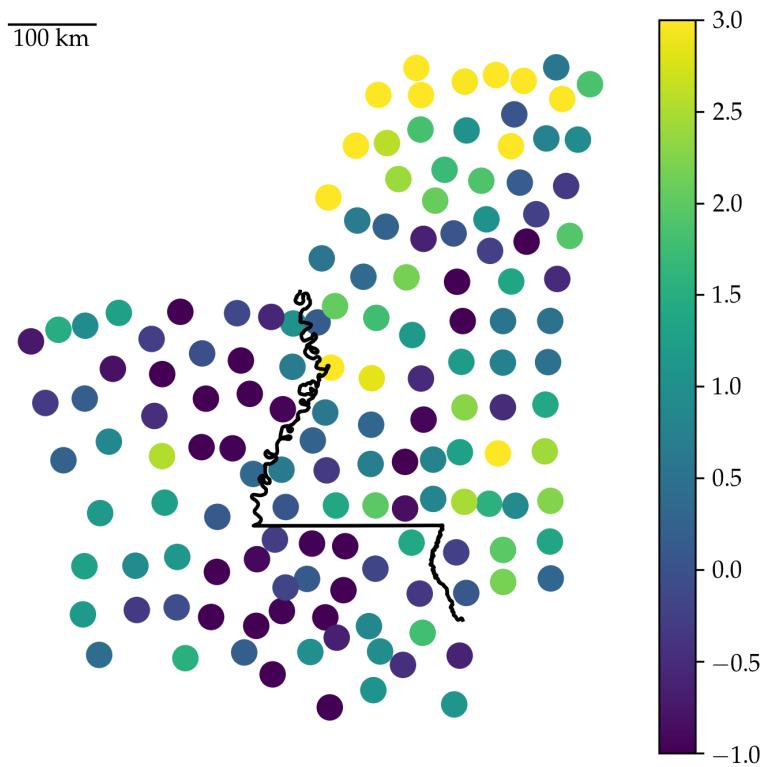


Figure 6: Set-up of the imaginary experiment in Louisiana and Mississippi. Each unit is at the centroid of a county. The colors indicated the observed outcomes in one draw of the simulation under  $\tau = 1.5$ . In this particular run, the p-values were 0.0016, 0.0018, and 0.0013 for the mLL,  $\chi^2$ , and inverse-variance test respectively.

Table 3: Power of marginal likelihood, chi-squared, and inverse-variance tests, with nominal significance of  $\alpha = 0.05$ , under null and alternative hypothesis for simulated outcomes at the centroids of Louisiana and Mississippi counties.

Test	Power under	
	$\tau = 0$	$\tau = 1.2$
Marginal log-likelihood	0.048	0.656
$\chi^2$	0.047	0.635
Inverse-variance $\Sigma^{-1}$	0.085	0.866
Bootstrap-calibrated $\Sigma^{-1}$	0.050	0.799
Analytically calibrated $\Sigma^{-1}$	0.051	0.801

hypothesis) and  $\tau = 1.2$  and significance level  $\alpha = 0.05$ .

We see that under the null, the  $\chi^2$  and likelihood ratio tests are valid (rejection of the null in 5% of simulations up to simulation error). This is enforced by the parametric bootstrap, which draws test statistics from the same null distribution to calibrate the tests. However, the p-values for the inverse-variance test are biased down, so that we will falsely reject the null 6.7% instead of 5% of the time. While unfortunate, this is unsurprising, since the inverse-variance test was derived heuristically rather than from a rigorous frequentist procedure.

After calibration, the hypothesis test based on the inverse-variance mean is valid, but retains higher power to detect the constant treatment effect than the mLL and  $\chi^2$  tests. This can lead to a paradox: we may reject the weak null hypothesis, but fail to reject the sharp null hypothesis (using the  $\chi^2$  or likelihood test), even though rejection of the weak null logically implies rejection of the sharp null. This paradox isn't specific to this setting, and is discussed in depth in the context of randomization-based inference by Ding (2014). To maximize power, we therefore recommend using the calibrated inverse-variance test in studies where the main interest is in the detection of an overall (average) increase or decrease in outcomes.

## 4.5 Placebo Tests

Gaussian process models are almost always misspecified. We do not believe that the Gaussian process with stationary squared exponential kernel is the true data-generating process, although we hope that the model is sufficiently flexible to represent reality well. Under misspecification, we should be skeptical of results that rely on the truth of the model specification. We therefore encourage practitioners to probe the validity of the above hypothesis tests by running a “placebo” test. A placebo test repeatedly applies the hypothesis test on data that are known to have zero treatment effect (a “placebo”), in order to verify that the returned p-values are uniformly distributed. In our spatial setting, we will use the treatment and control regions separately as placebo groups. Within each placebo group, we repeatedly draw an arbitrary

geographical border, creating new treatment and control groups. Here we drew lines that split the placebo units in half at a sequence of angles  $1^\circ, 2^\circ, 3^\circ, \dots, 180^\circ$  counter-clockwise from horizontal, each positioned so that half of the units fall on either side of the line in order to maximize power. Because the border was chosen arbitrarily by us, we should not expect there to be a discontinuous jump in outcomes at this border. We then apply the bootstrapped likelihood test procedure described above to this arbitrarily divided data, store the results, and hope to obtain a roughly uniform distribution of p-values. The resulting p-values will obviously be highly correlated, so we should only expect a very roughly uniform distribution (because of the small effective sample size), but at the very least, this procedure allows us to visually verify that the p-values are not blatantly biased.

## 5 Example: NYC School Districts

We illustrate the analysis of geographical regression discontinuity designs using house sales data from New York City. The city publishes information pertaining to property sales within the city in the last 12 months on a rolling basis. This includes the sale price, building class, and the address of the property. Public schools in the city are all part of the City School District of the City of New York, but the city-wide district is itself divided into 32 sub-districts. It is a common belief that school districts have an impact on real estate price, as parents are willing to pay more to live in districts with better schools. We therefore ask: can we measure a discontinuous jump in house prices across the borders separating school districts?

### 5.1 Preprocessing

In order to model the property sale prices, we need to obtain their locations. We geocode the address of each sale by merging the sales with NYC's Pluto database, which contains X and Y coordinates for each house, identified by its borough, zip code, block and lot. These coordinates are given in the EPSG:2263 projection in units of feet, which we also adopt. For addresses that do not find a match in Pluto, we use Google's geocoding API to obtain a latitude and longitude, which we then project to EPSG:2263.

We then filter the sales data, by removing sales (1) without a reported sale price, (2) outside of the family homes building class categories (one, two, and three family dwellings); (3) missing the square footage or other covariates; (3) without a location due to failed geocoding; (4) smaller than 100 sq ft, and (5) outliers with log price per square foot less than 3 or more than 8. We exclude condos and coops because only very few sales report square footage alongside the price.

We display the resulting dataset of 19,578 sales recorded between January 8, 2015 and December 7, 2016

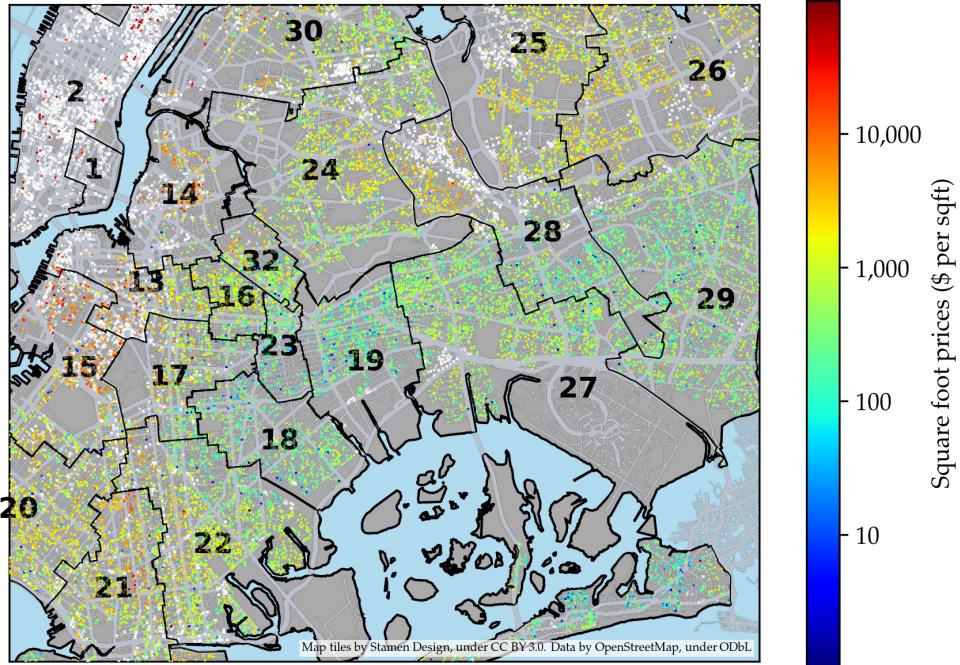


Figure 7: Map of property sales in New York City. Each dot is a sale, and its color indicates the price per square foot. White crosses indicate sales of properties with missing square footage, which are therefore excluded from the analysis. School district boundaries are shown, and each district is labeled by its number.

in Figure 7. We also show the 27,394 residential properties that have missing square footage information. Properties with missing square footage are almost all coops and condos, which explains the clustering of missing data in areas of higher density.

## 5.2 Model for Property Prices

The outcome of interest is price per square foot. As is commonly done in analyses of real estate prices, we take its logarithm to reduce the skew of the outcome. The complete model is then a Gaussian process within each of  $J_{\text{Distr}}$  districts over the geographical covariates  $s$ , super-imposed with a linear regression on the property covariates (provided by the  $L_{\text{BuildClass}}$  building categories encoded as dummy variables).

This model can be written:

$$\begin{aligned}
Y_i &= \mu_{Distr[i]} + \beta_{BuildClass[i]} + f_{Distr[i]}(s_i) + \epsilon_i \\
\epsilon_i &\sim \mathcal{N}(0, \sigma_y^2) \\
\mu_j &\sim \mathcal{N}(0, \sigma_\mu^2), j = 1, \dots, J_{Distr} \\
f_j &\sim GP(0, k(s, s')), j = 1, \dots, J_{Distr} \\
\beta_l &\sim \mathcal{N}(0, \sigma_\beta^2), l = 1, \dots, L_{BuildClass} \\
k(s, s') &= \sigma_{GP}^2 \exp\left\{-\frac{(s - s')^\top (s - s')}{2\ell^2}\right\}
\end{aligned} \tag{24}$$

A visual inspection of the house sales map in Figure 7 suggests examining the border between districts 19 and 27. We arbitrarily designate district 19 as the “treatment” area and district 27 as the “control” area. Importantly, the border between the two districts is also part of the border between Brooklyn and Queens, so we will not be able to attribute a difference in price solely to the causal effect of the school districts. This is an instance of what Keele and Titiunik (2015) term “compound treatments,” a frequent concern in GeoRDDs. Here we are therefore just *measuring* a discontinuity in the house prices at the district. Attributing the discontinuity to a particular cause (school district or borough) is not directly supported by the data.

Another concern is units sorting around the border, which would violate the identification assumptions for GeoRDDs. If people move across the border to live in a better school district, does this invalidate the analysis? We take the view that the unit of analysis here is the tract of land on which houses are built, rather than the residents themselves. If a district becomes more attractive, people may move to it, whereas land does not move but its price adjusts. A sale gives a snapshot of the price of the land, made more accurate by correcting in our model (24) for covariates that pertain to the building rather than land. Note that of course, the limited covariates provided by the data cannot fully capture the value of the building. For example, the wealthier residents who inhabit the more desirable school districts may also have more funds available to maintain and enhance their home, which will drive up the property’s resale value. Since it is not captured by the available covariates, this added value is folded into the treatment effect by our analysis.

The histogram in Figure 8 of outcomes  $Y$  in both districts also shows that marginally the house prices are very different. Our goal is to establish whether this difference is measurable at the border, and not merely an underlying trend that spans both districts.

We fit the hyperparameters  $\sigma_\beta$ ,  $\sigma_{GP}$ ,  $\ell$  and  $\sigma_\epsilon$  by optimizing the marginal log-likelihood of the data within neighboring school districts 18, 19, 23, 24, 25, 26, 27, 28, and 29. We hold  $\sigma_\mu$  fixed to 20 to give the

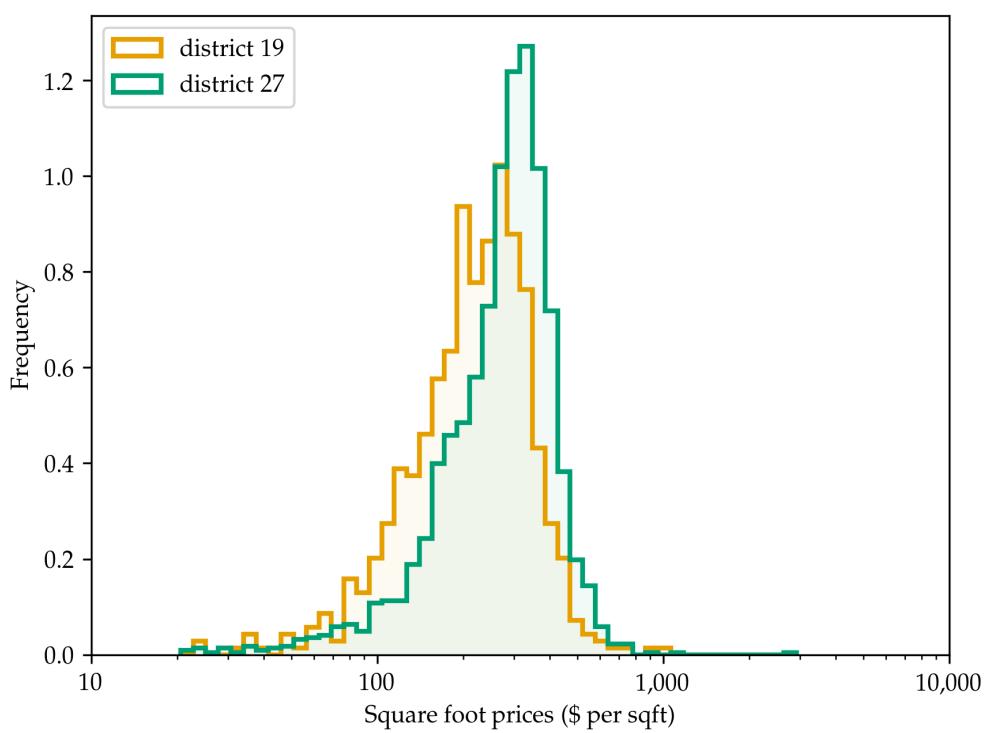


Figure 8: Histogram of log sale prices per square foot in NYC school districts 19 and 27.

district means  $\mu_j$ ; a fairly uninformative prior. The fitted hyperparameters were  $\widehat{\sigma}_\epsilon = 0.4020$ ,  $\widehat{\sigma}_{GP} = 0.1955$ ,  $\widehat{\sigma}_\beta = 0.1465$ , and  $\widehat{\ell} = 4482$  ft.

### 5.3 Cliff Face Estimator

We seek the treatment effect function  $\tau(\mathcal{B})$  between the two districts. We could proceed by computing the cliff face estimator with covariates (8). But to simplify the analysis as discussed in Section 2.4, we can instead obtain the posterior means of the  $\beta_{1j}$  and  $\beta_{2j}$  coefficients (following the procedure outlined in Appendix C, but extended to  $J_{Distr}$  rather than just two areas in accordance with (24)), extract the residuals  $\mathbf{Y}_T - \mathbf{D}_T \hat{\beta}$  and  $\mathbf{Y}_C - \mathbf{D}_C \hat{\beta}$ , which we then treat as the observed outcomes in a GeoRDD analysis with no non-spatial covariates. In this example, we find that the posterior variance of  $\beta$  is low, and therefore the two approaches yield very similar results, but conditioning on the estimate of  $\beta$  is computationally convenient. We therefore proceed with this two-step approach.

Following the inference procedure outlined in Section 2.2, we obtain the posterior distribution of the cliff height  $\tau(\mathcal{B})$  obtained at the sentinel locations. The cliff face is shown in Figure 9, and shows that  $\tau(\mathcal{B})$  is estimated as negative everywhere along the border, which corresponds to higher property prices in district 27. However, the credible envelope is fairly wide, especially in the southern section of the border, so we cannot visually rule out the null hypothesis that  $\tau(\mathcal{B}) = 0$ .

The “cliff face” can also be visualized directly in Figure 10 as the difference between the two log-price mean surfaces  $g(s)$ . The figure also gives a better sense of the spatial variation in prices captured by the model.

### 5.4 Average Log-Price Increase

The cliff face Figure 9 shows a negative treatment effect everywhere along the border, which can be averaged by the estimators we developed in Section 3. Our two recommended estimators, based on inverse-variance weighting and finite-population projection, yield ATE estimates of  $-0.19$  and  $-0.18$  respectively, which corresponds to an almost 20% increase in property prices going from district 19 to district 27. All ATE estimators from Section 3 applied to this setting are shown in Table 5.4. In this example the different estimators yield similar answers, as the border is fairly straight and short relative to the fitted lengthscale.

### 5.5 Significant Difference in Price?

The inverse-variance weighted mean treatment effect estimated suggests a significant treatment effect. But the posterior tail probability cannot be interpreted as a p-value. For this, we turn to the three tests developed

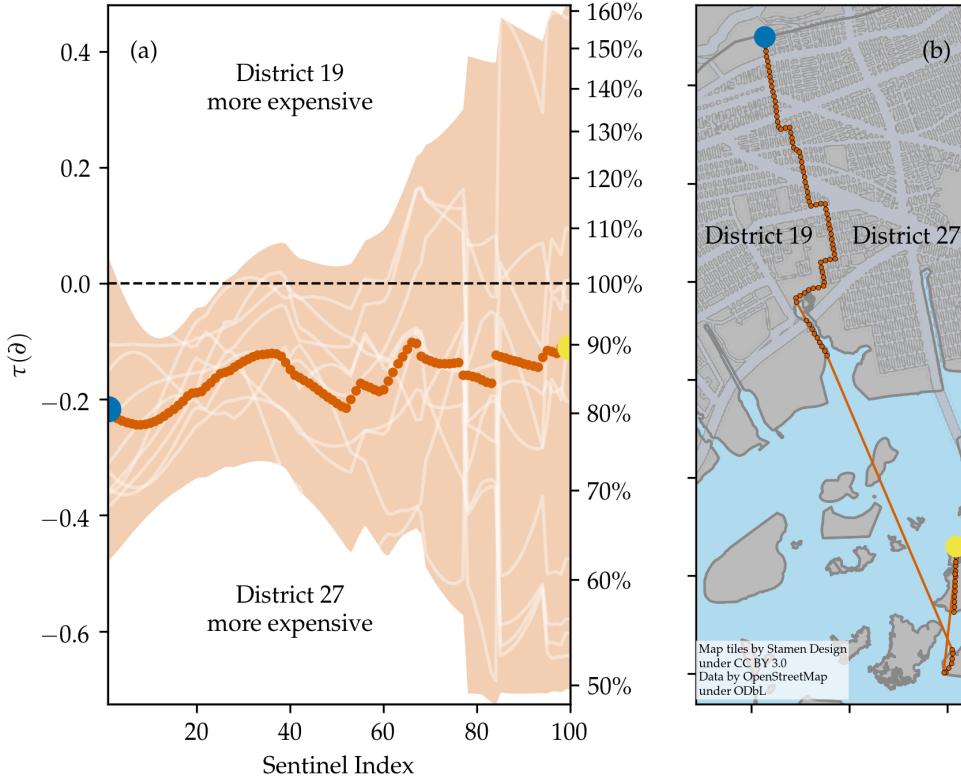


Figure 9: (a) Cliff face estimator for the school district effect on house prices per square foot between district 27 and district 19, with 95% credible envelope. The left y-axis is in the scale of log prices per square foot; positive values correspond to houses near the border being more expensive in district 19 than 27. The right y-axis shows the corresponding ratio of the price of a house near the border in district 19 over its price in district 27. A few draws from the posterior are shown in lighter color to show the posterior correlations between sentinels. Note the jumps from sentinels 77 to 78, and 84 to 85, which correspond to the school district border jumping from Long Island to islands in Jamaica Bay (Black Wall Island and then Rulers Bar Hassock). These islands are sparsely populated, and the posterior standard deviation is correspondingly much higher. (b) The map of sentinels, evenly spaced along the border between school districts 27 and 19. The northernmost sentinel—shown as a blue circle in both plots—has index 1, while the last sentinels—shown in yellow—is on Rulers Bar Hassock.

Estimand	Posterior		
	Mean	Standard Dev.	Tail Prob.
$\tau^{\text{UNIF}}$	-0.17	0.09	2.93%
$\tau^{\rho}$	-0.19	0.06	0.04%
$\tau^{\text{INV}}$	-0.19	0.06	0.03%
$\tau^{\text{PROJ}}$	-0.18	0.08	1.31%
$\tau^{\text{GEO}}$	-0.16	0.09	3.99%
$\tau^{\text{POP}}$	-0.18	0.06	0.08%

Table 4: Average difference in log price per square foot between school districts 19 and 27. For each ATE estimand, we show the mean and standard deviation of its posterior distribution, and the tail probability  $\mathbb{P}(\tau > 0 | Y, \hat{\beta}, \theta)$  of the average treatment being greater than zero. Negative ATEs correspond to district 27 being more expensive.

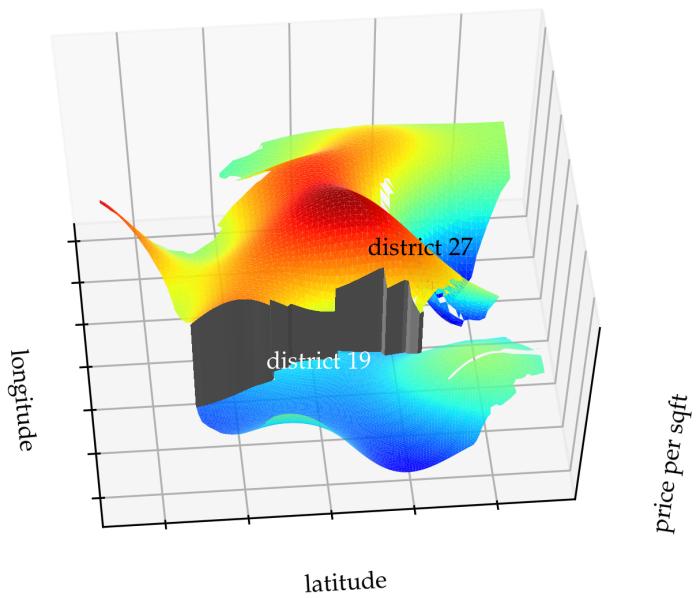


Figure 10: NYC surface plot viewed from the West. The grey cliff face connects the fitted price surfaces of districts 19 and 27, and has height given by (6) and shown in Figure 9.

Test	p-value
$\chi^2$ bootstrap	0.012
mLL bootstrap	0.002
$\tau^{\text{INV}}$ uncalibrated	0.0007
$\tau^{\text{INV}}$ calibrated	0.002

Table 5: Results of hypothesis tests for New York school district house prices example.

in Section 4. In applied settings, running multiple tests should be done with care, but as we are proposing this new methodology, we apply all three tests in order to gain insight into their differences. Their results are found in Table 5.5.

The three tests reject the null hypothesis that  $\tau(\mathcal{B}) = 0$  along the border between districts 19 and 27. This will not always be the case, as the calibrated inverse-variance test has higher power than the other two tests. The  $\chi^2$  test had the lowest power in the simulated example of Section 4.4, and here also returns the highest p-value.

To assess the validity of the three tests, we apply the placebo tests devised in Section 4.5. Within each district, we split the data in half by a line at angles  $1^\circ, 3^\circ, 5^\circ, 6^\circ, \dots, 179^\circ$ . Because these lines were drawn arbitrarily, we don't expect a discontinuous treatment effect between the two halves, and so we hope to see a uniform distribution of placebo p-values. However, these tests will be highly correlated, and so the low effective sample size could lead to some apparent departures from uniformity. There is in fact visible autocorrelation in the graphs of placebo p-values as a function of angle.

The mLL placebo p-values show a pronounced bias towards low values. This seems to confirm our concern that the marginal log-likelihood may be sensitive to features of the data other than the discontinuity at the border. In particular, model misspecification, which is a concern in spatial models, makes the interpretation of the mLL test unreliable. Based on this vulnerability, and its manifestation in this example, we do not recommend relying on the likelihood-ratio test.

The  $\chi^2$  test shows more robustness, with Figure 11(d) showing some negative bias in district 27, and some positive bias in district 19, which could simply be due to the low effective sample size. We therefore believe that the  $\chi^2$  test will continue to be reliable under misspecification. It is only due to its low power that we hesitate to recommend its use in applications where the treatment effect is expected to be fairly homogenous.

Lastly, the calibrated inverse-variance placebo p-values display no obvious bias, with Figure 11(f) close to uniformly distributed, and Figure 11(e) showing a lower auto-correlation than the mLL and  $\chi^2$  tests. The high power and robustness of the inverse-variance test make a strong case for its use in most applications.

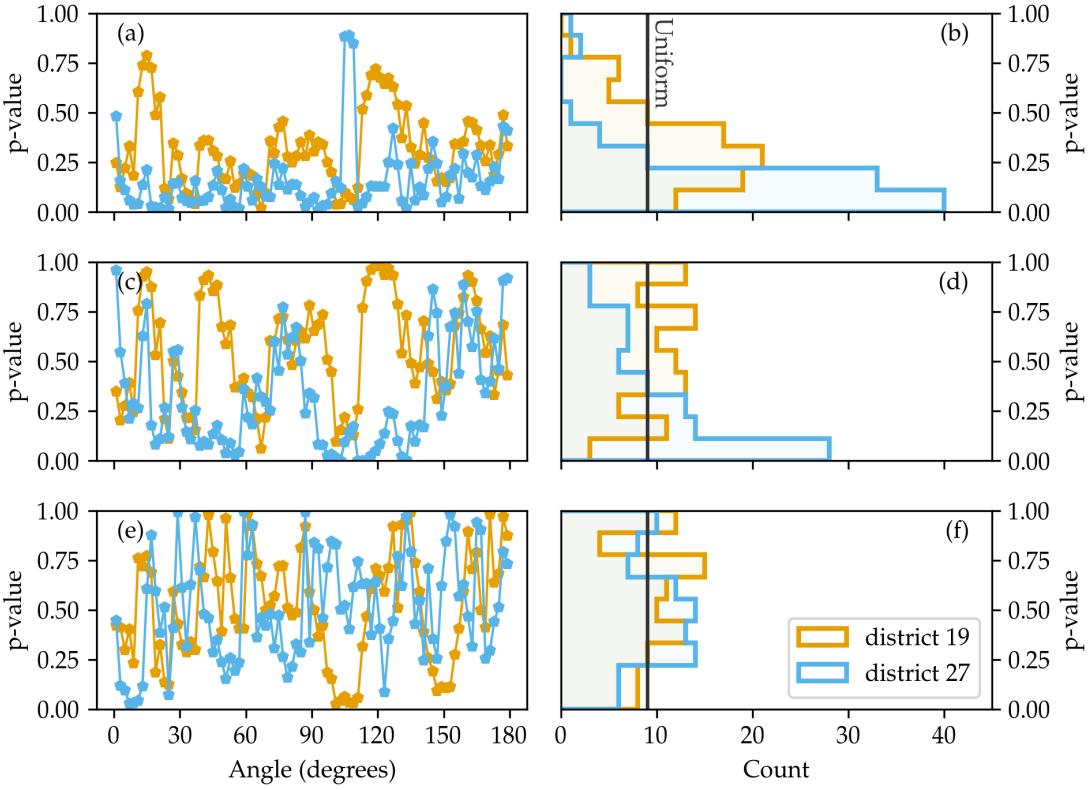


Figure 11: Placebo tests for significance tests applied to NYC school district house price example, applied within districts 19 and 27. The three rows respectively show results for the marginal log-likelihood bootstrap test, chi-squared bootstrap test, and calibrated inverse-variance test. The first column shows the placebo p-value as a function of the border angle, in order to visualize the auto-correlation of the placebo tests. The second column shows histograms of the placebo p-values, with the black vertical line indicating the uniform distribution for comparison.

## 6 Pairs of School Districts

The GeoRDD analysis can be repeated for each pair of adjacent districts. Figure 12 and Table 8 give an overview of the results by showing the posterior mean and standard deviation of the inverse variance ATE estimated at each border. Significant effects are found between many districts, but interpreting the results requires some caution. We have already mentioned the issue of compound treatments for borders between school districts that overlap with the border between boroughs. School districts 19, 32, and 14 are in Brooklyn, while districts 30, 24, and 27 are in Queens.

Some school districts are separated by parks (or other non-residential zones), for example districts 15 & 17 or 19 & 24, so that house sales do not extend all the way to the border on one or both sides. A significant treatment effect between these pairs cannot be interpreted as the detection of a discontinuity in prices at the border, let alone any kind of causal interpretation, but rather it means that the difference in prices between the two sides of the park exceeds the typical spatial variation of house prices expected over the same distance. This is not unsurprising, and one may speculate that physical barriers like parks, rivers, railways and major roads can separate neighborhoods with distinct character, demographics and thus house prices. This in turn challenges the stationarity assumption of the spatial model (3). The higher distance between data and the border also stretches the spatial model's ability to extrapolate, which makes it more vulnerable to model misspecification. And it violates the assumptions

22 . 140275816016985

Other pairs of district, like 13 & 14, 13 & 17, and 25 & 28 have clusters of missing data (condo sales with unknown square footage) near the border that cast doubt on the interpretation of the estimated effect. Nonetheless, significant effects are also found between pairs of school districts without issues due to compound treatments, physical barriers, or missing data. House prices increase going across the border from districts 16 to 13, 18 to 17, 24 to 30, 23 to 17, 25 to 26, 28 to 29, and 29 to 26. Overall, it seems that school district borders in Brooklyn and Queens can correspond to measurable jumps in house prices per square foot. The estimated size of this effect varies: zero or negligible in some cases, such as between districts 15, 20, 21, and 22; and quite pronounced in others, such as a 20% price increase from 29 to 26, or 22% from 18 to 17.

## 7 Conclusion

Geographic regression discontinuity designs (GeoRDDs) arise when a treatment is assigned to one region, but not to another adjacent region. For outcomes that vary spatially, a direct comparison of mean outcomes

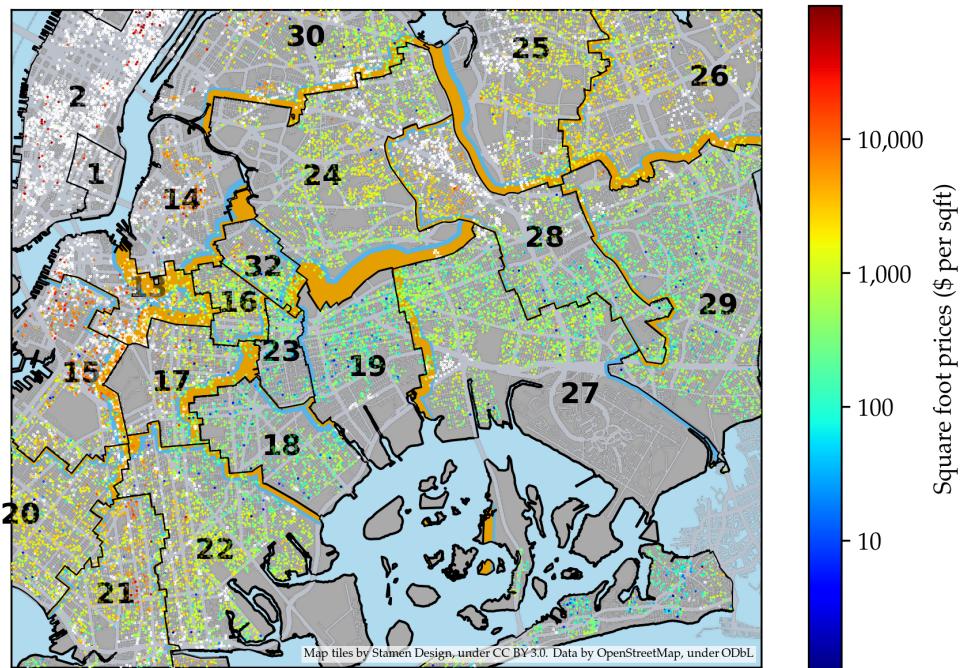


Figure 12: Pairwise estimates of the inverse variance ATE between adjacent districts. The thickness of the orange buffer adjacent to borders is proportional to the posterior mean of the inverse variance ATE, and the blue buffer beyond it is proportional to the posterior standard deviation of the ATE. The buffers are drawn on the side of the border that is estimated to have higher house prices.

between  $Y_T$  and  $Y_C$ , such as a t-test, is an invalid estimator of the treatment effect, as it is confounded by the spatial covariates. However, under smoothness assumptions, units adjacent to the border are comparable, and form a natural experiment. The same idea underpins causal interpretations of one-dimensional regression discontinuity designs (1D RDDs), where a single “forcing” variable controls the treatment assignment instead of a border separating two geographical regions. We use this similarity to motivate a general framework for the analysis of GeoRDDs. One-dimensional methods can be abstracted to three steps: (1) fit a smooth function on either side of the threshold; (2) extrapolate the functions to the threshold; and (3) take the difference of the two extrapolations to estimate the treatment effect at the threshold. For GeoRDDs we propose to (1) fit a smooth surface on either side of the border; (2) extrapolate the surfaces to the border; and (3) take the difference of the two extrapolations to estimate the treatment effect along the border.

Previous research has focused on extending methods developed for 1D RDDs to GeoRDDs. In applied settings, some have used the signed distance from the border as the forcing variable in a 1D RDD, but the resulting estimator is still spatially confounded. In this paper, our aim was to recognize the importance of the geographical aspect of the problem, and therefore draw from the spatial statistics literature, which brings a rich set of tools designed specifically to model and exploit spatial correlations to obtain more powerful inference. We therefore used Gaussian process regression, known as kriging in the spatial statistics literature, to fit the smooth surfaces to the outcomes in step (1) of our general framework. Our approach yields a multivariate normal posterior distribution of the treatment effect for a collection of “sentinel” locations along the border.

Defining an “average treatment effect” estimand turns out to have surprising pitfalls. Simply integrating the treatment effect uniformly along the border yields an estimand that is inefficient and undesirably sensitive to the topology of the border. More sophisticated estimands, summarized in Table 2, are robust to this effect, and use the information available in the data more efficiently.

There are multiple valid approaches to hypothesis testing against the null hypothesis of zero treatment effect along the border. We recommend the use of the calibrated inverse-variance test, derived from the posterior distribution of the inverse-variance ATE estimator. It has generally high statistical power and behaved well in placebo tests in the NYC empirical example.

We applied our method to a publicly available dataset of one year of New York City property sales, to examine whether school district cause difference in property prices. Focusing on the border between school districts 19 and 27, we estimated a roughly 20% average increase (inverse variance ATE) in house prices per square foot when crossing the border from district 19 to district 27. However, the border between these two districts is also the border between the NYC boroughs of Brooklyn and Queens, so we cannot

attribute this difference to the causal effect of the school districts. Other limitations apply to many pairs of districts. Parks, commercial zones, railways, and major roads can separate neighborhoods, keep data away from the borders, break the stationarity assumption of the spatial model, and increase the amount of extrapolation performed by the model, which casts doubt on the legitimacy of the estimated treatment effects. Missing data from condo sales which do not report square footage can also distort estimated effects. Some pairs of adjacent school district remain where a large effect was found, such as a the 22% increase in house prices per square footage when crossing the border from district 18 to 17. Overall, it seems that school district borders in Brooklyn and Queens are often accompanied by a discontinuity in house prices, but the causal attribution of this difference to the reputation of the school districts is often questionable due to the aforementioned geographical and political complications.

The main limitation of our approach to GeoRDDs is the reliance on modeling assumptions. We modeled the response surfaces as two independent Gaussian processes, with iid normal noise for each observation. As is common in spatial statistics, we use Gaussian processes as non-parametric smoothing devices used to capture spatial correlations, but do not think of them as truthful approximations to the data generating mechanism. We believe care must be taken not to lean heavily on this modeling assumption. In particular, we recommend that hypothesis tests always be accompanied by placebo tests: by applying the same procedure on data where no treatment was applied, we can verify that the test behaves appropriately under the null hypothesis despite any potential model misspecification. We also assumed a stationary covariance structure, with hyperparameters equal in the treatment and control regions, and in particular we chose the squared exponential kernel. This kernel makes smoothness assumptions that are often considered unrealistic in spatial settings, and so the Matérn covariance family is often recommended as a more robust alternative. The assumption of equal covariance parameters in the two areas can also be relaxed, by separately tuning the parameters within each area.

Because of the need to extrapolate the fitted processes a short distance to the border, our GeoRDD method may be vulnerable to the limitations of Gaussian processes when extrapolating. The distinction between interpolation and extrapolation of spatial models is explored in some depth in Stein (2012). We expect that methodological advances that improve the extrapolating behavior of Gaussian processes would also improve the robustness of our method. For example, Wilson and Adams (2013) develop spectral mixture (SM) covariance kernels with good extrapolating behavior, which could be applied beneficially to GeoRDDs. However, SM kernels are motivated by time series with some periodic or oscillatory behavior, which is more unusual in spatial applications, and may therefore not be as well-suited for use with GeoRDDs.

The use of Gaussian process regression to analyse GeoRDDs gives flexibility and extensibility to the

method. This presents many opportunities for future research, inspired by the past and future development of methods in spatial statistics and machine learning that are based on Gaussian processes. In spatial statistics, kriging has been used as the foundation for a plethora of spatial models, which may be adapted for the purposes of analyzing GeoRDDs. Banerjee et al. (2014) provides a good introduction to the richness of the spatial statistics field. For example, if the outcomes are binary, proportions, or counts, then binomial or Poisson likelihoods could be substituted for the iid normal likelihood used in this paper. Besides ATE and hypothesis testing, another potential question is whether the treatment effects are homogenous or heterogeneous. Hypothesis tests of homogeneity would be an interesting possible extension of our framework.

The framework and techniques of this paper could also be extended to spatio-temporal settings. If the treatment is only applied to the treatment region after a time  $t_*$ , one could envision a three-dimensional RDD consisting of the geographical border in the spatial dimensions, and a straight line through  $t_*$  in the temporal dimension. The only necessary modification to our approach is that the Gaussian process model would need to be augmented with a temporal component, for example with an anisotropic squared exponential covariance function. Longitudinal studies could also be handled by such an approach, with the addition of random intercepts for each unit. We leave spatio-temporal RDDs using Gaussian process models to future research.

## A Spatial Confounding of 1D RDD Applied to GeoRDD

Analysing GeoRDDs by using the signed distance from the border as a forcing variable in a 1D RDD can lead to spatial confounding. We demonstrate this with a simple artificial example, depicted in Figure 13.

Suppose we have units in a 2D square, with spatial coordinates  $s_1 \in [-1, 1]$ , and  $s_2 \in [-1, 1]$ , and with a horizontal border at  $s_2 = 0$  separating a treatment region from a control region. Let us assume the null hypothesis, with outcomes driven only by  $s_1$  (parallel to the border), given by  $Y_i = \alpha s_{1i} + \epsilon_i$ , where  $\epsilon_i$  is an iid noise term  $\epsilon_i \sim \mathcal{N}(0, \sigma_\epsilon^2)$ . Lastly, let us consider the situation where the density  $\rho(\mathbf{s})$  of units is different in each quadrant of the square:

$$\begin{aligned}
\rho(\mathbf{s}) &= 2\rho_0, \text{ where } s_1 < 0, s_2 > 0 && \text{(top left)} \\
\rho(\mathbf{s}) &= \rho_0, \text{ where } s_1 > 0, s_2 > 0 && \text{(top right)} \\
\rho(\mathbf{s}) &= 2\rho_0, \text{ where } s_1 > 0, s_2 < 0 && \text{(bottom right)} \\
\rho(\mathbf{s}) &= \rho_0, \text{ where } s_1 < 0, s_2 < 0 && \text{(bottom left)}
\end{aligned} \tag{25}$$

The projection RDD then considers a 1D RDD along  $s_2$ . The usual RDD estimand (1) can be obtained an-

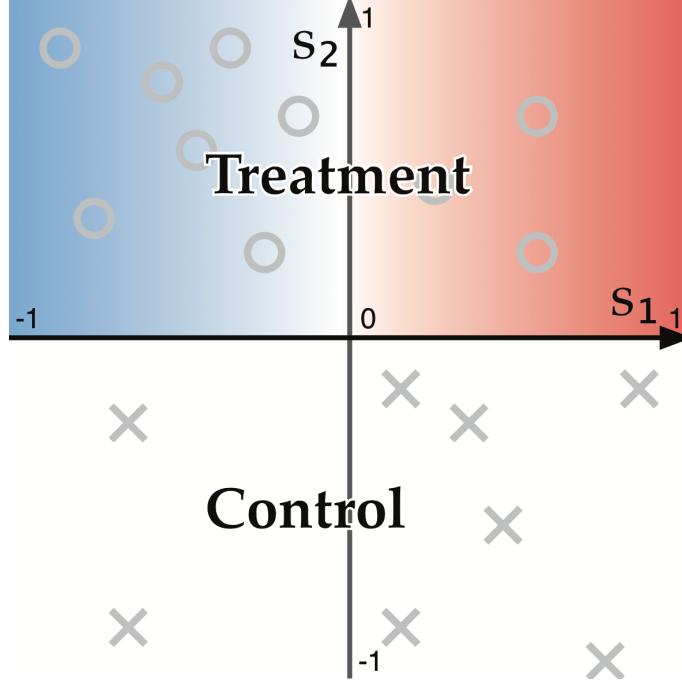


Figure 13: A theoretical example illustrating the susceptibility of the projected 1D RDD method to spatial confounding. The locations of treatment and control units are shown with circles and crosses respectively, separated by a border at  $s_2 = 0$ . Units are denser in the upper left and lower right quadrants. The treatment effect, depicted as a gradient from blue to red, increases linearly with  $s_1$ .

alytically, and equals  $\tau = \frac{-\alpha}{3}$ , despite assuming the null hypothesis. This is because  $s_1$  acts as a hidden confounder, whose distribution changes discontinuously at the border, which leads to bias and inconsistency in the projected 1D RDD estimate. In geographical settings, a discontinuous change in the density of units at the border is not unusual: for example a border could run alongside a park, or a small body of water, therefore with zero population density on one side of the border. A visual inspection of Figure 7 showing the locations of units in a New York City property sales dataset reveals many examples of this.

## B Covariances for Gaussian Process Model

All covariances below are conditional on the hyperparameters  $\theta \equiv (\ell, \sigma_{GP}, \sigma_\epsilon, \sigma_\mu)$ , omitted for concision.

$$\begin{aligned}
m_T, m_C &\sim \mathcal{N}(0, \sigma_\mu^2) \\
\text{Cov}(Y_{iT}, m_T) &= \sigma_\mu^2 \\
\text{Cov}(Y_{iC}, m_C) &= \sigma_\mu^2 \\
\text{Cov}(Y_{iT}, m_C) &= \text{Cov}(Y_{iC}, m_T) = 0 \\
\text{Cov}(Y_{iT}, f_T(s')) &= k(s_i, s') \\
\text{Cov}(Y_{iC}, f_C(s')) &= k(s_i, s') \\
\text{Cov}(Y_{iT}, f_C(s')) &= \text{Cov}(Y_{iC}, f_T(s')) = 0 \\
\text{Cov}(Y_{iT}, Y_{jT}) &= \text{Cov}(Y_{iC}, Y_{jC}) = \sigma_\mu^2 + k(s_i, s_j) + \delta_{ij} \sigma_\epsilon^2 \\
\text{Cov}(Y_{iT}, Y_{jC}) &= 0
\end{aligned} \tag{26}$$

We further define some shorthand notation.

Symbol	Size	$i,j$ th entry
$\mathbf{K}_{BB}$	$R \times R$	$\sigma_m^2 + k(b_i, b_j)$
$\mathbf{K}_{BT}$	$R \times n_T$	$\sigma_m^2 + k(b_i, s_{jT})$
$\mathbf{K}_{BC}$	$R \times n_C$	$\sigma_m^2 + k(b_i, s_{jC})$
$\mathbf{K}_{TT}$	$n_T \times n_T$	$\sigma_m^2 + k(s_{iT}, s_{jT})$
$\mathbf{K}_{CC}$	$n_C \times n_C$	$\sigma_m^2 + k(s_{iT}, s_{jC})$
$\Sigma_{TT}$	$n_T \times n_T$	$\sigma_m^2 + k(s_{iT}, s_{jT}) + \delta_{ij} \sigma_\epsilon^2$
$\Sigma_{CC}$	$n_C \times n_C$	$\sigma_m^2 + k(s_{iT}, s_{jC}) + \delta_{ij} \sigma_\epsilon^2$

## C Posterior Mean of $\beta$

We derive the posterior mean of the linear regression coefficients vector  $\beta$  for the model specified in (7).

$$\begin{aligned}
\mathbf{Y} &\equiv \begin{pmatrix} \mathbf{Y}_T \\ \mathbf{Y}_C \end{pmatrix} \\
\Sigma_{Y|\beta} &\equiv \text{Cov}(\mathbf{Y} | \beta) = \begin{bmatrix} \Sigma_{TT} & 0 \\ 0 & \Sigma_{CC} \end{bmatrix} \quad \text{conditional variance of } Y \\
\text{Cov}(\beta) &= \sigma_\beta^2 \mathbf{I}_p \quad \text{prior variance of } \beta \\
\Sigma_Y &\equiv \text{Cov}(\mathbf{Y}) = \Sigma_{Y|\beta} + \sigma_\beta^2 \mathbf{D}\mathbf{D}^\top \quad \text{prior variance of } Y \\
\hat{\beta} &= \sigma_\beta^2 \mathbf{D}^\top \Sigma_Y^{-1} \mathbf{Y} \quad \text{posterior mean of } \beta
\end{aligned} \tag{27}$$

The treatment and control residuals can then be obtained respectively as  $\mathbf{R}_T = \mathbf{Y}_T - \mathbf{D}_T \hat{\beta}$  and  $\mathbf{R}_C = \mathbf{Y}_C - \mathbf{D}_C \hat{\beta}$ . Conditionally on  $\hat{\beta}$ , the residuals  $\mathbf{R}_T$  and  $\mathbf{R}_C$  then have independent multivariate normal distributions with the same mean and covariance as  $\mathbf{Y}_T$  and  $\mathbf{Y}_C$  respectively.

## D Calibration of Inverse-variance Test

We seek to obtain a valid hypothesis test against the null hypothesis of zero treatment effect everywhere along the border by using the inverse-variance weighted ATE estimate obtained in Section 3.3 as a test statistic.

Under the parametric null hypothesis  $\mathcal{M}_0$ ,  $\mathbf{Y}_T$  and  $\mathbf{Y}_C$  are drawn from a single Gaussian process, with no discontinuity at the border. Their joint covariance is

$$\text{Cov} \left( \begin{pmatrix} \mathbf{Y}_T \\ \mathbf{Y}_C \end{pmatrix} \mid \mathcal{M}_0 \right) = \begin{bmatrix} \Sigma_{TT} & \mathbf{K}_{TC} \\ \mathbf{K}_{TC}^\top & \Sigma_{CC} \end{bmatrix} \tag{28}$$

where  $\mathbf{K}_{TC}$  is the  $n_T \times n_C$  matrix with  $ij^{\text{th}}$  entry equal to  $k(\mathbf{s}_{iT}, \mathbf{s}_{jC})$ . The predicted mean outcomes (5) at the sentinels  $\mu_{b_{1:R}|T}$  and  $\mu_{b_{1:R}|T}$  are obtained by left-multiplying  $\mathbf{Y}_T$  and  $\mathbf{Y}_C$  by matrices  $\mathbf{A}_T$  and  $\mathbf{A}_C$  (respectively) that are deterministic functions of the unit locations and the hyperparameters:

$$\begin{aligned}
\mathbf{A}_T &\equiv \mathbf{K}_{BT} \Sigma_{TT}^{-1}, \text{ and} \\
\mathbf{A}_C &\equiv \mathbf{K}_{BC} \Sigma_{CC}^{-1}.
\end{aligned} \tag{29}$$

Under  $\mathcal{M}_0$ , the joint distribution of  $\mu_{b_{1:R}|T}$  and  $\mu_{b_{1:R}|T}$  is consequently also multivariate normal with

mean zero and covariance given by:

$$\text{Cov}\left(\begin{pmatrix} \mathbf{A}_T Y_T \\ \mathbf{A}_C Y_C \end{pmatrix} \mid \mathcal{M}_0\right) = \begin{bmatrix} \mathbf{A}_T \Sigma_{TT} \mathbf{A}_T^\top & \mathbf{A}_T \mathbf{K}_{TC} \mathbf{A}_C^\top \\ (\mathbf{A}_T \mathbf{K}_{TC} \mathbf{A}_C^\top)^\top & \mathbf{A}_C \Sigma_{CC} \mathbf{A}_C^\top \end{bmatrix} \quad (30)$$

Continuing in this fashion, the cliff-face estimate  $\mu_{b_{1:R}|Y}$  (6) is yet another zero-mean multivariate normal with covariance given by:

$$\begin{aligned} \text{Cov}\left(\mu_{b_{1:R}|Y} \mid \mathcal{M}_0\right) &= \text{Cov}(\mathbf{A}_T Y_T - \mathbf{A}_C Y_C) \\ &= \mathbf{A}_T \Sigma_{TT} \mathbf{A}_T^\top + \mathbf{A}_C \Sigma_{CC} \mathbf{A}_C^\top - \mathbf{A}_T \mathbf{K}_{TC} \mathbf{A}_C^\top - (\mathbf{A}_T \mathbf{K}_{TC} \mathbf{A}_C^\top)^\top \end{aligned} \quad (31)$$

Weighted mean ATE estimators of the form defined in (10) are linear transformations of  $\mu_{b_{1:R}|Y}$  and so under  $\mathcal{M}_0$ , they are normally distributed with mean zero. For a weight function  $w_B(\mathbf{b})$ , its variance is given by

$$\begin{aligned} \text{var}\left(\mu_{\tau^w|Y} \mid \mathcal{M}_0\right) &= \text{Cov}\left(\frac{w_B(\mathbf{b}_{1:R})^\top \mu_{b_{1:R}|Y}}{\mathbf{1}_R^\top w_B(\mathbf{b}_{1:R})}\right) \\ &= \frac{w_B(\mathbf{b}_{1:R})^\top \text{Cov}(\mu_{b_{1:R}|Y}) w_B(\mathbf{b}_{1:R})}{(\mathbf{1}_R^\top w_B(\mathbf{b}_{1:R}))^2}. \end{aligned} \quad (32)$$

The p-value follows from treating the ATE estimate as a test statistic. Under the null hypothesis, the probability of  $\mu_{\tau^w|Y}$  exceeding in magnitude its observed value  $\mu_{\tau^w|Y}^{\text{obs}}$  is:

$$\mathbb{P}\left(\left|\mu_{\tau^w|Y}\right| > \left|\mu_{\tau^w|Y}^{\text{obs}}\right| \mid \mathcal{M}_0\right) = 2\Phi\left(-\frac{\left|\mu_{\tau^w|Y}^{\text{obs}}\right|}{\sqrt{\text{var}(\mu_{\tau^w|Y} \mid \mathcal{M}_0)}}\right). \quad (33)$$

The calibrated inverse-variance introduced in Section 4.3 test is the special case of this procedure where the weights are chosen to be  $w_B(\mathbf{b}_{1:R}) = \Sigma_{b_{1:R}|Y}^{-1} \mathbf{1}_R$ .

## E Tables

### E.1 Wiggly Border Simulation Results

$n_{\text{wiggles}}$	$\widehat{\tau^{\text{UNIF}}}$	$\tau^{\text{UNIF}}$	$\widehat{\tau^{\text{INV}}}$	$\tau^{\text{INV}}$	$\widehat{\tau^\rho}$	$\tau^\rho$	$\widehat{\tau^{\text{PROJ}}}$	$\tau^{\text{PROJ}}$	$\widehat{\tau^{\text{GEO}}}$	$\tau^{\text{GEO}}$	$\widehat{\tau^{\text{POP}}}$	$\tau^{\text{POP}}$
0	1.02 (0.14)	1.00	1.23 (0.10)	1.23	1.21 (0.10)	1.21	1.23 (0.10)	1.24	1.02 (0.14)	1.00	1.21 (0.10)	1.21
1	1.01 (0.13)	0.99	1.14 (0.09)	1.16	1.19 (0.10)	1.19	1.24 (0.10)	1.24	0.99 (0.13)	0.97	1.17 (0.10)	1.17
2	0.98 (0.13)	0.95	1.14 (0.09)	1.16	1.15 (0.10)	1.14	1.24 (0.10)	1.24	0.97 (0.13)	0.94	1.14 (0.10)	1.14
3	0.94 (0.13)	0.91	1.14 (0.09)	1.16	1.09 (0.10)	1.08	1.23 (0.10)	1.23	0.96 (0.13)	0.93	1.13 (0.10)	1.12
5	0.86 (0.13)	0.82	1.14 (0.09)	1.15	0.98 (0.11)	0.96	1.23 (0.10)	1.23	0.95 (0.13)	0.92	1.12 (0.10)	1.11
10	0.72 (0.14)	0.67	1.14 (0.09)	1.15	0.80 (0.13)	0.76	1.23 (0.10)	1.23	0.95 (0.13)	0.92	1.12 (0.10)	1.11
20	0.58 (0.15)	0.52	1.14 (0.09)	1.15	0.63 (0.14)	0.58	1.23 (0.10)	1.23	0.95 (0.13)	0.92	1.12 (0.10)	1.11
40	0.48 (0.16)	0.41	1.14 (0.09)	1.15	0.50 (0.16)	0.44	1.23 (0.10)	1.23	0.95 (0.13)	0.92	1.12 (0.10)	1.11
80	0.41 (0.17)	0.34	1.14 (0.09)	1.15	0.42 (0.17)	0.35	1.23 (0.10)	1.23	0.95 (0.13)	0.92	1.12 (0.10)	1.11
160	0.37 (0.18)	0.30	1.14 (0.09)	1.15	0.38 (0.18)	0.30	1.23 (0.10)	1.23	0.94 (0.13)	0.92	1.11 (0.10)	1.11
320	0.36 (0.18)	0.27	1.14 (0.09)	1.15	0.36 (0.18)	0.28	1.23 (0.10)	1.23	0.95 (0.13)	0.92	1.12 (0.10)	1.11
640	0.34 (0.18)	0.26	1.14 (0.09)	1.15	0.35 (0.18)	0.26	1.24 (0.10)	1.23	0.95 (0.13)	0.92	1.12 (0.10)	1.11
1000	0.35 (0.18)	0.26	1.15 (0.09)	1.15	0.35 (0.18)	0.26	1.24 (0.10)	1.23	0.95 (0.13)	0.92	1.12 (0.10)	1.11

Table 7: Posterior mean averaged over 10,000 simulations, posterior standard deviation and true value for each average treatment effect estimand as the wiggliness of the border is increased in the simulations of Section 3.7.

### E.2 NYC School District Estimated Treatment Effects

<b>13</b>	<b>14</b> : $-0.29 \pm 0.09$	<b>15</b> : $+0.03 \pm 0.07$	<b>16</b> : $-0.13 \pm 0.07$	<b>17</b> : $-0.26 \pm 0.08$				
<b>14</b>	<b>13</b> : $+0.29 \pm 0.09$	<b>16</b> : $-0.16 \pm 0.10$	<b>24</b> : $-0.38 \pm 0.15$	<b>32</b> : $-0.07 \pm 0.12$				
<b>15</b>	<b>13</b> : $-0.03 \pm 0.07$	<b>17</b> : $-0.18 \pm 0.10$	<b>20</b> : $+0.05 \pm 0.06$	<b>22</b> : $+0.24 \pm 0.11$				
<b>16</b>	<b>13</b> : $+0.13 \pm 0.07$	<b>14</b> : $+0.16 \pm 0.10$	<b>17</b> : $-0.04 \pm 0.07$	<b>23</b> : $-0.10 \pm 0.07$	<b>32</b> : $+0.05 \pm 0.06$			
<b>17</b>	<b>13</b> : $+0.26 \pm 0.08$	<b>15</b> : $+0.18 \pm 0.10$	<b>16</b> : $+0.04 \pm 0.07$	<b>18</b> : $-0.20 \pm 0.07$	<b>22</b> : $+0.06 \pm 0.07$	<b>23</b> : $-0.29 \pm 0.10$		
<b>18</b>	<b>17</b> : $+0.20 \pm 0.07$	<b>19</b> : $-0.06 \pm 0.12$	<b>22</b> : $+0.10 \pm 0.07$	<b>23</b> : $-0.03 \pm 0.09$				
<b>19</b>	<b>18</b> : $+0.06 \pm 0.12$	<b>23</b> : $-0.00 \pm 0.08$	<b>24</b> : $+0.39 \pm 0.11$	<b>27</b> : $+0.19 \pm 0.06$	<b>32</b> : $+0.27 \pm 0.12$			
<b>20</b>	<b>15</b> : $-0.05 \pm 0.06$	<b>21</b> : $+0.04 \pm 0.05$	<b>22</b> : $-0.11 \pm 0.08$					
<b>21</b>	<b>20</b> : $-0.04 \pm 0.05$	<b>22</b> : $-0.04 \pm 0.05$						
<b>22</b>	<b>15</b> : $-0.24 \pm 0.11$	<b>17</b> : $-0.06 \pm 0.07$	<b>18</b> : $-0.10 \pm 0.07$	<b>20</b> : $+0.11 \pm 0.08$	<b>21</b> : $+0.04 \pm 0.05$			
<b>23</b>	<b>16</b> : $+0.10 \pm 0.07$	<b>17</b> : $+0.29 \pm 0.10$	<b>18</b> : $+0.03 \pm 0.09$	<b>19</b> : $+0.00 \pm 0.08$	<b>32</b> : $-0.04 \pm 0.08$			
<b>24</b>	<b>14</b> : $+0.38 \pm 0.15$	<b>19</b> : $-0.39 \pm 0.11$	<b>25</b> : $+0.25 \pm 0.13$	<b>27</b> : $-0.22 \pm 0.10$	<b>28</b> : $+0.06 \pm 0.06$	<b>30</b> : $+0.14 \pm 0.05$	<b>32</b> : $+0.02 \pm 0.08$	
<b>25</b>	<b>24</b> : $-0.25 \pm 0.13$	<b>26</b> : $+0.08 \pm 0.04$	<b>28</b> : $-0.15 \pm 0.08$	<b>29</b> : $-0.06 \pm 0.10$	<b>30</b> : $-0.28 \pm 0.15$			
<b>26</b>	<b>25</b> : $-0.08 \pm 0.04$	<b>29</b> : $-0.18 \pm 0.05$						
<b>27</b>	<b>19</b> : $-0.19 \pm 0.06$	<b>24</b> : $+0.22 \pm 0.10$	<b>28</b> : $+0.04 \pm 0.04$	<b>29</b> : $-0.01 \pm 0.08$				
<b>28</b>	<b>24</b> : $-0.06 \pm 0.06$	<b>25</b> : $+0.15 \pm 0.08$	<b>27</b> : $-0.04 \pm 0.04$	<b>29</b> : $+0.09 \pm 0.04$				
<b>29</b>	<b>25</b> : $+0.06 \pm 0.10$	<b>26</b> : $+0.18 \pm 0.05$	<b>27</b> : $+0.01 \pm 0.08$	<b>28</b> : $-0.09 \pm 0.04$				
<b>30</b>	<b>24</b> : $-0.14 \pm 0.05$	<b>25</b> : $+0.28 \pm 0.15$						
<b>32</b>	<b>14</b> : $+0.07 \pm 0.12$	<b>16</b> : $-0.05 \pm 0.06$	<b>19</b> : $-0.27 \pm 0.12$	<b>23</b> : $+0.04 \pm 0.08$	<b>24</b> : $-0.02 \pm 0.08$			

Table 8: Estimated inverse variance ATE ( $\pm$  posterior standard deviation) for pairs of school districts in NYC. Each row contains ATEs estimated for one district compared to its neighbors. For example the first cell indicates an estimated average change difference log house prices per square foot going from district 13 to 14 of -0.29.

## References

- Antonelli, J., M. Cefalu, and L. Bornn, 2016: The positive effects of population-based preferential sampling in environmental epidemiology. *Biostatistics*, **17** (4), 764–778.
- Banerjee, S., B. P. Carlin, and A. E. Gelfand, 2014: *Hierarchical modeling and analysis for spatial data*. Crc Press.
- Branson, Z., M. Rischard, L. Bornn, and L. Miratrix, 2017: A nonparametric bayesian methodology for regression discontinuity designs. URL <https://arxiv.org/abs/1704.04858>, **1704.04858**.
- Chen, Y., A. Ebenstein, M. Greenstone, and H. Li, 2013: Evidence on the impact of sustained exposure to air pollution on life expectancy from china’s huai river policy. *Proceedings of the National Academy of Sciences*, **110** (32), 12 936–12 941.
- Cook, T. D., 2008: “waiting for life to arrive”: a history of the regression-discontinuity design in psychology, statistics and economics. *Journal of Econometrics*, **142** (2), 636–654.
- Crump, R. K., V. J. Hotz, G. W. Imbens, and O. A. Mitnik, 2009: Dealing with limited overlap in estimation of average treatment effects. *Biometrika*, **96** (1), 187–199.
- Ding, P., 2014: A paradox from randomization-based causal inference. URL <https://arxiv.org/abs/1402.0142>, **1402.0142**.
- Hahn, J., P. Todd, and W. Van der Klaauw, 2001: Identification and estimation of treatment effects with a regression-discontinuity design. *Econometrica*, **69** (1), 201–209.
- Imbens, G., and K. Kalyanaraman, 2012: Optimal bandwidth choice for the regression discontinuity estimator. *The Review of economic studies*, **79** (3), 933–959.
- Imbens, G., and T. Zajonc, 2011: Regression discontinuity design with multiple forcing variables. *Report, Harvard University.[972]*.
- Imbens, G. W., and T. Lemieux, 2008: Regression discontinuity designs: A guide to practice. *Journal of econometrics*, **142** (2), 615–635.
- Keele, L., S. Lorch, M. Passarella, D. Small, and R. Titiunik, 2017: *An Overview of Geographically Discontinuous Treatment Assignments with an Application to Children’s Health Insurance*, chap. 4, 147–194. Emerald Publishing Limited, doi:10.1108/S0731-905320170000038007, URL <http://www.emeraldinsight.com/doi/abs/10.1108/S0731-905320170000038007>, <http://www.emeraldinsight.com/doi/pdf/10.1108/S0731-905320170000038007>.

- Keele, L., R. Titiunik, and J. R. Zubizarreta, 2015: Enhancing a geographic regression discontinuity design through matching to estimate the effect of ballot initiatives on voter turnout. *Journal of the Royal Statistical Society: Series A (Statistics in Society)*, **178** (1), 223–239.
- Keele, L. J., and R. Titiunik, 2015: Geographic boundaries as regression discontinuities. *Political Analysis*, **23** (1), 127–155, doi:10.1093/pan/mpu014.
- Li, F., K. L. Morgan, and A. M. Zaslavsky, 2016: Balancing covariates via propensity score weighting. *Journal of the American Statistical Association*, **(just-accepted)**.
- MacDonald, J. M., J. Klick, and B. Grunwald, 2015: The effect of private police on crime: evidence from a geographic regression discontinuity design. *Journal of the Royal Statistical Society: Series A (Statistics in Society)*.
- Papay, J. P., J. B. Willett, and R. J. Murnane, 2011: Extending the regression-discontinuity approach to multiple assignment variables. *Journal of Econometrics*, **161** (2), 203–207.
- Rasmussen, C. E., and C. K. Williams, 2006: *Gaussian processes for machine learning*, Vol. 1. MIT press Cambridge.
- Rencher, A. C., 2003: *Methods of multivariate analysis*, Vol. 492. John Wiley & Sons.
- Stein, M. L., 2012: *Interpolation of spatial data: some theory for kriging*. Springer Science & Business Media.
- Thistlethwaite, D. L., and D. T. Campbell, 1960: Regression-discontinuity analysis: An alternative to the ex post facto experiment. *Journal of Educational psychology*, **51** (6), 309.
- Wilson, A., and R. Adams, 2013: Gaussian process kernels for pattern discovery and extrapolation. *Proceedings of the 30th International Conference on Machine Learning (ICML-13)*, 1067–1075.
- Zubizarreta, J. R., 2012: Using mixed integer programming for matching in an observational study of kidney failure after surgery. *Journal of the American Statistical Association*, **107** (500), 1360–1371.