ROYAL STATISTICAL SOCIETY
DATA | EVIDENCE | DECISIONS

Journal of the Statistics Society
**Series B**
Statistical Methodology

**B**

**Original Article**

# Elastic integrative analysis of randomised trial and real-world data for treatment heterogeneity estimation

## Shu Yang[1] 📙, Chenyin Gao[1] 📙, Donglin Zeng[2] and Xiaofei Wang[3]

[1]Department of Statistics, North Carolina State University, Raleigh, NC 27695, USA
[2]Department of Biostatistics, University of North Carolina at Chapel Hill, Chapel Hill, NC, USA
[3]Department of Biostatistics and Bioinformatics, Duke University, Durham, NC, USA

*Address for correspondence:* Shu Yang, Department of Statistics, North Carolina State University, North Carolina 27695, USA. Email: syang24@ncsu.edu

## Abstract

We propose a test-based elastic integrative analysis of the randomised trial and real-world data to estimate treatment effect heterogeneity with a vector of known effect modifiers. When the real-world data are not subject to bias, our approach combines the trial and real-world data for efficient estimation. Utilising the trial design, we construct a test to decide whether or not to use real-world data. We characterise the asymptotic distribution of the test-based estimator under local alternatives. We provide a data-adaptive procedure to select the test threshold that promises the smallest mean square error and an elastic confidence interval with a good finite-sample coverage property.

**Keywords:** counterfactual outcome, least favourable confidence interval, non-regularity, precision medicine, pre-test estimator, semiparametric efficiency

## 1 Introduction

Precision medicine (Hamburg & Collins, 2010), which aims at customising medical treatments to individual patient characteristics, has recently received lots of attention. A critical step toward precision medicine is to characterise the heterogeneity of treatment effect (HTE; Rothwell, 2005; Rothwell et al., 2005) entailing how patient characteristics are related to treatment effect. Randomised trials (RTs) are the gold-standard method for treatment effect evaluation because randomisation of treatment ensures that treatment groups are comparable and biases are minimised to the extent possible. However, due to high costs and eligibility criteria for recruiting patients, the trial sample is often small and limited in the patient diversity, which renders the trial underpowered to estimate the HTE and unable to estimate the HTE for specific patient characteristics. On the other hand, extensive real-world (RW) data are increasingly available for research purposes, such as electronic health records, claims databases, and disease registries, with much larger sample sizes and broader demographic and diversity than RT cohorts. Several national organisations (Norris et al., 2010) and regulatory agencies (Sherman et al., 2016) have recently advocated using RW data to have a faster and less costly drug discovery process. Indeed, big data provide unprecedented opportunities for new scientific discovery; however, they also present challenges with possible incomparability with RT data due to selection bias, unmeasured confounding, lack of concurrency, data quality, outcome validity, etc. (US Food and Drug Administration, 2019).

The motivating application is to evaluate adjuvant chemotherapy for resected non-small cell lung cancer (NSCLC) at early-stage disease. Adjuvant chemotherapy for resected NSCLC was shown to be effective in late-stage II and IIIA disease based on RTs (Le Chevalier, 2003).

However, the benefit of adjuvant chemotherapy in stage IB NSCLC disease is unclear. Cancer and Leukemia Group B (CALGB) 9633 is the only RT designed specifically for stage IB NSCLC (Strauss et al., 2008); however, it comprises about 300 patients, which was undersized to detect clinically meaningful improvements for adjuvant chemotherapy (Katz & Saad, 2009). '*Who can benefit from adjuvant chemotherapy with stage IB NSCLC?*' remains an important clinical question. An exploratory analysis of CALGB 9633 showed that patients with tumour size $\geq 4.0$ cm might benefit from adjuvant chemotherapy (Strauss et al., 2008). On the other hand, the National Cancer Database (NCDB) is a clinical oncology registry database that captures the information from approximately 75% of all newly diagnosed cancer patients in the USA. Our goal is to integrate the CALGB 9633 trial with a cohort selected under the same trial eligibility criteria from the NCDB. We expect that an integrated analysis of the CALGB 9633 and NCDB data can considerably improve the efficiency of the HTE estimation on adjuvant chemotherapy regarding tumour size over the RT-only analysis. Although such population-based disease registries provide rich information citing the real-world usage of adjuvant chemotherapy, the concern is the potential bias associated with RW data.

Many authors have proposed methods for generalising treatment effects from RTs to the target population, whose covariate distribution can be characterised by the RW data (Buchanan et al., 2018; Colnet et al., 2020; Lee, Yang, Dong, et al., 2022; Lee, Yang, & Wang, 2022; Zhao et al., 2019). When both RT and RW data provide covariate, treatment, and outcome information, there are two main approaches for integrative analysis: meta-analyses of summary statistics (e.g., Verde & Ohmann, 2015) and pooled patient data (Sobel et al., 2017). The major drawback of meta-analyses of the first kind is that they use only aggregated information and do not distinguish the roles of the RT and RW data, both having unique strengths and weaknesses. Meta-analyses of the second kind include all patients, but pooling the data from two sources breaks the randomisation of treatments and relies on causal inference methods to adjust for confounding bias (e.g., Prentice et al., 2005). More importantly, one cannot rule out possible unmeasured confounding in the RW data. In addition, most existing integrative methods focused on average treatment effects (ATEs) but not on HTEs, which lies at the heart of precision medicine.

To acknowledge the advantages of the RT and RW data, we propose an elastic algorithm for combining the RT and RW data for accurate and robust estimation of the HTE function with a vector of known effect modifiers. The primary identification assumptions underpinning our method are (i) the transportability of the HTE from the RT data to the target population and (ii) the strong ignorability of treatment assignment in the RT data. Transportability is a common assumption in the trial generalisability literature, which holds if the HTE function captures all the treatment effect modifiers, or the study sample is a random sample from the target population. The well-controlled trial design can also ensure the strong ignorability of treatment assignment. If the RW sample satisfies the parallel assumptions (i) and (ii), it is comparable to the RT sample in estimating the HTE. In this case, integrating the RW sample would increase the efficiency of HTE estimation. Toward this end, we use the semiparametric efficiency theory (Bickel et al., 1993; Robins, 1994) to derive the semiparametrically efficient integrative estimator of the HTE. However, due to many practical limitations, the RW sample may violate the desirable comparability assumption (i) or (ii). In this case, integrating the RW sample would lead to bias in HTE estimation. Utilising the design advantage of RTs, we derive a preliminary test statistic to gauge the comparability and reliability of the RW data and decide whether or not to use the RW data in an integrative analysis. Therefore, our test-based elastic integrative estimator uses the efficient combination strategy for estimation if the violation test is insignificant and retains only the RT data if the violation test is significant.

The proposed estimator belongs to pre-test estimation by construction (Giles & Giles, 1993) and is non-regular. We consider null, local, and fixed alternative hypotheses for the pre-testing, representing three scenarios when the comparability assumption required for the RW data is zero, weakly, and strongly violated, respectively. Notably, the fixed alternative formulates the bias of the RW score of the HTE parameter to be fixed, under which the pre-test statistic goes to infinity as the sample size increases. Thus, the inference under the fixed alternative cannot capture the finite-sample behaviour of the test and estimator well and lacks uniform validity. A common strategy to obtain uniform inference validity for non-regular estimators is considering the local alternative, which formulates the bias of the RW score to be in the $n^{-1/2}$ neighbourhood

of zero. The inference under the local alternative provides a better approximation of the finite-sample behaviour of the proposed estimator. Such strategies have been considered in designing trials for sample size/power calculation and in the weak instrument, partial identification, and classification literature (Cheng, 2008; Laber & Murphy, 2011; Staiger & Stock, 1997). Under the local alternative, when the testing distribution is non-degenerate, exact inference for pre-test estimation is complex because the estimator depends on the randomness of the test procedure. This issue cannot be solved by splitting the sample into two parts for testing and estimation separately (Toyoda & Wallace, 1979). The reason is that sample splitting cannot bypass the issue of the additional randomness due to pre-testing, and therefore the impact of pre-testing remains. Also, our test statistic and estimator are constructed based on the whole sample data. To consider the effect of pre-testing, we decompose the test-based elastic integrative estimator into orthogonal components; one is affected by the pre-testing, and the other is not. This step reveals the asymptotic distributions of the proposed estimator to be mixture distributions involving a truncated normal component with ellipsoid truncation and a normal component. Under this framework, we provide a data-adaptive procedure to select the threshold of the test statistic that promises the smallest mean square error (MSE) of the proposed estimator. Lastly, we propose an elastic procedure to construct confidence intervals (CIs), which are adaptive to the local and fixed alternative and have good finite-sample coverage properties.

This article is organised as follows. Section 2 introduces the basic setup, HTE, identification assumptions, and semiparametric efficient estimation. Section 3 establishes a test statistic for gauging the comparability of the RW data with the RT data, a test-based elastic integrative estimator, the asymptotic properties, and an elastic inference procedure. Section 4 presents a simulation study to evaluate the performance of the proposed estimator in terms of robustness and efficiency. Section 5 applies the proposed method to combined CALGB 9633 (RT) and NCDB (RW) data to characterise the HTE of adjuvant chemotherapy in patients with stage IB non-small cell lung cancer. We relegate technical details and all proofs to the Online supplementary material.

## 2 Basic set-up

### 2.1 Notation, the HTE, and two data sources

Let $A \in \{0, 1\}$ be the binary treatment, $Z$ a vector of pre-treatment covariates of interest with the first component being 1, $X$ a vector of auxiliary variables including $Z$, and $Y$ the outcome of interest. We consider $Y$ to be continuous or binary to fix ideas, although our framework can be extended to general-type outcomes, including the survival outcome. To define causal effects, we follow the potential outcomes framework (Neyman, 1923; Rubin, 1974). Under the Stable Unit of Treatment Value assumption, let $Y(a)$ be the potential outcome had the subject been given treatment $a$, for $a = 0, 1$. And, by the causal consistency assumption, the observed outcome is $Y = Y(A) = AY(1) + (1 - A)Y(0)$.

Based on the potential outcomes, the individual treatment effect is $Y(1) - Y(0)$, and $\tau(Z) = \mathbb{E}\{Y(1) - Y(0) \mid Z\}$ characterises the HTE. For a binary outcome, $\tau(Z)$ is also called the causal risk difference. In clinical settings, the parametric family of HTEs is desirable and has wide applications in precision medicine to discover optimal treatment regimes tailored to individual characteristics (Chakraborty & Moodie, 2013). We assume the HTE function to be

$$\tau(Z) = \tau_{\psi_0}(Z) = \mathbb{E}\{Y(1) - Y(0) \mid Z; \psi_0\}, \tag{1}$$

where $\psi_0 \in \mathbb{R}^p$ is a vector of unknown parameters and $p$ is fixed.

We illustrate the HTE function in the following examples.

**Example 1** (Shi et al., 2016; Tian et al., 2014). For a continuous outcome, a linear HTE function is $\tau_{\psi_0}(Z) = Z^{\mathsf{T}}\psi_0$, where each component of $\psi_0$ quantifies how the treatment effect varies over each $Z$.

**Example 2** (Richardson et al., 2017; Tian et al., 2014). For a binary outcome, an HTE function for the causal risk difference is $\tau_{\psi_0}(Z) = \{\exp(Z^{\mathsf{T}}\psi_0) - 1\}/\{\exp(Z^{\mathsf{T}}\psi_0) + 1\}$, ranging from −1 to 1.

To evaluate the effect of adjuvant chemotherapy, let $Y$ be the indication of cancer recurrence within one year of surgery. Consider the HTE function in Example 2 with $Z = (1, \text{age, tumour size})^T$ and $\psi_0 = (\psi_{0,0}, \psi_{0,1}, \psi_{0,2})^T$. This model entails that, on average, the treatment would increase or decrease the risk of cancer recurrence by $|\tau_{\psi_0}(Z)|$ had the patient received adjuvant chemotherapy, and the magnitude of increase depends on age and tumour size. If $Z^T\psi_0 < 0$, it indicates that the treatment is beneficial for this patient. Moreover, if $\psi_{0,1} < 0$ and $\psi_{0,2} < 0$, then older patients with larger tumour sizes would benefit more from adjuvant chemotherapy.

We consider two independent data sources: one from the RT study and the other from the RW study. Let $\delta = 1$ denote RT participation, and let $\delta = 0$ denote RW study participation. Let $V$ summarise the entire record of observed variables $(A, X, \delta, Y)$. The RT data consist of $\{V_i : i \in \mathcal{A}\}$ with sample size $m$, and the RW data consist of $\{V_i : i \in \mathcal{B}\}$ with sample size $n$, where $\mathcal{A}$ and $\mathcal{B}$ are sample index sets for the two data sources. Our setup requires the RT and RW samples to contain $Z$'s information but may include different sets of auxiliary information in $X$. The total sample size is $N = m + n$. Generally, $n$ is larger than $m$. In our asymptotic framework, we assume both $m$ and $n$ go to infinity, and $m/n \to \rho$, where $0 < \rho < 1$.

For simplicity of exposition, we use the following notations throughout the paper: $\mathbb{P}_N$ denotes the empirical measure over the combined RT and RW data, $M^{\otimes 2}$ denotes $MM^T$ for a vector or matrix $M$, $\mathbb{E}_a(\cdot)$ and $\mathbb{V}_a(\cdot)$ are the asymptotic expectation and variance of a random variable, $A_n \perp\!\!\!\perp B_n$ denotes $A_n$ is independent of $B_n$, $A_n \sim B_n$ denotes that $A_n$ follows the same distribution as $B_n$, and $A_n \overset{\cdot}{\sim} B_n$ denotes that $A_n$ and $B_n$ have the same asymptotic distribution as $n \to \infty$. Let $e_\delta(X) = \mathbb{P}(A = 1 \mid X, \delta)$ be the propensity score.

## 2.2 Identification of the HTE from the RT and RW data

The fundamental problem of causal inference is that $Y(0)$ and $Y(1)$ are not jointly observable. Therefore, the HTE is not identifiable without additional assumptions.

We view the RT sample as the gold standard for HTE estimation, satisfying the following assumption.

> **Assumption 1** (RT validity). (i) $\mathbb{E}\{Y(1) - Y(0) \mid X, \delta = 1\} = \tau(Z)$, and (ii) $Y(a) \perp\!\!\!\perp A \mid (X, \delta = 1)$ for $a \in \{0, 1\}$ and $0 < e_1(X) < 1$ for all $(X, \delta = 1)$.

Assumption 1(i) states that the HTE function is transportable from the RT sample to the target population. This assumption is a common assumption in the data integration literature. Stronger versions of Assumption 1(i) have also been considered in the literature, including the ignorability of study participation, i.e., $\{Y(0), Y(1)\} \perp\!\!\!\perp \delta \mid X$ (Buchanan et al., 2018; Stuart et al., 2011), or the mean exchangeability, i.e., $\mathbb{E}\{Y(a) \mid X, \delta\} = \mathbb{E}\{Y(a) \mid X\}$ for $a = 0, 1$ (Dahabreh et al., 2019). Assumption 1(i) holds if $Z$ captures the heterogeneity of effect modifiers or if the study sample is a random sample from the target population. Under the structural equation model framework, Pearl and Bareinboim (2011) provided graphical conditions for transportability. The graphical representation can aid the investigator in assessing the plausibility of Assumption 1(i). Assumption 1(ii) entails that treatment assignment in the RT study follows a randomisation mechanism based on the pre-treatment variables $X$, and all subjects have positive probabilities of receiving each treatment. Assumption 1(ii) holds by the design of complete randomisation of treatment, where the treatment is independent of the potential outcomes and covariates, i.e., $\{Y(a), X\} \perp\!\!\!\perp A \mid \delta = 1$. It also holds by the design of stratified block randomisation of treatment based on discrete $X$, where the treatment is independent of the potential outcomes within each stratum of $X$. The propensity score $e_1(X)$ is known by design.

We consider a parallel assumption for the RW sample, termed RW comparability.

> **Assumption 2** (RW comparability). (i) $\mathbb{E}\{Y(1) - Y(0) \mid X, \delta = 0\} = \tau(Z)$, and (ii) $Y(a) \perp\!\!\!\perp A \mid (X, \delta = 0)$ for $a \in \{0, 1\}$ and $0 < e_0(X) < 1$ for all $(X, \delta = 0)$.

Although Assumption 2 appears similar to Assumption 1, its implications differ substantively. Assumption 2(i) states that the HTE function is transportable from the RW sample to the target population. To make this assumption more plausible, one can use the same trial eligibility criteria

to select the RW sample to ensure a sufficient overlap of the RW covariate space with the RT sample. However, this assumption can be violated in various ways. For example, RT and RW studies may be conducted in different care settings (large academic medical centres versus smaller community hospitals), contexts (geography, policy-related, or socio-structural factors), or time frames. Each of these concerns can violate Assumption 2(i). In addition, due to the lack of control of treatment assignment in RW data, Assumption 2(ii) implies that the observed covariates $X$ capture all the confounding variables related to the treatment and outcome. This assumption may also be restrictive in practice. For example, in the NCDB cohort, the physicians or patients decided, based on experiences or preferences, whether patients received adjuvant chemotherapy after tumour resection. While the database captures many site-level and patient-level information, there may be unmeasured confounding variables that associate with the treatment selection and clinical outcome, e.g., financial status and accessibility to health care facilities.

By trial design, we assume Assumption 1 for the RT data holds throughout the paper; however, we regard Assumption 2 for the RW data as an idealistic assumption, which may be violated. If Assumption 2 holds, we will use a semiparametric efficient strategy to combine both data sources for optimal estimation. However, if Assumption 2 is violated, our proposed method will automatically detect the violation and retain only the RT data for estimation. In practice, it is important to identify a 'similar' RW sample to be integrated with the RT sample. Hernán and Robins (2016) provided a framework for using big real-world data to emulate a target trial when a randomised trial is unavailable. When selecting an RW sample, we can check the rubrics for the eligibility criteria that defines the target population, treatment definitions, assignment procedures, follow-up time, outcome, and effect contrast of interest, to increase the chance of successfully integrating the RW sample with the RT sample.

Unlike our focus on testing the comparability of the RW in HTE estimation, testing transportability alone may be of more importance in some contexts. Under Assumptions 1(ii) and 2(ii), i.e., the treatment ignorability holds, possible tests can be adopted to test $\mathbb{E}\{Y(1) - Y(0) \mid X, \delta = 1\} = \mathbb{E}\{Y(1) - Y(0) \mid X, \delta = 0\}$, e.g., the U-statistics-based test (Luedtke et al., 2019).

Under Assumptions 1 and 2, the following identification formula holds for the HTE:

$$\mathbb{E}\left\{\frac{AY}{e_\delta(X)} - \frac{(1-A)Y}{1 - e_\delta(X)}\middle| Z, \delta\right\} = \tau(Z). \tag{2}$$

The identification formula motivates regression analysis based on the modified outcome $A\{e_\delta(X)\}^{-1}Y - (1-A)\{1 - e_\delta(X)\}^{-1}Y$ to estimate the HTE. This approach involves the inverse of the treatment probability, and thus the resulting estimator may be unstable if some estimated treatment probabilities are close to zero or one. It calls for a principled way to construct improved estimators of the HTE. Rudolph and van der Laan (2017) derived the semiparametric efficiency score (SES) and bound for the average treatment effect. In the next subsection, we derive the SES of the HTE under Assumptions 1 and 2 that motivates improved estimators.

## 2.3 Semiparametric efficiency score

The semiparametric model consists of model (1) with the parameter of interest $\psi_0$ and the unspecified distribution. Assumptions 1 and 2 impose restrictions on $\psi_0$. To see this, define

$$H_\psi = Y - \tau_\psi(Z)A. \tag{3}$$

Intuitively, $H_{\psi_0}$ subtracts from the subject's observed outcome $Y$ the treatment effect of the subject's observed treatment $\tau_{\psi_0}(Z)A$, which mimics the potential outcome $Y(0)$. Formally, following Robins (1994), we can show that $\mathbb{E}(H_{\psi_0} \mid A, X, \delta) = \mathbb{E}\{Y(0) \mid A, X, \delta\}$. Therefore, by Assumptions 1 and 2, $\psi_0$ must satisfy the restriction:

$$\mathbb{E}(H_{\psi_0} \mid A, X, \delta) = \mathbb{E}(H_{\psi_0} \mid X, \delta). \tag{4}$$

For simplicity of exposition, denote

$$\mathbb{E}(H_{\psi_0} \mid X, \delta) = \mu_\delta(X), \quad \mathbb{V}(H_{\psi_0} \mid X, \delta) = \sigma_\delta^2(X),$$

where $\mu_\delta(X)$ is the outcome mean function and $\sigma_\delta^2(X)$ is the outcome variance function. By viewing $(X, \delta)$ jointly as the set of confounders, we invoke the SES of the structural nested mean model in Robins (1994). We further make a simplifying assumption that

$$\mathbb{E}(H_{\psi_0}^2 \mid A, X, \delta) = \mathbb{E}(H_{\psi_0}^2 \mid X, \delta), \tag{5}$$

which is a natural extension of (4). This assumption allows us to derive the SES of $\psi_0$ as

$$S_{\psi_0}(V) = q^*(X, \delta)\{H_{\psi_0} - \mu_\delta(X)\}\{A - e_\delta(X)\}, \quad q^*(X, \delta) = \left\{\partial\tau_{\psi_0}(Z)/\partial\psi\right\}\left\{\sigma_\delta^2(X)\right\}^{-1}, \tag{6}$$

which separates the term with the outcome, i.e., $H_{\psi_0} - \mu_\delta(X)$, and the term with the treatment, i.e., $A - e_\delta(X)$. This feature relaxes model assumptions of the nuisance functions while retaining root-$n$ consistency in the estimation of $\psi_0$; see Section 2.4. Even without the simplifying assumption in (5), by the mean independence property in (4), we can verify that

$$\mathbb{E}\{S_{\psi_0}(V)\} = \mathbb{E}[q^*(X, \delta)\mathbb{E}\{H_{\psi_0} - \mu_\delta(X) \mid X, \delta\} \times \mathbb{E}\{A - e_\delta(X) \mid X, \delta\}] = 0.$$

Therefore, if (5) holds, $S_{\psi_0}(V)$ is the SES of $\psi_0$; if (5) does not hold, $S_{\psi_0}(V)$ is unbiased and permits robust estimation. We provide examples to elucidate the SES below before delving into robust estimation in the following subsection.

**Example 3** For a continuous outcome and the HTE function given in Example 1, the SES of $\psi_0$ is

$$S_{\psi_0}(V) = Z\left\{\sigma_\delta^2(X)\right\}^{-1}\{H_{\psi_0} - \mu_\delta(X)\}\{A - e_\delta(X)\}.$$

For a binary outcome and the HTE function given in Example 2, the SES of $\psi_0$ is

$$S_{\psi_0}(V) = Z\frac{2\exp(Z^{\mathrm{T}}\psi_0)}{\{\exp(Z^{\mathrm{T}}\psi_0) + 1\}^2}[\mu_\delta(X)\{1 - \mu_\delta(X)\}]^{-1}\{H_{\psi_0} - \mu_\delta(X)\}\{A - e_\delta(X)\}.$$

**Remark 1** (Comparison with other doubly robust approaches). The identification formula (2) motivates the inverse probability weighted (IPW)-adjusted regression. However, IPW is known to be inefficient and sensitive to model misspecification of the propensity score. Alternatively, Kennedy (2020) proposed a pseudo-outcome regression approach using augmented IPW (AIPW) pseudo-outcomes that leverages weighting and outcome mean functions and improves the performance of IPW-adjusted regression. The doubly robust loss function for the treatment contrast or blip function in Luedtke and van der Laan (2016) also exploits weighting and outcome mean functions. Both IPW and AIPW use weighting to remove confounding biases; differently, the SES in (6) uses the mean independence of $H_{\psi_0} - \mu_\delta(X)$ and $A - e_\delta(X)$ to construct unbiased estimating equations. The simulation study in Online Supplementary Material, Section S4.1 shows that the SES approach outperforms the AIPW-adjusted approach when the propensity score can be close to zero or one.

## 2.4 From SES to robust estimation

In principle, an efficient estimator for $\psi_0$ can be obtained by solving $\mathbb{P}_N S_{\text{eff},\psi}(V) = 0$. However, $S_{\text{eff},\psi}$ depends on the unknown distribution through $e_0(X)$, $\mu_\delta(X)$, and $\sigma_\delta^2(X)$, and thus solving $\mathbb{P}_N S_{\text{eff},\psi}(V) = 0$ is infeasible. Nevertheless, the state-of-art causal inference literature suggests that estimators constructed based on SES are robust to approximation errors using machine learning methods, the so-called rate double robustness; see, e.g., Chernozhukov et al. (2018) and Rotnitzky et al. (2019).

In order to obtain a robust estimator with good efficiency properties, we consider approximating the unknown functions using non-parametric or machine learning methods. In summary, our algorithm for the estimation of $\psi_0$ proceeds as follows.

Step 1. Obtain an estimator of $e_0(X)$ using non-parametric or machine learning methods, denoted by $\widehat{e}_0(X)$, based on $\{(A_i, X_i, \delta_i = 0) : i \in \mathcal{B}\}$.

Step 2. Obtain a preliminary estimator $\widehat{\psi}_{\text{p}}$ by solving $\sum_{i \in \mathcal{A}} [q^*(X_i, \delta_i)\{A_i - e_1(X_i)\}H_{\psi,i}] = 0$, based on $\{(A_i, X_i, Y_i, \delta_i = 1) : i \in \mathcal{A}\}$.

Step 3. Obtain the estimators of $\mu_1(X)$ and $\mu_0(X)$ using non-parametric or machine learning methods, denoted by $\widehat{\mu}_1(X)$ and $\widehat{\mu}_0(X)$, based on $\{(A_i, X_i, H_{\widehat{\psi}_{\text{p}},i}, \delta_i = 1) : i \in \mathcal{A}\}$ and $\{(A_i, X_i, H_{\widehat{\psi}_{\text{p}},i}, \delta_i = 0) : i \in \mathcal{B}\}$, respectively.

Step 4. Let $\widehat{S}_{\text{eff},\psi}(V)$ be $S_{\text{eff},\psi}(V)$ with the unknown quantities replaced by the estimated parametric models in Steps 1 and 3. Obtain the efficient integrative estimator $\widehat{\psi}_{\text{eff}}$ by solving

$$\mathbb{P}_N \widehat{S}_\psi(V) = 0. \tag{7}$$

The estimator $\widehat{\psi}_{\text{eff}}$ depends on the approximation of nuisance functions. To establish the asymptotic properties of $\widehat{\psi}_{\text{eff}}$, we provide the regularity conditions.

**Assumption 3** (i) $\|\widehat{e}_0(X) - e_0(X)\| = o_{\mathbb{P}}(1)$ and $\|\widehat{\mu}_\delta(X) - \mu_\delta(X)\| = o_{\mathbb{P}}(1)$; (ii) $\|\widehat{e}_0(X) - e_0(X)\| \times \|\widehat{\mu}_\delta(X) - \mu_\delta(X)\| = o_{\mathbb{P}}(n^{-1/2})$; and (iii) additional regularity conditions in Online Supplementary Material, Assumption S1.

Assumption 3 is typical regularity conditions for Z-estimation or M-estimation (van der Vaart, 2000). Assumption 3(i) states that we require the posited models to be consistent for the two nuisance functions. Assumption 3(ii) states that the *combined rate* of convergence of the posited models is $o_{\mathbb{P}}(n^{-1/2})$. Online Supplementary Material, Assumption S1 regularises the complexity of the functional space. Importantly, these conditions ensure $\widehat{\psi}_{\text{eff}}$ retains the parametric-rate consistency, allowing flexible data-adaptive models and not restricting to stringent parametric models.

**Theorem 1** Suppose Assumptions 1–3 hold. Then, $\widehat{\psi}_{\text{eff}}$ is root-$n$ consistent for $\psi_0$ and asymptotically normal.

Theorem 1 implies that asymptotically, $\widehat{\psi}_{\text{eff}}$ can be viewed as the solution to $\mathbb{P}_N S_\psi(V) = 0$ when the nuisance functions are known. Therefore, for consistent variance estimation of $\widehat{\psi}_{\text{eff}}$, we can use the standard sandwich formula (Stefanski & Boos, 2002) or the perturbation-based resampling (Hu & Kalbfleisch, 2000), treating the nuisance functions to be known.

## 3 Test-based elastic integrative analysis

A major concern for integrating the RT and RW data lies in the possibly poor quality of the RW data. Then, combining the RT and RW data into an integrative analysis would lead to a biased HTE estimator. This section addresses the critical challenge of preventing any biases present in the RW data from leaking into the proposed estimator.

## 3.1 Detection of the RW incompatibility

We consider all assumptions in Theorem 1 hold except that Assumption 2 may be violated. We derive a test that detects the violation of this crucial assumption for using the RW data. For simplicity, we denote the SES based solely on the RT or RW data as

$$S_{\mathrm{rt},\psi}(V) = \delta S_\psi(V), \quad S_{\mathrm{rw},\psi}(V) = (1-\delta)S_\psi(V),$$

respectively. Moreover, let $\widehat{S}_{\mathrm{rt},\psi}(V)$ and $\widehat{S}_{\mathrm{rw},\psi}(V)$ be $S_{\mathrm{rt},\psi}(V)$ and $S_{\mathrm{rw},\psi}(V)$ with the nuisance functions replaced by their estimates, and let $\mathcal{I}_{\mathrm{rt}} = \mathbb{E}\{S_{\mathrm{rt},\psi_0}(V)^{\otimes 2} \mid \delta = 1\}$ and $\mathcal{I}_{\mathrm{rw}} = \mathbb{E}\{S_{\mathrm{rw},\psi_0}(V)^{\otimes 2} \mid \delta = 0\}$ be Fisher information matrices.

We now formulate the null hypothesis $H_0$ for the case when Assumption 2 holds and fixed and local alternatives $H_a$ and $H_{a,n}$ for the case when Assumption 2 is violated:

$H_0$ (Null) $\mathbb{E}\{S_{\mathrm{rw},\psi_0}(V)\} = 0$.
$H_a$ (Fixed alternative) $\mathbb{E}\{S_{\mathrm{rw},\psi_0}(V)\} = \eta_{\mathrm{fix}}$ , where $\eta_{\mathrm{fix}}$ is a $p$-vector of constants with at least one non-zero component.
$H_{a,n}$ (Local alternative) $\mathbb{E}\{S_{\mathrm{rw},\psi_0}(V)\} = n^{-1/2}\eta$ , where $\eta$ is a $p$-vector of constants with at least one non-zero component.

Considering the fixed alternative is common to establish asymptotic properties of standard estimators and tests; however, the local alternative is useful to study finite-sample properties and regularity of non-standard estimators and tests. In finite samples, the violation of Assumption 2 may be weak; e.g., there exists a hidden confounder in the RW data, but the association between the hidden confounder and the outcome or the treatment is small. In such cases, the test statistic can be small or moderate. The fixed alternative formulates the bias of the RW score to be fixed, implying that the test statistic goes to infinity with the sample size. Consequently, the fixed alternative inference cannot capture the finite-sample behaviour well in the cases of weak violation and does not have uniform validity. That is, there exist scenarios where the finite-sample coverage probability from standard inference is far from the nominal level for any sample size. The local alternative asymptotics is a common approach to obtaining uniform inference validity for non-regular estimators. In the local alternative $H_{a,n}$, the bias of $S_{\mathrm{rw},\psi_0}(V)$ may be small as quantified by $n^{-1/2}\eta$. The values of $\eta$ represent different tracks that the bias of $S_{\mathrm{rw},\psi_0}(V)$ follows to converge to zero. We will show that the test statistic is $O_\mathbb{P}(1)$, thus better capturing the finite-sample behaviour in the weak violation cases. The local alternative encompasses the null and fixed alternative as special cases by considering different values of $\eta$. In particular, $H_0$ corresponds to $H_{a,n}$ with $\eta = 0$. Also, $H_a$ corresponds to $H_{a,n}$ with $\eta = \pm\infty$; hence, considering $H_a$ alone is not informative about the finite-sample behaviours of the proposed test and estimator.

We detect biases in the RW data based on the following two key insights. First, we obtain an initial estimator $\widehat{\psi}_{\mathrm{rt}}$ by solving the estimating equation based solely on the RT data, $\sum_{i \in \mathcal{A}} \widehat{S}_{\mathrm{rt},\psi}(V_i) = 0$. It is important to emphasise that the propensity score in the RT $e_1(X)$ is known by design and, therefore, $\widehat{\psi}_{\mathrm{rt}}$ is always consistent. Second, if Assumption 2 holds for the RW data, $S_{\mathrm{rw},\psi_0}(V)$ is unbiased, but $S_{\mathrm{rw},\psi_0}(V)$ is no longer unbiased if it is violated. Therefore, large values of $n^{-1/2} \sum_{i \in \mathcal{B}} \widehat{S}_{\mathrm{rw},\hat{\psi}_{\mathrm{rt}}}(V_i)$ provide evidence of the violation of Assumption 2.

To detect the violation of Assumption 2 for using the RW data, we construct the test statistic

$$T = \left\{ n^{-1/2} \sum_{i \in \mathcal{B}} \widehat{S}_{\mathrm{rw},\hat{\psi}_{\mathrm{rt}}}(V_i) \right\}^{\mathrm{T}} \widehat{\Sigma}_{SS}^{-1} \left\{ n^{-1/2} \sum_{i \in \mathcal{B}} \widehat{S}_{\mathrm{rw},\hat{\psi}_{\mathrm{rt}}}(V_i) \right\}, \tag{8}$$

where $\Sigma_{SS} = \Gamma^{\mathrm{T}}\mathcal{I}_{\mathrm{rt}}\Gamma + \mathcal{I}_{\mathrm{rw}}$ is the asymptotic variance of $n^{-1/2} \sum_{i \in \mathcal{B}} \widehat{S}_{\mathrm{rw},\hat{\psi}_{\mathrm{rt}}}(V_i)$, $\Gamma = \mathcal{I}_{\mathrm{rt}}^{-1}\mathcal{I}_{\mathrm{rw}}\rho^{-1/2}$, and $\widehat{\Sigma}_{SS}$ is a consistent estimator for $\Sigma_{SS}$. The test statistic $T$ measures the distance between $n^{-1/2} \sum_{i \in \mathcal{B}} S_{\mathrm{rw},\hat{\psi}_{\mathrm{rt}}}(V_i)$ and zero. If the idealistic assumption holds, we expect $T$ to be small. By

the standard asymptotic theory, we show in the Online supplementary material that under $H_0$, $T \dot\sim \chi_p^2$, a Chi-square distribution with degrees of freedom $p$, as $n \to \infty$. This result serves to detect the violation of the assumption required for the RW data.

## 3.2 Elastic integration

Let $c_\gamma = \chi_{p,\gamma}^2$ be the $100(1 - \gamma)$th percentile of $\chi_p^2$. For a small $\gamma$, if $T \geq c_\gamma$, there is strong evidence to reject $H_0$ for the RW data; i.e., there is a detectable bias for the RW data estimator. In this case, we would only use the RT data for estimation. On the other hand, if $T < c_\gamma$, there is no strong evidence that the RW data estimator is biased; therefore, we would combine both the RT and RW data for optimal estimation. Our strategy leads to the elastic integrative estimator $\widehat\psi_{\text{elas}}$ solving

$$\sum_{i \in \mathcal{A} \cup \mathcal{B}} \left\{ \delta_i \widehat S_\psi(V_i) + \mathbf{1}(T < c_\gamma)(1 - \delta_i) \widehat S_\psi(V_i) \right\} = 0. \tag{9}$$

The choice of $\gamma$ involves the bias-variance trade-off. On the one hand, under $H_0$, the acceptance probability of integrating the RW data is $\mathbb{P}(T < c_\gamma) = 1 - \gamma$. Therefore, for a relatively large sample size, we will accept good-quality RW data with probability $1 - \gamma$ and reject good-quality RW data with type I error $\gamma$. Hence, a small $\gamma$ is desirable; similarly, for $H_{a,n}$ with small $\eta$. On the other hand, under $H_{a,n}$ with large $\eta$, the reverse is true, and hence a large $\gamma$ is desirable.

To formally investigate the trade-off, we characterise the asymptotic distributions of the elastic integrative estimator $\widehat\psi_{\text{elas}}$ under the null, fixed, and local alternatives. We do not discuss the trivial cases when $\gamma = 0$ and $1$, corresponding to $\widehat\psi_{\text{elas}} = \widehat\psi_{\text{rt}}$ or $\widehat\psi_{\text{eff}}$. With $\gamma \in (0, 1)$, $\widehat\psi_{\text{elas}}$ mixes two distributions, namely, $\widehat\psi_{\text{rt}} \mid (T \geq c_\gamma)$ and $\widehat\psi_{\text{eff}} \mid (T < c_\gamma)$. Each distribution can be non-standard because the estimators and test are constructed based on the same data and, therefore, may be asymptotically dependent.

To characterise those non-standard distributions, we decompose this task into three steps. First, by the standard asymptotic theory, it follows that $T \dot\sim \mathcal{Z}_1^{\mathsf{T}} \mathcal{Z}_1$, where $\mathcal{Z}_1$ is a standard $p$-variate normal random vector, $n^{1/2}(\widehat\psi_{\text{rt}} - \psi_0) \dot\sim \mathcal{N}_{\text{rt}}$, and $n^{1/2}(\widehat\psi_{\text{eff}} - \psi_0) \dot\sim \mathcal{N}_{\text{eff}}$, where $\mathcal{N}_{\text{rt}}$ and $\mathcal{N}_{\text{eff}}$ are some $p$-variate normal random vectors with variances $V_{\text{rt}} = (\rho \mathcal{I}_{\text{rt}})^{-1}$ and $V_{\text{eff}} = (\rho \mathcal{I}_{\text{rt}} + \mathcal{I}_{\text{rw}})^{-1}$, respectively.

Second, we find another standard $p$-variate normal random vector $\mathcal{Z}_2$ that is independent of $\mathcal{Z}_1$, and decompose the normal distributions $\mathcal{N}_{\text{rt}}$ and $\mathcal{N}_{\text{eff}}$ into two orthogonal components: i) one corresponds to $\mathcal{Z}_1$ and ii) the other one corresponds to $\mathcal{Z}_2$. Importantly, component i) would be affected by the test constraints induced by $\mathcal{Z}_1^{\mathsf{T}} \mathcal{Z}_1$, but component ii) would not be affected. For $\mathcal{N}_{\text{eff}}$, we show that it is fully represented by $\mathcal{Z}_2$ as $\mathcal{N}_{\text{eff}} = -V_{\text{eff}}^{1/2} \mathcal{Z}_2$. Therefore, its distribution is not affected by $\mathcal{Z}_1^{\mathsf{T}} \mathcal{Z}_1 < c_\gamma$; that is,

$$\mathcal{N}_{\text{eff}} \mid (\mathcal{Z}_1^{\mathsf{T}} \mathcal{Z}_1 < c_\gamma) \sim -V_{\text{eff}}^{1/2} \mathcal{Z}_2.$$
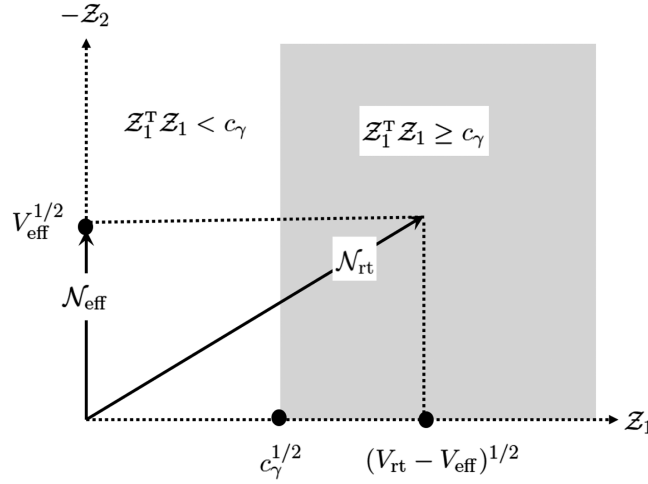
For $\mathcal{N}_{\text{rt}}$, we show that $\mathcal{N}_{\text{rt}} = V_{\text{rt-eff}}^{1/2} \mathcal{Z}_1 - V_{\text{eff}}^{1/2} \mathcal{Z}_2$ with $V_{\text{rt-eff}} = V_{\text{rt}} - V_{\text{eff}}$. Due to the independence between $\mathcal{Z}_1$ and $\mathcal{Z}_2$, $\mathcal{N}_{\text{rt}} \mid (\mathcal{Z}_1^{\mathsf{T}} \mathcal{Z}_1 \geq c_\gamma)$ is a mixture distribution

$$\mathcal{N}_{\text{rt}} \mid (\mathcal{Z}_1^{\mathsf{T}} \mathcal{Z}_1 \geq c_\gamma) \sim V_{\text{rt-eff}}^{1/2} \mathcal{Z}_{c_\gamma}^{\mathsf{t}} - V_{\text{eff}}^{1/2} \mathcal{Z}_2,$$

mixing a non-normal component, where $\mathcal{Z}_c^{\mathsf{t}}$ represents the truncated normal distribution $\mathcal{Z}_1 \mid (\mathcal{Z}_1^{\mathsf{T}} \mathcal{Z}_1 \geq c)$, and a normal component. For illustration, Figure 1 demonstrates the geometry of the decomposition of distributions with scalar variables.

Third, we formally characterise the distribution of $\mathcal{Z}_c^{\mathsf{t}}$, a multivariate normal distribution with ellipsoid truncation (Li et al., 2018; Tallis, 1963). This step enables us to quantify the asymptotic bias and variance of the proposed estimator; see Section 3.3.

Let $F_p(\cdot)$ be the cumulative distribution function (CDF) of a $\chi_p^2$ random variable, and $F_p(\cdot; \lambda)$ be the CDF of a $\chi_p^2(\lambda)$ random variable, where $\chi_p^2$ and $\chi_p^2(\lambda)$ are the central Chi-square distribution and

- $\mathcal{N}_{\mathrm{rt}} = V_{\mathrm{rt\text{-}eff}}^{1/2}\mathcal{Z}_1 - V_{\mathrm{eff}}^{1/2}\mathcal{Z}_2$ and $\mathcal{N}_{\mathrm{rt}} \mid (\mathcal{Z}_1^{\mathrm{T}}\mathcal{Z}_1 \geq c_\gamma) \sim V_{\mathrm{rt\text{-}eff}}^{1/2}\mathcal{Z}_1 \mid (\mathcal{Z}_1^{\mathrm{T}}\mathcal{Z}_1 \geq c_\gamma) - V_{\mathrm{eff}}^{1/2}\mathcal{Z}_2$

- $\mathcal{N}_{\mathrm{eff}} = -V_{\mathrm{eff}}^{1/2}\mathcal{Z}_2$ and $\mathcal{N}_{\mathrm{eff}} \mid (\mathcal{Z}_1^{\mathrm{T}}\mathcal{Z}_1 < c_\gamma) \sim \mathcal{N}_{\mathrm{eff}}$

**Figure 1.** Representation of the normal distributions $\mathcal{N}_{\mathrm{rt}}$ and $\mathcal{N}_{\mathrm{eff}}$ based on $\mathcal{Z}_1$ and $\mathcal{Z}_2$ with $p = 1$.

the non-central Chi-square distribution with the non-centrality parameter $\lambda$, respectively. Theorem 2 summarises the asymptotic distribution of $\widehat{\psi}_{\mathrm{elas}}$.

**Theorem 2** Suppose assumptions in Theorem 1 hold except that Assumption 2 may be violated. Let $\mathcal{Z}_1$ and $\mathcal{Z}_2$ be independent normal random vectors with mean $\mu_1 = \Sigma_{SS}^{-1/2}\eta$ and $\mu_2 = V_{\mathrm{eff}}^{1/2}\eta$, respectively, and covariance $I_{p\times p}$. Let $\mathcal{Z}_c^{\mathrm{t}}$ be the truncated normal distribution $\mathcal{Z}_1 \mid (\mathcal{Z}_1^{\mathrm{T}}\mathcal{Z}_1 \geq c)$. Let the elastic integrative estimator $\widehat{\psi}_{\mathrm{elas}}$ be obtained by solving (9). Then, $n^{1/2}(\widehat{\psi}_{\mathrm{elas}} - \psi_0)$ has a limiting mixture distribution

$$\mathcal{M}(\gamma; \eta) = \begin{cases} \mathcal{M}_1(\gamma; \eta) = V_{\mathrm{rt-eff}}^{1/2}\mathcal{Z}_{c_\gamma}^{\mathrm{t}} - V_{\mathrm{eff}}^{1/2}\mathcal{Z}_2, & \text{w.p. } \xi, \\ \mathcal{M}_2(\eta) = -V_{\mathrm{eff}}^{1/2}\mathcal{Z}_2, & \text{w.p. } 1 - \xi, \end{cases} \tag{10}$$

(a) Under $H_0$, $\mu_1 = \mu_2 = 0$ and $\xi = 1 - F_p(c_\gamma) = \gamma$.
(b) Under $H_a$, $\mu_1 = \mu_2 = \pm\infty$ and $\xi = 1$; i.e., (10) reduces to a normal distribution with mean 0 and variance $V_{\mathrm{rt}}$.
(c) Under $H_{a,n}$, $\mu_1 = \Sigma_{SS}^{-1/2}\eta$, $\mu_2 = V_{\mathrm{eff}}^{1/2}\eta$ with $\eta \in \mathbb{R}^p$, and $\xi = 1 - F_p(c_\gamma; \lambda)$, where $\lambda = \eta^{\mathrm{T}}\Sigma_{SS}^{-1}\eta$.

In Theorem 2, $\mathcal{M}(\gamma; \eta)$ in (10) is a general characterisation of the asymptotic distribution of $n^{1/2}(\widehat{\psi}_{\mathrm{elas}} - \psi_0)$. It implies different asymptotic behaviours of $n^{1/2}(\widehat{\psi}_{\mathrm{elas}} - \psi_0)$ depending on whether Assumption 2 is strongly, weakly, or not violated. First, $H_a$ corresponds to the situation where Assumption 2 is strongly violated. Under $H_a$, $T$ rejects the RW data (i.e., $\mathcal{Z}_1^{\mathrm{T}}\mathcal{Z}_1 \geq c_\gamma$ holds) with probability converging to one, $\mathcal{Z}_{c_\gamma}^{\mathrm{t}}$ becomes $\mathcal{Z}_1$, and $\mathcal{M}(\gamma; \eta = \pm\infty)$ becomes $V_{\mathrm{rt-eff}}^{1/2}\mathcal{Z}_1 - V_{\mathrm{eff}}^{1/2}\mathcal{Z}_2$, a normal distribution with mean 0 and variance $V_{\mathrm{rt}}$. As expected, under $H_a$, $n^{1/2}(\widehat{\psi}_{\mathrm{elas}} - \psi_0)$ is asymptotically normal and regular. Second, $H_0$ and $H_{a,n}$ correspond to the situations when Assumption 2 is not and weakly violated, respectively. Under $H_0$ and $H_{a,n}$, $T$ has positive probabilities of accepting and rejecting the RW data, $\widehat{\psi}_{\mathrm{elas}}$ switches between $\widehat{\psi}_{\mathrm{eff}}$ and $\widehat{\psi}_{\mathrm{rt}}$, and $n^{1/2}(\widehat{\psi}_{\mathrm{elas}} - \psi_0)$ follows a limiting mixing distribution $\mathcal{M}(\gamma; \eta)$, indexed by $\eta$.

Although the exact form of $\mathcal{M}(\gamma; \eta)$ is complicated, the entire distribution and summary statistics such as mean, variance, and quantiles can be simulated by rejective sampling. Importantly, under $H_0$ and $H_{a,n}$, $n^{1/2}(\widehat{\psi}_{\text{elas}} - \psi_0)$ is non-normal and non-regular. The non-regularity is determined by the local parameter $\eta$, which entails that the asymptotic distribution of $n^{1/2}(\widehat{\psi}_{\text{elas}} - \psi_0)$ may change abruptly when $H_0$ is slightly violated. It is worth emphasising that the local asymptotics provides a better approach to demonstrate the finite-sample properties of the test and estimators than the fixed asymptotics does.

### 3.3 Asymptotic bias and MSE

Based on Theorem 2, it is essential to understand the asymptotic behaviours of $\mathcal{Z}_c^{\text{t}}$ and the truncated multivariate normal distribution in general. Toward that end, we derive the moment generating functions (MGFs) of such distributions in the Online supplementary material, which shed light on the moments of $n^{1/2}(\widehat{\psi}_{\text{elas}} - \psi_0)$.

Corollary 1 provides the analytical formula of the asymptotic bias and MSE of $n^{1/2}(\widehat{\psi}_{\text{elas}} - \psi_0)$.

> **Corollary 1**    Suppose assumptions in Theorem 1 hold except that Assumption 2 may be violated.
>
> (a) Under $H_0$, the bias and MSE of $n^{1/2}(\widehat{\psi}_{\text{elas}} - \psi_0)$ are bias = 0 and mse $= V_{\text{eff}} + V_{\text{rt}-\text{eff}}\{1 - F_{p+2}(c_\gamma)\}$.
>
> (b) Under $H_a$, the bias and MSE of $n^{1/2}(\widehat{\psi}_{\text{elas}} - \psi_0)$ are bias = 0 and mse $= V_{\text{rt}}$.
>
> (c) Under $H_{a,n}$, the bias and MSE of $n^{1/2}(\widehat{\psi}_{\text{elas}} - \psi_0)$ are
>
> $$\text{bias}(\gamma, \eta) = -V_{\text{eff}}\eta F_{p+2}(c_\gamma; \lambda), \tag{11}$$
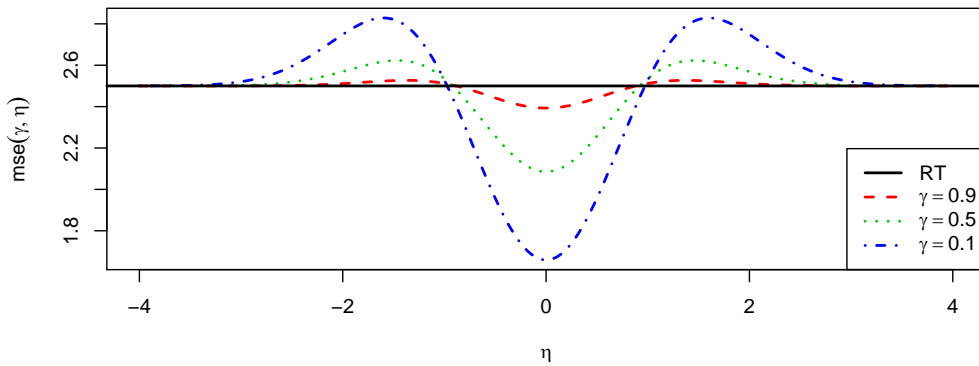>
> and
>
> $$\begin{aligned} \text{mse}(\gamma, \eta) = V_{\text{eff}} &+ V_{\text{rt}-\text{eff}}\{1 - F_{p+2}(c_\gamma; \lambda)\} \\ &+ (V_{\text{eff}}\eta)^{\otimes 2}\{2F_{p+2}(c_\gamma; \lambda) - F_{p+4}(c_\gamma; \lambda)\} \end{aligned} \tag{12}$$
>
> with $\lambda = \eta^{\text{T}}\Sigma_{SS}^{-1}\eta$.

Corollary 1 enables us to demonstrate the potential advantages and disadvantages of $\widehat{\psi}_{\text{elas}}$ compared with $\widehat{\psi}_{\text{rt}}$ and $\widehat{\psi}_{\text{eff}}$ under different scenarios. To illustrate, we consider the case of a scalar $\psi_0$, $V_{\text{eff}} = 1$, $V_{\text{rt}} = 2.5$, and $\Sigma_{SS} = 0.5$. Figure 2 shows mse$(\gamma, \eta)$ as a function of $\eta$ by varying $\gamma \in \{0.9, 0.5, 0.1\}$ compared to $\widehat{\psi}_{\text{rt}}$. For a given $\gamma \in (0, 1)$, when $\eta$ is small, $\widehat{\psi}_{\text{elas}}$ is more efficient than $\widehat{\psi}_{\text{rt}}$; and when $\eta$ increases, the MSE of $\widehat{\psi}_{\text{elas}}$ increases, exceeds, and gradually returns to the MSE of $\widehat{\psi}_{\text{rt}}$. This phenomenon reveals the super-efficiency (related to the problem of non-regularity) of $\widehat{\psi}_{\text{elas}}$ at small values of $\eta$ at the cost of the MSE inflation for some $\eta$ values. LeCam (1953) obtained an earlier result of super-efficiency for the famous Hodges estimator. Also, $\widehat{\psi}_{\text{elas}}$ with a smaller $\gamma$ achieves a larger deduction of the MSE at small values of $\eta$ but also more considerable inflation of the MSE at big values of $\eta$ compared to $\widehat{\psi}_{\text{rt}}$, and vice versa. This observation motivates our adaptive selection of $\gamma$ in Section 3.5 to produce an elastic integrative estimator with small bias and mean squared error for a possible value of $\eta$. Also, super-efficiency and non-regularity are the root causes for the standard asymptotic inference to fail, which motivates the proposed elastic confidence intervals to provide uniformly valid confidence intervals (Section 3.4); however, they can be conservative at certain parameter values when the sample size is small (Section 4).

> **Remark 2**    (Sample splitting and cross fitting). Sample splitting and cross fitting are helpful tactics to simplify asymptotic analyses by removing the dependence between nuisance parameter estimation and primary parameter estimation (Chernozhukov et al., 2018; Kennedy, 2020). To apply sample splitting to our context, one can divide the sample into two parts for testing and estimation

**Figure 2.** Illustration of the super-efficiency of $\widehat{\psi}_{\text{elas}}$ in terms of mse($\gamma, \eta$) as a function of $\eta$ by varying $\gamma \in \{0.9(\text{dashed}), 0.5(\text{dotted}), 0.1(\text{dotdash})\}$ compared to $\widehat{\psi}_{\text{rt}}$.

separately. While sample splitting and cross fitting are beneficial in theoretical development, they may come with expenses of heavier computation and fewer data for estimating different components. Thus, we do not use sampling splitting or cross fitting as a device to establish the theoretical properties of the proposed pre-test estimator. Without sample splitting, the test and estimators are intimately related, requiring careful decompositions of the estimators into components that are asymptotically dependent and independent of the test statistic, as shown in our three steps toward Theorem 2. Also, sample splitting cannot resolve the non-regularity issue of the pre-test estimator (Toyoda & Wallace, 1979). This is because sample splitting cannot bypass additional randomness due to pre-testing. Thus, the impact of pre-testing and superefficiency remains an issue; see the simulation study in Online Supplementary Material, Section S4.6.

**Remark 3** (Soft thresholding to mitigate the non-regularity). The proposed elastic integrative estimator involves an indicator function to make a binary decision to include or exclude the RW data from analysis. The indicator function serves as hard thresholding. To alleviate the non-regularity issue and refine the proposed estimator, one may use soft thresholding by imposing the smoothness of the indicator function. For example, similar to Yang and Ding (2018), one can use a smooth weight function $\Phi_\epsilon(c_\gamma - T)$ to replace $I(T < c_\gamma)$, where $\Phi_\epsilon(z)$ is the normal cumulative distribution with zero mean and variance $\epsilon^2$. As $\epsilon \to 0$, $\Phi_\epsilon(c_\gamma - T)$ becomes closer to $I(T < c_\gamma)$. Also, as suggested by a reviewer, one can weigh the RW data based on the $p$-value from the test, i.e., $1 - F_p(T)$. A small $p$-value indicates a large bias in the RW data, and we should give the RW data less weight. Conversely, a large $p$-value suggests a small bias, and we should provide the RW data with more weight. The third idea is to create bootstrap replications of the elastic integrative estimator and obtain the average of the bootstrap replications to impose smoothness. Chakraborty et al. (2010) showed in simulation that soft-thresholding reduces the non-regularity of Q-learner in the dynamic treatment regime literature; however, they also provided a caveat that soft-thresholding cannot eliminate the non-regularity. Heuristically, the standard inference under the fixed alternative still provides poor finite sample coverage properties. Therefore, one still requires the local alternative asymptotics to derive inference procedures with uniform validity as we did for the hard thresholding estimator. We will leave this topic for future research.

### 3.4 Inference

The non-parametric bootstrap method provides consistent inference in many cases of regular estimators. However, this feature prevents using the non-parametric bootstrap inference for $\widehat{\psi}_{\text{elas}}$ because the indicator function of the preliminary test in (9) renders $\widehat{\psi}_{\text{elas}}$ a non-smooth and non-regular estimator (Shao, 1994). We formally show in the Online supplementary material the inconsistency of the nonparametric bootstrap inference for $\widehat{\psi}_{\text{elas}}$. Alternatively, Laber and Murphy (2011) proposed an adaptive confidence interval for the test error in classification, a non-regular statistics, by bootstrapping the upper and lower bounds of the test error. In this article, we propose an adaptive procedure for robust inference of $\psi_0$ accommodating the strength of violation of Assumption 2 in finite samples.

Let $e_k$ be a $p$-vector of zeros except that the $k$th component is one, and let $e_k^{\mathrm{T}} \psi_0$ be the $k$th component of $\psi_0$, for $k = 1, \ldots, p$. Because the asymptotic distribution of $n^{1/2} e_k^{\mathrm{T}} (\widehat{\psi}_{\text{elas}} - \psi_0)$ is different under the local and fixed alternatives, we propose different strategies for constructing CIs: under $H_{a,n}$, the asymptotics is non-standard, we construct a least favourable CI that guarantees good coverage properties uniformly over possible values of the local parameter; under $H_a$, the asymptotics is standard, we construct the usual Wald CI based on the normal limiting distribution.

First, under $H_{a,n}$, we rewrite $\mathcal{M}(\gamma; \eta)$ in (10) as $\mathcal{D}^{\text{NR}}(\mu_1) + \mathcal{D}^{\text{R}}$, where $\mathcal{D}^{\text{NR}}(\mu_1) = V_{\text{rt-eff}}^{1/2} \mathcal{Z}_1 \mathbf{1}(\mathcal{Z}_1^{\mathrm{T}} \mathcal{Z}_1 \geq c_\gamma)$ is the non-regular component with $\mathcal{Z}_1$ having mean $\mu_1$, $\mathcal{D}^{\text{R}} = -V_{\text{eff}}^{1/2} \mathcal{Z}_2$ is the regular component, and $\mathcal{D}^{\text{NR}}(\mu_1)$ and $\mathcal{D}^{\text{R}}$ are independent. For a fixed $\mu_1$, let $\widehat{Q}_{k,a}(\mu_1)$ be the approximated $100\alpha$th quantile of $\mathcal{D}^{\text{NR}}(\mu_1) + \mathcal{D}^{\text{R}}$, which can be obtained by rejective sampling. We can construct a $(1-\alpha)100\%$ confidence interval of $n^{1/2} e_k^{\mathrm{T}} (\widehat{\psi}_{\text{elas}} - \psi_0)$ as $[\widehat{Q}_{k,\alpha/2}(\mu_1), \widehat{Q}_{k,1-\alpha/2}(\mu_1)]$. Different CIs are required for different values of $\mu_1$. To accommodate different possible values of $\mu_1$, one solution is to construct the least favourable CI by taking the infimum of the lower bound of the CI $\widehat{Q}_{k,\alpha/2}(\mu_1)$ and the supremum of the upper bound of the CI $\widehat{Q}_{k,1-\alpha/2}(\mu_1)$ over all possible values of $\mu_1$. However, the range of $\mu_1$ can be vast, rendering the least favourable CI non-informative. We identify the plausible values of $\mu_1$ following a multivariate normal distribution with mean $n^{-1/2} \widehat{\Sigma}_{SS}^{-1/2} \sum_{i \in \mathcal{B}} \widehat{S}_{\text{rw},\hat{\psi}_{\text{rt}}}(V_i)$ and variance $I_{p \times p}$. Let $\widetilde{\alpha} = 1 - (1-\alpha)^{1/2}$, such that $(1 - \widetilde{\alpha})^2 = 1 - \alpha$ and let $\mathcal{B}_{1-\tilde{\alpha}}^{\mathrm{N}}$ be a $1 - \widetilde{\alpha}$ bounded region of a standard $p$-variate normal distribution. Then,

$$\mathcal{B}_{1-\tilde{\alpha}} = \left\{ \mu_1 : \left\{ n^{-1/2} \widehat{\Sigma}_{SS}^{-1/2} \sum_{i \in \mathcal{B}} \widehat{S}_{\text{rw},\hat{\psi}_{\text{rt}}}(V_i) - \mu_1 \right\} \in \mathcal{B}_{1-\tilde{\alpha}}^{\mathrm{N}} \right\}$$

is a bounded region of $\mu_1$ with asymptotic probability $1 - \widetilde{\alpha}$. We construct the $(1-\alpha)100\%$ least favourable CI for $n^{1/2} e_k^{\mathrm{T}} (\widehat{\psi}_{\text{elas}} - \psi_0)$ as $[\inf_{\mu_1 \in \mathcal{B}_{1-\tilde{\alpha}}} \widehat{Q}_{k,\tilde{\alpha}/2}(\mu_1), \sup_{\mu_1 \in \mathcal{B}_{1-\tilde{\alpha}}} \widehat{Q}_{k,1-\tilde{\alpha}/2}(\mu_1)]$. Here, using the wider $(1 - \widetilde{\alpha})100\%$ quantile range of $\widehat{Q}_k(\mu_1)$ instead of the $(1-\alpha)$ quantile range is necessary to guarantee the coverage of $(1-\alpha)$ due to ignoring other possible values of $\mu_1$ outside $\mathcal{B}_{1-\tilde{\alpha}}$.

Second, under $H_a$, Assumption 2 is strongly violated. As shown in Theorem 2, $n^{1/2} e_k^{\mathrm{T}} (\widehat{\psi}_{\text{elas}} - \psi_0)$ is regular and asymptotically normal, denoted by $\mathcal{M}(\gamma; \pm\infty, \pm\infty)$. Therefore, a $(1-\alpha)100\%$ confidence interval of $n^{1/2} e_k^{\mathrm{T}} (\widehat{\psi}_{\text{elas}} - \psi_0)$ can be constructed based on the $100\alpha/2$- and $100(1-\alpha/2)$th quantiles of the normal distribution $\mathcal{M}(\gamma; \pm\infty, \pm\infty)$, denoted by $[\widehat{Q}_{k,\alpha/2}(\pm\infty), \widehat{Q}_{k,1-\alpha/2}(\pm\infty)]$.

Finally, because the least favourable CI may be unnecessarily wide under $H_a$, we require a strategy to distinguish between $H_{a,n}$ corresponding to finite values of $\mu_1$ and $H_a$ corresponding to $\mu_1 = \pm\infty$. To do this, we use the test statistic $T$. Under $H_{a,n}$, $T = O_{\mathbb{P}}(1)$; while under $H_a$, $T = \infty$. Therefore, we specify a sequence of thresholds $\{\kappa_n : n \geq 1\}$ that diverges to infinity as $n \to \infty$ and compare $T$ to $\kappa_n$. Many choices of $\kappa_n$ can be considered, e.g., $\kappa_n = (\log n)^{1/2}$, which is similar to the BIC criterion (Andrews & Soares, 2010; Cheng, 2008). If $T \leq \kappa_n$, we choose the local alternative strategy to construct the least favourable CI, and if $T > \kappa_n$, we choose the fixed alternative

strategy to construct a normal CI, leading to an elastic CI

$$\text{ECI}_{k,1-\alpha} = \begin{cases} [\inf_{\mu_1 \in \mathcal{B}_{1-\tilde{\alpha}}} \widehat{Q}_{k,\tilde{\alpha}/2}(\mu_1), \ \sup_{\mu_1 \in \mathcal{B}_{1-\tilde{\alpha}}} \widehat{Q}_{k,1-\tilde{\alpha}/2}(\mu_1)], & \text{if } T \leq \kappa_n, \\ [\widehat{Q}_{k,\alpha/2}(\pm\infty), \ \widehat{Q}_{k,1-\alpha/2}(\pm\infty)], & \text{if } T > \kappa_n. \end{cases} \tag{13}$$

**Theorem 3** Suppose assumptions in Theorem 1 hold except that Assumption 2 may be violated. The asymptotic coverage rate of the elastic CI of $n^{1/2}e_k^{\mathrm{T}}(\widehat{\psi}_{\text{elas}} - \psi_0)$ in (13) satisfies

$$\lim_{n \to \infty} \mathbb{P}\{n^{1/2}e_k^{\mathrm{T}}(\widehat{\psi}_{\text{elas}} - \psi_0) \in \text{ECI}_{k,1-\alpha}\} \geq 1 - \alpha,$$

and the equality holds under $H_a$.

### 3.5 Adaptive selection of $\gamma$

The selection of $\gamma$ involves the bas-variance trade-off and therefore is important to determine the MSE of $\widehat{\psi}_{\text{elas}}$. Corollary 1 indicates that under $H_{a,n}$, the MSE of $\widehat{\psi}_{\text{elas}}$ in (12) involves two terms: Term 1 is $V_{\text{eff}} + V_{\text{rt-eff}}\{1 - F_{p+2}(c_\gamma; \lambda)\}$, and Term 2 involves $(V_{\text{eff}}\eta)^{\otimes 2}$. If $\eta$ is small, the MSE is dominated by Term 1, which can be made small if we select a small $\gamma$; while if $\eta$ is large, the MSE is dominated by Term 2, which can be made small if we select a large $\gamma$.

The above observation motivates an adaptive selection of $\gamma$. We propose to estimate $\eta$ by $\widehat{\eta} = n^{-1/2}\sum_{i \in \mathcal{B}} \widehat{S}_{\text{rw},\widehat{\psi}_{\text{rt}}}(V_i)$ and select $\gamma$ that minimises $\text{mse}(\gamma; \widehat{\eta})$, where $\text{mse}(\gamma; \eta)$ is given by (12) or approximated by rejective sampling. In practice, we can specify a grid of values from 0 to 1 for $\gamma$, denoted by $\mathcal{G}$, simulate the distribution of $\mathcal{M}(\gamma; \widehat{\eta})$ for all $\gamma \in \mathcal{G}$, and finally choose $\gamma$ to be the one in $\mathcal{G}$ that minimises the MSE of $\mathcal{M}(\gamma; \widehat{\eta})$. As corroborated by simulation, the selection strategy is effective in the sense that when the signal of violation is weak, the selected value of $\gamma$ is small and when the signal of violation is strong, the selected value of $\gamma$ is large.

## 4 Simulation study

We evaluate the finite sample performance of the proposed elastic estimator via simulation for robustness against unmeasured confounding and adaptive inference. Specifically, we compare the RT estimator, the efficient combining estimator, and the elastic estimator under settings that vary the strength of unmeasured confounding in the RW data. We also carry out simulation under a setting when the transportability assumption is violated in the RW data; see Online Supplementary Material, Section S4.3 in the supplementary material.

We first generate populations of size $10^5$. For each population, we generate the covariate $X = (1, X_1, X_2, X_3)^{\mathrm{T}}$, where $X_j \sim \text{Normal}(1, 1)$ for $j = 1, 2, 3$, and the treatment effect modifier is $Z = (1, X_1, X_2)^{\mathrm{T}}$. We generate $Y(a)$ by

$$Y(a) \mid X = \mu(X) + a \times \tau(Z) + \epsilon(a), \quad \epsilon(a) \sim \text{Normal}(0, 1), \tag{14}$$
$$\mu(X) = X_1 + X_2 + X_3, \quad \tau(Z) = \psi_0 + \psi_1 X_1 + \psi_2 X_2,$$

for $a = 0, 1$. Throughout the simulation, we fix $\psi_0$ to be zero and consider two cases for $(\psi_1, \psi_2)$: a) zero effect modification $(\psi_1, \psi_2) = (0, 0)$ and b) nonzero effect modification $(\psi_1, \psi_2) = (1, 1)$.

We then generate two samples from the target population. We generate the RT selection indicator by $\delta \mid X \sim \text{Bernoulli}\{\pi_\delta(X)\}$, where $\text{logit}\{\pi_\delta(X)\} = -4.5 - 2X_1 - 2X_2$. Under this selection mechanism, the selection rate is around 0.6%, which results in $m \approx 620$ RT subjects. We also take a random sample of size $n \in \{2000, 5000\}$ from the population to form an RW sample. In the RT sample, the treatment assignment is $A \mid X, \delta = 1 \sim \text{Bernoulli}\{e_1(X)\}$, where $e_1(X) = 0.5$. In the RW sample, $A \mid X, \delta = 0 \sim \text{Bernoulli}\{e_0(X)\}$, where $\text{logit}\{e_0(X)\} = \alpha - X_1 - X_2 - bX_3$ with adaptively chosen $\alpha$ to ensure the mean of $e_0(X)$ to be around 0.5. In addition, we vary $b$

to indicate the different strengths of unmeasured confounding in the analysis (violation of Assumption 2). The observed outcome $Y$ in both samples is $Y = AY(1) + (1 - A)Y(0)$.

To assess the robustness of the elastic integrative estimator against unmeasured confounding, we consider the omission of $X_3$ in all estimators, resulting in unmeasured confounding in the RW data. The strength of unmeasured confounding is indexed by $b$ in (14); high values of $b$ indicate strong levels of unmeasured confounding and vice versa. We specify the range of $b$ by 10 values in an irregular grid from 0 to 2 {0, 0.11, 0.23, 0.34, 0.46, 0.57, 0.69, 0.80, 1, 2}, which places more emphasis on the scenarios where Assumption 2 is weakly violated. We compare the following estimators for the HTE parameter $\psi$:
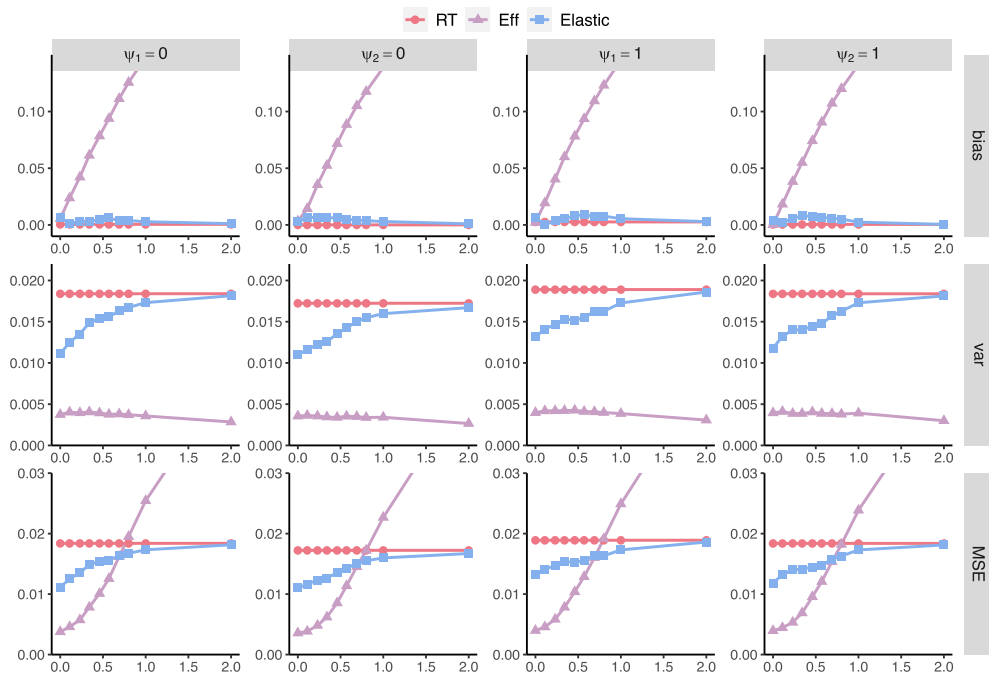
(a) RT $\widehat{\psi}_{\mathrm{rt}}$: the efficient estimator based only on the RT data solving (9) with $\mathbf{1}(T < c_\gamma) \equiv 0$;
(b) Eff $\widehat{\psi}_{\mathrm{eff}}$: the efficient integrative estimator solving (9) with $\mathbf{1}(T < c_\gamma) \equiv 1$;
(c) Elastic $\widehat{\psi}_{\mathrm{elas}}$: the proposed elastic integrative estimator solving (9) with adaptive selection of $\gamma$.

For all estimators, we estimate the propensity score function by a logistic sieve model with the power series $X$, $X^2$ and their two-way interactions (omitting $X_3$) and the outcome mean functions by linear sieve models with the power series $X$, $X^2$ and their two-way interactions (omitting $X_3$). If higher-order series is specified, it is necessary to select the series to balance the bias and variance in estimating the nuisance functions, such as using the penalised estimating equation approach (Lee, Yang, Dong, et al., 2022). The CIs are constructed for $\widehat{\psi}_{\mathrm{aipw}}$, $\widehat{\psi}_{\mathrm{rt}}$ and $\widehat{\psi}_{\mathrm{eff}}$ based on the perturbation-based resampling with the replication size 100 and for $\widehat{\psi}_{\mathrm{elas}}$ based on the elastic approach with $\kappa_n = (\log n)^{1/2}$. Sensitivity analysis shows that the coverage rates and widths of the CIs stay close with $\kappa_n = 0.5(\log n)^{1/2}$ (Online Supplementary Material, Section S4.4).

Figure 3 presents the plots of Monte Carlo biases, variances, and MSEs of estimators based on 2000 simulated datasets with numerical results reported in Online Supplementary Material, Table S3. Table 1 reports the coverage rates and widths of 95% CIs. The RT estimator $\widehat{\psi}_{\mathrm{rt}}$ is unbiased across different scenarios, and the coverage rates are close to the nominal level. However, $\widehat{\psi}_{\mathrm{rt}}$ has larger variances than other integrative estimators due to the small RT sample size. The efficient integrative estimator $\widehat{\psi}_{\mathrm{eff}}$ gains efficiency over $\widehat{\psi}_{\mathrm{rt}}$ by leveraging the large sample size of the RW data. However, the bias of $\widehat{\psi}_{\mathrm{eff}}$ increases as $b$ increases. Thus, $\widehat{\psi}_{\mathrm{eff}}$ has smaller MSEs than $\widehat{\psi}_{\mathrm{rt}}$ for small values of $b$ but larger MSEs for large values of $b$. The coverage rates of the CIs for $\widehat{\psi}_{\mathrm{eff}}$ deviate away from the nominal level as $b$ increases. This can lead to an uncontrolled false discovery of important treatment effect modifiers (see the case of zero effect modification with $\psi_1 = \psi_2 = 0$). The elastic integrative estimator $\widehat{\psi}_{\mathrm{elas}}$ with the adaptive selection of $\gamma$ reduces $\widehat{\psi}_{\mathrm{eff}}$'s biases across all scenarios regardless of the strength of unmeasured confounding. The challenging scenarios are indexed by $b$ around 0.44 and 0.67, where the small biases of $\widehat{\psi}_{\mathrm{elas}}$ occur. In these scenarios, the pretesting (built in the elastic estimator) has difficulty in detecting the RW sample's biases. However, $\widehat{\psi}_{\mathrm{elas}}$ with an adaptive selection of $\gamma$ achieves the smallest MSE among all estimators across all scenarios (Figure 3 and Online Supplementary Material, Table S3).

To inspect the performance of the proposed data-adaptive selection strategy, Online Supplementary Material, Table S8 reports Monte Carlo averages and standard deviations of the selected values for the local parameter $\eta$, the threshold $c_\gamma$, and the proportion of combining the RT and RW samples. As expected, $\widehat{\eta}$ increases as $b$ increases, indicating increased biases in the RW sample. The selected $\gamma$ increases (as a result, the proportion of combining the RT and RW samples decreases) as $b$ increases, which shows the proposed adaptive selection strategy is effective. To compare the adaptive selection strategy with the fixed threshold strategy, a simulation study in Online Supplementary Material, Section S4.5 shows that the elastic integrative estimator $\widehat{\psi}_{\mathrm{elas}}$ with a fixed threshold can have increased biases compared to a data-adaptive selected threshold.

The coverage rates of the ECIs for $\widehat{\psi}_{\mathrm{elas}}$ are close to the nominal level for all settings with different values of $b$. The ECIs are narrower than the CIs for $\widehat{\psi}_{\mathrm{rt}}$ when $b$ is small ($b \leq 0.46$ for $\psi_1 = \psi_2 = 0$ and $b \leq 0.34$ for $\psi_1 = \psi_2 = 1$), are wider than the CIs for $\widehat{\psi}_{\mathrm{rt}}$ when $b$ increases, and become close to the CIs for $\widehat{\psi}_{\mathrm{rt}}$ when $b$ reaches 1 or larger. However, the conservativity of the ECIs reduces as $n$ increases, and the ECIs can perform at least as well as the CIs for $\widehat{\psi}_{\mathrm{rt}}$ for any $b$ (see Online Supplementary Material, Table S6 for $n = 5000$).

**Figure 3.** Summary statistics plots of estimators of $(\psi_1, \psi_2)$ with respect to the strength of unmeasured confounding labelled by 'b'. In each plot, the three estimators $\widehat{\psi}_{\text{rt}}$, $\widehat{\psi}_{\text{eff}}$, and $\widehat{\psi}_{\text{elas}}$ are labelled by 'RT', 'Eff', and 'Elastic'. Each row of the plots corresponds to a different metrics: 'bias' for bias, 'var' for variance, 'MSE' for mean square error; each column of the plots corresponds to one component of $(\psi_1, \psi_2)$ in the two cases: $\psi_1 = 0$, $\psi_2 = 0$, $\psi_1 = 1$, and $\psi_2 = 1$ with $n = 2000$.

## 5 An application

We illustrate the potential benefit of the proposed elastic estimator to evaluate the effect of adjuvant chemotherapy for early-stage resected non-small cell lung cancer (NSCLC) using the CALGB 9633 data and a large clinical oncology database, the NCDB. In CALGB 9633, we include 319 patients, with 163 randomly assigned to observation $(A = 0)$ and 156 randomly assigned to chemotherapy $(A = 1)$. The NCDB cohort is selected based on the same patient eligibility criteria as the CALGB 9633 trial; see Online Supplementary Material, Section S5. The comparable NCDB sample includes 15,166 patients diagnosed with NSCLC between 2004 and 2016 in stage IB disease, with 10,903 on observation and 4,316 receiving chemotherapy after surgery. The numbers of treated and controls are relatively balanced in the CALGB 9633 trial, while they are unbalanced in the NCDB sample. We include five covariates in the analysis: gender $(1 = \text{male}, 0 = \text{female})$, age, the indicator for histology $(1 = \text{squamous}, 0 = \text{non} - \text{squamous})$, race $(1 = \text{white}, 0 = \text{non-white})$, and tumour size in centimetre. The outcome is the overall survival within three years after the surgery, i.e., $Y = 1$ if died due to all causes and $Y = 0$ otherwise. We are interested in estimating the HTE of adjuvant chemotherapy over observation after resection for the patient population with the same set of eligibility criteria as that of CALGB 9633.

Table 2 reports the covariate means by sample and treatment group. Due to treatment randomisation, covariates are balanced between the treated and the control in the CALGB 9633 trial sample. While due to a lack of treatment randomisation, covariates are relatively unbalanced in the NCDB sample. Older patients with histology and smaller tumours are likely to choose a conservative treatment on observation. Moreover, we cannot rule out the possibility of unmeasured confounders in the NCDB sample.

We assume a linear HTE function with tumour size as the treatment effect modifier. We compare the same set of estimators and variance estimators considered in the simulation study and the efficient estimator applied to the real-world NCDB cohort, denoted by $\widehat{\psi}_{\text{rw}}$. Table 3 reports

**Table 1.** Simulation results for coverage rates and widths of 95% confidence intervals for $\widehat{\psi}_{rt}$, $\widehat{\psi}_{eff}$, and $\widehat{\psi}_{elas}$ (labelled as 'RT', 'Eff', and 'Elastic') in the two cases: zero effect modification $\psi_1 = \psi_2 = 0$ (left) and nonzero effect modification $\psi_1 = \psi_2 = 1$ (right) with $n = 2000$; the slightly wider ECIs for $\widehat{\psi}_{eff}$ (than CIs for $\widehat{\psi}_{rt}$) are bolded
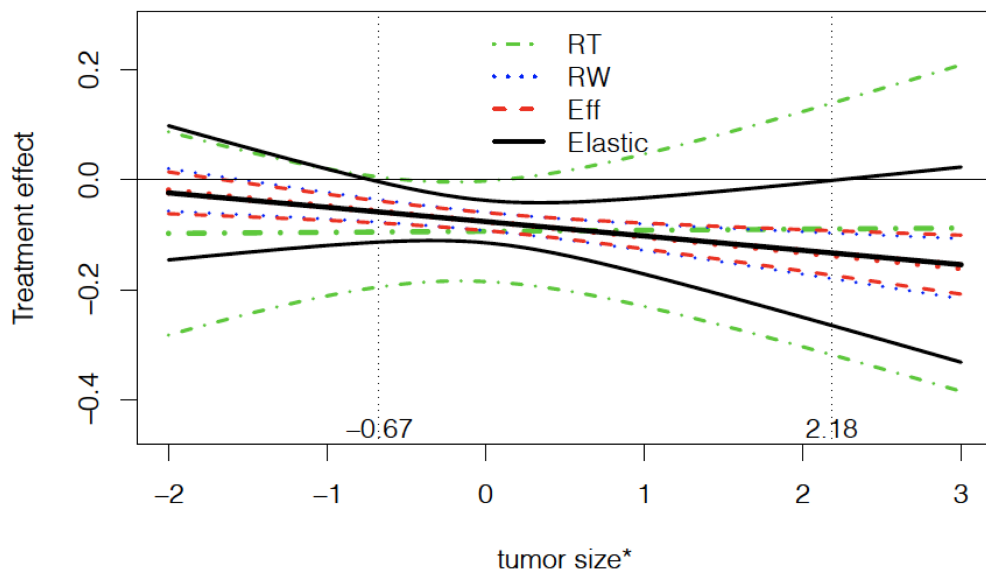
| | Case 1: zero effect modification | | | | | | Case 2: nonzero effect modification | | | | | |
| | RT | | Eff | | Elastic | | RT | | Eff | | Elastic | |
| $b$ | $\psi_1 = 0$ | $\psi_2 = 0$ | $\psi_1 = 0$ | $\psi_2 = 0$ | $\psi_1 = 0$ | $\psi_2 = 0$ | $\psi_1 = 1$ | $\psi_2 = 1$ | $\psi_1 = 1$ | $\psi_2 = 1$ | $\psi_1 = 1$ | $\psi_2 = 1$ |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| **Coverage Rate (%)** | | | | | | | | | | | | |
| 0 | 94.1 | 94.1 | 93.8 | 93.7 | 92.7 | 92.5 | 94.3 | 93.8 | 95.0 | 94.2 | 92.7 | 92.5 |
| 0.11 | 94.1 | 94.1 | 92.2 | 92.7 | 93.2 | 92.8 | 94.3 | 93.8 | 93.3 | 92.9 | 92.9 | 92.7 |
| 0.23 | 94.1 | 94.0 | 88.5 | 89.8 | 92.8 | 92.8 | 94.3 | 93.8 | 89.8 | 89.0 | 93.3 | 92.7 |
| 0.34 | 94.1 | 94.0 | 83.2 | 84.5 | 94.0 | 93.8 | 94.3 | 93.8 | 84.9 | 83.5 | 94.4 | 93.5 |
| 0.46 | 94.1 | 94.0 | 74.7 | 76.3 | 94.5 | 94.5 | 94.3 | 93.8 | 76.8 | 75.8 | 94.5 | 94.4 |
| 0.57 | 94.1 | 94.0 | 66.4 | 66.1 | 95.5 | 95.2 | 94.3 | 93.8 | 67.2 | 66.8 | 95.5 | 94.8 |
| 0.69 | 94.1 | 94.1 | 56.1 | 56.3 | 95.5 | 95.8 | 94.3 | 93.8 | 56.8 | 55.9 | 95.3 | 94.6 |
| 0.8 | 94.1 | 94.0 | 46.3 | 46.8 | 95.5 | 95.6 | 94.3 | 93.8 | 46.5 | 45.2 | 95.3 | 95.0 |
| 1 | 94.1 | 94.0 | 31.5 | 31.1 | 95.5 | 95.0 | 94.3 | 93.8 | 30.9 | 29.4 | 95.5 | 94.9 |
| 2 | 94.1 | 94.0 | 2.9 | 3.6 | 94.3 | 94.4 | 94.3 | 93.8 | 2.6 | 3.0 | 94.7 | 94.2 |
| **Width (×10⁻³)** | | | | | | | | | | | | |
| 0 | 528 | 528 | 243 | 242 | 472 | 473 | 529 | 530 | 243 | 243 | 472 | 474 |
| 0.11 | 527 | 528 | 242 | 242 | 488 | 487 | 529 | 530 | 242 | 243 | 479 | 480 |
| 0.23 | 527 | 528 | 241 | 242 | 496 | 497 | 529 | 530 | 241 | 242 | 498 | 500 |
| 0.34 | 528 | 528 | 241 | 241 | 516 | 516 | 529 | 530 | 241 | 242 | 511 | 514 |
| 0.46 | 528 | 528 | 239 | 240 | **530** | **530** | 529 | 530 | 240 | 240 | 524 | 526 |
| 0.57 | 528 | 528 | 238 | 238 | **535** | **535** | 529 | 530 | 238 | 239 | **530** | **532** |
| 0.69 | 528 | 528 | 235 | 236 | **534** | **534** | 529 | 530 | 236 | 236 | **529** | **531** |
| 0.8 | 528 | 528 | 233 | 234 | **532** | **532** | 529 | 530 | 233 | 234 | **530** | **532** |
| 1 | 528 | 528 | 229 | 230 | **529** | **529** | 529 | 530 | 229 | 230 | **530** | **532** |
| 2 | 528 | 528 | 207 | 208 | 527 | 527 | 529 | 530 | 208 | 209 | 528 | 530 |

**Table 2.** Covariate means with standard errors in parentheses by sample and treatment group in the CALGB 9633 trial and NCDB samples

| | *A* | *N* | Age (years) | tumour size (cm) | Male (%) | Squamous (%) | White (%) |
|---|---|---|---|---|---|---|---|
| RT: | 0, 1 | 319 | 60.8 (9.62) | 4.60 (2.08) | 63.9 | 39.8 | 89.3 |
| CALGB 9633 | 1 | 156 | 60.6 (10) | 4.62 (2.09) | 64.1 | 40.4 | 90.4 |
| | 0 | 163 | 61.1 (9.25) | 4.57 (2.07) | 63.8 | 39.3 | 88.3 |
| RW: | 0, 1 | 15,166 | 67.9 (10.2) | 4.82 (1.71) | 54.6 | 39.1 | 89.6 |
| NCDB | 1 | 4,263 | 63.9 (9.23) | 5.19 (1.79) | 54.3 | 35.6 | 88.6 |
| | 0 | 10,903 | 69.4 (10.1) | 4.67 (1.65) | 54.8 | 40.5 | 90.0 |

**Table 3.** Point estimate, standard error, and 95% Wald confidence interval of the causal risk difference between adjuvant chemotherapy and observation based on the CALGB 9633 trial sample and the NCDB sample: tumour size* = (tumour size − 4.82)/1.72

| | Intercept ($\psi_{0,1}$) | | | tumour size* ($\psi_{0,2}$) | | |
|---|---|---|---|---|---|---|
| | Est. | S.E. | C.I. | Est. | S.E. | C.I. |
| RT | −0.094 | 0.054 | (−0.202, 0.015) | 0.002 | 0.055 | (−0.107, 0.111) |
| RW | −0.076 | 0.0085 | (−0.093, −0.059) | −0.029 | 0.009 | (−0.046, −0.011) |
| Eff | −0.076 | 0.0083 | (−0.093, −0.059) | −0.026 | 0.009 | (−0.043, −0.009) |
| Elastic | −0.076 | 0.0196 | (−0.115, −0.037) | −0.026 | 0.029 | (−0.084, 0.032) |



**Figure 4.** Estimated treatment effect as a function of the (standardised) tumour size along with the 95% Wald confidence intervals: tumour size* = (tumour size − 4.82)/1.72, RT, RW, and Eff are the efficient estimator applied to the RT, RW, and combined sample, respectively, and Elastic is the proposed elastic combining estimator.

the results. Figure 4 shows the estimated treatment effect as a function of the standardised tumour size. Due to the limited sample size of the trial sample, all components in $\widehat{\psi}_{\mathrm{rt}}$ are not significant. Due to the large sample size of the NCDB sample, $\widehat{\psi}_{\mathrm{rw}}$ and $\widehat{\psi}_{\mathrm{eff}}$ are close and reveal that adjuvant chemotherapy significantly reduced cancer recurrence within three years after the surgery. Patients with larger tumour sizes benefit more from adjuvant chemotherapy. However, this finding may be subject to possible biases of the NCDB sample. In the proposed elastic integrative analysis, the test statistic is $T = 1.9$; there is no strong evidence that the NCDB presents hidden confounding in our analysis. As a result, the elastic integrative estimator $\widehat{\psi}_{\mathrm{elas}}$ remains the same as $\widehat{\psi}_{\mathrm{eff}}$. In reflection of the pre-testing procedure, the estimated standard error of $\widehat{\psi}_{\mathrm{elas}}$ is larger than $\widehat{\psi}_{\mathrm{eff}}$'s. From Figure 4, patients with tumour sizes in $[4.82 + 1.72 \times (-0.67), 4.82 + 1.72 \times (2.18)] = [3.67, 8.57]$ significantly benefit from adjuvant chemotherapy in improving overall survival within three years after the surgery.

## 6 Concluding remarks

The proposed elastic estimator integrates 'high-quality small data' with 'big data' to simultaneously leverage small but carefully controlled unbiased experiments and massive but possibly biased RW datasets for HTEs. Most causal inference methods require the no unmeasured confounding assumption. However, this assumption may not hold for the RW data due to the uncontrolled, real-world data collection mechanism and is unverifiable based only on the RW data. Utilising the design advantage of RTs, we can gauge the reliability of the RW data and decide whether or not to use RW data in an integrative analysis.

The key assumptions underpinning our framework are the structural HTE model, i.e., Model (1), HTE transportability, and no unmeasured confounding. In practice, RTs usually consider much narrower populations than seen in the real world. Improving the generalisability or external validity of RT findings has been an important research topic in the data integration literature (e.g., Cole & Stuart, 2010; Lee, Yang, Dong, et al., 2022; Rudolph & van der Laan, 2017). Besides Assumption 1(i), the positivity of trial participation or the overlap of the covariate distribution between the RT and RW samples is required in the problem of generalisability. We emphasise that although, formally, we do not require the overlap assumption between the RT and RW samples, its violation renders Model (1) and transportability vulnerable. When transporting from the narrow RT sample to the broader RW sample, the reliable information of treatment effects for the non-overlapping region essentially hinges on the extrapolation from the RT sample. If there is no strong prior knowledge, Model (1) and transportability may not hold. In this case, the RT estimate and the RW estimate of the HTE can be inconsistent due to model misspecification even when there are no unmeasured confounders. See a simulation study in Online Supplementary Material, Section S4.3. The inconsistency of the RW estimator with the RT estimator may reflect violation of either transportability (e.g., due to model misspecification) or unmeasured confounding. Some practical strategies (e.g., matching) can be implemented to select an RW sample with sufficient overlap with the RT sample to improve their comparability and the chance of successfully integrating the information from two separate sources; see Online Supplementary Material, Section S5.2.

The elastic integrative estimator gains efficiency over the RT-only estimator by integrating the reliable RW data and also automatically detecting bias in the RW data and gears to the RT data. However, the proposed estimator is non-regular and belongs to pre-test estimation by construction (Giles & Giles, 1993). To demonstrate the non-regularity issue, we characterise the distribution of the elastic integrative estimator under local alternatives, which better approximates the finite-sample behaviours. Moreover, we provide a data-adaptive selection of the threshold in the testing procedure, which guarantees small MSEs of the estimator. Nonetheless, fixing the threshold may not control bias well under $H_{a,n}$; see a simulation study in Online Supplementary Material, Section S4.5. If the investigator prefers small biases in the elastic combining estimator, we recommend setting the lower bounds of a grid for selecting $\gamma$. Although the elastic confidence intervals demonstrate good coverage properties in our simulation under all hypotheses $H_0$, $H_{a,n}$, and $H_a$, an open problem remains for the post-selection inference after a data-adaptive selection of the threshold in the testing procedure, which will be rigorously analysed theoretically and empirically in the future study.

The proposed framework can also be extended to individualised treatment regime learning (Chu et al., 2022; Wu & Yang, 2021, 2022) and the data integration problem of combining probability

and non-probability samples (Yang & Kim, 2020; Yang et al., 2019, 2021). However, an additional complication arises due to the mixed design-based and super-population inference framework, which will be overcome in future research.

*Conflict of interest:* None declared.

## Funding

## Data availability

The authors have access to the data from CALGB 9633 through Alliance Statistics and Data Center, and the NCDB database through the approved license from American College of Surgeons to Duke University School of Medicine. De-identified individual patient data were used for the application, and the results of the analysis were summarized in the manuscript. The readers may contact Alliance data sharing working group to request an access to CALGB 9633 data and initiate an application process with American College of Surgeons to gain the access to the NCDB database.

## Supplementary material

Supplementary material are available at *Journal of the Royal Statistical Society: Series B* online.

## References

Andrews D. W., & Soares G. (2010). Inference for parameters defined by moment inequalities using generalized moment selection. *Econometrica*, 78(1), 119–157. https://doi.org/10.3982/ECTA7502

Bickel P. J., Klaassen C., Ritov Y., & Wellner J. (1993). *Efficient and adaptive inference in semiparametric models*. Johns Hopkins University Press.

Buchanan A. L., Hudgens M. G., Cole S. R., Mollan K. R., Sax P. E., Daar E. S., Adimora A. A., Eron J. J., & Mugavero M. J. (2018). Generalizing evidence from randomized trials using inverse probability of sampling weights. *Journal of the Royal Statistical Society: Series A (Statistics in Society)*, 181(4), 1193– 1209. https://doi.org/10.1111/rssa.12357

Chakraborty B., & Moodie E. E. (2013). *Statistical methods for dynamic treatment regimes*. Springer.

Chakraborty B., Murphy S., & Strecher V. (2010). Inference for non-regular parameters in optimal dynamic treatment regimes. *Statistical Methods in Medical Research*, 19(3), 317–343. https://doi.org/10.1177/0962280209105013

Cheng X. (2008). Robust confidence intervals in nonlinear regression under weak identification, *Manuscript, Department of Economics, Yale University*.

Chernozhukov V., Chetverikov D., Demirer M., Duflo E., Hansen C., Newey W., & Robins J. (2018). Double/debiased machine learning for treatment and structural parameters. *The Econometrics Journal*, 21(1), C1–C68. https://doi.org/10.1111/ectj.12097

Chu J., Lu W., & Yang S. (2022). 'Targeted optimal treatment regime learning using summary statistics', arXiv, arXiv:2201.06229, preprint: not peer reviewed.

Cole S. R., & Stuart E. A. (2010). Generalizing evidence from randomized clinical trials to target populations: The ACTG 320 trial. *American Journal of Epidemiology*, 172(1), 107–115. https://doi.org/10.1093/aje/kwq084

Colnet B., Mayer I., Chen G., Dieng A., Li R., Varoquaux G., Vert J.-P., Josse J., & Yang S. (2020). 'Causal inference methods for combining randomized trials and observational studies: a review', arXiv, arXiv:2011.08047, preprint: not peer reviewed.

Dahabreh I. J., Robertson S. E., Tchetgen E. J., Stuart E. A., & Hernán M. A. (2019). Generalizing causal inferences from individuals in randomized trials to all trial-eligible individuals. *Biometrics*, 75(2), 685–694. https://doi.org/10.1111/biom.13009

Giles J. A., & Giles D. E. A. (1993). Pre-test estimation and testing in econometrics: Recent developments. *Journal of Economic Surveys*, 7(2), 145–197. https://doi.org/10.1111/j.1467-6419.1993.tb00163.x

Hamburg M. A., & Collins F. S. (2010). The path to personalized medicine. *New England Journal of Medicine*, 363(4), 301–304. https://doi.org/10.1056/NEJMp1006304

Hernán M. A., & Robins J. M. (2016). Using big data to emulate a target trial when a randomized trial is not available. *American Journal of Epidemiology*, 183(8), 758–764. https://doi.org/10.1093/aje/kwv254

Hu F., & Kalbfleisch J. D. (2000). The estimating function bootstrap. *Canadian Journal of Statistics*, 28(3), 449–481. https://doi.org/10.2307/3315958

Katz A., & Saad E. D. (2009). CALGB 9633: An underpowered trial with a methodologically questionable conclusion. *Journal of Clinical Oncology*, 27(13), 2300–2301. https://doi.org/10.1200/JCO.2008.21.1565

Kennedy E. H. (2020). 'Optimal doubly robust estimation of heterogeneous causal effects', arXiv, arXiv:2004.14497, preprint: not peer reviewed.

Laber E. B., & Murphy S. A. (2011). Adaptive confidence intervals for the test error in classification. *Journal of the American Statistical Association*, 106(495), 904–913. https://doi.org/10.1198/jasa.2010.tm10053

LeCam L. (1953). On some asymptotic properties of maximum likelihood estimates and related Bayes estimates. *University of California Publications in Statistics*, 1, 277–330.

Le Chevalier T. (2003). Results of the Randomized International Adjuvant Lung Cancer Trial (IALT): Cisplatin-based chemotherapy (CT) vs no CT in 1867 patients with resected non-small cell lung cancer (NSCLC). *Lung Cancer*, 21, 238–238. https://doi.org/10.1016/S0169-5002(03)91656-4

Lee D., Yang S., Dong L., Wang X., Zeng D., & Cai J. (2022). Improving trial generalizability using observational studies. *Biometrics*. https://doi.org/10.1111/biom.13609

Lee D., Yang S., & Wang X. (2022). 'Generalizable survival analysis of randomized controlled trials with observational studies', arXiv, arXiv:2201.06595, preprint: not peer reviewed.

Li X., Ding P., & Rubin D. B. (2018). Asymptotic theory of rerandomization in treatment–control experiments. *Proceedings of the National Academy of Sciences*, 115(37),9157–9162. https://doi.org/10.1073/pnas.1808191115

Luedtke A., Carone M., & van der Laan M. J. (2019). 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. https://doi.org/10.1111/rssb.12299

Luedtke A. R., & van der Laan M. J. (2016). Super-learning of an optimal dynamic treatment rule. *The International Journal of Biostatistics*, 12(1), 305–332. https://doi.org/10.1515/ijb-2015-0052

Neyman J. (1923). Sur les applications de la thar des probabilités aux experiences Agaricales: Essay de principle. English translation of excerpts by Dabrowska, D. and Speed, T.. *Statistical Science*, 5, 465–472.

Norris S., Atkins D., Bruening W., Fox S., Johnson E., Kane R., Morton S. C., Oremus M., Ospina M., Randhawa G., Schoelles K., Shekelle P., & Viswanathan M. (2010). Selecting observational studies for comparing medical interventions. In *Methods guide for effectiveness and comparative effectiveness reviews [Internet]*. Agency for Healthcare Research and Quality.

Pearl J., & Bareinboim E. (2011). Transportability of causal and statistical relations: A formal approach. In *Data Mining Workshops (ICDMW), 2011 IEEE 11th International Conference on* (pp. 540–547). IEEE.

Prentice R. L., Langer R., Stefanick M. L., Howard B. V., Pettinger M., Anderson G., Barad D., Curb J. D., Kotchen J., Kuller L., Limacher M., & Wactawski-Wende J. (2005). 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. https://doi.org/10.1093/aje/kwi223

Richardson T. S., Robins J. M., & Wang L. (2017). On modeling and estimation for the relative risk and risk difference. *Journal of the American Statistical Association*, 112(519), 1121–1130. https://doi.org/10.1080/01621459.2016.1192546

Robins J. M. (1994). Correcting for non-compliance in randomized trials using structural nested mean models. *Communications in Statistics-Theory and Methods*, 23(8), 2379–2412. https://doi.org/10.1080/03610929408831393

Rothwell P. M. (2005). Subgroup analysis in randomised controlled trials: Importance, indications, and interpretation. *The Lancet*, 365(9454), 176–186. https://doi.org/10.1016/S0140-6736(05)17709-5

Rothwell P. M., Mehta Z., Howard S. C., Gutnikov S. A., & Warlow C. P. (2005). From subgroups to individuals: General principles and the example of carotid endarterectomy. *The Lancet*, 365(9455), 256–265. https://doi.org/10.1016/S0140-6736(05)70156-2

Rotnitzky A., Smucler E., & Robins J. M. (2019). 'Characterization of parameters with a mixed bias property', arXiv, arXiv:1904.03725, preprint: not peer reviewed.

Rubin D. B. (1974). Estimating causal effects of treatments in randomized and nonrandomized studies. *Journal of Educational Psychology*, 66(5), 688–701. https://doi.org/10.1037/h0037350

Rudolph K. E., & van der Laan M. J. (2017). Robust estimation of encouragement design intervention effects transported across sites. *Journal of the Royal Statistical Society: Series B (Statistical Methodology)*, 79(5), 1509–1525. https://doi.org/10.1111/rssb.12213

Shao J. (1994). Bootstrap sample size in nonregular cases. *Proceedings of the American Mathematical Society*, 122(4), 1251–1262. https://doi.org/10.1090/S0002-9939-1994-1227529-8

Sherman R. E., Anderson S. A., Dal Pan G. J., Gray G. W., Gross T., Hunter N. L., LaVange L., Marinac-Dabic D., Marks P. W., Robb M. A., Shuren J., Temple R., Woodcock J., Yue L. Q., & Califf R. M. (2016). Real-world evidence—what is it and what can it tell us. *New England Journal of Medicine*, *375*(23), 2293–2297. https://doi.org/10.1056/NEJMsb1609216

Shi C., Song R., & Lu W. (2016). Robust learning for optimal treatment decision with np-dimensionality. *Electronic Journal of Statistics*, *10*(2), 2894–2921. https://doi.org/10.1214/16-EJS1178

Sobel M., Madigan D., & Wang W. (2017). Causal inference for meta-analysis and multi-level data structures, with application to randomized studies of Vioxx. *Psychometrika*, *82*(2), 459–474. https://doi.org/10.1007/s11336-016-9507-z

Staiger D., & Stock J. H. (1997). Instrumental variables regression with weak instruments. *Econometrica*, *65*(3), 557–586. https://doi.org/10.2307/2171753

Stefanski L. A., & Boos D. D. (2002). The calculus of M-estimation. *The American Statistician*, *56*(1), 29–38. https://doi.org/10.1198/000313002753631330

Strauss G. M., Herndon J. E., Maddaus M. A., Johnstone D. W., Johnson E. A., Harpole D. H., Gillenwater H. H., Watson D. M., Sugarbaker D. J., Schilsky R. L., Vokes E. E., & Green M. R. (2008). Adjuvant paclitaxel plus carboplatin compared with observation in stage IB non–small-cell lung cancer: CALGB 9633 with the Cancer and Leukemia Group B, Radiation Therapy Oncology Group, and North Central Cancer Treatment Group Study Groups. *Journal of Clinical Oncology*, *26*(31), 5043–5051. https://doi.org/10.1200/JCO.2008.16.4855

Stuart E. A., Cole S. R., Bradshaw C. P., & Leaf P. J. (2011). The use of propensity scores to assess the generalizability of results from randomized trials. *Journal of the Royal Statistical Society: Series A (Statistics in Society)*, *174*(2), 369–386. https://doi.org/10.1111/j.1467-985X.2010.00673.x

Tallis G. M. (1963). Elliptical and radial truncation in normal populations. *The Annals of Mathematical Statistics*, *34*(3), 940–944. https://doi.org/10.1214/aoms/1177704016

Tian L., Alizadeh A. A., Gentles A. J., & Tibshirani R. (2014). A simple method for estimating interactions between a treatment and a large number of covariates. *Journal of the American Statistical Association*, *109*(508), 1517–1532. https://doi.org/10.1080/01621459.2014.951443

Toyoda T., & Wallace T. D. (1979). Pre-testing on part of the data. *Journal of Econometrics*, *10*(1), 119–123. https://doi.org/10.1016/0304-4076(79)90071-X

US Food and Drug Administration (2019). Rare diseases: Natural history studies for drug development, *https://www.fda.gov/media/122425/ (accessed 1 May 2022)*.

van der Vaart A. W. (2000). *Asymptotic statistics*. Cambridge University Press.

Verde P. E., & Ohmann C. (2015). Combining randomized and non-randomized evidence in clinical research: A review of methods and applications. *Research Synthesis Methods*, *6*(1), 45–62. https://doi.org/10.1002/jrsm.1122

Wu L., & Yang S. (2021). 'Transfer learning of individualized treatment rules from experimental to real-world data', arXiv, arXiv:2108.08415, preprint: not peer reviewed.

Wu L., & Yang S. (2022). Integrative *r*-learner of heterogeneous treatment effects combining experimental and observational studies. In *Proceedings of Machine Learning Research* (Vol. 140, pp. 1–S5).

Yang S., & Ding P. (2018). Asymptotic inference of causal effects with observational studies trimmed by the estimated propensity scores. *Biometrika*, *105*(2), 487–493. https://doi.org/10.1093/biomet/asy008

Yang S., & Kim J. K. (2020). Statistical data integration in survey sampling: A review. *Japanese Journal of Statistics and Data Science*, *3*(2), 625–650. https://doi.org/10.1007/s42081-020-00093-w

Yang S., Kim J. K., & Hwang, Y. (2021). Integration of data from probability surveys and big found data for finite population inference using mass imputation. *Survey Methodology*, *47*(1), 29–58.

Yang S., Kim J. K., & Song R. (2019). Doubly robust inference when combining probability and non-probability samples with high-dimensional data. *Journal of the Royal Statistical Society, Series B*, *82*(2), 445–465. https://doi.org/10.1111/rssb.12354

Zhao Y.-Q., Zeng D., Tangen C. M., & Leblanc M. L. (2019). Robustifying trial-derived optimal treatment rules for a target population. *Electronic Journal of Statistics*, *13*, 1717–1743. https://doi.org/10.1214/19-EJS1540