
Detecting critical treatment effect bias in small subgroups

Piersilvio de Bartolomeis¹ Javier Abad¹ Konstantin Donhauser¹ Fanny Yang¹

Abstract

Randomized trials are considered the gold standard for making informed decisions in medicine, yet they often lack generalizability to the patient populations in clinical practice. Observational studies, on the other hand, cover a broader patient population but are prone to various biases. Thus, before using an observational study for decision-making, it is crucial to *benchmark* its treatment effect estimates against those derived from a randomized trial. We propose a novel strategy to benchmark observational studies beyond the average treatment effect. First, we design a statistical test for the null hypothesis that the treatment effects estimated from the two studies, conditioned on a set of relevant features, differ up to some tolerance. We then estimate an asymptotically valid lower bound on the maximum bias strength for any subgroup in the observational study. Finally, we validate our benchmarking strategy in a real-world setting and show that it leads to conclusions that align with established medical knowledge.

1. Introduction

Randomized trials have traditionally been the gold standard for informed decision-making in medicine, as they allow for unbiased estimation of treatment effects under mild assumptions. However, there is often a significant discrepancy between the patients observed in clinical practice and those enrolled in randomized trials, limiting the generalizability of the trial results (Rothwell, 2005; Duma et al., 2018). To address this issue, the U.S. Food and Drug Administration advocates for using observational data, as it is usually more representative of the patient population in clinical practice (Platt et al., 2018; Klonoff, 2020). Yet, a major caveat to this recommendation is that several sources of bias, including hidden confounding, can compromise the

¹Department of Computer Science, ETH Zurich. Correspondence to: Piersilvio de Bartolomeis <piersilvio.debartolomeis@inf.ethz.ch>.

causal conclusions drawn from observational data.

In light of the inherent limitations of randomized and observational data, it has become a popular strategy to *benchmark* observational studies against existing randomized trials to assess their quality (Dahabreh et al., 2020; Forbes & Dahabreh, 2020). The main idea behind this approach is first to emulate the procedures adopted in the randomized trial within the observational study; see e.g. (Hernán & Robins, 2016) for a detailed explanation. Then, the treatment effect estimates from the observational data are compared with those from the randomized data. If the estimates are similar, we may be willing to trust the observational study for patient populations where the randomized data is insufficient.

To support the benchmarking framework, several works propose statistical tests that compare treatment effect estimates between randomized and observational data (Viele et al., 2014; Hussain et al., 2023; De Bartolomeis et al., 2024; Yang et al., 2023; Demirel et al., 2024). In particular, two properties have been identified as essential for effective benchmarking of observational studies: *tolerance* and *granularity*. Tolerance allows the acceptance of studies with negligible bias that does not impact decision-making, thereby significantly reducing false rejections in real-world settings where some bias is expected. Granularity, on the other hand, allows the detection of bias on small subgroups or individuals that would otherwise go unnoticed.

In this work, we design a statistical test for the null hypothesis that treatment effects differ up to some tolerance value when conditioned on a relevant subset of features. Our test is the first, to our knowledge, to satisfy granularity and tolerance. Further, we use our test to estimate an asymptotically valid lower bound on the maximum bias strength for any individual. Finally, we show that our lower bound leads to conclusions that align with established medical knowledge.

2. Problem setting

We have access to two datasets: D_{rct} of size n_{rct} from a randomized trial (rct) and D_{os} of size n_{os} from an observational study (os), containing tuples $Z := (X, Y, T)$ of covariates $X \in \mathbb{R}^d$, bounded observed outcome $Y \in \mathbb{R}$, and treatment assignment variable $T \in \{0, 1\}$. We assume that the data is drawn i.i.d from the distributions \mathbb{P}^{rct} and \mathbb{P}^{os}

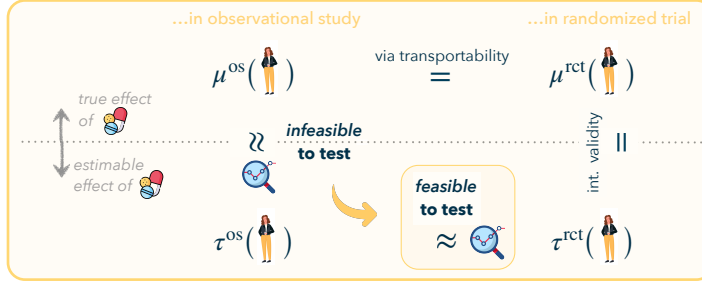


Figure 1. High-level illustration of our approach. We want to test if the bias in the observational study, i.e. $\mu^{\text{os}} - \tau^{\text{os}}$, is contained within a tolerance range. However, the true treatment effect μ^{os} is not identifiable, and instead, we test the bias between the treatment effects estimated from the two studies, i.e. $\tau^{\text{os}} - \tau^{\text{rc}}$.

that are marginal distributions of the respective full distribution $\mathbb{P}_{\text{full}}^{\diamond}$ over $(X, U, Y(0), Y(1), Y, T)$ for $\diamond \in \{\text{rc}, \text{os}\}$. In particular, the full distribution also includes randomness over a vector of unobserved covariates $U \in \mathbb{R}^k$ and potential outcomes $(Y(0), Y(1)) \in \mathbb{R}^2$. We further assume that the support of the randomized trial is included in the support of the observational study, i.e. $\text{supp}(\mathbb{P}_X^{\text{rc}}) \subseteq \text{supp}(\mathbb{P}_X^{\text{os}})$.

Treatment effect estimation A crucial quantity to estimate for decision-making in many domains is the conditional average treatment effect (CATE):

$$\mu^{\diamond}(x) := \mathbb{E}_{\mathbb{P}_{\text{full}}^{\diamond}} [Y(1) - Y(0) \mid X = x],$$

for $\diamond \in \{\text{rc}, \text{os}\}$ and $\mathcal{X} \subseteq \text{supp}(\mathbb{P}_X^{\text{rc}})$. Unfortunately, we cannot estimate the CATE from the observed data as we never observe the potential outcomes. Instead, we can estimate the regression function defined by

$$\tau^{\diamond}(x) = \mathbb{E}_{\mathbb{P}^{\diamond}} [Y \mid T = 1, X = x] - \mathbb{E}_{\mathbb{P}^{\diamond}} [Y \mid T = 0, X = x],$$

for $\diamond \in \{\text{rc}, \text{os}\}$. For the treatment effect in the randomized trial, we observe that $\tau^{\text{rc}}(x) = \mu^{\text{rc}}(x)$ holds for all $x \in \mathcal{X}$, under the assumption of internal validity outlined below.

Assumption 2.1 (Internal validity). *The data-generating process of the randomized trial satisfies*

- (i) $Y = Y(T)$ $\mathbb{P}_{\text{full}}^{\text{rc}}$ - almost surely.
- (ii) $T \perp\!\!\!\perp (Y(1), Y(0))$.
- (iii) $\mathbb{P}_{\text{full}}^{\text{rc}}(T = 1 \mid X, U) = \pi \in (0, 1)$.

In particular, Assumption 2.1 is expected to hold in a completely randomized experiment, and thus, μ^{rc} can be estimated from the observed data under mild assumptions (Rubin, 1978). On the other hand, we cannot estimate μ^{os} from the observed data due to hidden confounding or other sources of bias in the observational study, i.e. we cannot rule out the existence of $x \in \mathcal{X}$ such that $\tau^{\text{os}}(x) \neq \mu^{\text{os}}(x)$. Therefore, it is crucial to benchmark the observational study before using the estimate of τ^{os} for any downstream task.

2.1. Null hypothesis

Our goal is to test if the bias in the observational study, defined as $\delta^*(x) := \tau^{\text{os}}(x) - \mu^{\text{os}}(x)$ for all $x \in \mathcal{X}$, is contained within a tolerance range. However, the bias δ^* is not estimable from the data. Instead, we can test the bias $\tilde{\delta}(x) := \tau^{\text{os}}(x) - \tau^{\text{rc}}(x)$, which is equivalent to δ^* under internal validity and transportability, i.e. $\mu^{\text{os}}(x) = \mu^{\text{rc}}(x)$ for all $x \in \mathcal{X}$ (see Figure 1). We would like to test if the bias $\tilde{\delta}$ between the two studies is contained within a tolerance range (requires tolerance) across all patient subgroups (requires granularity). Hence, we will now introduce a null hypothesis that allows for both tolerance and granularity.

To do so, we define two bounded tolerance functions $\tau_{\pm}^{\text{os}} : \mathcal{X} \rightarrow \mathbb{R}$ that capture how much the estimated treatment effects can differ between studies and satisfy $\tau_{-}^{\text{os}}(x) \leq \tau^{\text{os}}(x) \leq \tau_{+}^{\text{os}}(x)$ for all $x \in \mathcal{X}$. Further, we define the patient subgroups via a subset of features $X^{\mathcal{J}}$, corresponding to the covariates with indices $\mathcal{J} \subseteq \{1, \dots, d\}$. We can then introduce our null hypothesis, given by

$$H_0 : \mathbb{E}_{\mathbb{P}^{\text{rc}}} [\tau^{\text{rc}}(X) \mid X^{\mathcal{J}}] \in [\mathbb{E}_{\mathbb{P}^{\text{rc}}} [\tau_{-}^{\text{os}}(X) \mid X^{\mathcal{J}}], \mathbb{E}_{\mathbb{P}^{\text{rc}}} [\tau_{+}^{\text{os}}(X) \mid X^{\mathcal{J}}]], \quad \mathbb{P}_{X^{\mathcal{J}}}^{\text{rc}} - \text{a.s.} \quad (1)$$

Discussion of our null hypothesis We provide several remarks on the null hypothesis in Equation (1). First, we satisfy tolerance by testing if $\tau^{\text{rc}}(x)$ is contained (in probability) in an interval around $\tau^{\text{os}}(x)$, for all $x \in \mathcal{X}$. Second, we can satisfy granularity by choosing an appropriate subset \mathcal{J} : When $|\mathcal{J}| = d$, we detect bias at the individual level, thereby satisfying the strictest definition of granularity. On the other hand, when $|\mathcal{J}| = 0$, we test if the average treatment effects are equal, thus potentially ignoring bias in small subgroups and individuals. Third, we test if the treatment effects are equal (up to tolerance) on the support of the randomized trial since we cannot extrapolate outside the support of \mathbb{P}_X^{rc} without further assumptions.

Example: User-specified tolerance A natural choice for the tolerance functions is to add (respectively subtract) a

user-specified function $\delta(x) \geq 0$, that is

$$\tau_{\pm}^{\text{os}}(x) = \tau^{\text{os}}(x) \pm \delta(x), \quad \text{for all } x \in \mathcal{X}.$$

The function δ can incorporate all sources of bias in the observational study, such as unobserved confounding and non-adherence to treatment assignments. For instance, we can test whether $\|\tilde{\delta}\|_{L^\infty(\mathbb{P}_{X^{\text{rct}}})}$ is larger than a critical value $\delta_{\text{CT}} \in \mathbb{R}$ by choosing $\tau_{\pm}^{\text{os}}(x) = \tau^{\text{os}}(x) \pm \delta_{\text{CT}}$.

3. Methodology

In this section, we rewrite the null hypothesis from Equation (1) in terms of a *signal* function that captures the bias between τ^{os} and τ^{rct} . Then, we propose an oracle test statistic assuming that the tolerance functions τ_{\pm}^{os} are known. Finally, we provide asymptotic guarantees for the finite-sample test statistic where the tolerance functions are estimated.

3.1. Null hypothesis using signal function

We first observe that, for some tolerance functions τ_{\pm}^{os} , Equation (1) is equivalent to stating that there exists a function $g : \mathbb{R}^{|\mathcal{J}|} \rightarrow [0, 1]$ such that $\tau_g^{\text{os}}(X) := g(X^{\mathcal{J}}) \tau_+^{\text{os}}(X) + (1 - g(X^{\mathcal{J}})) \tau_-^{\text{os}}(X)$ satisfies

$$\mathbb{E}_{\mathbb{P}^{\text{rct}}} [\tau^{\text{rct}}(X) | X^{\mathcal{J}}] = \mathbb{E}_{\mathbb{P}^{\text{rct}}} [\tau_g^{\text{os}}(X) | X^{\mathcal{J}}], \quad \mathbb{P}_{X^{\mathcal{J}}}^{\text{rct}} \text{-a.s.}$$

We test a slightly more restrictive hypothesis by assuming that g lies in a sufficiently rich function class \mathcal{G} :

$$H_0^{\mathcal{G}} : \mathbb{E}_{\mathbb{P}^{\text{rct}}} [\tau^{\text{rct}}(X) | X^{\mathcal{J}}] = \mathbb{E}_{\mathbb{P}^{\text{rct}}} [\tau_{g^*}^{\text{os}}(X) | X^{\mathcal{J}}], \\ \text{for some } g^* \in \mathcal{G}, \quad \mathbb{P}_{X^{\mathcal{J}}}^{\text{rct}} \text{-a.s.}$$

In practice, one can either restrict \mathcal{G} to a particular function class if domain knowledge is available or use neural networks as general function approximations.

We can then rewrite the null hypothesis above using a *signal* function that captures the bias between the estimates from observational and randomized data. Recall that $Z = (X, Y, T)$ is the vector of observed variables, we define

$$\psi_g(Z) = Y \left(\frac{T}{\pi} - \frac{1-T}{1-\pi} \right) - \tau_g^{\text{os}}(X)$$

and finally arrive at the null hypothesis

$$H_0^{\mathcal{G}} : \mathbb{E}_{\mathbb{P}^{\text{rct}}} [\psi_{g^*}(Z) | X^{\mathcal{J}}] = 0, \quad (2) \\ \text{for some } g^* \in \mathcal{G}, \quad \mathbb{P}_{X^{\mathcal{J}}}^{\text{rct}} \text{-a.s.}$$

At first glance, testing the null hypothesis in Equation (2) may seem equivalent to testing equality of conditional means (Delgado, 1993; Neumeyer & Dette, 2003; Racine et al., 2006; Luedtke et al., 2019; Muandet et al., 2020); however, we remark that this equivalence holds only if the function g^* is known, and to our knowledge, the scenario where g^* is unknown has not been previously explored.

3.2. Oracle test statistic

We now derive a kernelized test statistic for the null hypothesis in Equation (2). First, we observe that the hypothesis $H_0^{\mathcal{G}}$ implies an infinite set of unconditional moment constraints, i.e. for any $g \in \mathcal{G}$, it holds that

$$\mathbb{E}_{\mathbb{P}^{\text{rct}}} [\psi_g(Z) | X^{\mathcal{J}}] = 0, \quad \mathbb{P}_{X^{\mathcal{J}}}^{\text{rct}} \text{-a.s.} \implies \\ \mathbb{E}_{\mathbb{P}^{\text{rct}}} [\psi_g(Z) f(X^{\mathcal{J}})] = 0, \quad \text{for all measurable } f.$$

Therefore, the validity of testing the RHS would carry over to the validity of testing $H_0^{\mathcal{G}}$. However, testing the RHS of the implication above for all measurable functions is infeasible. Instead, we can restrict f to be in a reproducing kernel Hilbert space (RKHS). The problem then becomes more tractable since it holds that

$$\mathbb{H}^2(\psi_g) := \left(\sup_{\|f\|_{\mathcal{F}} \leq 1} \mathbb{E}_{\mathbb{P}^{\text{rct}}} [\psi_g(Z) f(X^{\mathcal{J}})] \right)^2 \quad (3) \\ = \|\mathbb{E}_{\mathbb{P}^{\text{rct}}} [\psi_g(Z) k(X^{\mathcal{J}}, \cdot)]\|_{\mathcal{F}}^2 \\ = \mathbb{E}_{\mathbb{P}^{\text{rct}}} [\psi_g(Z) k(X^{\mathcal{J}}, \tilde{X}^{\mathcal{J}}) \psi_g(\tilde{Z})],$$

where k is a uniformly bounded reproducing kernel corresponding to an RKHS \mathcal{F} , and \tilde{Z} is an independent copy of Z following the same distribution. In particular, the null hypothesis $H_0^{\mathcal{G}}$ implies that $\mathbb{H}^2(\psi_g) = 0$ for $g = g^*$, and thus we can construct a valid test based on $\mathbb{H}^2(\psi_{g^*})$.

A valid test statistic Given i.i.d. samples Z_i from \mathbb{P}^{rct} , an unbiased empirical estimate of $\mathbb{H}^2(\psi_g)$ is the cross U-statistic (Kim & Ramdas, 2024), defined as

$$\hat{\mathbb{H}}^2(\psi_g) := \frac{2}{n_{\text{rct}}} \sum_{i=1}^{n_{\text{rct}}/2} h(Z_i; \psi_g), \quad \text{with} \\ h(Z_i; \psi_g) := \frac{2}{n_{\text{rct}}} \sum_{j=n_{\text{rct}}/2+1}^{n_{\text{rct}}} \psi_g(Z_i) k(X_i^{\mathcal{J}}, X_j^{\mathcal{J}}) \psi_g(Z_j),$$

for all $g \in \mathcal{G}$. The main advantage of the cross U-statistic is that, for $g = g^*$, it is asymptotically normal under the null hypothesis $H_0^{\mathcal{G}}$ and weak regularity assumptions (see Theorem 3.1), i.e. as $n_{\text{rct}} \rightarrow \infty$ it holds that

$$\sqrt{\frac{n_{\text{rct}}}{2}} \frac{\hat{\mathbb{H}}^2(\psi_{g^*})}{\hat{\sigma}(\hat{\mathbb{H}}^2(\psi_{g^*}))} \rightarrow \mathcal{N}(0, 1),$$

where $\hat{\sigma}(\hat{\mathbb{H}}^2(\psi_{g^*}))$ is the finite sample estimate of the variance term defined as

$$\sigma^2(\hat{\mathbb{H}}^2(\psi_g)) := \mathbb{E}_{\mathbb{P}^{\text{rct}}} [(h(Z; \psi_g) - \mathbb{E}_{\mathbb{P}^{\text{rct}}} [h(Z; \psi_g)])^2],$$

for all $g \in \mathcal{G}$. Observe that under the assumption that $g^* \in \mathcal{G}$, we have

$$\mathbb{H}_{\text{OPT}}^2 := \min_{g \in \mathcal{G}} \left| \sqrt{\frac{n_{\text{rct}}}{2}} \frac{\hat{\mathbb{H}}^2(\psi_g)}{\hat{\sigma}(\hat{\mathbb{H}}^2(\psi_g))} \right| \leq \left| \sqrt{\frac{n_{\text{rct}}}{2}} \frac{\hat{\mathbb{H}}^2(\psi_{g^*})}{\hat{\sigma}(\hat{\mathbb{H}}^2(\psi_{g^*}))} \right| \quad (4)$$

Therefore, we can achieve validity by comparing $\mathbb{H}_{\text{OPT}}^2$ with the quantiles of the half-normal distribution.

3.3. Theoretical guarantees

Since, in practice, we do not have access to the signal function ψ_g , we define the finite-sample analogous as

$$\hat{\psi}_g(Z) = Y \left(\frac{T}{\pi} - \frac{1-T}{1-\pi} \right) - \hat{\tau}_g^{\text{os}}(X), \quad \text{where}$$

$$\hat{\tau}_g^{\text{os}}(X) := g(X^{\mathcal{J}}) \hat{\tau}_+^{\text{os}}(X) + (1-g(X^{\mathcal{J}})) \hat{\tau}_-^{\text{os}}(X),$$

and $\hat{\tau}_{\pm}^{\text{os}}$ is a consistent estimate of τ_{\pm}^{os} that uses only the observational data D_{os} . We can then define our finite-sample test statistic as

$$\hat{\mathbb{H}}_{\text{OPT}}^2 := \min_{g \in \mathcal{G}} \left| \sqrt{\frac{n_{\text{rct}}}{2}} \frac{\hat{\mathbb{H}}^2(\hat{\psi}_g)}{\hat{\sigma}(\hat{\mathbb{H}}^2(\hat{\psi}_g))} \right|,$$

and the testing function $\hat{\phi}(\alpha) := \mathbb{I} \left\{ \hat{\mathbb{H}}_{\text{OPT}}^2 \geq z_{1-\alpha} \right\}$, where z_{α} is the α -quantile of the half-normal distribution. Below, we provide sufficient conditions for $\hat{\phi}$ to be an asymptotically valid test.

Theorem 3.1 (Validity of the test). *We make the following assumptions:*

- (i) *The variance term is non-zero, i.e.*
 $\mathbb{E}_{\text{Prct}} \left[\psi_{g^*}^2(Z) k^2(X^{\mathcal{J}}, \tilde{X}^{\mathcal{J}}) \psi_{g^*}^2(\tilde{Z}) \right] > 0.$
- (ii) *$\hat{\tau}_{\pm}^{\text{os}}$ satisfy $\|\tau_{\pm}^{\text{os}} - \hat{\tau}_{\pm}^{\text{os}}\|_{L^2(\text{Prct})} = O_{\text{Pos}} \left(\frac{1}{\sqrt{n_{\text{os}}}} \right)$, and it holds that $\lim_{n_{\text{rct}}, n_{\text{os}} \rightarrow \infty} n_{\text{rct}}/n_{\text{os}} = 0.$*

Then, we have that

$$\sqrt{\frac{n_{\text{rct}}}{2}} \frac{\hat{\mathbb{H}}^2(\hat{\psi}_{g^*})}{\hat{\sigma}(\hat{\mathbb{H}}^2(\hat{\psi}_{g^*}))} \rightarrow \mathcal{N}(0, 1), \quad \text{as } n_{\text{rct}}, n_{\text{os}} \rightarrow \infty.$$

Hence, $\hat{\phi}(\alpha)$ is a valid asymptotic test at level α for the null hypothesis $H_0^{\mathcal{G}}$ from Equation (2).

Power of the test While Theorem 3.1 only shows asymptotic validity, we further present guarantees for the asymptotic power of the test in Appendix A.2. In particular, in Theorem A.1, we show that under the alternative hypothesis

$$H_A^{\mathcal{G}} : \inf_{g \in \mathcal{G}} \sup_{\|f\|_{\mathcal{F}} \leq 1} \mathbb{E}_{\text{Prct}} [\psi_g(Z) f(X^{\mathcal{J}})] > 0,$$

the test statistic $\hat{\mathbb{H}}_{\text{OPT}}^2$ in Equation (4) grows at the typical rate of order $\sqrt{n_{\text{rct}}}$ for a fixed function class \mathcal{G} .

3.4. Benchmarking the observational study

Given the theoretical results in this section, we can now introduce our strategy to benchmark observational studies. More concretely, we choose as tolerance functions $\tau_{\pm}^{\text{os}}(X) = \tau^{\text{os}}(X) \pm \delta$, for some constant $\delta \in \mathbb{R}^+$, and we define a data-dependent lower bound on the bias as

$$\hat{\delta}_{\text{LB}} := \inf_{\delta} \{ \delta : \hat{\phi}(\alpha) = 0 \}, \quad (5)$$

where $\hat{\phi}$ depends implicitly on δ via the tolerance functions and we fix $\mathcal{J} = \{1, \dots, d\}$. Then, under the assumptions in Theorem 3.1, it holds that

$$\mathbb{P} \left(\tilde{\delta} \geq \hat{\delta}_{\text{LB}} \right) \geq 1 - \alpha + o_{\mathbb{P}}(1).$$

Crucially, to benchmark the observational study, we propose to compare the lower bound on the bias against a critical value, e.g. the minimum bias strength that would explain away the estimated treatment effect in a subgroup of interest. If the lower bound is greater than the critical value, we discard the conclusions drawn from the observational study. In Section 5, we demonstrate that our strategy yields conclusions consistent with epidemiological knowledge using real-world data from the Women's Health Initiative.

4. Semi-synthetic experiments

4.1. Experimental setting

Dataset We evaluate our testing procedure on a semi-synthetic dataset derived from a real-world randomized trial: Hillstrom's MineThatData Email dataset (Hillstrom, 2008). By default, we use 80% of the full dataset as the os and the remaining 20% as the rct.

Bias model We consider three different models for the bias between studies, given by $\delta^*(x) = \mu^{\text{os}}(x) - \tau^{\text{os}}(x)$, for all $x \in \mathcal{X}$. In Scenario 1, we consider a single subgroup with a constant bias of $\delta^* = 60$, while the rest of os remains unbiased. In Scenario 2 (Figure 5a), we add biases of varying magnitudes across 12 subgroups where the largest bias is $\delta^* = 60$, and it affects only 12% of the observational dataset. Finally, in Scenario 3 (Figure 5b), we model the bias as a quadratic polynomial.

User-defined tolerance and baselines We refer to the testing function proposed in this paper as $\hat{\phi}^{\text{CATE}}$, and we instantiate it using constant upper and lower bounds for the tolerance function: $\tau_{\pm}^{\text{os}}(X) = \tau^{\text{os}}(X) \pm \delta$ for some constant $\delta \in \mathbb{R}^+$. We compare our test against $\hat{\phi}^{\text{ATE}}$, which is a slight modification¹ of the test with tolerance proposed in (De Bartolomeis et al., 2024).

¹ $\hat{\phi}^{\text{ATE}}$ is a t-test for the null hypothesis that average treatment effects between the studies differ at most δ .

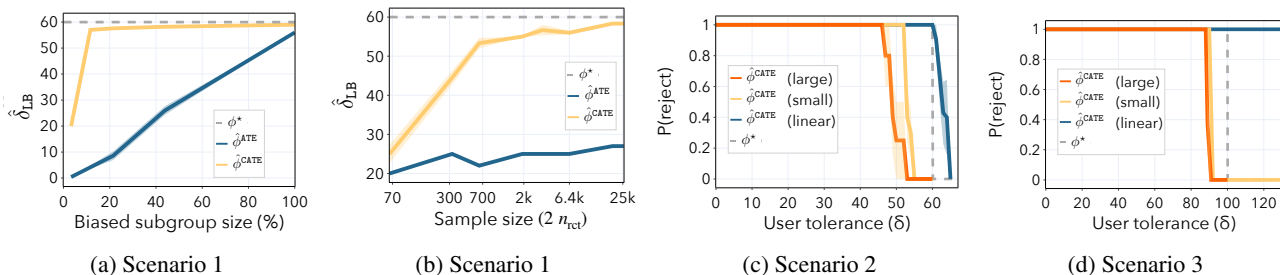


Figure 2. For all the plots: the significance level is set at $\alpha = 0.05$, ϕ^* denotes the oracle test, which rejects for $\delta < \delta^*$. (a-b) Scenario 1, comprising a single subgroup with a constant bias $\delta^* = 60$: we plot the bias lower bound $\hat{\delta}_{LB}$ as a function of (a) the biased subgroup percentage w.r.t. total sample size and (b) the randomized trial sample size. (c-d) Probability of rejection for different function classes \mathcal{G} as a function of the user-specified tolerance δ for (c) Scenario 2 (Figure 5a) based on 12 subgroups with different biases and (d) Scenario 3 (Figure 5b) based on a quadratic polynomial bias. We report mean and standard error over 5 runs.

4.2. Experimental results

We first study the effect of the biased subgroup size (Figure 2a) and the randomized trial sample size (Figure 2b) on the lower bounds $\hat{\delta}_{LB}$ obtained from our test $\hat{\phi}^{CATE}$ and the baseline $\hat{\phi}^{ATE}$. Next, we assess the validity and power of our test $\hat{\phi}^{CATE}$ in two more complex settings: Scenario 2 (Figure 2c) and Scenario 3 (Figure 2d).

Effect of biased subgroup and rct sample sizes Figure 2a shows that our test yields an average lower bound $\hat{\delta}_{LB}$ smaller and close to the true maximum bias δ^* . This implies that the test remains valid and exhibits significant power, even when the biased subgroup represents roughly 14% of the observational dataset. In contrast, $\hat{\phi}^{ATE}$ experiences a significant drop in power as the proportion of biased data points decreases. Such behavior is expected since $\hat{\phi}^{ATE}$ only tests for the difference of averages, and it cannot detect bias in small subgroups, i.e. it is not granular. In Figure 2b, we add a constant bias of 60 to 44% of the observational data points and study the effect of the randomized trial sample size. While our test suffers more than $\hat{\phi}^{ATE}$ from a decrease in the sample size due to the use of kernels, it always yields higher power, even in the very small sample size regime with 70 data points.

Validity and power in complex scenarios Figure 2c and Figure 2d show the validity and power of our testing procedure for Scenario 2 (Figure 5a) and Scenario 3 (Figure 5b), respectively. In both scenarios, if we use a neural network to approximate the bias function, our test remains valid and shows very high power since it rejects the null hypothesis at values of δ close to the true bias δ^* .

Effect of misspecified function class Notably, when g is modeled with a linear function, our test loses its validity, rejecting values of δ that are larger than the true bias. Such behavior is expected as the chosen function class \mathcal{G} lacks the

complexity necessary to capture the true bias model. Nevertheless, we observe that the *small* network with one hidden layer is already sufficient. Further, significantly increasing the complexity – the *large* network has approximately 45 times more parameters than the *small* one – still yields high power. Therefore, we recommend practitioners to be conservative in their choice of function class to ensure validity, even if it might come at the potential cost of some power. Although we cannot guarantee convergence to a global optimum, given the non-convexity of the problem for complex function classes, we show that the optimization procedure is stable and reaches the same solution in Appendix B.2.

5. Real-world experiments

In this section, we provide a concrete application of the benchmarking framework using the Women’s Health Initiative (WHI) study. We show how tolerance and granularity are necessary for effective benchmarking.

5.1. The WHI controversy

The WHI study included a randomized trial and an observational study that investigated the use of hormone therapy (HT) for preventing common sources of mortality among postmenopausal women, including cardiovascular disease, cancer, and fractures (Anderson et al., 2003).

To HT, or not to HT The initial results of the WHI study in 2002 led to fear and confusion regarding the use of hormone therapy (HT) after menopause, resulting in a dramatic reduction in prescriptions for HT around the world. Although in 2002, it was stated that HT increases the risk of coronary heart disease (CHD) for all women, subsequent studies clearly showed that younger women close to menopause can benefit from HT. Further, subsequent randomized trials have continued demonstrating the benefits of HT when started early in young women close to

Table 1. The significance level is set at $\alpha = 0.05$. $\hat{\delta}_{CT}$ is the amount of bias that would explain away the positive effect of HT in young women close to menopause. $\hat{\delta}_{LB}$ is the maximum bias detected in the observational study. $\hat{\phi}_{\delta=0}^{ATE}$ and $\hat{\phi}_{\delta=0}^{CATE}$ denote the respective tests without tolerance, i.e. when the tolerance function is set at $\delta = 0$.

Statistical tests	$\hat{\phi}^{CATE}$	$\hat{\phi}^{ATE}$	$\hat{\phi}_{\delta=0}^{CATE}$	$\hat{\phi}_{\delta=0}^{ATE}$
$\hat{\delta}_{CT}$	0.32	0.32	0.32	0.32
$\hat{\delta}_{LB}$	0.25	0.11	X	X
Reject the study	0	0	1	1

menopause (Hodis et al., 2016; Taylor et al., 2017). To date, the consensus among epidemiologists is that hormone therapy reduces the risk of CHD in women aged less than 60 years and within 10 years of menopause; see e.g. the current guidelines for hormone therapy (Lee et al., 2020).

Limitations of the WHI randomized trial The main issue with the randomized trial is that younger women’s cardiac events are relatively rare. Indeed, not only would it have been prohibitively expensive to conduct a randomized trial exclusively in younger women, but it would have also taken many years to accumulate enough events to reach statistical significance. Hence, the trial lacked enough events to reach statistical significance on the subgroup of interest. On the other hand, the average treatment effect over all the patients suggested an increase in CHD risk because the majority of cardiac events came from older women, and epidemiologists concluded that HT is harmful to all women.

Benchmarking can help! The natural question is, thus, if benchmarking the observational study could have prevented such a turn of events. Indeed, this is the perfect setting to test our methodology, as we would like to ask the question:

Is the bias in the observational study enough to explain away the benefits of HT in young women close to menopause?

In what follows, we show that answering such a question requires a statistical test that offers tolerance. Further, even though we cannot demonstrate that granularity is necessary in this concrete example², we stress that it is equally important in practice.

5.2. Experimental results

Linking back to our question of interest, we demonstrate how our method can provide a correct answer, i.e. one that

²To do so, we would need to know a small biased subgroup in the observational study and show that only the tests with granularity detect the bias. Unfortunately, we are unaware of subgroups that were found to be biased in the WHI study.

aligns with the epidemiology literature. A natural way to do so is to first estimate from the available data the amount of bias that would explain away the treatment effect on the group of interest, defined as

$$\hat{\delta}_{CT} := \left| \mathbb{E}_{\mathbb{P}^{OS}} [\tau^{OS}(X) \mid X \in G] \right|.$$

In essence, the critical value quantifies the minimum strength of bias for which positive and negative values of treatment effect are reasonable, thereby invalidating the observational study results³. In our example, the group G is defined as young women (age ≤ 60) who are close to menopause (≤ 10 years).

Similarly to the semi-synthetic experiments, we instantiate the tolerance functions using constant upper and lower bounds, i.e. $\tau_{\pm}^{OS}(X) = \tau^{OS}(X) \pm \delta$ for some constant $\delta \in \mathbb{R}^+$. We compute the lower bound $\hat{\delta}_{LB}$ on the maximum amount of treatment effect bias in the observational study, as defined in Equation (5). We remark that this quantity can be computed only for tests that allow some tolerance. Then, our decision-making procedure will flag the observational study as invalid if $\hat{\delta}_{LB} \geq \hat{\delta}_{CT}$.

Experimental details We consider a binary-valued outcome: the presence of coronary heart disease within the follow-up period. We choose as covariates X the basic adjustment variables used in many existing analyses, and we further limit patients to those who were not current users of HT at the time of enrolment, as the duration of HT use has been found to have a substantial impact on treatment effects (Prentice et al., 2005; Vandembroucke, 2009). We refer to Appendix C.2 for complete experimental details.

We present evidence that our procedure can yield the conclusions established in the epidemiological literature. It avoids issuing false alarms when the bias is negligible (tolerance) and detects a larger amount of bias, as it is more powerful than tests based on average treatment effect (granularity).

Results In Table 1, we show the result for all the statistical tests on the WHI study. First, we observe that both tests that allow for tolerance correctly do not flag the study, while $\hat{\phi}_{\delta=0}^{CATE}$ and $\hat{\phi}_{\delta=0}^{ATE}$ do. This difference shows the importance of tolerance for distinguishing between small and large amounts of bias. Second, we observe that the lower bound on the bias is larger for the test with granularity $\hat{\phi}^{CATE}$. Such behavior is expected and shows the importance of granularity to detect bias that would otherwise go unnoticed using the test without any granularity $\hat{\phi}^{ATE}$.

³Note that other choices for the critical value are possible, and practitioners should determine the most appropriate one given the specific context.

Acknowledgements

PDB was supported by the Hasler Foundation grant number 21050. JA was supported by the ETH AI Center. KD was supported by the ETH AI Center and the ETH Foundations of Data Science.

References

- Anderson, G., Manson, J., Wallace, R., Lund, B., Hall, D., Davis, S., Shumaker, S., Wang, C.-Y., Stein, E., and Prentice, R. Implementation of the Women’s Health Initiative study design. *Annals of Epidemiology*, 13(9): S5–S17, 2003.
- Dahabreh, I., Robins, J., and Hernán, M. Benchmarking observational methods by comparing randomized trials and their emulations. *Epidemiology*, 31(5):614–619, 2020.
- De Bartolomeis, P., Abad, J., Donhauser, K., and Yang, F. Hidden yet quantifiable: A lower bound for confounding strength using randomized trials. *International Conference on Artificial Intelligence and Statistics*, 2024.
- Delgado, M. Testing the equality of nonparametric regression curves. *Statistics & Probability Letters*, 17(3):199–204, 1993.
- Demirel, I., De Brouwer, E., Hussain, Z., Oberst, M., Philipakis, A., and Sontag, D. Benchmarking observational studies with experimental data under right-censoring. *International Conference on Artificial Intelligence and Statistics*, 2024.
- Duma, N., Vera Aguilera, J., Paludo, J., Haddox, C., Gonzalez Velez, M., Wang, Y., Leventakos, K., Hubbard, J., Mansfield, A., Go, R., et al. Representation of minorities and women in oncology clinical trials: review of the past 14 years. *Journal of Oncology Practice*, 14(1):e1–e10, 2018.
- Forbes, S. and Dahabreh, I. Benchmarking observational analyses against randomized trials: a review of studies assessing propensity score methods. *Journal of General Internal Medicine*, 35:1396–1404, 2020.
- Franklin, J., Glynn, R., Martin, D., and Schneeweiss, S. Evaluating the use of nonrandomized real-world data analyses for regulatory decision making. *Clinical Pharmacology & Therapeutics*, 105(4):867–877, 2019.
- He, Z., Tang, X., Yang, X., Guo, Y., George, T., Charness, N., Quan Hem, K. B., Hogan, W., and Bian, J. Clinical trial generalizability assessment in the big data era: a review. *Clinical and Translational Science*, 13(4):675–684, 2020.
- Hernán, M. and Robins, J. Using big data to emulate a target trial when a randomized trial is not available. *American Journal of Epidemiology*, 183(8):758–764, 2016.
- Hillstrom, K. The MineThatData e-mail analytics and data mining challenge, 2008.
- Hodis, H., Mack, W., Henderson, V., Shoupe, D., Budoff, M., Hwang-Levine, J., Li, Y., Feng, M., Dustin, L., Kono, N., et al. Vascular effects of early versus late postmenopausal treatment with estradiol. *New England Journal of Medicine*, 374(13):1221–1231, 2016.
- Huskova, M. and Janssen, P. Consistency of the generalized bootstrap for degenerate U-statistics. *The Annals of Statistics*, pp. 1811–1823, 1993.
- Hussain, Z., Shih, M.-C., Oberst, M., Demirel, I., and Sontag, D. Falsification of internal and external validity in observational studies via conditional moment restrictions. *International Conference on Artificial Intelligence and Statistics*, 2023.
- Kennedy, E. Towards optimal doubly robust estimation of heterogeneous causal effects. *Electronic Journal of Statistics*, 17(2):3008–3049, 2023.
- Kim, I. and Ramdas, A. Dimension-agnostic inference using cross U-statistics. *Bernoulli*, 30(1):683–711, 2024.
- Klonoff, D. The new FDA real-world evidence program to support development of drugs and biologics. *Journal of Diabetes Science and Technology*, 14(2):345–349, 2020.
- Lee, S. R., Cho, M. K., Cho, Y. J., Chun, S., Hong, S.-H., Hwang, K. R., Jeon, G.-H., Joo, J. K., Kim, S. K., Lee, D. O., et al. The 2020 menopausal hormone therapy guidelines. *Journal of Menopausal Medicine*, 26(2):69, 2020.
- Luedtke, A., Carone, M., and van der Laan, M. An omnibus non-parametric test of equality in distribution for unknown functions. *Journal of the Royal Statistical Society Series B: Statistical Methodology*, 81(1):75–99, 2019.
- Muandet, K., Jitkrittum, W., and Kübler, J. Kernel conditional moment test via maximum moment restriction. *Conference on Uncertainty in Artificial Intelligence*, 2020.
- Neumeyer, N. and Dette, H. Nonparametric comparison of regression curves: an empirical process approach. *The Annals of Statistics*, 31(3):880–920, 2003.
- Platt, R., Brown, J., Robb, M., McClellan, M., Ball, R., Nguyen, M., and Sherman, R. The FDA Sentinel Initiative—an evolving national resource. *New England Journal of Medicine*, 379(22):2091–2093, 2018.

- Prentice, R., Langer, R., Stefanick, M., Howard, B., Pettinger, M., Anderson, G., Barad, D., Curb, D., Kotchen, J., Kuller, L., et al. Combined postmenopausal hormone therapy and cardiovascular disease: toward resolving the discrepancy between observational studies and the Women's Health Initiative clinical trial. *American Journal of Epidemiology*, 162(5):404–414, 2005.
- Racine, J., Hart, J., and Li, Q. Testing the significance of categorical predictor variables in nonparametric regression models. *Econometric Reviews*, 25(4):523–544, 2006.
- Rothwell, P. External validity of randomised controlled trials: “to whom do the results of this trial apply?”. *The Lancet*, 365(9453):82–93, 2005.
- Rubin, D. Bayesian inference for causal effects: The role of randomization. *The Annals of Statistics*, pp. 34–58, 1978.
- Schurman, B. The framework for FDA's real-world evidence program. *Applied Clinical Trials*, 28(4), 2019.
- Serfling, R. *Approximation Theorems of Mathematical Statistics*. Wiley, 1980.
- Taylor, H., Tal, A., Pal, L., Li, F., Black, D., Brinton, E., Budoff, M., Cedars, M., Du, W., Hodis, H., et al. Effects of oral vs transdermal estrogen therapy on sexual function in early postmenopause: ancillary study of the Kronos Early Estrogen Prevention Study (KEEPS). *JAMA Internal Medicine*, 177(10):1471–1479, 2017.
- Vandenbroucke, J. The HRT controversy: observational studies and RCTs fall in line. *The Lancet*, 373(9671): 1233–1235, 2009.
- Viele, K., Berry, S., Neuenschwander, B., Amzal, B., Chen, F., Enas, N., Hobbs, B., Ibrahim, J., Kinnersley, N., Lindborg, S., et al. Use of historical control data for assessing treatment effects in clinical trials. *Pharmaceutical Statistics*, 13(1):41–54, 2014.
- Yang, S., Gao, C., Zeng, D., and Wang, X. Elastic integrative analysis of randomised trial and real-world data for treatment heterogeneity estimation. *Journal of the Royal Statistical Society Series B: Statistical Methodology*, 85(3):575–596, 04 2023.

Appendices

The following appendices provide deferred proofs, experiment details, and ablation studies.

A. Methodology

For the sake of clarity, we write $n := n_{\text{rct}}/2$, $\mathbb{P} := \mathbb{P}^{\text{rct}}$ and $\mathbb{E}[\cdot] := \mathbb{E}_{\mathbb{P}^{\text{rct}}}[\cdot]$ throughout this section.

A.1. Proof of Theorem 3.1

We begin with the simple observation that

$$\min_{g \in \mathcal{G}} \left| \frac{\sqrt{n} \hat{\mathbb{H}}^2(\hat{\psi}_g)}{\hat{\sigma}(\hat{\mathbb{H}}^2(\hat{\psi}_g))} \right| \leq \left| \frac{\sqrt{n} \hat{\mathbb{H}}^2(\hat{\psi}_{g^*})}{\hat{\sigma}(\hat{\mathbb{H}}^2(\hat{\psi}_{g^*}))} \right|,$$

which holds under the assumption that $g^* \in \mathcal{G}$. Thus, asymptotic validity of $\hat{\phi}$ follows when showing that the RHS converges in distribution to an absolute normal distribution.

The key ingredient to prove this statement is to show that the following convergence in probability holds for all fixed n :

$$\hat{\mathbb{H}}^2(\hat{\psi}_{g^*}) \rightarrow \hat{\mathbb{H}}^2(\psi_{g^*}) \quad \text{and} \quad \hat{\sigma}^2(\hat{\mathbb{H}}^2(\hat{\psi}_{g^*})) \rightarrow \hat{\sigma}^2(\hat{\mathbb{H}}^2(\psi_{g^*})), \quad \text{as } n_{\text{os}} \rightarrow \infty. \quad (6)$$

Then, under Assumption (i), we can apply Theorem 4.2 from Kim & Ramdas (2024) to show that

$$\frac{\sqrt{n} \hat{\mathbb{H}}^2(\psi_{g^*})}{\hat{\sigma}(\hat{\mathbb{H}}^2(\psi_{g^*}))} \rightarrow \mathcal{N}(0, 1), \quad \text{as } n \rightarrow \infty.$$

Moreover, as a consequence of Equation (12) and (57) in the proof of Theorem 4.2 from Kim & Ramdas (2024), we have that

$$\frac{1}{n \hat{\sigma}^2(\hat{\mathbb{H}}^2(\psi_{g^*}))} = O_{\mathbb{P}}(1).$$

Thus, when applying Slutsky's Theorem we have that

$$\frac{\sqrt{n} \hat{\mathbb{H}}^2(\hat{\psi}_{g^*})}{\hat{\sigma}(\hat{\mathbb{H}}^2(\hat{\psi}_{g^*}))} \rightarrow \mathcal{N}(0, 1), \quad \text{as } n \rightarrow \infty \text{ and } n_{\text{os}} \rightarrow \infty,$$

and the statement in Theorem 3.1 follows. It now remains to prove Equation 6.

PROOF OF STATEMENT IN EQUATION 6

We begin by defining the error term

$$\Delta := \hat{\psi}_{g^*}(Z) - \psi_{g^*}(Z) = g^*(X^{\mathcal{J}}) (\hat{\tau}_+^{\text{os}}(X) - \tau_+^{\text{os}}(X)) + (1 - g^*(X^{\mathcal{J}})) (\hat{\tau}_-^{\text{os}}(X) - \tau_-^{\text{os}}(X)),$$

and we denote with Δ_i the i.i.d. samples from \mathbb{P} . We restate the definition of the mean and variance terms here. Formally, we split the dataset D_{rct} equally into two folds, \mathcal{I}_1 and \mathcal{I}_2 , of size n and obtain

$$n \hat{\mathbb{H}}^2(\hat{\psi}_{g^*}) = \frac{1}{\sqrt{n}} \sum_{i \in \mathcal{I}_1} \hat{\psi}_i \frac{1}{\sqrt{n}} \sum_{j \in \mathcal{I}_2} k(X_i^{\mathcal{J}}, X_j^{\mathcal{J}}) \hat{\psi}_j, \quad (7)$$

$$n \hat{\sigma}^2(\hat{\mathbb{H}}^2(\hat{\psi}_{g^*})) = \frac{1}{n} \sum_{i \in \mathcal{I}_1} \hat{\psi}_i^2 \left(\frac{1}{\sqrt{n}} \sum_{j \in \mathcal{I}_2} k(X_i^{\mathcal{J}}, X_j^{\mathcal{J}}) \hat{\psi}_j \right)^2 - \left(\sqrt{n} \hat{\mathbb{H}}^2(\hat{\psi}_{g^*}) \right)^2, \quad (8)$$

where we use the shorthand $\hat{\psi}_i := \hat{\psi}_{g^*}(Z_i)$ and $\psi_i := \psi_{g^*}(Z_i)$.

Preliminary step: bounds for the error term Δ_i By Assumption (ii) in Theorem 3.1, we have that

$$\mathbb{E} [\Delta^2] =: \|\Delta\|_{L_2(\mathbb{P})}^2 \leq 2\|\hat{\tau}_+^{\text{os}} - \tau_+^{\text{os}}\|_{L_2(\mathbb{P})}^2 + 2\|\hat{\tau}_-^{\text{os}} - \tau_-^{\text{os}}\|_{L_2(\mathbb{P})}^2 = o_{\mathbb{P}^{\text{os}}} \left(\frac{1}{n} \right), \quad (9)$$

where the probability \mathbb{P}^{os} is over the dataset \mathcal{D}^{os} used to train $\hat{\tau}_{\pm}^{\text{os}}$. We further define

$$\tau_2(X^{\mathcal{J}}) := \frac{1}{\sqrt{n}} \sum_{j \in \mathcal{I}_2} k(X_i^{\mathcal{J}}, X_j^{\mathcal{J}}) \Delta_j \quad \text{and} \quad \tau_1(X^{\mathcal{J}}) := \frac{1}{\sqrt{n}} \sum_{i \in \mathcal{I}_1} k(X_i^{\mathcal{J}}, X_j^{\mathcal{J}}) \Delta_i.$$

We will repeatedly make use of the following bound, which holds analogously for τ_1 :

$$\sup_{X^{\mathcal{J}}} |\mathbb{E} [\tau_2(X^{\mathcal{J}}) | X^{\mathcal{J}}]| = \sup_{X^{\mathcal{J}}} \left| \sqrt{n} \mathbb{E} [k(X^{\mathcal{J}}, \tilde{X}^{\mathcal{J}}) \tilde{\Delta} | X^{\mathcal{J}}] \right| \lesssim \sqrt{n} \sqrt{\mathbb{E} [\Delta^2]},$$

where $\tilde{X}^{\mathcal{J}}$ and $\tilde{\Delta}$ are i.i.d. copies of $X^{\mathcal{J}}$ and Δ , and in the last inequality, we used Cauchy-Schwartz together with the fact that the kernel is uniformly bounded. We will further use that

$$\begin{aligned} \sup_{X^{\mathcal{J}}} \left[\mathbb{E} \left[(\tau_2(X^{\mathcal{J}}))^2 | X^{\mathcal{J}} \right] \right] &= \sup_{X^{\mathcal{J}}} \left[\mathbb{E} \left[\frac{1}{n} \sum_{j \in \mathcal{I}_2} k(X^{\mathcal{J}}, X_j^{\mathcal{J}})^2 \Delta_j^2 + \frac{1}{n} \sum_{\substack{j, j' \in \mathcal{I}_2 \\ j \neq j'}} k(X^{\mathcal{J}}, X_j^{\mathcal{J}}) k(X^{\mathcal{J}}, X_{j'}^{\mathcal{J}}) \Delta_j \Delta_{j'} \middle| X^{\mathcal{J}} \right] \right] \\ &\lesssim \mathbb{E} [\Delta^2] + \sup_{X^{\mathcal{J}}} \left[\frac{n(n-1)}{n} \left(\mathbb{E} [k(X^{\mathcal{J}}, \tilde{X}^{\mathcal{J}}) \tilde{\Delta} | X^{\mathcal{J}}] \right)^2 \right] \\ &\leq \mathbb{E} [\Delta^2] + (n-1) \sup_{X^{\mathcal{J}}} \left[\mathbb{E} [k^2(X^{\mathcal{J}}, \tilde{X}^{\mathcal{J}}) | X^{\mathcal{J}}] \mathbb{E} [\Delta^2] \right] = o_{\mathbb{P}^{\text{os}}}(1), \end{aligned}$$

where we used Cauchy-Schwartz again. We are now ready to show the convergences in Equation (6).

Term 1: controlling $\hat{\mathbb{H}}^2(\hat{\psi}_{g^*})$ We first control the mean term $\hat{\mathbb{H}}^2(\hat{\psi}_{g^*})$. Since n is held fixed, it is equivalent to show that

$$\left| n\hat{\mathbb{H}}^2(\hat{\psi}_{g^*}) - n\hat{\mathbb{H}}^2(\psi_{g^*}) \right| = o_{\mathbb{P}^{\text{os}}}(1). \quad (10)$$

We decompose the difference into the following three terms:

$$n\hat{\mathbb{H}}^2(\hat{\psi}_{g^*}) - n\hat{\mathbb{H}}^2(\psi_{g^*}) = \underbrace{\frac{1}{\sqrt{n}} \sum_{i \in \mathcal{I}_1} \psi_i \tau_2(X_i^{\mathcal{J}})}_{=: T_1} + \underbrace{\frac{1}{\sqrt{n}} \sum_{j \in \mathcal{I}_2} \psi_j \tau_1(X_j^{\mathcal{J}})}_{=: T_2} + \underbrace{\frac{1}{\sqrt{n}} \sum_{i \in \mathcal{I}_1} \Delta_i \tau_2(X_i^{\mathcal{J}})}_{=: T_3}. \quad (11)$$

To control the first two terms, we note that under the null hypothesis in Equation (2), it holds that $\mathbb{E} [\psi_{g^*} | X^{\mathcal{J}} = x] = 0$, for all $x \in \text{supp} \left(\mathbb{P}_{X^{\mathcal{J}}}^{\text{rect}} \right)$. Thus, it suffices to show that the variance goes to zero:

$$\begin{aligned} \text{Var} \left[\frac{1}{\sqrt{n}} \sum_{i \in \mathcal{I}_1} \psi_i \tau_2(X_i^{\mathcal{J}}) \right] &= \mathbb{E} \left[\left(\frac{1}{\sqrt{n}} \sum_{i \in \mathcal{I}_1} \psi_i \tau_2(X_i^{\mathcal{J}}) \right)^2 \right] = \mathbb{E} \left[\mathbb{E} [\psi^2 | X^{\mathcal{J}}] (\tau_2(X^{\mathcal{J}}))^2 \right] \\ &\lesssim \mathbb{E} [(\tau_2(X^{\mathcal{J}}))^2] = o_{\mathbb{P}^{\text{os}}}(1), \end{aligned}$$

where we used that the conditional second moment $\mathbb{E} [\psi_{g^*}^2 | X^{\mathcal{J}}]$ is uniformly bounded, since the outcome Y and the tolerance function τ_{\pm}^{os} are both bounded. Further, the same argument also applies when swapping \mathcal{I}_1 with \mathcal{I}_2 , we thus can conclude from Chebyshev's inequality that

$$|T_1| = o_{\mathbb{P}^{\text{os}}}(1) \quad \text{and} \quad |T_2| = o_{\mathbb{P}^{\text{os}}}(1).$$

Next, we bound the last term T_3 . We first consider the mean of T_3 :

$$\begin{aligned}\mathbb{E}[T_3] &= \sqrt{n} \mathbb{E}[\Delta \tau_2(X^{\mathcal{J}})] = \sqrt{n} \mathbb{E}[\Delta \mathbb{E}[\tau_2(X^{\mathcal{J}})|X^{\mathcal{J}}]] \leq \sup_{X^{\mathcal{J}}} [|\mathbb{E}[\tau_2(X^{\mathcal{J}})|X^{\mathcal{J}}]|] \mathbb{E}[\sqrt{n}|\Delta|] \\ &\leq \sup_{X^{\mathcal{J}}} [|\mathbb{E}[\tau_2(X^{\mathcal{J}})|X^{\mathcal{J}}]|] \sqrt{n} \sqrt{\mathbb{E}[\Delta^2]} \\ &= o_{\mathbb{P}^{\text{pos}}}(1).\end{aligned}$$

Then, we consider the variance of T_3 :

$$\mathbb{E}[T_3^2] = \mathbb{E}[\Delta^2 (\tau_2(X^{\mathcal{J}}))^2] = \mathbb{E}[\Delta^2 \mathbb{E}[(\tau_2(X^{\mathcal{J}}))^2 | X^{\mathcal{J}}]] \leq \sup_{X^{\mathcal{J}}} [\mathbb{E}[(\tau_2(X^{\mathcal{J}}))^2 | X^{\mathcal{J}}]] \mathbb{E}[\Delta^2] = o_{\mathbb{P}^{\text{pos}}}(1).$$

Thus, we can conclude that $|T_3| = o_{\mathbb{P}}(1)$, and the equality in Equation (10) follows.

Term 2: controlling $\hat{\sigma}^2(\hat{\mathbb{H}}^2(\hat{\psi}_{g^*}))$ As a second step, we control the variance term $\hat{\sigma}^2(\hat{\mathbb{H}}^2(\hat{\psi}_{g^*}))$. Our goal is again to show that

$$\left| n\hat{\sigma}^2(\hat{\mathbb{H}}^2(\hat{\psi}_{g^*})) - n\hat{\sigma}^2(\hat{\mathbb{H}}^2(\psi_{g^*})) \right| = o_{\mathbb{P}^{\text{pos}}}(1).$$

Given the results from the previous paragraph in Equation (10), we note that it suffices to show that

$$\left| \frac{1}{n} \sum_{i \in \mathcal{I}_1} \hat{\psi}_i^2 \left(\frac{1}{\sqrt{n}} \sum_{j \in \mathcal{I}_2} k(X_i^{\mathcal{J}}, X_j^{\mathcal{J}}) \hat{\psi}_j \right)^2 - \frac{1}{n} \sum_{i \in \mathcal{I}_1} \psi_i^2 \left(\frac{1}{\sqrt{n}} \sum_{j \in \mathcal{I}_2} k(X_i^{\mathcal{J}}, X_j^{\mathcal{J}}) \psi_j \right)^2 \right| = o_{\mathbb{P}^{\text{pos}}}(1). \quad (12)$$

We begin again by decomposing the difference of the two terms on the LHS into the following six terms:

$$\begin{aligned}&= \underbrace{\frac{1}{n} \sum_{i \in \mathcal{I}_1} \Delta_i^2 \left(\frac{1}{\sqrt{n}} \sum_{j \in \mathcal{I}_2} k(X_i^{\mathcal{J}}, X_j^{\mathcal{J}}) (\psi_j + \Delta_j) \right)^2}_{=: T_1} + \underbrace{\frac{1}{n} \sum_{i \in \mathcal{I}_1} (\psi_i + \Delta_i)^2 (\tau_2(X_i^{\mathcal{J}}))^2}_{=: T_2} - \underbrace{\frac{1}{n} \sum_{i \in \mathcal{I}_1} \Delta_i^2 (\tau_2(X_i^{\mathcal{J}}))^2}_{=: T_3} \\ &+ \underbrace{\frac{2}{n} \sum_{i \in \mathcal{I}_1} \psi_i^2 \left(\frac{1}{\sqrt{n}} \sum_{j \in \mathcal{I}_2} k(X_i^{\mathcal{J}}, X_j^{\mathcal{J}}) \psi_j \right) \tau_2(X_i^{\mathcal{J}})}_{=: T_4} + \underbrace{\frac{4}{n} \sum_{i \in \mathcal{I}_1} \psi_i \Delta_i \left(\frac{1}{\sqrt{n}} \sum_{j \in \mathcal{I}_2} k(X_i^{\mathcal{J}}, X_j^{\mathcal{J}}) \psi_j \right) \tau_2(X_i^{\mathcal{J}})}_{=: T_5} \\ &+ \underbrace{\frac{2}{n} \sum_{i \in \mathcal{I}_1} \psi_i \Delta_i \left(\frac{1}{\sqrt{n}} \sum_{j \in \mathcal{I}_2} k(X_i^{\mathcal{J}}, X_j^{\mathcal{J}}) \psi_j \right)^2}_{=: T_6}.\end{aligned}$$

We now show that $\forall i \in [1, \dots, 6]$, $|T_i| = o_{\mathbb{P}^{\text{pos}}}(1)$.

Controlling T_1 : Since the term is non-negative, it suffices to show that the expectation $\mathbb{E}[T_1] = o_{\mathbb{P}^{\text{pos}}}(1)$ and then apply

Markov's inequality. More formally, we have

$$\begin{aligned}
 \mathbb{E}[T_1] &= \mathbb{E} \left[\Delta^2 \left(\frac{1}{\sqrt{n}} \sum_{j \in \mathcal{I}_2} k(X^{\mathcal{J}}, X_j^{\mathcal{J}}) (\psi_j + \Delta_j) \right)^2 \right] \\
 &= \mathbb{E} \left[\Delta^2 \left[\frac{1}{n} \sum_{j \in \mathcal{I}_2} k^2(X^{\mathcal{J}}, X_j^{\mathcal{J}}) \psi_j^2 + \left(\frac{1}{\sqrt{n}} \sum_{j \in \mathcal{I}_2} k(X^{\mathcal{J}}, X_j^{\mathcal{J}}) \Delta_j \right)^2 + \frac{1}{n} \sum_{j \in \mathcal{I}_2} k^2(X^{\mathcal{J}}, X_j^{\mathcal{J}}) \psi_j \Delta_j \right] \right] \\
 &= \mathbb{E} \left[\Delta^2 k(X^{\mathcal{J}}, \tilde{X}^{\mathcal{J}})^2 [\tilde{\psi}^2 + \tilde{\psi} \tilde{\Delta}] \right] + \mathbb{E} \left[\Delta^2 (\tau_2(X^{\mathcal{J}}))^2 \right] \\
 &\lesssim \mathbb{E}[\Delta^2] \left[\mathbb{E}[\psi^2] + \sqrt{\mathbb{E}[\psi^2] \mathbb{E}[\Delta^2]} + \sup_{X^{\mathcal{J}}} \mathbb{E}[(\tau_2(X^{\mathcal{J}}))^2 | X^{\mathcal{J}}] \right] \\
 &= o_{\mathbb{P}^{\text{pos}}}(1),
 \end{aligned}$$

where in the second equality we use again $\mathbb{E}[\psi_{g^*} | X^{\mathcal{J}} = x] = 0$, for all $x \in \text{supp}(\mathbb{P}_{X^{\mathcal{J}}}^{\text{prct}})$.

Controlling T_2 and T_3 : We can again upper-bound the expectation and apply Markov's inequality. We have

$$\mathbb{E}[T_2] \leq \mathbb{E} \left[(2\psi^2 + 2\Delta^2) \mathbb{E}[(\tau_2(X^{\mathcal{J}}))^2 | X^{\mathcal{J}}] \right] = \sup_{X^{\mathcal{J}}} \left[\mathbb{E}[(\tau_2(X^{\mathcal{J}}))^2 | X^{\mathcal{J}}] \right] (2\mathbb{E}[\psi^2] + 2\mathbb{E}[\Delta^2]) = o_{\mathbb{P}^{\text{pos}}}(1),$$

and thus it also follows that $\mathbb{E}[T_3] = o_{\mathbb{P}^{\text{pos}}}(1)$.

Controlling T_4 , T_5 and T_6 : We note that the expectations $\mathbb{E}[T_4] = 0$, $\mathbb{E}[T_5] = 0$ and $\mathbb{E}[T_6] = 0$ are all zero. Thus, we can bound the terms in probability by showing that the respective variances converge to zero and then applying Chebyshev's inequality. We first upper-bound the variance of T_4 :

$$\begin{aligned}
 \text{Var}[T_4] &= \text{Var} \left[\frac{2}{n} \sum_{i \in \mathcal{I}_1} \psi_i^2 \left(\frac{1}{\sqrt{n}} \sum_{j \in \mathcal{I}_2} k(X_i^{\mathcal{J}}, X_j^{\mathcal{J}}) \psi_j \right) \tau_2(X_i^{\mathcal{J}}) \right] \\
 &= \frac{4}{n} \mathbb{E} \left[\psi^4 \left(\frac{1}{\sqrt{n}} \sum_{j \in \mathcal{I}_2} k(X^{\mathcal{J}}, X_j^{\mathcal{J}}) \psi_j \right)^2 (\tau_2(X^{\mathcal{J}}))^2 \right] \\
 &= 4\mathbb{E} \left[\psi^4 \left(\frac{1}{n} \sum_{j \in \mathcal{I}_2} k(X^{\mathcal{J}}, X_j^{\mathcal{J}}) \psi_j \right)^2 (\tau_2(X^{\mathcal{J}}))^2 \right] \\
 &\lesssim \sup_{X^{\mathcal{J}}} \mathbb{E}[(\tau_2(X^{\mathcal{J}}))^2 | X^{\mathcal{J}}] \\
 &= o_{\mathbb{P}^{\text{pos}}}(1),
 \end{aligned}$$

where we use the fact that both the kernel k and the fourth conditional moment of $\psi_{g^*} | X^{\mathcal{J}}$ are almost surely upper bounded by a constant. Next, we bound the variance of T_5 :

$$\begin{aligned}
 \text{Var}[T_5] &= \text{Var} \left[\frac{2}{n} \sum_{i \in \mathcal{I}_1} \psi_i \Delta_i \left(\frac{1}{\sqrt{n}} \sum_{j \in \mathcal{I}_2} k(X_i^{\mathcal{J}}, X_j^{\mathcal{J}}) \psi_j \right) \tau_2(X_i^{\mathcal{J}}) \right] \\
 &= 4\mathbb{E} \left[\psi^2 \Delta^2 \left(\frac{1}{n} \sum_{j \in \mathcal{I}_2} k(X^{\mathcal{J}}, X_j^{\mathcal{J}}) \psi_j \right)^2 (\tau_2(X^{\mathcal{J}}))^2 \right] \\
 &= o_{\mathbb{P}^{\text{pos}}}(1).
 \end{aligned}$$

Finally, we upper-bound the variance of the term T_6 :

$$\begin{aligned} \text{Var}[T_6] &= \text{Var} \left[\frac{2}{n} \sum_{i \in \mathcal{I}_1} \psi_i \Delta_i \left(\frac{1}{\sqrt{n}} \sum_{j \in \mathcal{I}_2} k(X_i^{\mathcal{J}}, X_j^{\mathcal{J}}) \psi_j \right) \right]^2 \\ &= \frac{4}{n} \mathbb{E} \left[\psi^2 \Delta^2 \left(\frac{1}{n^2} \sum_{j \in \mathcal{I}_2} k(X^{\mathcal{J}}, X_j^{\mathcal{J}})^4 \psi_j^4 + \frac{6}{n^2} \sum_{j, j' \in \mathcal{I}_2; j \neq j'} k(X^{\mathcal{J}}, X_j^{\mathcal{J}})^2 k(X^{\mathcal{J}}, X_{j'}^{\mathcal{J}})^2 \psi_j^2 \psi_{j'}^2 \right) \right] \\ &= o_{\mathbb{P}^{\text{os}}}(1). \end{aligned}$$

As a result, we conclude that $|T_4| = o_{\mathbb{P}^{\text{os}}}(1)$, $|T_5| = o_{\mathbb{P}^{\text{os}}}(1)$, and $|T_6| = o_{\mathbb{P}^{\text{os}}}(1)$.

Discussion of assumptions Assumption (i) is mild and applies to very general settings, e.g. it is satisfied when Y is a non-deterministic random variable. Assumption (ii) is stronger and generally only expected to hold when $n_{\text{os}} \gg n_{\text{rct}}$ and the support of the randomized control trial is contained in the support of the observational study, i.e. $\text{supp}(\mathbb{P}_X^{\text{rct}}) \subseteq \text{supp}(\mathbb{P}_X^{\text{os}})$. These two conditions are realistic in our setting, as they align with the standard design of observational studies (Franklin et al., 2019; Schurman, 2019; He et al., 2020). Further, we remark that previous works either assume oracle access to the functions τ_{\pm}^{os} (Hussain et al., 2023; Demirel et al., 2024) or impose similar assumptions on the rates (De Bartolomeis et al., 2024).

Why not a classic U-statistic? We remark that it is not clear how to test the null hypothesis $H_0^{\mathcal{G}}$ using the classic U-statistic (Serfling, 1980), as done in previous works (see e.g. (Hussain et al., 2023; Demirel et al., 2024)). The main challenge is that under the null hypothesis $\mathbb{H}^2(\psi_{g^*}) = 0$, the U-statistic converges in distribution to a weighted χ^2 -statistic. However, estimating the quantiles (needed for a valid test) of this asymptotic distribution via bootstrapping requires knowing the function g^* (Huskova & Janssen, 1993). In contrast, our test statistic $\mathbb{H}_{\text{OPT}}^2$ is bounded by a valid asymptotic pivot, i.e. a function of the data and the unknown function g^* whose asymptotic distribution does not depend on g^* . Hence, we can compute the quantiles of the RHS in Equation (4) and construct an asymptotically valid test.

A.2. Power of the test

We first discuss a simple result on the power of the test from Equation (4) in rejecting an alternative hypothesis. While Hussain et al. (2023); Muandet et al. (2020) show asymptotic normality for the kernel conditional moment test statistics \mathbb{M}^2 under the alternative hypothesis that $\mathbb{M}^2 \neq 0$, the same result does not directly apply to our test statistic. However, as we show in the following theorem, our test statistic grows at a rate \sqrt{n} . For the sake of clarity, we only prove the result for the oracle test statistic $\mathbb{H}_{\text{OPT}}^2$ computed from ψ_g . Nevertheless, we remark that our result can be easily extended to the empirical test statistics via the same argument used in the proof of Theorem 3.1.

Theorem A.1. *Assume that for every $\epsilon > 0$, the function class \mathcal{G} has a finite ℓ_{∞} -norm covering number. Then, we can lower-bound the test statistic, in probability as $n \rightarrow \infty$, by*

$$\mathbb{H}_{\text{OPT}}^2 = \min_{g \in \mathcal{G}} \left| \frac{\sqrt{n} \hat{\mathbb{H}}^2(\psi_g)}{\hat{\sigma}(\hat{\mathbb{H}}^2(\psi_g))} \right| \gtrsim \sqrt{n} \left(\inf_{g \in \mathcal{G}} \sup_{\|f\|_{\mathcal{F}} \leq 1} \mathbb{E}[\psi_g(Z) f(X^{\mathcal{J}})] \right)^2,$$

where we use \gtrsim to hide universal constants not depending on n .

Thus, under the alternative hypothesis

$$H_A^{\mathcal{G}} : \inf_{g \in \mathcal{G}} \sup_{\|f\|_{\mathcal{F}} \leq 1} \mathbb{E}_{\text{prct}}[\psi_g(Z) f(X^{\mathcal{J}})] > 0,$$

the RHS grows at a rate \sqrt{n} , which implies that our test has an asymptotic power of one (note that the same rate is achieved by existing conditional moment tests (Hussain et al., 2023; Muandet et al., 2020)).

Proof of Theorem A.1 Let us define

$$T := \inf_{g \in \mathcal{G}} \sup_{\|f\|_{\mathcal{F}} \leq 1} \mathbb{E}[\psi_g(Z) f(V)],$$

and note that if $T = 0$ the result follows trivially. Thus, we may assume that $T > 0$ is some constant independent of n . Additionally, observe that ψ_g is uniformly bounded since the outcome Y is a bounded random variable. Therefore, the variance term $\hat{\sigma}^2(\hat{\mathbb{H}}^2(\psi_g))$ is also uniformly bounded, and it suffices to show that $\hat{\mathbb{H}}^2(\psi_g) = \Omega_{\mathbb{P}}(1)$ is lower bounded in probability.

Controlling $\hat{\mathbb{H}}^2(\psi_g)$ First, recall that our test statistic is given by

$$\hat{\mathbb{H}}^2(\psi_g) = \frac{1}{n^2} \sum_{i=1}^n \sum_{j=n+1}^{2n} \psi_g(Z_i) k(X_i^{\mathcal{J}}, X_j^{\mathcal{J}}) \psi_g(Z_j). \quad (13)$$

Further, for all $g \in \mathcal{G}$, it holds that

$$\mathbb{E} \left[\hat{\mathbb{H}}^2(\psi_g) \right] = \mathbb{E} \left[\psi_g(Z) k(X^{\mathcal{J}}, \tilde{X}^{\mathcal{J}}) \psi_g(\tilde{Z}) \right] \geq \inf_{g \in \mathcal{G}} \mathbb{E} \left[\psi_g(Z) k(X^{\mathcal{J}}, \tilde{X}^{\mathcal{J}}) \psi_g(\tilde{Z}) \right] = T^2,$$

where \tilde{Z} is an independent copy of Z following the same distribution, and the last equality follows from Equation (3). Thus, it suffices to show that the following inequality holds with probability one as $n \rightarrow \infty$

$$\sup_{g \in \mathcal{G}} \left| \hat{\mathbb{H}}^2(\psi_g) - \mathbb{E} \left[\hat{\mathbb{H}}^2(\psi_g) \right] \right| \leq \frac{T^2}{2}.$$

We use a simple ϵ -net argument to show this result. Let \mathcal{G}_ϵ be the epsilon net in ℓ_∞ distance of balls with radii ϵ . Then, since ψ_g is uniformly bounded, it holds that for all Z and $g \in \mathcal{G}$,

$$\inf_{\tilde{g} \in \mathcal{G}_\epsilon} |\psi_g(Z) - \psi_{\tilde{g}}(Z)| \lesssim \epsilon.$$

Thus, from the definition of $\hat{\mathbb{H}}^2(\psi_g)$ in Equation (13) it follows that we can choose a constant $\epsilon > 0$, such that the following inequality holds almost surely,

$$\sup_{g \in \mathcal{G}} \inf_{\tilde{g} \in \mathcal{G}_\epsilon} \left| \hat{\mathbb{H}}^2(\psi_g) - \hat{\mathbb{H}}^2(\psi_{\tilde{g}}) \right| \leq \frac{T^2}{4}. \quad (14)$$

Then, for any constant $c > 0$, it holds that

$$\begin{aligned} & \mathbb{P} \left(\sup_{g \in \mathcal{G}} \left| \hat{\mathbb{H}}^2(\psi_g) - \mathbb{E} \left[\hat{\mathbb{H}}^2(\psi_g) \right] \right| \leq \frac{T^2}{2} \right) \\ &= \mathbb{P} \left(\sup_{g \in \mathcal{G}} \inf_{\tilde{g} \in \mathcal{G}_\epsilon} \left| \hat{\mathbb{H}}^2(\psi_g) - \hat{\mathbb{H}}^2(\psi_{\tilde{g}}) + \hat{\mathbb{H}}^2(\psi_{\tilde{g}}) - \mathbb{E}[\hat{\mathbb{H}}^2(\psi_{\tilde{g}})] + \mathbb{E}[\hat{\mathbb{H}}^2(\psi_{\tilde{g}})] - \mathbb{E}[\hat{\mathbb{H}}^2(\psi_g)] \right| \right) \\ &\stackrel{(i)}{\geq} \mathbb{P} \left(\sup_{\tilde{g} \in \mathcal{G}_\epsilon} \left| \hat{\mathbb{H}}^2(\psi_{\tilde{g}}) - \mathbb{E} \left[\hat{\mathbb{H}}^2(\psi_{\tilde{g}}) \right] \right| \leq 2c \right) \\ &\stackrel{(ii)}{\geq} 1 - \sum_{\tilde{g} \in \mathcal{G}_\epsilon} \underbrace{\mathbb{P} \left(\left| \hat{\mathbb{H}}^2(\psi_{\tilde{g}}) - \mathbb{E} \left[\hat{\mathbb{H}}^2(\psi_{\tilde{g}}) \right] \right| \geq 2c \right)}_{\xrightarrow{n \rightarrow \infty} 0 \text{ (L.L.N.)}}, \end{aligned}$$

where (i) follows from applying the inequality in Equation (14) and (ii) follows since, by assumption, for every fixed $\epsilon > 0$ the cover $|\mathcal{G}_\epsilon| < \infty$ is constant as a function of n .



Figure 3. For all the plots: the significance level is set at $\alpha = 0.05$, and the bias model is from Scenario 2. (a) Effect of varying the feature set X^J on the average lower bound $\hat{\delta}_{LB}$, illustrating the trade-off between feature set size and the power of the test. ϕ^* represents the oracle test, which rejects for $\delta < \delta^*$. The highest power is achieved when the feature set size $|X^J| = 3$, including only the relevant features to model the bias. We average runs over 5 seeds and report the standard error. (b) Evolution of the test statistic with respect to the training epochs using the small neural network. We set the user tolerance to $\delta = 58$, close to the maximum true bias $\delta^* = 60$. The dashed red line represents the α -quantile of the absolute normal distribution.

B. Additional experiments

B.1. Ablation study of the feature subset X^J

In Scenario 2 from Figure 5a, we introduced constant bias in the subgroups resulting from different combinations of the features `newbie`, `mens` and `channel`, with a maximum true bias $\delta^* = 60$. Figure 3a shows the effect of the selected feature set X^J on the average lower bound $\hat{\delta}_{LB}$ for the bias model from Scenario 2. When $|X^J| = 3$, we select the features that capture the bias between `rct` and `os` datasets (`newbie`, `mens`, `channel`), and hence we achieve the highest power. Intuitively, if the feature set is smaller, some of the bias averages out, and the test loses power. On the other hand, when increasing the feature set, the test loses power due to the curse of dimensionality, being particularly severe with smaller sample sizes. After $|X^J| = 6$ (i.e. $X^J = X$), we add redundant features sampled from a standard normal distribution $\mathcal{N}(0, 1)$.

B.2. Convergence of the optimization procedure

We provide evidence that our testing procedure is reliable, meaning that the optimizer consistently reaches the same solution for the bias model from Scenario 2 and the small neural network model. Recall that, given the non-convex nature of the optimization problem, we cannot guarantee convergence to the true global minimum g^* . Figure 3b shows the test statistic as a function of the training epoch under different random network initializations. We observe that the test statistic consistently reaches the same minimum and that the optimization stabilizes after 10000 epochs.

B.3. Interpretability of the testing procedure

Similar to the test proposed by Hussain et al. (2023), our testing procedure outputs a “witness function” that enables practitioners to identify the most biased subgroups within the observational dataset. Additionally, our witness function provides insights into the bias strength and direction for each subgroup. This is achieved by minimizing the objective in Equation (4), where we learn the bias function \hat{g} . If the function class \mathcal{G} is sufficiently rich to model the bias structure, and the optimizer converges to the global minima, we expect \hat{g} to be a good approximation of g^* .

To interpret this bias function, we observe that $\hat{g}(X) \in [0, 1]$ interpolates between the tolerance bounds $\tau_-^{\text{os}}(X)$ and $\tau_+^{\text{os}}(X)$; therefore, values close to zero indicate a negative bias of magnitude close to user-tolerance δ , while values close to one indicate the same for positive bias. Hence, we can estimate the subgroup bias as

$$\text{bias}(G) = \hat{\delta}_{LB} \left(\frac{2}{|G|} \sum_{X_i \in G} \hat{g}(X_i) - 1 \right), \quad (15)$$

where G represents the subgroup of interest. In Figure 4, we illustrate how practitioners could use the witness function for Scenario 2, where the categorical nature of the features defines subgroups. We compare the estimated bias with the ground

truth and observe that our estimates closely align with the true bias model. In scenarios where subgroups are not predefined, a practitioner can select the bottom or top 10% of witness function values, as suggested by Hussain et al. (2023).

However, it is important to note that, unlike the approach by Hussain et al. (2023), we do not have guarantees for the correctness of the witness function, i.e. we cannot guarantee that $\hat{g} \rightarrow g^*$. Therefore, any claims based on it should be approached cautiously and contrasted with domain expertise.

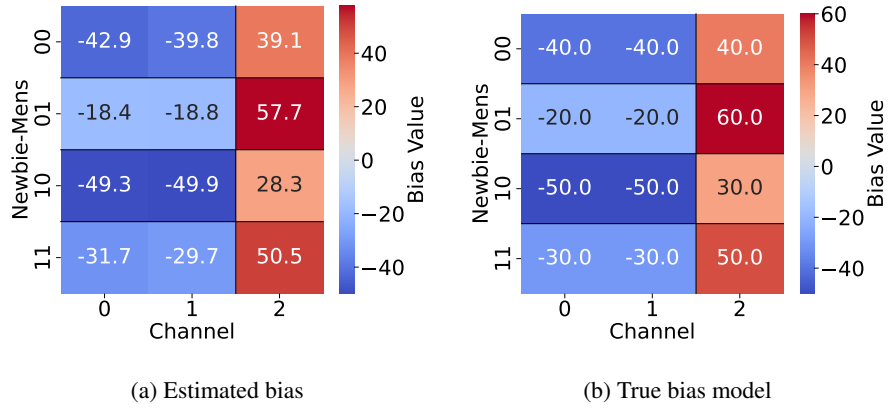


Figure 4. Comparison between the estimated and true bias models for Scenario 2. Our estimates of the bias from Equation (15) closely align with the true bias. We run the test with a random seed, using the same hyperparameters as in our experimental evaluation, and set the user tolerance to $\delta = 57$.

C. Experimental details

C.1. Hillstrom’s MineThatData

Hillstrom’s MineThatData Email dataset (Hillstrom, 2008) is a large-scale, real-world randomized trial that contains records of 64,000 customers who made purchases online within the last twelve months. They were part of an email campaign designed to assess the effectiveness of different campaign strategies. Two treatment groups, “Men’s” and “Women’s” email campaigns, and a control group were established, with treatments randomly assigned. Our analysis primarily focuses on a combined treatment group, which constitutes approximately 66% of the dataset. Although the original dataset presents various outcomes, including binary indicators of customers visiting or purchasing in the days after the campaign, we focus on the dollars spent in the two weeks post-campaign. The dataset provides data on annual spending (`history`), merchandise type (`mens` and `womens`), geographical location (`zip code`), newcomer status (`newbie`), and purchasing avenues (`channel`). We, therefore, discard features describing the history segment (`history segment`) and recency of the last purchase (`recency`). Since the average treatment effect is close to zero, we add a constant shift of 30 to all treated individuals, allowing us more flexibility to introduce bias. We normalize continuous features and one-hot-encode categorical features, resulting in a 13-dimensional dataset. By default, we use 80% of the full dataset as the observational study (`os`), and the remaining 20% as the randomized controlled trial (`rct`).

We fit the propensity score using logistic regression with default hyperparameters from `scikit-learn`. We train a `Random Forest Classifier` for the selection score (`rct` or `os`), also with default hyperparameters from `scikit-learn`. Finally, we estimate the CATE functions using the doubly-robust learner from Kennedy (2023), instantiating `Random Forest Regressors` for the potential outcome functions and the pseudo-outcome regression, fixing hyperparameters to 300 `tree` estimators with a maximum `depth` of 6.

Bias models We illustrate the bias model for Scenario 2 and Scenario 3 in Figure 5. For scenario 3, we sample the coefficient for the polynomial bias model in Figure 5a from a normal distribution $\mathcal{N}(0, 0.01^2)$.

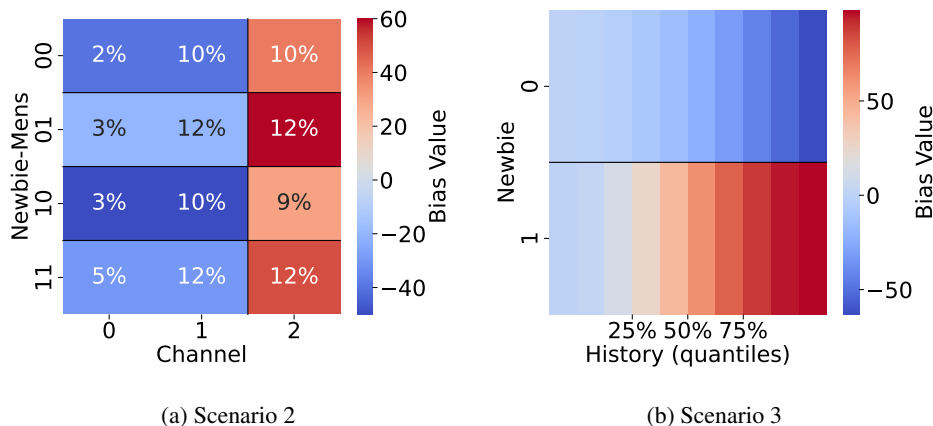


Figure 5. Heatmap visualizations of the bias for (a) Scenario 2 based on 12 subgroups with different biases (the numbers in the cells represent the percentage w.r.t. the full observational dataset), and (b) Scenario 3 based on a quadratic polynomial bias.

Implementation We use the Laplacian kernel with a scale of 1.0 to compute our test statistic $\hat{\phi}^{\text{CATE}}$. We perform gradient descent for 6000 epochs using the Adam optimizer from the JAX-based library `optax` with its default hyperparameters and record the smallest test statistic. As function class \mathcal{G} , we consider linear functions and two multilayer perceptrons (MLPs), one *small* and one *large*, with hidden layer widths of 10 and 100-50-10-5 neurons, respectively. For the linear function and the small MLP, we set the learning rate to 0.1, and for the large MLP, we set it to 0.01. For the test $\hat{\phi}^{\text{ATE}}$, we use 500 bootstrap samples to estimate the variance of the test statistic.

C.2. Women’s Health Initiative

The Women’s Health Initiative (WHI) is a long-term national health study that has focused on strategies for preventing the major causes of death, disability, and frailty in older women, specifically heart disease, cancer, and osteoporotic fractures.

This multi-million dollar, 20+ year project, sponsored by the National Institutes of Health (NIH) and the National Heart, Lung, and Blood Institute (NHLBI), initially enrolled 161,808 women aged 50-79 between 1993 and 1998. The WHI was one of the most definitive, far-reaching clinical trials of post-menopausal women’s health ever undertaken in the US.

The WHI had two major parts: a randomized trial and an observational study. The randomized trial enrolled 68,132 women in trials testing three prevention strategies. Eligible women could choose to enroll in one, two, or three of the trial components.

- A Hormone Therapy Trial (HT) that examined the effects of combined hormones or estrogen alone on the prevention of heart disease and osteoporotic fractures and associated risk for breast cancer.
- A Dietary Modification Trial (DM) that evaluated the effect of a low-fat and high-fruit, vegetable, and grain diet on preventing breast and colorectal cancers and heart disease.
- A Calcium and Vitamin D Trial (CaD) that evaluated the effect of calcium and vitamin D supplementation on preventing osteoporotic fractures and colorectal cancer.

The Observational Study (OS) examines the relationship between lifestyle, health risk factors, and disease outcomes. This component involves tracking the medical events and health habits of 93,676 women. Recruitment for the observational study was completed in 1998, and participants have been followed since.

We use observational study and randomized trial data from the Women’s Health Initiative (WHI) to assess our method in a real-world scenario. We use the Hormone Therapy (HT) trial as the RCT in our analysis ($n_{\text{rct}} = 16,608$), run on postmenopausal women aged 50-79 years with an intact uterus. The trial investigated the effect of hormone therapy on several types of cancers, cardiovascular events, and fractures, measuring the “time-to-event” for each outcome. In the WHI setup, the observational study component was run in parallel, and outcomes were tracked similarly to those of the RCT.

Data preprocessing We binarize a composite outcome, where $Y = 1$ if coronary heart disease was observed in the first seven years of follow-up, and $Y = 0$ otherwise. To establish treatment and control groups in the observational study, we use questionnaire data in which participants confirm or deny usage of combination hormones (i.e. both estrogen and progesterone) in the first three years. Using this procedure, we end up with a total of $n_{\text{os}} = 33,511$ patients. Finally, we restrict the set of covariates used to those that are measured in both the RCT and the observational study. In particular, we use as covariates only those measured in both the RCT and observational study, and we further restrict them to those identified as significant in epidemiological literature, such as in (Prentice et al., 2005). Specifically, the covariates in our analysis are: AGE, ETHNIC_White, BMI, SMOKING_Past_Smoker, SMOKING_Current_Smoker, EDUC_x_College_graduate_or_Baccalaureate Degree, EDUC_x_Some_post_graduate_or_professional, MENO, PHYSFUN. The data used is available on BIOLINCC.

Experimental details We use a gaussian kernel with `bandwidth = 1.0`. The set of features for the granularity of the test is chosen to be $J = \{\text{AGE}, \text{MENO}\}$. We use a logistic regression model for both the outcome model and propensity score (default hyperparameters in `scikit-learn` were used). We train a neural (1 hidden layer and 10 neurons) network with Adam, with a `learning rate` of 0.01 for 500 epochs. We repeat the optimization for 10 seeds with different initializations to ensure that we converge.