



Distribution-Free Assessment of Population Overlap in Observational Studies

Lihua Lei (Stanford), Alexander D'Amour (Google Brain), Peng Ding (UC Berkeley), Avi Feller (UC Berkeley), and Jasjeet Sekhon (Yale)

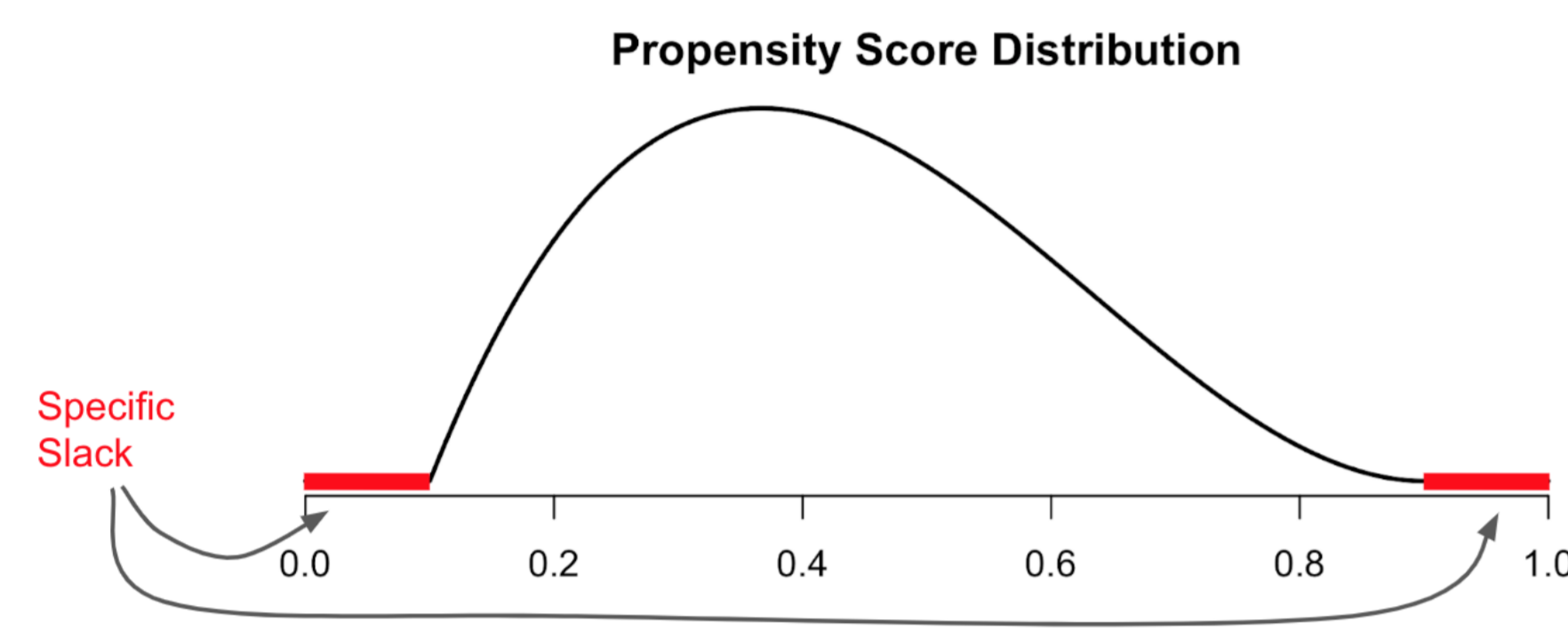
Population overlap in observational studies

Setting: binary treatment T , baseline covariates X (arbitrary),
 $(T_i, X_i) \stackrel{i.i.d.}{\sim} (T, X)$ (the **only assumption!**)

Strict overlap condition:

$$\exists \mathcal{O}_0 > 0, \quad \mathcal{O}_0 \leq \underbrace{e(X)}_{\text{propensity score}} \triangleq \mathbb{P}(T=1 | X) \leq 1 - \mathcal{O}_0, \quad \text{a.s.}$$

Population overlap slack: $\mathcal{O}^* = \inf_x \min\{e(x), 1 - e(x)\}$



- Strict overlap condition $\iff \mathcal{O}^* \geq \mathcal{O}_0$
- $n\mathcal{O}^*$ is the **effective sample size** without outcome restriction (Hong et al. '18)
- $1 - n\mathcal{O}^* / \min\{n_1, n_0\}$ measures the **relative efficiency loss** compared to an RCT
- In practice, high $\mathcal{O}^* \implies$ stability of doubly robust estimators

Current approaches for assessing overlap:

- Informal comparisons or plug-in estimates based on estimated propensity scores
 - useful but lack of statistical guarantees
 - “sample overlap” \neq population overlap
 - sensitive to model mis-specification or finite sample errors
- Standard two-sample test: testing the wrong null

$$H_0 : \mathbb{P}(X | T=1) = \mathbb{P}(X | T=0) \implies H_0 : e(X) \equiv e_0$$

Major challenge: \mathcal{O}^* is irregular (extreme of an unknown function)

O-value

Definition. $\hat{\mathcal{O}}$ is an O-value if it is an upper confidence bound of \mathcal{O}^* , i.e.

$$\mathbb{P}(\mathcal{O}^* \leq \hat{\mathcal{O}}) \geq 1 - \alpha$$

Analogous to p-value:

- A small $\hat{\mathcal{O}}$ provides strong evidence against overlap
- A large $\hat{\mathcal{O}}$ does not necessarily imply sufficient overlap

Some practical implications:

- Strict overlap condition as a composite null hypothesis: reject if $\hat{\mathcal{O}} < \mathcal{O}_0$
- $(1 - n\hat{\mathcal{O}} / \min\{n_1, n_0\})_+$ estimates **efficiency loss** caused by the imbalance
- Assessing if trimming (say, at 0.1 and 0.9) is successful by comparing $\hat{\mathcal{O}}$ with 0.1
- Comparing different matches based on $\hat{\mathcal{O}}$

Main (and perhaps surprising) contribution:

We develop **distribution-free** O-values that are valid **in finite samples!**

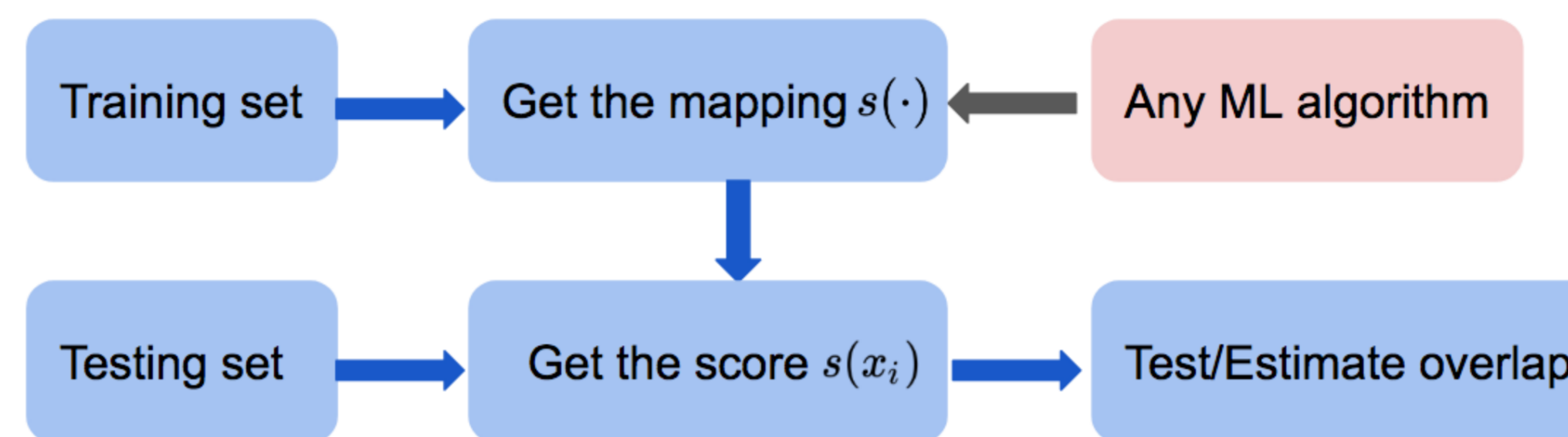
Step I: covariate standardization

Key observation: **overlap is preserved under transformation of X**

$$\mathcal{O}_0 \leq e(x) \leq 1 - \mathcal{O}_0 \implies \mathcal{O}_0 \leq \mathbb{P}(T=1 | s(X)=s) \leq 1 - \mathcal{O}_0$$

for any fixed function $s(\cdot)$

- Let \mathcal{O}_s^* be the population overlap slack for $(T, s(X))$
- $\mathcal{O}^* = \mathcal{O}_e^*$ (Rosenbaum and Rubin, '83)
- **Data splitting** guarantees that $\hat{e}(\cdot) \perp\!\!\!\perp$ (second half of data)
- $\mathcal{O}^* \leq \mathcal{O}_{\hat{e}}^*$ always holds even if \hat{e} is bad; tight if \hat{e} is good



From now on, we assume that

$$(S_i, T_i) \text{ are i.i.d. with } S_i = \hat{e}(X_i) \in [0, 1]$$

Goal: construct upper confidence bounds on \mathcal{O}_e^* (w/ **standardized covariates**)

Step II: careful balance check

Key observation: **overlap \iff bounded likelihood ratio**

$$b_{\min}(\mathcal{O}^*; \pi) \leq \frac{dP_{S|T=1}(s)}{dP_{S|T=0}(s)} \triangleq \frac{dP_1}{dP_0}(s) \leq b_{\max}(\mathcal{O}^*; \pi), \quad \forall s \in [0, 1],$$

$$\text{where } \pi = \mathbb{P}(T=1), \quad b_{\min}(\mathcal{O}^*; \pi) = \frac{\mathcal{O}^*}{1 - \mathcal{O}^*} \frac{\pi}{1 - \pi}, \quad b_{\max}(\mathcal{O}^*; \pi) = \frac{1 - \mathcal{O}^*}{\mathcal{O}^*} \frac{\pi}{1 - \pi}$$

Intuition: **larger $\mathcal{O}^* \implies$ smaller discrepancy between P_1 and P_0**

A generic strategy:

- Find an estimable “discrepancy” $\Delta(P_0, P_1)$ and $B_{\Delta}(\mathcal{O}) \downarrow \mathcal{O}$

$$\Delta(P_0, P_1) \leq B_{\Delta}(\mathcal{O}^*) \quad (\text{population property})$$
- Compute a lower confidence bound on $\Delta(P_0, P_1)$

$$\mathbb{P}(\hat{\Delta}^- \leq \Delta(P_0, P_1)) \geq 1 - \alpha \quad (\text{sample property})$$
- $\hat{\mathcal{O}} = B_{\Delta}^{-1}(\hat{\Delta}^-)$ is a valid O-value:

$$\mathbb{P}(\hat{\mathcal{O}} \geq \mathcal{O}^*) = \mathbb{P}(\hat{\Delta}^- \leq B_{\Delta}(\mathcal{O}^*)) \geq \mathbb{P}(\hat{\Delta}^- \leq \Delta(P_0, P_1)) \geq 1 - \alpha$$

Summary of DiM/DiT/DiR/CE O-values

	Δ	$B_{\Delta}(\mathcal{O}^*)$	$\hat{\Delta}^-$	Two-sample test analogy
DiM	T-stat.	χ^2 -divergence	Hedged capital bound <small>Waudby-Smith-Ramdas ('20)</small>	t-test
DiT	LR	Simple algebra	Line-crossing <small>Dempster ('59)</small> Simes' inequality <small>Sarkar ('98)</small>	Kolmogorov-Smirnov test
DiR	AUC	Generalized Neyman-Pearson	Hybrid bound for U-stat <small>Bates-Candès-Lei-Romano-Sesia ('21)</small>	Wilcoxon rank-sum test
CE	class. error	Formula of Bayes risk	Same as DiT O-values	Classification-based test

Example: (simplified) DiM O-value

Population property D'Amour-Ding-Feller-Lei-Sekhon ('17)

Theorem. Let μ_t, σ_t^2 be the mean and variance of P_t . Then

$$(\Delta(P_0, P_1) =) T_0 \triangleq \frac{|\mu_1 - \mu_0|}{\sigma_0} \leq \sqrt{(1 - b_{\min}(\mathcal{O}^*; \pi))(b_{\max}(\mathcal{O}^*; \pi) - 1)} (= B_{\Delta}(\mathcal{O}^*)),$$

Sample property Maurer-Pontil ('09) (constants are not good enough!)

Proposition (Empirical Bernstein's inequality). Let $Z_1, \dots, Z_n \in [0, 1]$ be i.i.d. with $\mathbb{E} Z_i = \mu$, $\text{Var}(Z_i) = \sigma^2$. Then with probability $1 - \delta$,

$$|\hat{\mu} - \mu| \leq \hat{\sigma} \sqrt{\frac{2 \log(\frac{3}{\delta})}{n} + \frac{7 \log(\frac{3}{\delta})}{3(n-1)}}, \quad \sigma - \hat{\sigma} \leq \sqrt{\frac{2 \log(\frac{3}{\delta})}{n-1}}.$$

\implies a lower confidence bound on T_0 (with Bonferroni correction over (μ_1, μ_0, σ_0)):

$$\hat{T}_0^- = \frac{\hat{\mu}_1^- - \hat{\mu}_0^+}{\hat{\sigma}_0^+}$$

(Simplified) DiM O-value

$\mathcal{C}^\pi \leftarrow (1 - \alpha/2)$ -CI for π , $\hat{T}_0^- \leftarrow (1 - \alpha/2)$ -lower confidence bound on T_0

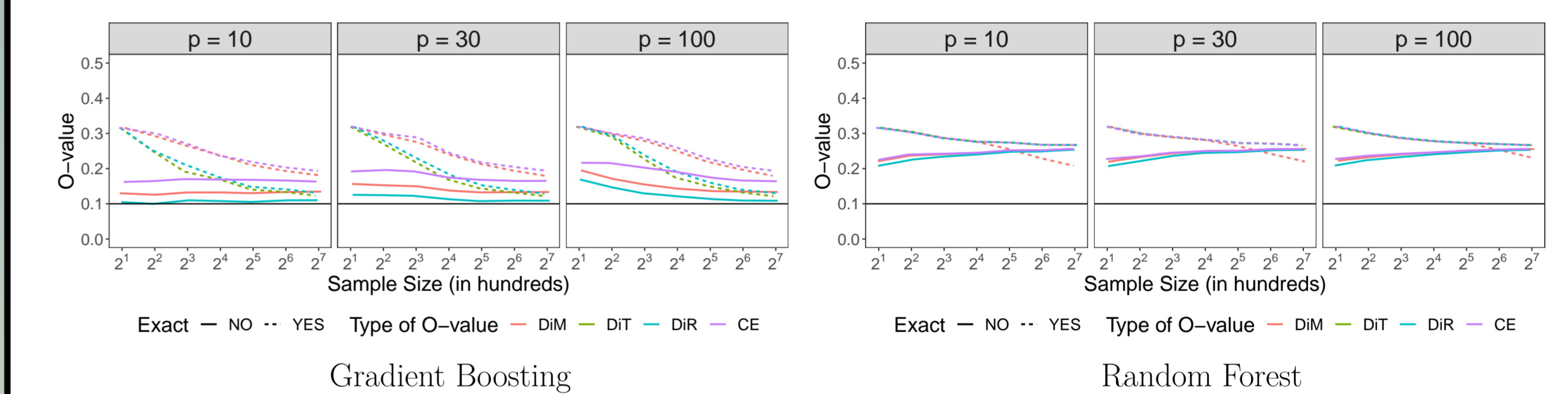
$$\hat{\mathcal{O}}_{\text{DiM}} = \sup_{\pi \in \mathcal{C}^\pi} \frac{1}{2} - \frac{1}{\sqrt{4 - \frac{\pi(1-\pi)}{\pi^2(\hat{T}_0^-)^2 + 1}}}$$

Theorem. With solely the i.i.d. assumption, $\mathbb{P}(\mathcal{O}^* \leq \hat{\mathcal{O}}_{\text{DiM}}) \geq 1 - \alpha$

Comparisons of O-values

An illustrative simulation study

- $X \sim N(0, I_p)$ with $p \in \{10, 30, 100\}$
- $e(x) = f(x^T \beta)$, $f(y) = \begin{cases} 0.1 & (y < c) \\ 0.9 & (y > c) \end{cases}$
- β sparse; c is chosen such that $\mathbb{P}(e(X) = 0.1) = 0.8$



Practical recommendation based on extensive numerical experiments

- Algorithm to estimate propensity scores: **gradient boosting**
- Type of O-value: **DiT**

O-values for Lalonde data

- National Supported Work Demonstration program Lalonde ('86)
- Treatment group has $n_1 = 185$ units
- 7 control groups: 6 from observational studies, 1 from an RCT
- Apply gradient boosting for DiT O-values

	CPS			PSID			RCT		
	n_0	$\hat{\mathcal{O}}$	\hat{L}	n_0	$\hat{\mathcal{O}}$	\hat{L}	n_0	$\hat{\mathcal{O}}$	\hat{L}
Raw	15992	0.003	77%	2490	0.018	75%	260	0.483	0%
V2	2369	0.021	71%	253	0.234	45%			
V3	429	0.143	53%	128	0.313	23%			