



Sharp Closed-Form Bounds for Interference Bias in Linear ATT Estimators

Joao Alipio-Correa

Department of Political Science & Department of Statistics, University of Pittsburgh



Overview

Setting.

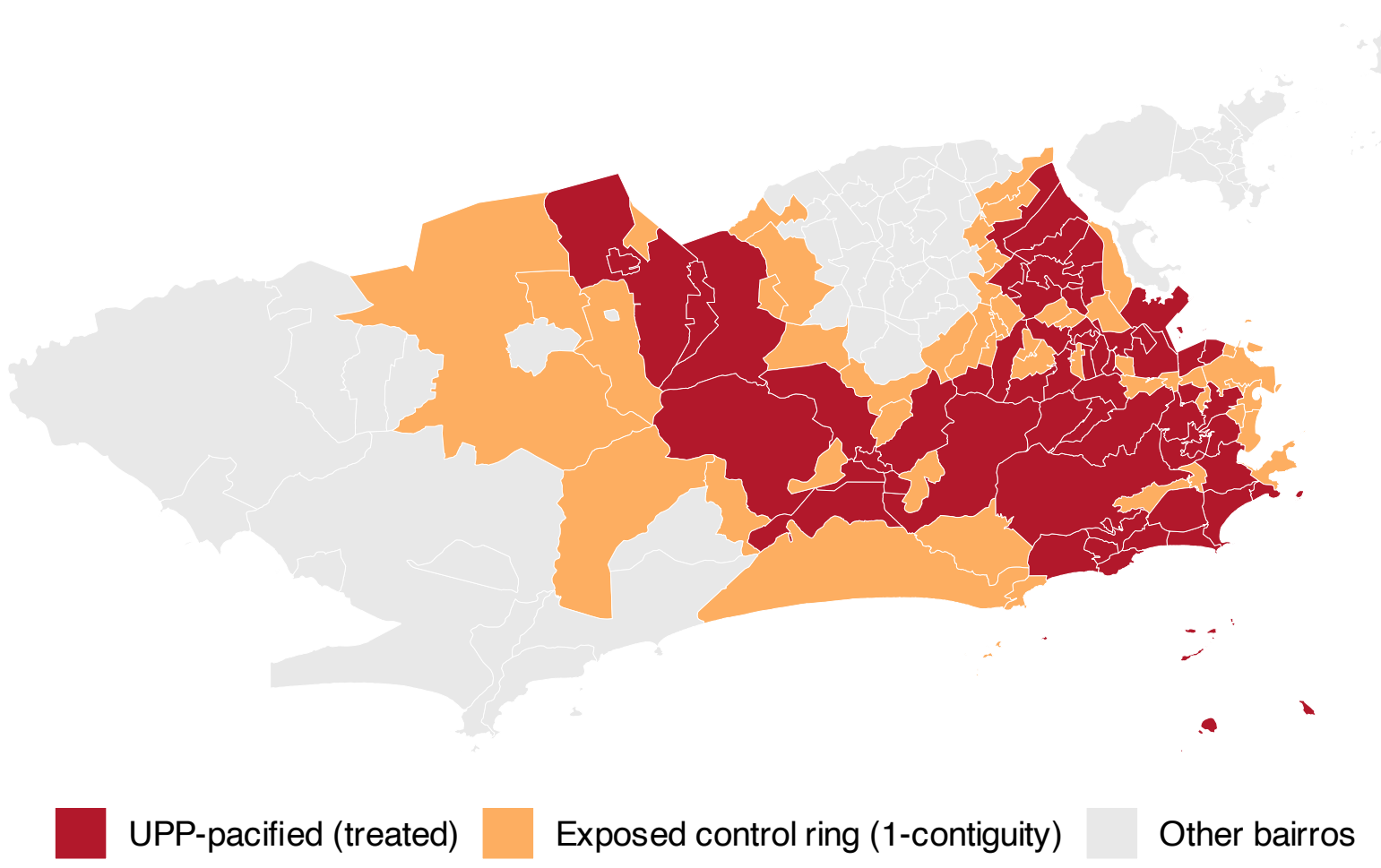
- Social-science theories constantly entertain interference: policies diffuse, information travels, interventions re-shape nearby units. Comparison groups are rarely beyond a treatment's reach.
- Every linear ATT estimator: $\hat{\tau}(w) = \Delta \bar{Y}_T - \sum_j w_j \Delta Y_j^{obs}$. DiD, staggered adoption, synthetic control differ only in w .
- Spillovers reach nearby controls \rightarrow the counterfactual is partially treated \rightarrow a bias $B(w)$ built from unobserved counterfactuals.

Contribution.

- Detect interference (doubly robust $\hat{\psi}$).
- Bound its damage (sharp closed-form bounds, confidence intervals, sensitivity). **An effect bounded against interference, with estimator unchanged.**
- Prove the identified mean of the missing counterfactuals is enough to sharply bound the bias of any linear estimator.
- The width of the bounds is governed by how unevenly the estimator weights exposed controls:** uniform weights point-identify; concentration widens the set.

Motivating problem: police pacification in Rio

2008–2014: Rio de Janeiro installs Pacifying Police Units (UPPs) in 55 communities. Outcome: the **left bloc's vote share** across 156 neighborhoods. A control neighborhood bordering a pacified community shares its streets, commerce, and policing. Its observed change embeds part of the intervention.



Control	E_j	ΔY_j^{obs}	$\Delta Y_j(0,0)$
connected to a UPP	1	✓	missing
distant	0	✓	$= \Delta Y_j^{obs}$

Each control reveals one of two potential changes:

$$\Delta Y_j^{obs} = E_j \Delta Y_j(0,1) + (1 - E_j) \Delta Y_j(0,0)$$

Theorem 1.2 (Contamination decomposition).

$$\hat{\tau}(w) = \tau^{SUTVA}(w) - B(w)$$

$$B(w) = \sum_{j \in C_1} w_j [\Delta Y_j^{obs} - \Delta Y_j(0,0)]$$

Only exposed controls enter. Nonnegative spillovers $\Rightarrow B(w) \geq 0$: the effect is understated.

How large is the interference bias $B(w)$, and how tightly can it be bounded from the data?

- The proposal: identify and recover the mean of the missing vector $\{\Delta Y_j(0,0)\}_{j \in C_1}$.

Assumptions

Notation guide

$\Delta Y_j(a,e)$	outcome change under own treatment a , exposure e
E_j	$= 1$ if a treated unit lies within distance r of j — computed from the data, never assumed
$C_0, C_1; n_i; \rho$	unexposed / exposed controls; $ C_1 ; n_1/ C $
$\delta_j = \Delta Y_j(0,1) - \Delta Y_j(0,0)$	unit-level spillover
$\psi = \mathbb{E}[\delta_j A_j = 0]$	average spillover on controls
$d_j = \Delta Y_j^{obs} - a$	spillover capacity (a : support floor)

Assumptions

- Exposure ignorability.** $\Delta Y_j(0,0) \perp\!\!\!\perp E_j | X_j, \{X_k\}_{k \in N_j}, A_j = 0$. *Sitting near a treated unit is as good as random, given own and neighbors' covariates.*
- Treatment ignorability.** $A_j \perp\!\!\!\perp \Delta Y_j(a,e) | X_j$. *Adoption unconfounded given covariates.*
- Independent assignments.** $A_i \perp\!\!\!\perp A_j | X_i, X_j$. *With A2, delivers A1.*
- Positivity.** $\eta < \Pr(E_j = 1 | \cdot) < 1 - \eta$. *Every control could sit near or far.*
- Overlap.** Exposed controls have unexposed counterparts with the same covariates.
- Feasibility.** $0 \leq \psi \leq \bar{d}$. *Average spillover within average capacity; checked from data.*

- A1–A5: unexposed neighborhoods reveal, on average, what exposed ones would have done.

From detection to partial identification

Detection estimand: $\psi = \mathbb{E}[\Delta Y_j(0,1) - \Delta Y_j(0,0) | A_j = 0]$, the ATT of exposure on controls, identified under A1–A5. Doubly robust estimation:

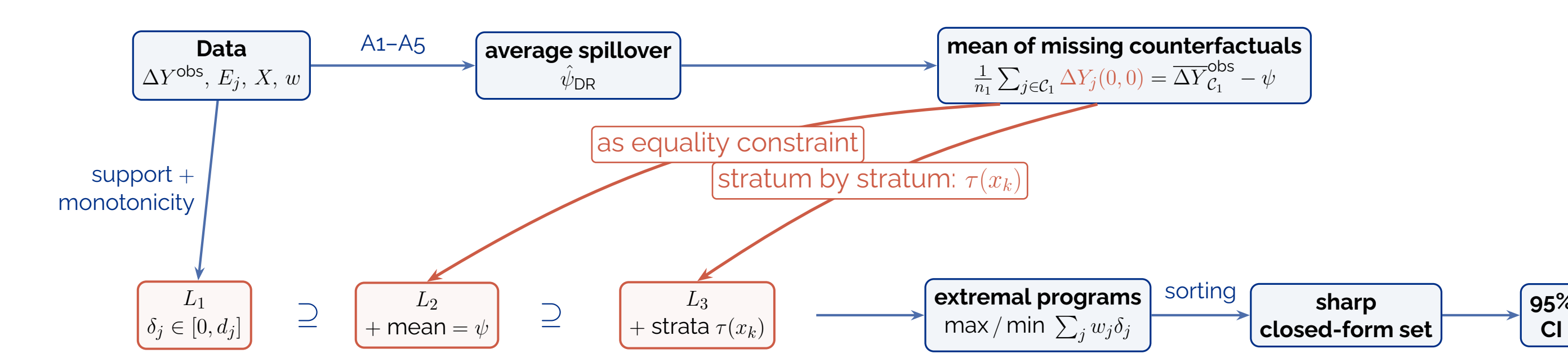
$$\hat{\psi}_{DR} = \frac{1}{n_1} \sum_{j \in C_1} [\Delta Y_j - \hat{\mu}_0(X_j)] - \frac{1}{n_1} \sum_{j \in C_0} \frac{\hat{\pi}_j}{1 - \hat{\pi}_j} [\Delta Y_j - \hat{\mu}_0(X_j)]$$

outcome-model residuals (exposed) propensity-reweighted correction (unexposed)

$\hat{\mu}_0(x)$: outcome model on unexposed controls. $\pi_j = 1 - \prod_{k \in N_j} (1 - e(X_k))$: exposure propensity, from neighbors' adoption probabilities.

Properties: doubly robust — either nuisance correct suffices; \sqrt{n} -normal; semiparametric efficient.

$\hat{\psi}_{DR}$ identifies the mean of the missing counterfactuals: $\frac{1}{n_1} \sum_{j \in C_1} \Delta Y_j(0,0) = \Delta \bar{Y}_{C_1}^{obs} - \psi$. Not the end product: the detected average enters the bounding problem as an equality constraint on the missing vector.



- Nested requirements, nested sets: $\mathcal{F}_3 \subseteq \mathcal{F}_2 \subseteq \mathcal{F}_1$. Each addition shrinks the bounds.

Sharp bounds by sorting

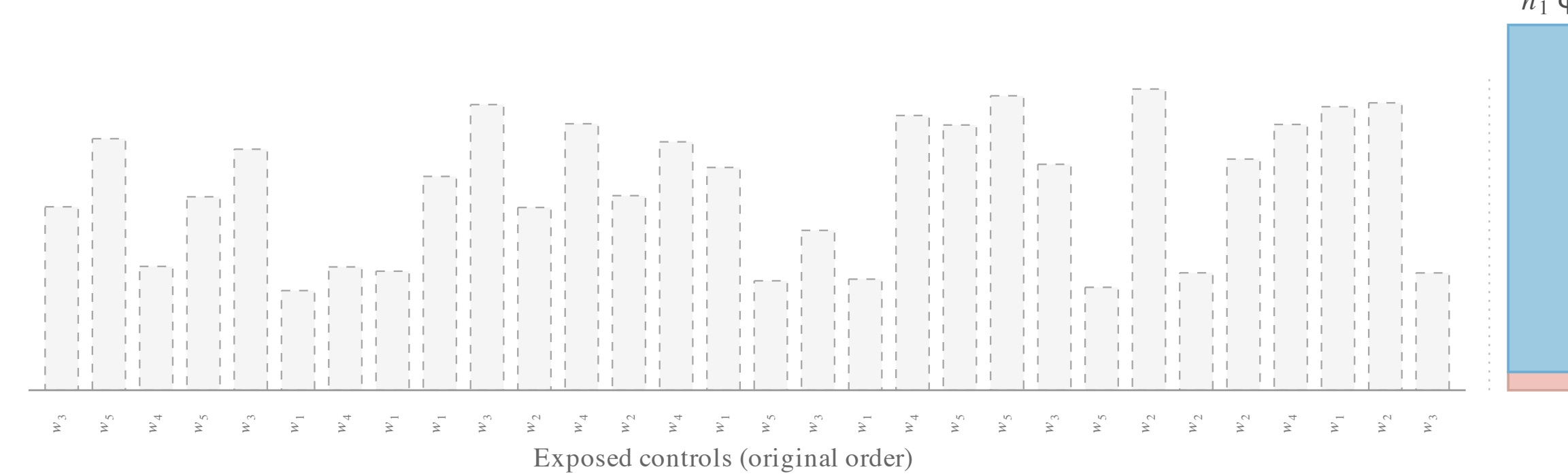
The L_2 extremes solve

$$B_2^U(w) = \max_{\delta} \sum_{j \in C_1} w_j \delta_j \quad \text{s.t.} \quad \underbrace{0 \leq \delta_j \leq d_j}_{\text{capacity}}, \quad \underbrace{\frac{1}{n_1} \sum_j \delta_j = \psi}_{\text{detection}}$$

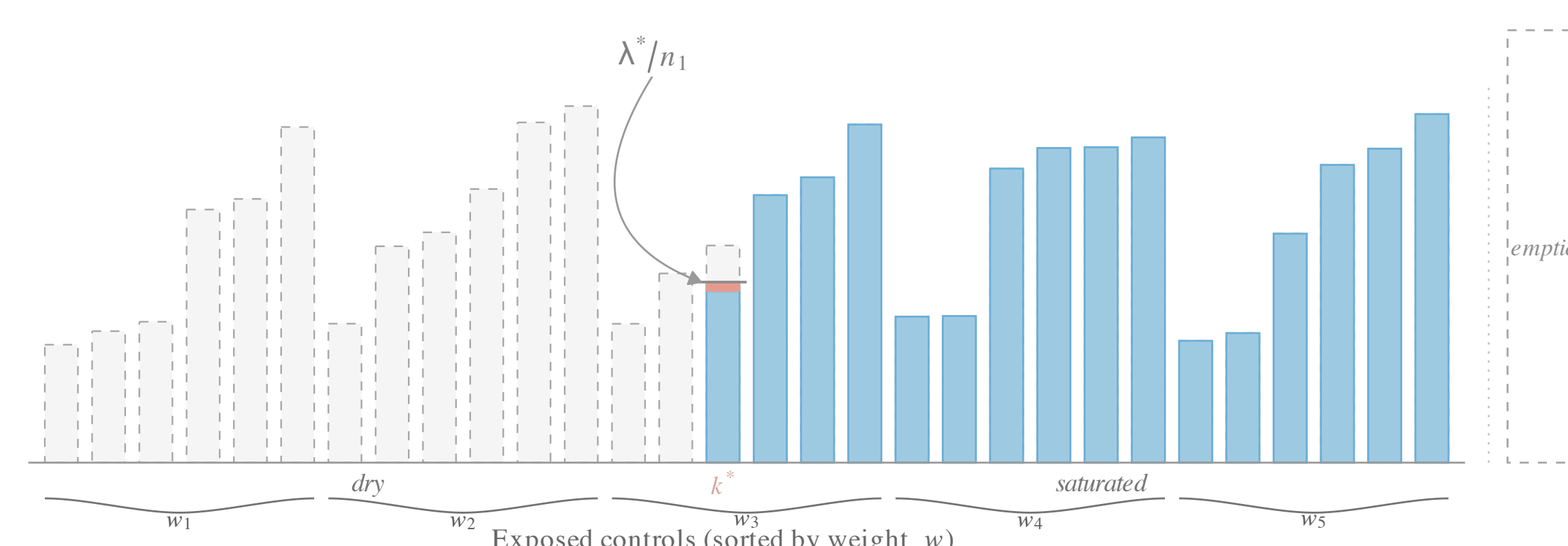
Theorem 3.8 (Sorting). Sort weights ascending; the maximizing allocation fills capacities from the highest weight down until the budget $n_1 \psi$ is spent; one threshold unit k^* absorbs the residual. **Theorem 3.9** evaluates it in closed form:

$$B_2^U(w) = \sum_{j > k^*} w_{(j)} d_{(j)} + w_{(k^*)} [n_1 \psi - C(k^* + 1)], \quad C(k^* + 1) < n_1 \psi \leq C(k^*)$$

Unsorted capacities



Adversary's optimal allocation



Low-weight units stay *dry*, high-weight units *saturate*, the threshold unit takes the residual. The lower bound mirrors the allocation from the lowest weight up.

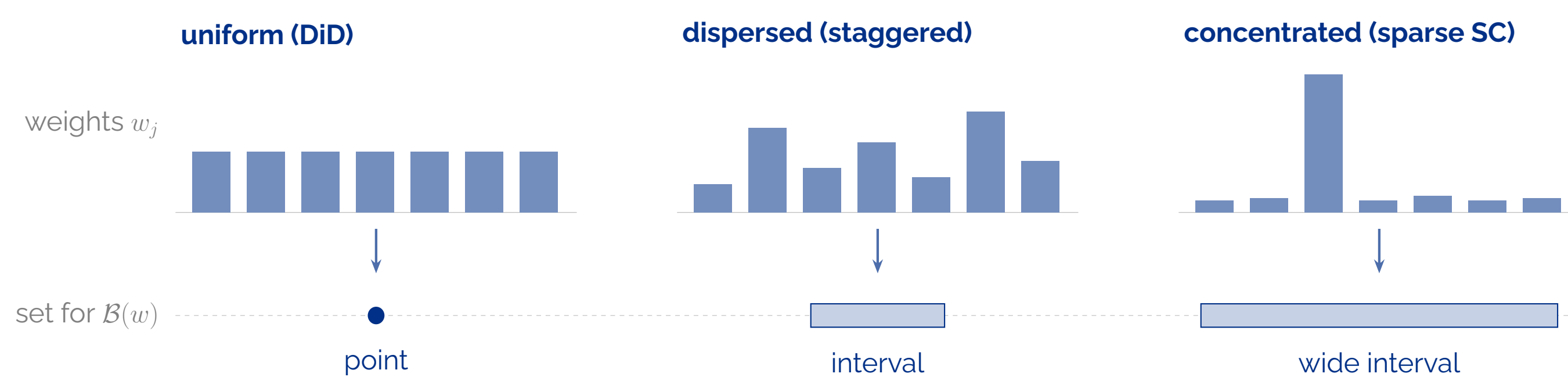
Sharp: no spillover configuration consistent with the data can escape these bounds. LP duality certifies the bounds and an explicit data-generating process attains them: sharpness is certified twice (computation: one sort, $O(n_1 \log n_1)$).

Width of the identification region

Theorem 3.10.

$$W(w) = B_2^U - B_2^L = \sum_{j=1}^{n_1} \tilde{w}_{(j)} \Delta_{(j)} \leq \underbrace{2 \max_j |\tilde{w}_{(j)}|}_{\text{weight dispersion}} \cdot \underbrace{n_1 \min(\psi, \bar{d} - \psi)}_{\text{budget flexibility}}$$

with $\tilde{w}_{(j)} = w_{(j)} - \bar{w}$ the centered sorted weights and $\Delta_{(j)}$ the gap between the two extreme allocations. Zero width \leftarrow uniform weights: then $B = \rho \hat{\psi}$ and detection alone de-contaminates. Concentrating weight widens the set.

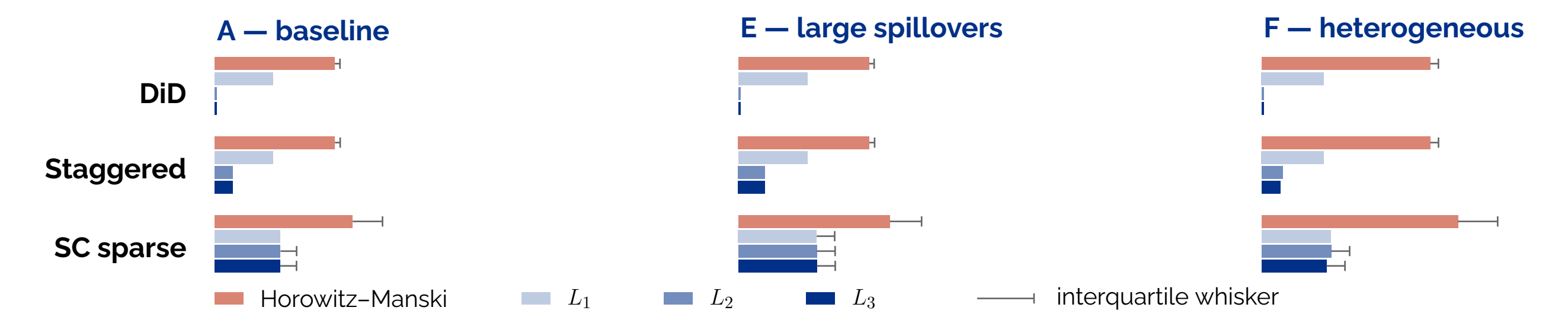


Same data, same $\hat{\psi}$: the weight profile alone sets the width.

- The ranking of widths is set by the weighting strategy before any data arrive. The evidence below tests it.

Simulation evidence

11 DGP configurations \times 4 weight schemes (HHI 0.0002 \rightarrow 0.028) \times 5,000 replications; median widths at $N = 5000$, interquartile whiskers.



- Coverage** — the probability that the reported interval contains the true contamination $B(w)$. The Imbens–Manski interval covers it, and every value of the identified set, with probability $\geq 95\%$: observed ≥ 0.95 in every cell.
 - Ignore interference and the naive 95% CI covers the true ATT in **0% of replications**.
- Median contamination: 25% of the ATT at baseline, 50% under large spillovers. Stratification (Config F): a further 8–12%.

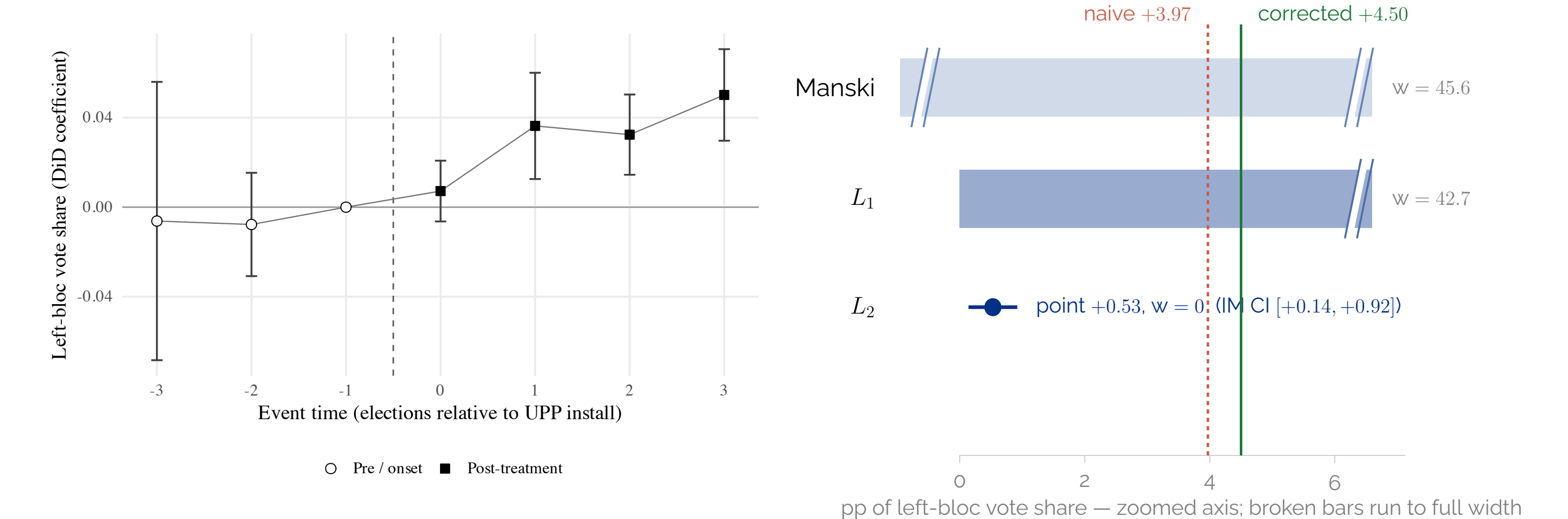
Compared with Horowitz–Manski: they use one worst-case capacity for every unit, $W_{exp}(b-a)$, and detection plays no role. Here, unit-specific capacities $d_j = \Delta Y_j^{obs} - a$ are strictly tighter, and the identified mean has no HM analog: HM $\rightarrow L_2$ removes 100 / 85.2 / 49.8% of the width (DiD / staggered / sparse SC, baseline).

Estimation and inference

- Fit $\hat{\mu}_0$ on unexposed controls; $\hat{\pi}_j$ from neighbors' adoption propensities (cross-fit if using ML).
- Compute $\hat{\psi}_{DR}$; check feasibility $0 \leq \hat{\psi} \leq \bar{d}$.
- Sort exposed controls by weight; thresholds: closed-form bounds (L_3 : per stratum).
- Influence-function SEs through the piecewise-linear bound map; Imbens–Manski interval.
- Repeat over r ; report the sensitivity curve.

Application A — UPP Rio: uniform weights, a point

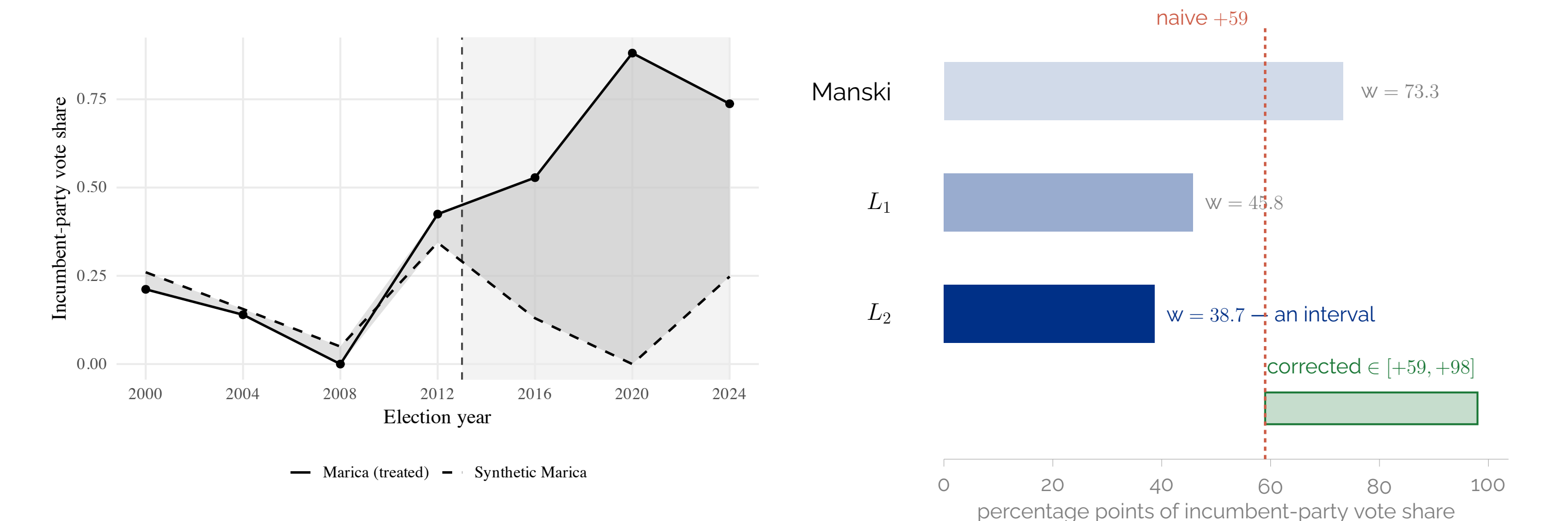
Outcome: **left-bloc vote share**. 156 neighborhoods, 55 pacified, staggered event study; exposure = bordering a pacified community ($\rho = 0.40$); pre-trends flat ($\rho = 0.84, 0.51$).



Pacification moved the left-bloc vote by 4.5 points, not 4.0: spillover into bordering neighborhoods hid 13% of the effect, and the correction is a point whose CI excludes zero.

Application B — Maricá: concentrated weights, an interval

Outcome: **incumbent-party vote share**. Synthetic control after the municipal basic income, close pre-treatment fit; donor municipalities within ~ 75 km exposed ($\rho = 0.24$).



The naive +59 points is a floor: exposed donors were lifted, deflating the synthetic counterfactual, so the true effect lies in [+59, +98]. A second outcome (vote concentration, $t = 3.0$) repeats the pattern: contamination follows the design, not the outcome.

Advice to applied researchers: Keep your estimator. Report the point estimate, the detected $\hat{\psi}$, the sharp de-contaminated set, and the sensitivity curve in r .