
Estimating individual treatment effects under unobserved confounding using binary instruments

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 Estimating individual treatment effects (ITEs) from observational data is relevant
2 in many fields such as personalized medicine. However, in practice, the treatment
3 assignment is usually confounded by unobserved variables and thus introduces
4 bias. A remedy to remove the bias is the use of instrumental variables (IVs). Such
5 settings are widespread in medicine (e.g., trials where compliance is used as binary
6 IV). In this paper, we propose a novel, multiply robust machine learning framework,
7 called MRIV, for estimating ITEs using binary IVs and thus yield an unbiased ITE
8 estimator. Different from previous work for binary IVs, our framework estimates
9 the ITE directly via a pseudo outcome regression. (1) We provide a theoretical
10 analysis where we show that our framework yields multiply robust convergence
11 rates: our ITE estimator achieves fast convergence even if several nuisance esti-
12 mators converge slowly. (2) We further show that our framework asymptotically
13 outperforms state-of-the-art plug-in IV methods for ITE estimation. (3) We build
14 upon our theoretical results and propose a tailored deep neural network architecture
15 called MRIV-Net for ITE estimation using binary IVs. Across various compu-
16 tational experiments, we demonstrate empirically that our MRIV-Net achieves
17 state-of-the-art performance. To the best of our knowledge, our MRIV is the first
18 multiply robust machine learning framework tailored to estimating ITEs in the
19 binary IV setting.

20 1 Introduction

21 Individual treatment effects (ITEs) are relevant across many disciplines such as marketing [41] and
22 personalized medicine [51]. Knowledge about ITEs provides insights into the heterogeneity of
23 treatment effects, and thus help in potentially better treatment decisions.

24 Many recent works that use machine learning to estimate ITEs are based on the assumption of
25 unconfoundedness [1, 15, 27, 36, 42], **In practice, however, this assumption is often violated because**
26 **it is common that some confounders are not reported in the data.** Typical examples are race,
27 income, gender, or the socioeconomic status of patients, which are not stored in medical files. If the
28 confounding is sufficiently strong, standard methods for estimating ITEs suffer from confounding
29 bias [31], which may lead to inferior treatment decisions.

30 To handle unobserved confounders, instrumental variables (IVs) can be leveraged to relax the
31 assumption of unconfoundedness and still compute reliable ITE estimates. IV methods were originally
32 developed in economics [48], but, only recently, there is a growing interest in combining IV methods
33 with machine learning (see Sec. 3). Importantly, IV methods outperform classical ITE estimators
34 if a sufficient amount of confounding is not observed [17]. We thus aim at estimating ITEs from
35 observational data under unobserved confounding using IVs.

36 In this paper, we consider the setting where a single binary instrument is available. This setting is
 37 widespread in personalized medicine (and other applications such as marketing or public policy)
 38 [9]. In fact, the setting is encountered in essentially all observational or randomized studies with
 39 observed non-compliance [19]. As an example, consider a randomized controlled trial (RCT), where
 40 treatments are randomly assigned to patients and their outcomes are observed. Due to some potentially
 41 unobserved confounders (e.g., income, education), some patients refuse to take the treatment initially
 42 assigned to them. Here, the treatment assignment serves as a binary IV. Moreover, such RCTs have
 43 been widely used by public decision-makers, e.g., to analyze the effect of health insurance on health
 44 outcome (see the so-called *Oregon health insurance experiment*) [16] or the effect of military service
 45 on lifetime earnings [2].

46 We propose a novel machine learning framework (called MRIV) for estimating ITEs using binary IVs.
 47 **Our framework takes an initial ITE estimator and nuisance parameter estimators as input to perform a**
 48 **pseudo-outcome regression. Importantly, our framework uses a multiply robust parametrization of**
 49 **the efficient influence function as pseudo outcome.**

50 We provide a theoretical analysis, **where we use tools from [22] to show that our framework achieves**
 51 **a multiply robust convergence rate**, i.e., our MRIV converges with a fast rate even if several nuisance
 52 parameters converge slowly. We further show that, compared to existing plug-in IV methods, the
 53 performance of our framework is asymptotically superior. Finally, we leverage our framework and,
 54 on top of it, build a tailored deep neural network called MRIV-Net.

55 **Differences to existing literature:** Our framework is **multiply robust**¹, i.e., it is consistent in
 56 the union of three different model specifications. This is different from existing methods for ITE
 57 estimation using IVs, which are only doubly robust (e.g., Syrgkanis et al. [40]) or plug-in estimators
 58 [5, 19].

59 **Contributions:**² (1) We propose a novel, multiply robust machine learning
 60 framework (called MRIV) to learn the ITE using the binary IV setting. To the best of our knowledge,
 61 ours is the first that is multiply robust, i.e., consistent in the union of three model specifications. **For comparison,**
 62 **existing works for ITE estimation are only double robust [45, 40].** (2) We
 63 prove that MRIV achieves a multiply robust convergence rate. **This is**
 64 **different to methods for IV settings which are only doubly robust, such**
 65 **as [40].** We further show that our MRIV is asymptotically superior to existing
 66 plug-in estimators. (3) We propose a tailored deep neural network,
 67 called MRIV-Net, which builds upon our framework to estimate ITEs.
 68 We demonstrate that MRIV-Net achieves state-of-the-art performance.

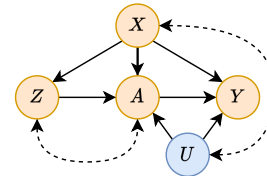


Figure 1: Underlying causal graph. The instrument Z has a direct influence on the treatment A , but does not have a direct effect on the outcome Y . Note that we allow for unobserved confounders for both Z – A (dashed line) and A – Y (given by U). **Our setting is general in that U can be correlated or uncorrelated with the observed confounders X .**

70 2 Problem setup

71 **Data generating process:** We observe data $\mathcal{D} = (x_i, z_i, a_i, y_i)_{i=1}^n$ con-
 72 sisting of $n \in \mathbb{N}$ observations of the tuple (X, Z, A, Y) . Here, $X \in \mathcal{X}$
 73 are observed confounders, $Z \in \{0, 1\}$ is a binary instrument, $A \in \{0, 1\}$
 74 is a binary treatment, and $Y \in \mathbb{R}$ is an outcome of interest. Furthermore,
 75 we assume the existence of unobserved confounders $U \in \mathcal{U}$, which affect
 76 both the treatment A and the outcome Y . The causal graph is shown in
 77 Fig. 1.

78 **Applicability:** Our proposed framework is widely applicable in practice, namely to all settings with
 79 the above data generating process. This includes both (1) observational data and (2) RCTs with
 80 non-compliance. For (1), observational data is commonly encountered in, e.g., personalized medicine.
 81 Here, modeling treatments as binary variables is consistent with previous literature on causal effect
 82 estimation and standard in medical practice [33]. For (2), our setting is further encountered in RCTs
 83 when the instrument Z is a randomized treatment assignment but individuals do not comply with
 84 their treatment assignment. Such RCTs have been extensively used by public decision-makers, e.g.,

¹For a detailed introduction to multiple robustness and its importance in treatment effect estimation, we refer to [46], Section 4.5.

²Codes are in the supplementary materials. Codes are also available at <https://anonymous.4open.science/r/MRIV-Net-0AC4> (Upon acceptance, we replace the link and point to a public GitHub repository).

85 to analyze the effect of health insurance on health outcome [16] or the effect of military service on
86 lifetime earnings [2].

87 We build upon the potential outcomes framework [34] for modeling causal effects. Let $Y(a, z)$
88 denote the potential outcome that would have been observed under $A = a$ and $Z = z$. Following
89 previous literature on IV estimation [45], we impose the following standard IV assumptions on the
90 data generating process.

91 **Assumption 1** (Standard IV assumptions [45, 47]). We assume: (1) *Exclusion*: $Y(a, z) = Y(a)$ for
92 all $a, z \in \{0, 1\}$, i.e., the instrument has no direct effect on the patient outcome; (2) *Independence*:
93 $Z \perp\!\!\!\perp U \mid X$; (3) *Relevance*: $Z \not\perp\!\!\!\perp A \mid X$, (iv) *The model includes all A–Y confounder*: $Y(a) \perp\!\!\!\perp$
94 $(A, Z) \mid (X, U)$ for all $a \in \{0, 1\}$.

95 **Assumption 1 is standard for IV methods and fulfilled in practical settings where IV methods**
96 **are applied [2, 4, 19].** Note that Assumption 1 does not prohibit the existence of unobserved Z –
97 A confounders. On the contrary, it merely prohibits the existence of unobserved confounders
98 that affect all Z , A , and Y simultaneously, as it is standard in IV settings [47]. A practical and
99 widespread example where Assumption 1 is satisfied are randomized controlled trials (RCTs) with
100 non-compliance [19]. Here, the treatment assignment Z is randomized, but the actual relationship
101 between treatment A and outcome Y may still be confounded. For instance, in the *Oregon health*
102 *insurance experiment* [16], people were given access to health insurance (Z) by a lottery with aim
103 to study the effect of health insurance (A) on health outcome (Y) [16]. Here, non-compliance
104 information is observed because the lottery winners needed to sign up for health insurance.

105 **Objective:** In this paper, we are interested in estimating the *individual treatment effect* (ITE)

$$\tau(x) = \mathbb{E}[Y(1) - Y(0) \mid X = x]. \quad (1)$$

106 If there is no unobserved confounding ($U = \emptyset$), the ITE is identifiable from observational data under
107 mild positivity assumptions [36]. However, in practice, it is often unlikely that all confounders are
108 observable. To account for this, we leverage the instrument Z to identify the ITE. We state the
109 following assumption for identifiability.

110 **Assumption 2** (Identifiability of the ITE [45]). **At least one of the following two statements holds**
111 **true:** (1) $\mathbb{E}[A \mid Z = 1, X, U] - \mathbb{E}[A \mid Z = 0, X, U] = \mathbb{E}[A \mid Z = 1, X] - \mathbb{E}[A \mid Z = 0, X]$; or
112 (2) $\mathbb{E}[Y(1) - Y(0) \mid X, U] = \mathbb{E}[Y(1) - Y(0) \mid X]$.

113 **Example:** Assumption 1 holds is when the function $f(a, X, U) = \mathbb{E}[Y(a) \mid X, U]$ is additive with
114 respect to a and U , e.g., $f(a, X, U) = g(a, X) + h(U)$ for measurable functions h and g .

115 Under Assumptions 1 and 2, the ITE is identifiable [45]. It can be written as

$$\tau(x) = \frac{\mu_1^Y(x) - \mu_0^Y(x)}{\mu_1^A(x) - \mu_0^A(x)} = \frac{\delta_Y(x)}{\delta_A(x)}, \quad (2)$$

116 where $\mu_i^Y(x) = \mathbb{E}[Y \mid Z = i, X = x]$ and $\mu_i^A(x) = \mathbb{E}[A \mid Z = i, X = x]$. Even if Assumption 2
117 does not hold, the quantity on the right-hand side of Eq. (2) still allows for interpretation. If no
118 unobserved Z – A confounders exist, it can be interpreted as conditional version of the *local average*
119 *treatment effect* (LATE) [19, 5] under a monotonicity assumption. Furthermore, under a no-current-
120 treatment-value-interaction assumption, it can be interpreted as conditional *treatment effect on the*
121 *treated* (ETT) [45].³ This has an important implication for our results: If Assumption 2 does not
122 hold in practice, our estimates still provide conditional LATE or ETT estimates under the respective
123 assumptions because they are based on Eq. (2). If Assumption 2 does hold, all three – i.e., ITE,
124 conditional LATE, and ETT – coincide [45].

125 3 Related work

126 **ITE methods without unconfoundedness:** Various machine learning methods for estimating ITEs
127 *without* unobserved confounding have been proposed in recent literature [1, 15, 25, 27, 36, 42, 52,

³The conditional LATE measures the ITE for individuals which are part of the complier subpopulation, i.e., the subpopulation for whom $A(Z = 1) > A(Z = 0)$. The conditional ETT measures the ITE for treated individuals.

53]. To remove plug-in bias, the DR-learner performs a second stage regression on the uncentered influence function of the average treatment effect [22, 14]. However, under unobserved confounding, all of these methods are biased (see Appendix). As a result, this hampers their performance in our setting.

ITE methods for unobserved confounding: There is a rich literature for causal effect estimation under unobserved confounding. Methods include deconfounding methods [46, 7, 18], proxy learning methods [13, 49], causal sensitivity analysis [21, 20], and IV methods. IV methods address the problem of unobserved confounding by exploiting the variance in treatment and outcome induced by the instruments. Traditionally, two-stage least squares (2SLS) has been used for estimating causal effects [48, 4]. 2SLS was originally developed in economics, and follows a two-stage procedure: it performs a first stage regression of treatment A on the instrument Z , and then uses the fitted values for a second stage regression to predict the outcome Y . Several nonparametric methods have been developed in econometric to generalize 2SLS in order to account for non-linearities within the data [28, 44], yet these are limited to low-dimensional settings.

Only recently, machine learning has been integrated into IV methods. These are: [37] and [50] generalize 2SLS by learning complex feature maps using kernel methods and deep learning, respectively. [17] adopts a two-stage neural network architecture that performs the first stage via conditional density estimation. [6] and [40] leverage moment conditions for IV estimation. However, the aforementioned methods are not specifically designed for the binary IV setting but, rather, for multiple IVs or treatment scenarios. In particular, they impose stronger assumptions such as additive confounding in order to identify the ITE. Note that additive confounding is a special case of when our Assumption 2 holds. Moreover, they are not multiply robust: Even though doubly robust IV methods have been proposed (e.g., Syrgkanis et al. [40]), these methods are not consistent in the union of more than two model specifications [45]. We provide more details below.

Doubly robust IV methods: Doubly robust estimators are commonly used in causal inference as they allow for consistent estimation under model misspecification and fast convergence rates [22]. Recently, they also have been adopted for IV settings: [23] proposes a pseudo regression estimator for the local average treatment effect using continuous instruments, which has been extended to individual effects by [35]. Furthermore, [38] uses a doubly robust approach to estimate average compiler parameters. Finally, Ogburn et al. [29] and Syrgkanis et al. [40] propose doubly robust ITE estimators in the IV setting which both rely on doubly robust parametrizations of the uncentered efficient influence function [30]. However, these estimators are not multiply robust in the sense that they are consistent in the union of more than two model specifications [45].

Multiply robust IV methods: Multiply robust estimators for IV settings have been proposed only for average treatment effects (ATEs) [45] and optimal treatment regimes [12] but not for ITEs. In particular, Wang et al. [45] derive a multiply robust parametrization of the efficient influence function for the ATE. However, there exists no similar approach for ITE estimation (see Table 1).

We provide a detailed, technical comparison of existing methods and our framework in Appendix E.

Binary IVs: In the binary IV setting, current methods proceed by estimating $\mu_i^Y(x)$ and $\mu_i^A(x)$ separately, before plugging them in Eq. 2 [19, 3, 5]. As result, these suffer from plug-in bias and do *not* offer robustness properties.

Research gap: To the best of our knowledge, there exists no method for ITE estimation under unobserved confounding that is *multiply robust*. To fill this gap, we propose MRIV: a *multiply robust* machine learning framework tailored to the binary IV setting. For this, we build upon the approach by Kennedy [22] to derive robust convergence rates, yet this approach has not been adapted to IV settings, which is our contribution.

4 MRIV for estimating ITEs using binary instruments

In the following, we present our MRIV framework for estimating ITEs under unobserved confounding (Sec. 4.1). We then derive an asymptotic convergence rate for MRIV (Sec. 4.2) and finally use our framework to develop a tailored deep neural network called MRIV-Net (Sec. 4.4).

Table 1: Key methods for causal effect estimation with IVs. This paper: Multiply robustness for ITEs.

Robustness	Estimand	ATE	ITE
Doubly robust		Okui et al. [30]	Syrgkanis et al. [40]
Multiply robust		Wang et al. [45]	MRIV (ours)

181 **4.1 Framework**

182 **Motivation:** A naïve approach to estimate the ITE is to leverage the identification result in Eq. (2).
 183 Assuming that we have estimated the nuisance components $\hat{\mu}_i^Y$ and $\hat{\mu}_i^A$ for $i \in \{0, 1\}$, we can simply
 184 plug them into Eq. (2) to obtain the so-called (plug-in) Wald estimator $\hat{\tau}_W(x)$ [43].

185 However, in practice, the true ITE curve $\tau(x)$ is often simpler (e.g., smoother, more sparse) than
 186 its complements $\mu_i^Y(x)$ or $\mu_i^A(x)$ [25]. In this case, $\hat{\tau}_W(x)$ is inefficient because it models all
 187 components separately, and, to address this, our proposed framework estimates τ directly using a
 188 pseudo outcome regression.

189 **Overview:** We now propose MRIV. MRIV is a two-stage meta learner that takes any base method for
 190 ITE estimation as input. For instance, the base ssssmethod could be the Wald estimator from Eq. (2),
 191 any other IV method such as 2SLS, or a deep neural network (as we propose in our MRIV-Net later
 192 in Sec. 4.4). In Stage 1, MRIV produces nuisance estimators $\hat{\mu}_0^Y(x)$, $\hat{\mu}_0^A(x)$, $\hat{\delta}_A(x)$, and $\hat{\pi}(x)$, where
 193 $\hat{\pi}(x)$ is an estimator of the propensity score $\pi(x) = \mathbb{P}(Z = 1 | X = x)$. In Stage 2, MRIV estimates
 194 $\tau(x)$ directly using a pseudo outcome \hat{Y}_{MR} as a regression target.

195 Given an arbitrary initial ITE estimator $\hat{\tau}_{init}(x)$ and nuisance estimates $\hat{\mu}_0^Y(x)$, $\hat{\mu}_0^A(x)$, $\hat{\delta}_A(x)$, and
 196 $\hat{\pi}(x)$, we define the pseudo outcome

$$\hat{Y}_{MR} = \left(\frac{Z - (1 - Z)}{\hat{\delta}_A(X)} \right) \left(\frac{Y - (\hat{\mu}_0^Y(X) + \hat{\tau}_{init}(X)(A - \hat{\mu}_0^A(X)))}{Z \hat{\pi}(X) + (1 - Z)(1 - \hat{\pi}(X))} \right) + \hat{\tau}_{init}(X). \quad (3)$$

197

198 The pseudo outcome \hat{Y}_{MR} in Eq. (3) is a multiply robust parameterization of the (uncentered) efficient
 199 influence function for the average treatment effect $\mathbb{E}_X[\tau(X)]$ (see the derivation in [45]). The initial
 200 estimator $\hat{\tau}_{init}(X)$ is corrected by a weighted difference of the observed outcome Y and the term
 201 $\hat{\mu}_0^Y(X) + \hat{\tau}_{init}(X)(A - \hat{\mu}_0^A(X))$. Individuals X with small $\hat{\delta}_A(X)$ (large estimated compliance) or
 202 small/large $\hat{\pi}(X)$ (i.e., low/high probability of receiving treatment Z) receive a larger correction.

203 Once we have obtained the pseudo outcome \hat{Y}_{MR} , we regress it on X to obtain the Stage 2 MRIV
 204 estimator $\hat{\tau}_{MRIV}(x)$ for $\tau(x)$. The pseudocode for MRIV is given in Algorithm 1. MRIV can be
 205 interpreted as a way to remove plug-in bias from $\hat{\tau}_{init}(x)$ via the efficient influence function [14]

Algorithm 1: MRIV

Input: data (X, Z, A, Y) , initial ITE estimator $\hat{\tau}_{init}(x)$
 // Stage 1: Estimate nuisance components
 $\hat{\pi}(x) \leftarrow \mathbb{E}[Z | X = x]$, $\hat{\mu}_0^Y(x) \leftarrow \mathbb{E}[Y | X = x, Z = 0]$, $\hat{\mu}_0^A(x) \leftarrow \mathbb{E}[A | X = x, Z = 0]$
 $\hat{\delta}_A(x) \leftarrow \mathbb{E}[A | X = x, Z = 1] - \mathbb{E}[A | X = x, Z = 0]$
 // Stage 2: Pseudo outcome regression
 $\hat{Y}_{MR} \leftarrow \left(\frac{Z - (1 - Z)}{\hat{\delta}_A(X)} \right) \left(\frac{Y - A \hat{\tau}_{init}(X) - \hat{\mu}_0^Y(X) + \hat{\mu}_0^A(X) \hat{\tau}_{init}(X)}{Z \hat{\pi}(X) + (1 - Z)(1 - \hat{\pi}(X))} \right) + \hat{\tau}_{init}(X)$
 $\hat{\tau}_{MRIV}(x) \leftarrow \mathbb{E}[\hat{Y}_{MR} | X = x]$

207 Using the fact that \hat{Y}_{MR} is a multiply robust parametrization of the efficient influence function, we
 208 derive a multiply robustness property of $\hat{\tau}_{MRIV}(x)$.

209 **Theorem 1** (multiply robustness property). *Let $\hat{\mu}_0^Y(x)$, $\hat{\mu}_0^A(x)$, $\hat{\delta}_A(x)$, $\hat{\pi}(x)$, and $\hat{\tau}_{init}(x)$ denote*
 210 *estimators of $\mu_0^Y(x)$, $\mu_0^A(x)$, $\delta_A(x)$, $\pi(x)$, and $\tau(x)$, respectively. Then, for all $x \in \mathcal{X}$, it holds that*
 211 *$\mathbb{E}[\hat{Y}_{MR} | X = x] = \tau(x)$, if least one of the following conditions is satisfied: (1) $\hat{\mu}_0^Y = \mu_0^Y$, $\hat{\mu}_0^A = \mu_0^A$,*
 212 *$\hat{\delta}_A = \delta_A$, and $\hat{\tau}_{init} = \tau$; or (2) $\hat{\pi} = \pi$ and $\hat{\delta}_A = \delta_A$; or (3) $\hat{\pi} = \pi$ and $\hat{\tau}_{init} = \tau$.*

213 **Theorem 1** implies that $\hat{\tau}_{MRIV}(x)$ is consistent for $\tau(x)$ if either condition (1), (2), or (3) holds.
 214 As a result, our MRIV framework is *multiply robust* in the sense that our estimator, $\hat{\tau}_{MRIV}(x)$, is
 215 consistent in the union of three different model specifications. Importantly, this is different from
 216 *doubly robust* estimators which are only consistent in the union of two model specifications [45].

217 **Example:** We illustrate the robustness under model specification (2) in an example. Let $\hat{\mu}_0^Y(x) =$
 218 $\hat{\mu}_0^A(x) = \hat{\tau}_{init}(x) = 0$ be misspecified and let $\hat{\pi} = \pi$ and $\hat{\delta}_A = \delta_A$ be correctly specified. It
 219 follows $\mathbb{E}[\hat{Y}_{MR} | X = x] = \frac{1}{\delta_A(X)} \mathbb{E} \left[\frac{ZY - (1 - Z)Y}{Z\pi(x) + (1 - Z)(1 - \pi(x))} | X = x \right] = \frac{\mu_1^Y(x) - \mu_0^Y(x)}{\delta_A(X)} = \tau(x)$. This
 220 justifies the pseudo-outcome regression in last step of MRIV.

221 Our MRIV is directly applicable to RCTs with non-compliance: Then, the treatment assignment is
 222 randomized and the propensity score $\pi(x)$ is known. Our MRIV framework can be thus adopted
 223 by plugging in the known $\pi(x)$ into the pseudo outcome in Eq. (3). Moreover, $\hat{\tau}_{\text{MRIV}}(x)$ is already
 224 consistent if either $\hat{\tau}_{\text{init}}(x)$ or $\hat{\delta}_A(x)$ are.

225 4.2 Theoretical analysis

226 In the following, we derive the asymptotic convergence rate of MRIV under smoothness assumptions.
 227 For this, we define s -smooth functions as functions contained in the Hölder class $\mathcal{H}(s)$, associated
 228 with Stone’s minimax rate [39] of $n^{-2s/(2s+p)}$, where p is the dimension of \mathcal{X} .

229 **Assumption 3** (Smoothness). We assume that (1) the nuisance components $\mu_i^Y(\cdot)$ are α -smooth,
 230 $\mu_i^A(\cdot)$ and $\delta_A(\cdot)$ are β -smooth, and $\pi(\cdot)$ is δ -smooth; (2) all nuisance components are estimated with
 231 their respective minimax rate of $n^{-\frac{2k}{2k+p}}$, where $k \in \{\alpha, \beta, \delta\}$; and (3) the oracle ITE $\tau(\cdot)$ is γ -smooth
 232 and the initial ITE estimator $\hat{\tau}_{\text{init}}$ converges with rate $r_\tau(n)$.

233 Assumption 3 for smoothness provides us with a way to quantify the difficulty of the underlying
 234 nonparametric regression problems. Similar assumptions have been imposed for asymptotic analysis
 235 of previous ITE estimators in [22, 15]. They can be replaced with other assumptions such as
 236 assumptions on the level of sparsity of the ITE components. We also provide an asymptotic analysis
 237 under sparsity assumptions (see Appendix B).

238 We additionally impose the following boundedness assumptions on the the underlying data generating
 239 process and estimators.

240 **Assumption 4** (Boundedness). We assume that there exist constants $C, \rho, \tilde{\rho}, \epsilon, K > 0$ such that for
 241 all $x \in \mathcal{X}$ it holds that: (1) $|\mu_i^Y(x)| \leq C$; (2) $|\delta_A(x)| = |\mu_1^A(x) - \mu_0^A(x)| \geq \rho$ and $|\hat{\delta}_A(x)| \geq \tilde{\rho}$;
 242 (3) $\epsilon \leq \hat{\pi}(x) \leq 1 - \epsilon$; and (4) $|\hat{\tau}_{\text{init}}(x)| \leq K$.

243 Assumptions 4.1, 4.3, and 4.4 are standard and in line with previous works on theoretical analyses
 244 of ITE estimators [15, 22]. Assumption 4.2 ensures that both the oracle ITE and the estimator are
 245 bounded. Violations of Assumption 4.2 may occur when working with so-called “weak” instruments,
 246 which are IVs that are only weakly correlated with the treatment. Using IV methods with weak
 247 instruments should generally be avoided [26]. However, in many applications such as RCTs with
 248 non-compliance, weak instruments are unlikely to occur as patients’ decisions to follow the treatment
 249 are generally correlated with the initial treatment assignments.

250 We state now our main theoretical result: an upper bound on the oracle risk of the MRIV estimator.
 251 **To derive our bound, we leverage the sample splitting approach from [22]. The approach in [22] has**
 252 **been initially used to analyze the DR-learner for ITE estimation under unconfoundedness and allows**
 253 **for the derivation of robust convergence rates. It has later been adapted to several other meta learners**
 254 **[15], yet not for IV methods.**

255 **Theorem 2** (Oracle upper bound under sample splitting). *Let \mathcal{D}_ℓ for $\ell \in \{1, 2, 3\}$ be independent*
 256 *samples of size n . Let $\hat{\tau}_{\text{init}}(x)$, $\hat{\mu}_0^Y(x)$, and $\hat{\mu}_0^A(x)$ be trained on \mathcal{D}_1 , and let $\hat{\delta}_A(x)$ and $\hat{\pi}(x)$ be*
 257 *trained on \mathcal{D}_2 . We denote \hat{Y}_{MR} as the pseudo outcome from Eq. (3) and Y_0 as the corresponding*
 258 *oracle. Let $\hat{\tau}_{\text{MRIV}}(x) = \hat{\mathbb{E}}_n[\hat{Y}_{\text{MR}} | X = x]$ and $\tilde{\tau}_{\text{MRIV}}(x) = \hat{\mathbb{E}}_n[Y_0 | X = x]$ denote the (oracle)*
 259 *pseudo outcome regression on \mathcal{D}_3 for some generic estimator $\hat{\mathbb{E}}_n[\cdot | X = x]$ of $\mathbb{E}[\cdot | X = x]$.*

260 *We assume that the second-stage estimator $\hat{\mathbb{E}}_n$ yields the minimax rate $n^{-\frac{2\gamma}{2\gamma+p}}$ and satisfies the fol-*
 261 *lowing two assumptions from Kennedy [22]: (1) $\hat{\mathbb{E}}_n[W + c | X = x] = \hat{\mathbb{E}}_n[W | X = x] + c$*
 262 *for any random W and constant c and (2) if $\mathbb{E}[W | X = x] = \mathbb{E}[V | X = x]$, then*
 263 $\mathbb{E} \left[\left(\hat{\mathbb{E}}_n[W | X = x] - \mathbb{E}[W | X = x] \right)^2 \right] \asymp \mathbb{E} \left[\left(\hat{\mathbb{E}}_n[V | X = x] - \mathbb{E}[V | X = x] \right)^2 \right]$. *Then, the*
 264 *oracle risk is upper bounded by*

$$\mathbb{E} \left[\left(\hat{\tau}_{\text{MRIV}}(x) - \tau(x) \right)^2 \right] \lesssim n^{-\frac{2\gamma}{2\gamma+p}} + r_\tau(n) \left(n^{-\frac{2\beta}{2\beta+p}} + n^{-\frac{2\delta}{2\delta+p}} \right) + n^{-2\left(\frac{\alpha}{2\alpha+p} + \frac{\delta}{2\delta+p}\right)} + n^{-2\left(\frac{\beta}{2\beta+p} + \frac{\delta}{2\delta+p}\right)}.$$

265 *Proof.* See Appendix A. □

266 Recall that the first summand of the lower bound in Eq. (2) is the minimax rate for the oracle ITE
 267 $\tau(x)$ which cannot be improved upon. Hence, for a fast convergence rate of $\hat{\tau}_{\text{MRIV}}(x)$, it is sufficient
 268 if either: (1) $r_\tau(n)$ decreases fast and δ is large; (2) $r_\tau(n)$ decreases fast and α and β are large;
 269 or (3) all α , β , and δ are large. This is in line with the multiply robustness property of MRIV and
 270 means that MRIV achieves a fast rate of convergence even if the initial estimator or several nuisance
 271 estimators converge slowly.

272 From the bound in Eq. (2), it follows that $\hat{\tau}_{\text{MRIV}}(x)$ improves on the convergence rate of the initial
 273 ITE estimator $\hat{\tau}_{\text{init}}(x)$ if its rate $r_\tau(n)$ is lower bounded by

$$r_\tau(n) \gtrsim n^{\frac{-2\gamma}{2\gamma+p}} + n^{-2\left(\frac{\alpha}{2\alpha+p} + \frac{\delta}{2\delta+p}\right)} + n^{-2\left(\frac{\beta}{2\beta+p} + \frac{\delta}{2\delta+p}\right)}. \quad (4)$$

274 Hence, our MRIV estimator is more likely to improve on the initial estimator for large α , β , and δ ,
 275 i.e., if the nuisance components are smooth. Note that it is sufficient if either (1) *only* the propensity
 276 score $\pi(x)$ is relatively smooth (large δ) or (2) that *all* other nuisance components are (large α and
 277 β). In fact, this is widely fulfilled in practice. For example, the former is fulfilled for RCTs with
 278 non-compliance, where $\pi(x)$ is often some known, fixed number $p \in (0, 1)$. Hence, for RCTs with
 279 non-compliance, MRIV should (at least asymptotically) improve the performance of most estimators.

280 4.3 MRIV vs. Wald estimator

281 In the following, we compare $\hat{\tau}_{\text{MRIV}}(x)$ to the Wald estimator $\hat{\tau}_{\text{W}}(x)$. First, we derive corresponding
 282 upper bound under smoothness.

283 **Theorem 3** (Wald oracle upper bound). *Given estimators $\hat{\mu}_i^Y(x)$ and $\hat{\mu}_i^A(x)$. Let $\hat{\delta}_A(x) = \hat{\mu}_1^A(x) -$
 284 $\hat{\mu}_0^A(x)$ satisfy Assumption 4. Then, the oracle risk of the Wald estimator $\hat{\tau}_{\text{W}}(x)$ is bounded by*

$$\mathbb{E} [(\hat{\tau}_{\text{W}}(x) - \tau(x))^2] \lesssim n^{-\frac{2\alpha}{2\alpha+p}} + n^{-\frac{2\beta}{2\beta+p}}. \quad (5)$$

285 *Proof.* See Appendix A. □

286 We now consider the MRIV estimator $\hat{\tau}_{\text{MRIV}}(x)$ with $\hat{\tau}_{\text{init}} = \hat{\tau}_{\text{W}}(x)$, i.e., initialized with the Wald
 287 estimator (under sample splitting). Plugging the Wald rate from Eq. (5) into the Eq. (2) yields

$$\mathbb{E} [(\hat{\tau}_{\text{MRIV}}(x) - \tau(x))^2] \lesssim n^{\frac{-2\gamma}{2\gamma+p}} + n^{\frac{-4\beta}{2\beta+p}} + n^{-2\left(\frac{\alpha}{2\alpha+p} + \frac{\beta}{2\beta+p}\right)} + n^{-2\left(\frac{\delta}{2\delta+p} + \frac{\alpha}{2\alpha+p}\right)} + n^{-2\left(\frac{\delta}{2\delta+p} + \frac{\beta}{2\beta+p}\right)}. \quad (6)$$

288 For $\alpha = \beta = \delta$, the rates of $\hat{\tau}_{\text{MRIV}}(x)$ and $\hat{\tau}_{\text{W}}(x)$ reduce to

$$\mathbb{E} [(\hat{\tau}_{\text{MRIV}}(x) - \tau(x))^2] \lesssim n^{\frac{-2\gamma}{2\gamma+p}} + n^{\frac{-4\alpha}{2\alpha+p}} \quad \text{and} \quad \mathbb{E} [(\hat{\tau}_{\text{W}}(x) - \tau(x))^2] \lesssim n^{\frac{-2\alpha}{2\alpha+p}}. \quad (7)$$

289 Hence, $\hat{\tau}_{\text{MRIV}}(x)$ outperforms $\hat{\tau}_{\text{W}}(x)$ asymptotically for $\gamma > \alpha$, i.e., when the ITE $\tau(x)$ is smoother
 290 than its components, which is usually the case in practice [25]. For $\gamma = \alpha$, the rates of both estimators
 291 coincide. Hence, we should expect MRIV to improve on the Wald estimator in real-world settings
 292 with sufficiently large sample size.

293 4.4 MRIV-Net

294 Based on our MRIV framework, we develop a tailored deep neural network called MRIV-Net for ITE
 295 estimation using IVs. Our MRIV-Net produces both an initial ITE estimator $\hat{\tau}_{\text{init}}(x)$ and nuisance
 296 estimators $\hat{\mu}_0^Y(x)$, $\hat{\mu}_0^A(x)$, $\hat{\delta}_A(x)$, and $\hat{\pi}(x)$.

297 For MRIV-Net, we choose deep neural networks for the nuisance components due to their predictive
 298 power and their ability to learn complex shared representations for several nuisance components.
 299 Sharing representations between nuisance components has been exploited previously for ITE estima-
 300 tion, yet only under unconfoundedness [36, 15]. Building shared representations is more efficient in
 301 finite sample regimes than estimating all nuisance components separately as they usually share some
 302 common structure.

303 In MRIV-Net, not all nuisance components should share a representation. Recall that, in
 304 Theorem 2, we assumed that (1) $\hat{\tau}_{\text{init}}(x)$, $\hat{\mu}_0^Y(x)$, and $\hat{\mu}_0^A(x)$; and (2) $\hat{\delta}_A(x)$ and $\hat{\pi}(x)$
 305 are trained on two independent samples in order to derive the upper bound on the oracle
 306 risk. Hence, we propose to build two separate representations Φ_1 and Φ_2 , so that (i) Φ_1
 307 is used to learn $\hat{\tau}_{\text{init}}(x)$, $\hat{\mu}_0^Y(x)$, and $\hat{\mu}_0^A(x)$, and (ii) Φ_2 is used to learn $\hat{\delta}_A(x)$ and $\hat{\pi}(x)$.

308 This ensures that the nuisance estimators (1) share minimal information
 309 with nuisance estimators (2) even though they are estimated on the same
 310 data. Intuitively, this should lead to a faster decay of the oracle upper
 311 bound (cf. [15]).

312 The architecture of MRIV-Net is shown in Fig. 2. MRIV-Net takes the
 313 observed covariates X as input to build the two representations Φ_1 and
 314 Φ_2 . The first representation Φ_1 is used to output estimates $\hat{\mu}_1^Y(x)$, $\hat{\mu}_0^Y(x)$,
 315 $\hat{\mu}_1^A(x)$, and $\hat{\mu}_0^A(x)$ of the ITE components. The second representation
 316 Φ_2 is used to output estimates $\tilde{\mu}_1^A(x)$, $\tilde{\mu}_0^A(x)$, and $\hat{\pi}(x)$. MRIV-Net is
 317 trained by minimizing an overall loss

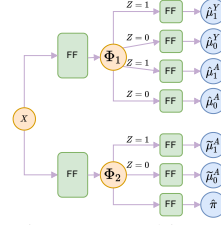


Figure 2: Architecture of MRIV-Net.

$$\mathcal{L}(\theta) = \sum_{i=1}^n \left[(\hat{\mu}_{z_i}^Y(x_i) - y_i)^2 + \text{BCE}(\hat{\mu}_{z_i}^A(x_i), a_i) + \text{BCE}(\tilde{\mu}_{z_i}^A(x_i), a_i) + \text{BCE}(\hat{\pi}(x_i), z_i) \right], \quad (8)$$

318 where θ denotes the neural network parameters and BCE is the binary cross entropy loss. After
 319 training MRIV-Net, we obtain the $\hat{\tau}_{\text{init}}(x) = \frac{\hat{\mu}_1^Y(x) - \hat{\mu}_0^Y(x)}{\hat{\mu}_1^A(x) - \hat{\mu}_0^A(x)}$ and obtain the nuisance estimators $\hat{\mu}_0^Y(x)$,
 320 $\hat{\mu}_0^A(x)$, $\hat{\delta}_A(x) = \tilde{\mu}_1^A(x) - \tilde{\mu}_0^A(x)$ and $\hat{\pi}(x)$. Then, we perform, we perform the pseudo regression
 321 (Stage 2) of MRIV to obtain $\hat{\tau}_{\text{MRIV}}(x)$.

322 **Implementation:** We use PyTorch Lightning for our implementation and train MRIV-Net with
 323 the Adam optimizer [24]. Details on the network architecture and hyperparameter tuning are in
 324 Appendix G. We perform both the training of MRIV-Net and the pseudo outcome regression on
 325 the full training data. Needless to say, MRIV-Net can be easily adopted for sample splitting or
 326 cross-fitting procedures as in [10], namely, by learning separate networks for each representation
 327 Φ_1 and Φ_2 . However, in our experiments, we do not use sample splitting or cross-fitting, as this can
 328 affect the performance in finite sample regimes. Of note, our choice is consistent with previous work
 329 [15].

330 5 Computational experiments

331 5.1 Simulated data

332 In causal inference literature, it is common practice to use simulated data for performance evaluations
 333 [8, 15, 17]. Simulated data offers the crucial benefit that it provides ground-truth information on the
 334 counterfactual outcomes and thus allows for direct benchmarking against the oracle ITE.

335 **Data generation:** We generate simulated data by sampling the oracle ITE $\tau(x)$ and the nuisance
 336 components $\mu_i^Y(x)$, $\mu_i^A(x)$, and $\pi(x)$ from Gaussian process priors. Using Gaussian processes has
 337 the following advantages: (1) It allows for a fair method comparison, as there is no need to explicitly
 338 specify the nuisance components, which could lead to unwanted inductive biases favoring a specific
 339 method; (2) the sampled nuisance components are non-linear and thus resemble real-world scenarios
 340 where machine learning methods would be applied; and, (3) by sampling from the prior induced by
 341 the Matérn kernel [32], we can control the smoothness of the nuisance components, which allows
 342 us to confirm our theoretical results from Sec. 4.2. For a detailed description of our data generating
 343 process, we refer to Appendix C.

344 **Baselines:** We compare our MRIV-Net with the following state-of-the-art baselines: (1) ITE methods
 345 for unconfoundedness: **TARNet** [36] and TARNet combined with the **DR-learner** [22]; (2) general
 346 IV methods: **2SLS** [48], kernel IV (**KIV**) [37], **DFIV** [50], **DeepIV** [17], **DeepGMM** [6], **DMLIV**
 347 [40], and DMLIV combined with **DRIV** (as described in [40]); (3) the (plug-in) Wald estimator using
 348 **linear models** and Bayesian additive regression trees (**BART**) [11]. Of note, the DR-learner assumes
 349 unconfoundedness, which is why we only combine it TARNet in our experiments. Implementation
 350 details regarding baselines and nuisance parameter estimation are in Appendix E. **Note that many of**
 351 **the baselines do not directly aim at ITE estimation but rather at counterfactual outcome prediction.**
 352 **We nevertheless use these methods as baselines and, for this, obtain the ITE by taking the difference**
 353 **between the predictions of the factual and counterfactual outcomes.**

354 **Performance evaluation:** For all experiments, we use a 80/20 split as training/test set. We calculate
 355 the root mean squared errors (RMSE) between the ITE estimates and the oracle ITE on the test set.

356 We report the mean RMSE and the standard deviation over five data sets generated from random
 357 seeds.

358 **Results:** Table 2 shows the results for
 359 all baselines. Here, the DR-learner does
 360 not improve the performance of TAR-
 361 Net, which is reasonable as both the
 362 DR-learner and TARNet assume uncon-
 363 foundedness and are thus biased in our
 364 setting. Our MRIV-Net outperforms all
 365 baselines. Our MRIV-Net also achieves
 366 a smaller standard deviation. For addi-
 367 tional results, we refer to Appendix H.

368 We further compare the performance of
 369 two different meta-learner frameworks –
 370 DRIV [40] and our MRIV– across differ-
 371 ent base methods. The nuisance param-
 372 eters are estimated using feed forward neural networks (DRIV) or TARNets with either binary or
 373 continuous outputs (MRIV). The results are in Table 3. Our MRIV improves over the variant without
 374 any meta-learner framework across all base methods (both in terms of RMSE and standard deviation).

Table 3: Base model with different meta-learners (i.e., none, DRIV, and our MRIV).

Base methods	n = 3000			n = 5000			n = 8000		
	None	DRIV	MRIV (ours)	None	DRIV	MRIV (ours)	None	DRIV	MRIV (ours)
(1) STANDARD ITE									
TARNet [36]	0.76 ± 0.14	0.31 ± 0.05	0.34 ± 0.13	0.70 ± 0.12	0.17 ± 0.06	0.17 ± 0.05	0.69 ± 0.17	0.21 ± 0.04	0.16 ± 0.04
(2) GENERAL IV									
2SLS [47]	1.22 ± 0.23	0.40 ± 0.11	0.31 ± 0.08	0.79 ± 0.37	0.17 ± 0.09	0.19 ± 0.05	1.12 ± 0.29	0.21 ± 0.05	0.16 ± 0.02
KIV [37]	1.54 ± 0.53	0.40 ± 0.10	0.39 ± 0.11	1.18 ± 1.14	0.20 ± 0.08	0.17 ± 0.06	3.80 ± 4.71	0.31 ± 0.18	0.28 ± 0.19
DFIV [50]	0.43 ± 0.11	0.26 ± 0.05	0.27 ± 0.07	0.40 ± 0.21	0.18 ± 0.09	0.16 ± 0.04	0.46 ± 0.54	0.21 ± 0.06	0.18 ± 0.05
DeepIV [17]	0.96 ± 0.30	0.27 ± 0.03	0.26 ± 0.05	0.28 ± 0.09	0.18 ± 0.08	0.18 ± 0.05	0.23 ± 0.04	0.21 ± 0.03	0.16 ± 0.03
DeepGMM [6]	0.95 ± 0.38	0.40 ± 0.15	0.36 ± 0.13	0.37 ± 0.09	0.24 ± 0.12	0.16 ± 0.05	0.42 ± 0.14	0.21 ± 0.03	0.17 ± 0.03
DMLIV [40]	1.92 ± 0.71	0.41 ± 0.12	0.37 ± 0.11	0.92 ± 0.41	0.22 ± 0.05	0.16 ± 0.05	1.14 ± 0.24	0.21 ± 0.06	0.18 ± 0.05
(3) WALD ESTIMATOR [43]									
Linear	1.06 ± 0.63	0.42 ± 0.15	0.38 ± 0.14	0.62 ± 0.22	0.19 ± 0.09	0.25 ± 0.09	0.81 ± 0.34	0.19 ± 0.09	0.18 ± 0.04
BART	0.95 ± 0.30	0.48 ± 0.14	0.46 ± 0.12	0.63 ± 0.33	0.26 ± 0.13	0.20 ± 0.07	0.88 ± 0.28	0.31 ± 0.08	0.29 ± 0.04
MRIV-Netw network only (ours)	0.39 ± 0.13	0.35 ± 0.12	0.26 ± 0.11	0.31 ± 0.04	0.19 ± 0.13	0.15 ± 0.03	0.26 ± 0.06	0.18 ± 0.08	0.13 ± 0.03

Reported: RMSE (mean ± standard deviation). Lower = better (best improvement over none meta-learner in bold)

375 Furthermore, MRIV is clearly superior
 376 over DRIV. This demonstrates the effec-
 377 tiveness of our MRIV across different
 378 base methods (note: MRIV with an ar-
 379 bitrary base model is typically superior
 380 to DRIV with our custom network from
 381 above). MRIV-Net is overall best. We
 382 also performed additional experiments where we used cross-fitting approaches for both meta-learners
 383 (see Appendix I).

384 **Ablation study:** Table 4 compares different variants of our MRIV-Net. These are: (1) MRIV but
 385 network only; (2) MRIV-Net with a single representation for all nuisance estimators; and (3) our
 386 MRIV-Net from above. We observe that MRIV-Net is best. This justifies our proposed network
 387 architecture for MRIV-Net. Hence, combing the result from above, our performance gain must be
 388 attributed to both our framework and the architecture of our deep neural network.

389 **Robustness checks for
 390 unobserved confounding
 391 and smoothness:** Here, we
 392 demonstrate the importance
 393 of handling unobserved
 394 confounding (as we do in our
 395 MRIV framework). For this,
 396 Fig. 3 plots the results for
 397 our MRIV-Net vs. standard
 398 ITE without customization

399 for confounding (i.e., TARNet with and without the DR-learner) over over different levels of
 400 unobserved confounding. The RMSE of both TARNet variants increase almost linearly with

Table 2: Performance comparison: our MRIV-Net vs. existing baselines.

Method	n = 3000	n = 5000	n = 8000
(1) STANDARD ITE			
TARNet [36]	0.76 ± 0.14	0.70 ± 0.12	0.69 ± 0.17
TARNet + DR [36, 22]	0.78 ± 0.10	0.66 ± 0.09	0.70 ± 0.10
(2) GENERAL IV			
2SLS [47]	1.22 ± 0.23	0.79 ± 0.37	1.12 ± 0.29
KIV [37]	1.54 ± 0.53	1.18 ± 1.14	3.80 ± 4.71
DFIV [50]	0.43 ± 0.11	0.40 ± 0.21	0.46 ± 0.54
DeepIV [17]	0.96 ± 0.30	0.28 ± 0.09	0.23 ± 0.04
DeepGMM [6]	0.95 ± 0.38	0.37 ± 0.09	0.42 ± 0.14
DMLIV [40]	1.92 ± 0.71	0.92 ± 0.41	1.14 ± 0.24
DMLIV + DRIV [40]	0.41 ± 0.12	0.22 ± 0.04	0.21 ± 0.06
(3) WALD ESTIMATOR [43]			
Linear	1.06 ± 0.63	0.62 ± 0.22	0.81 ± 0.34
BART	0.95 ± 0.30	0.63 ± 0.33	0.88 ± 0.28
MRIV-Net (ours)	0.26 ± 0.11	0.15 ± 0.03	0.13 ± 0.03

Reported: RMSE for base methods (mean ± standard deviation). Lower = better (best in bold)

Table 4: Ablation study.

Method	n = 3000	n = 5000	n = 8000
MRIV-Netw network only	0.39 ± 0.13	0.31 ± 0.04	0.26 ± 0.06
MRIV-Netw single repr.	0.28 ± 0.12	0.21 ± 0.04	0.32 ± 0.10
MRIV-Net (ours)	0.26 ± 0.11	0.15 ± 0.03	0.13 ± 0.03

Reported: RMSE (mean ± standard deviation). Lower = better (best in bold)

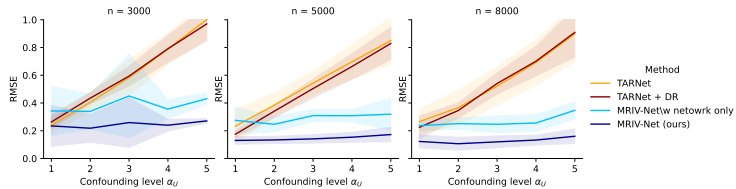


Figure 3: Results over different levels of confounding α_U . Shaded area shows standard deviation.

401 increasing confounding. In contrast, the RMSE of our MRIV-Net only marginally. Even for low
 402 confounding regimes, our MRIV-Net performs competitively.

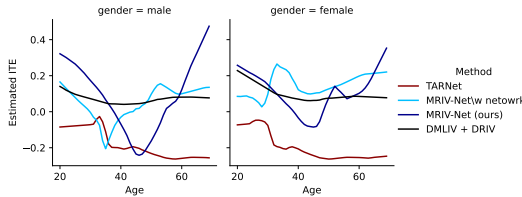
403 Fig. 4 varies the smoothness level. This is given by α of $\mu_i^Y(\cdot)$ (controlled by the Matérn kernel
 404 prior). Here, the performance decreases for the baselines, i.e., DeepIV and our network without
 405 MRIV framework. In contrast, the performance of our MRIV-Net remains robust and outperforms
 406 the baselines. This confirms our theoretical results from above. It thus indicates that our MRIV
 407 framework works best when the oracle ITE $\tau(x)$ is smoother than the baseline methods.

408 5.2 Case study with real-world data

409 **Setting:** We demonstrate effectiveness of our framework using
 410 a case study with real-world, medical data. Here, we use medi-
 411 cal data from the so-called *Oregon health insurance experiment*
 412 (OHIE) [16]. It provides data for an RCT with non-compliance:
 413 In 2008, $\sim 30,000$ low-income, uninsured adults in Oregon were
 414 offered participation in a health insurance program by a lottery.
 415 Individuals whose names were drawn could decide to sign up
 416 for health insurance. After a period of 12 months, in-person
 417 interviews took place to evaluate the health condition of the
 418 respective participant.
 419

420 In our analysis, the lottery assignment is the instrument Z , the decision to sign up for health insurance
 421 is treatment A , and an overall health score is the outcome Y . We also include five covariates X (age,
 422 gender, language, the number of emergency visits before the experiment, and the number of people
 423 the individual signed up with). It is important to include the latter in our analysis as it is the only
 424 variable influencing the propensity score. For details, we refer to Appendix D. We first estimate the
 425 ITE function and then report the treatment effect heterogeneity w.r.t. age and gender, while fixing
 426 the other covariates (i.e., we consider the English-speaking subpopulation with one emergency visit
 427 that signed up alone). We repeat the same procedure for our neural network architecture without the
 428 MRIV-Net framework and TARNet. The results are in Fig. 5.

429 **Results:** Our MRIV-Net estimates larger causal effects for an older age. In
 430 contrast, TARNet does not estimate positive ITEs even for an older age.



431 Figure 5: Results on real-world medical data.

432 Even though we cannot evaluate the estimation
 433 quality on real-world data, our estimates seem
 434 reasonable in light of the medical literature: the
 435 benefit of health insurance should increase with
 436 older age. This showcases that TARNet may
 437 suffer from bias induced by unobserved con-
 438 founders. We also report the results for DRIV
 439 with DMLIV as base method, and observe that
 440 in contrast to MRIV-Net, the corresponding ITE
 441 does not vary much between ages. Interestingly,
 442 both our MRIV-Net estimate a somewhat smaller ITE for middle ages (around 30–50 yrs). One
 443 explanation might be that individual in this age group are more likely to have stable jobs and, thus, are
 444 also more likely to be able to afford medical care, decreasing the direct effect of health insurance on
 445 individuals health. In sum, the findings from our case study are of direct relevance for decision-makers
 in public health [19], and highlight the practical value of our framework.

446 6 Conclusion

447 In this paper, we propose MRIV-Net: a novel ITE estimator based on a deep neural network.
 448 Importantly, our estimator is consistent in the union of three models specifications and, therefore,
 449 multiply robust. This is a crucial difference to existing works: previously, existing ITE estimators
 450 (such es DRIV from Syrgkanis et al. [40]) were only doubly robust. We show both theoretically and
 451 empirically that MRIV-Net is state-of-the-art for estimating ITEs using binary IVs. For future work,
 452 it would be interesting to derive finite sample results for MRIV-Net, because our theoretical analysis
 453 is purely asymptotic. Furthermore, one could develop multiply robust estimators for other IV settings
 454 (e.g., multiple or continuous instruments and treatments).

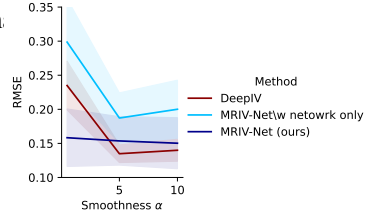


Figure 4: Results over different levels of smoothness α of $\mu_i^Y(\cdot)$, sample size $n = 8000$. Larger $\alpha =$ smoother. Shaded areas show standard deviation.

455 References

- 456 [1] Ahmed M. Alaa and Mihaela van der Schaar. “Bayesian inference of individualized treatment
457 effects using multi-task Gaussian processes”. In: *NeurIPS*. 2017.
- 458 [2] Joshua D. Angrist. “Lifetime earnings and the vietnam era draft lotter: Evidence from social
459 security administrative records”. In: *The American Economic Review* 80.3 (1990), pp. 313–336.
- 460 [3] Joshua D. Angrist, Guido W. Imbens, and Donald B. Rubin. “Identification of causal effects
461 using instrumental variables”. In: *Journal of the American Statistical Association* 91.434
462 (1996), pp. 444–455.
- 463 [4] Joshua D. Angrist and Alan B. Krueger. “Does compulsory school attendance affect schooling
464 and earnings?” In: *The Quarterly Journal of Economics* 106.4 (1991), pp. 979–1014.
- 465 [5] Falco J. Bargagli-Stoffi, Kristof de Witte, and Giorgio Gnecco. “Heterogeneous causal effects
466 with imperfect compliance: A Bayesian machine learning approach”. In: *Annals of Applied
467 Statistics* (2021).
- 468 [6] Andrew Bennett, Nathan Kallus, and Tobias Schnabel. “Deep generalized method of moments
469 for instrumental variable analysis”. In: *NeurIPS*. 2019.
- 470 [7] Ioana Bica, Ahmed M. Alaa, and Mihaela van der Schaar. “Time series deconfounder: Esti-
471 mating treatment effects over time in the presence of hidden confounders”. In: *ICML*. 2020.
- 472 [8] Ioana Bica et al. “Estimating counterfactual treatment outcomes over time through adversarially
473 balanced representations”. In: *ICLR*. 2020.
- 474 [9] Howard S. Bloom et al. “The benefits and costs of JTPA title II-A programs: Key Findings
475 from the National Job Training Partnership Act Study”. In: *Journal of Human Resources* 32.32
476 (1997), pp. 549–586.
- 477 [10] Victor Chernozhukov et al. “Double/debiased machine learning for treatment and structural
478 parameters”. In: *The Econometrics Journal* 21.1 (2018), pp. C1–C68.
- 479 [11] Hugh A. Chipman, Edward I. George, and Robert E. McCulloch. “BART: Bayesian additive
480 regression trees”. In: *The Annals of Applied Statistics* 4.1 (2010), pp. 266–298.
- 481 [12] Yifan Cui and Eric Tchetgen Tchetgen. “A semiparametric instrumental variable approach
482 to optimal treatment regimes under endogeneity”. In: *Journal of the American Statistical
483 Association* 116.553 (2021), pp. 126–137.
- 484 [13] Yifan Cui et al. “Semiparametric proximal causal inference”. In: *arXiv preprint* (2020).
- 485 [14] Alicia Curth, Ahmed M. Alaa, and Mihaela van der Schaar. “Estimating structural target
486 functions using machine learning and influence functions”. In: *arXiv preprint* (2020).
- 487 [15] Alicia Curth and Mihaela van der Schaar. “Nonparametric estimation of heterogeneous treat-
488 ment effects: From theory to learning Algorithms”. In: *AISTATS*. 2021.
- 489 [16] Amy Finkelstein et al. “The oregon health insurance experiment: Evidence from the first year”.
490 In: *The Quarterly Journal of Economics* 127.3 (2012), pp. 1057–1106.
- 491 [17] Jason Hartford et al. “Deep IV: A flexible approach for counterfactual prediction”. In: *ICML*.
492 2017.
- 493 [18] Tobias Hatt and Stefan Feuerriegel. “Sequential deconfounding for causal inference with
494 unobserved confounders”. In: *arXiv preprint* (2021).
- 495 [19] Guido W. Imbens and Joshua D. Angrist. “Identification and estimation of local average
496 treatment effects”. In: *Econometrica* 62.2 (1994), pp. 467–475.
- 497 [20] Andrew Jesson et al. “Quantifying ignorance in individual-level causal-effect estimates under
498 hidden confounding”. In: *ICML*. 2021.
- 499 [21] Nathan Kallus, Xiaojie Mao, and Angela Zhou. “Interval estimation of individual-level causal
500 effects under unobserved confounding”. In: *AISTATS*. 2019.
- 501 [22] Edward H. Kennedy. “Optimal doubly robust estimation of heterogeneous causal effects”. In:
502 *arXiv preprint* (2020).
- 503 [23] Edward H. Kennedy, Scott A. Lorch, and Dylan S. Small. “Robust causal inference with
504 continuous instruments using the local instrumental variable curve”. In: *Journal of the Royal
505 Statistical Society: Series B* 81.1 (2019), pp. 121–143.
- 506 [24] Diederik P. Kingma and Jimmy Ba. “Adam: A method for stochastic optimization”. In: *ICLR*.
507 2015.

- 508 [25] Sören R. Künzel et al. “Metalearners for estimating heterogeneous treatment effects using
509 machine learning”. In: *Proceedings of the National Academy of Sciences (PNAS)* 116.10
510 (2019), pp. 4156–4165.
- 511 [26] Chunxiao Li, Cynthia Rudin, and Tyler H. McCormick. “Rethinking nonlinear instrumental
512 variable models through prediction validity”. In: *Journal of Machine Learning Research* 23
513 (2022), pp. 1–55.
- 514 [27] Bryan Lim, Ahmed M. Alaa, and Mihaela van der Schaar. “Forecasting treatment responses
515 over time using recurrent marginal structural networks”. In: *NeurIPS*. 2018.
- 516 [28] Whitney K. Newey and James L. Powell. “Instrumental variable estimation of nonparametric
517 models”. In: *Econometrica* 71.5 (2003), pp. 1565–1578.
- 518 [29] Elizabeth L. Ogburn, Andrea Rotnitzky, and James M. Robins. “Doubly robust estimation of
519 the local average treatment effect curve”. In: *Journal of the Royal Statistical Society: Series B*
520 77.2 (2015), pp. 373–396.
- 521 [30] Ryo Okui et al. “Doubly robust instrumental variable regression”. In: *Statistica Sinica* 22.1
522 (2012), pp. 173–205.
- 523 [31] Judea Pearl. *Causality*. New York City: Cambridge University Press, 2009.
- 524 [32] Carl Edward Rasmussen and Christopher K. I. Williams. *Gaussian processes for machine
525 learning*. 3. print. Adaptive computation and machine learning. Cambridge, Mass.: MIT Press,
526 2008.
- 527 [33] James M. Robins, Miguel A. Hernán, and Babette Brumback. “Marginal structural models and
528 causal inference in epidemiology”. In: *Epidemiology* 11.5 (2000), pp. 550–560.
- 529 [34] Donald B. Rubin. “Estimating causal effects of treatments in randomized and nonrandomized
530 studies”. In: *Journal of Educational Psychology* 66.5 (1974), pp. 688–701.
- 531 [35] Vira Semenova and Victor Chernozhukov. “Debiased machine learning of conditional average
532 treatment effects and other causal functions”. In: *The Econometrics Journal* 24.2 (2021),
533 pp. 264–289.
- 534 [36] Uri Shalit, Fredrik D. Johansson, and David Sontag. “Estimating individual treatment effect:
535 Generalization bounds and algorithms”. In: *ICML*. 2017.
- 536 [37] Rahul Singh, Maneesh Sahani, and Arthur Gretton. “Kernel instrumental variable regression”.
537 In: *NeurIPS*. 2019.
- 538 [38] Rahul Singh and Liyang Sun. “Double robustness for complier parameters and a semiparamet-
539 ric test for complier characteristics”. In: *arXiv preprint* ().
- 540 [39] Charles J. Stone. “Optimal rates of convergence for nonparametric estimators”. In: *Annals of
541 Statistics* 8.6 (1980).
- 542 [40] Vasilis Syrgkanis et al. “Machine learning estimation of heterogeneous treatment effects with
543 instruments”. In: *NeurIPS*. 2019.
- 544 [41] Hal R. Varian. “Causal inference in economics and marketing”. In: *Proceedings of the National
545 Academy of Sciences (PNAS)* 113.27 (2016), pp. 7310–7315.
- 546 [42] Stefan Wager and Susan Athey. “Estimation and inference of heterogeneous treatment effects
547 using random forests”. In: *Journal of the American Statistical Association* 113.523 (2018),
548 pp. 1228–1242.
- 549 [43] Abraham Wald. “The fitting of straight lines if both variables are subject to error”. In: *Annals
550 of Mathematical Statistics* 11.3 (1940), pp. 284–300.
- 551 [44] Guihua Wang, Jun Li, and Wallace J. Hopp. “An instrumental variable forest approach for
552 detecting heterogeneous treatment effects in observational studies”. In: *Management Science*
553 (2021).
- 554 [45] Linbo Wang and Eric J. Tchetgen Tchetgen. “Bounded, efficient and multiply robust estimation
555 of average treatment effects using instrumental variables”. In: *Journal of the Royal Statistical
556 Society: Series B* 80.3 (2018), pp. 531–550.
- 557 [46] Yixin Wang and David M. Blei. “The blessings of multiple causes”. In: *Journal of the American
558 Statistical Association* 114.528 (2019), pp. 1574–1596.
- 559 [47] Jeffrey M. Wooldridge. *Introductory Econometrics: A modern approach*. Routledge, 2013.
- 560 [48] Phillip G. Wright. *The tariff on animal and vegetable oils*. New York: Macmillan, 1928.
- 561 [49] Liyuan Xu, Heishiro Kanagawa, and Arthur Gretton. “Deep proxy causal learning and its
562 application to confounded bandid policy evaluation”. In: *NeurIPS*. 2021.

- 563 [50] Liyuan Xu et al. “Learning deep features in instrumental variable regression”. In: *ICLR*. 2021.
- 564 [51] Azam M. Yazdani and Eric Boerwinkle. “Causal inference in the age of decision medicine”.
565 In: *Journal of Data Mining in Genomics & Proteomics* 6.1 (2015).
- 566 [52] Jinsung Yoon, James Jordon, and Mihaela van der Schaar. “GANITE: Estimation of individu-
567 alized treatment effects using generative adversarial nets”. In: *ICLR*. 2018.
- 568 [53] Yao Zhang, Alexis Bellot, and Mihaela van der Schaar. “Learning overlapping representations
569 for the estimation of individualized treatment effects”. In: *AISTATS*. 2020.

570 Checklist

- 571 1. For all authors...
- 572 (a) Do the main claims made in the abstract and introduction accurately reflect the paper’s
573 contributions and scope? [Yes] See here
- 574 (b) Did you describe the limitations of your work? [Yes] See Sec. 5.2.
- 575 (c) Did you discuss any potential negative societal impacts of your work? [No]
- 576 (d) Have you read the ethics review guidelines and ensured that your paper conforms to
577 them? [Yes]
- 578 2. If you are including theoretical results...
- 579 (a) Did you state the full set of assumptions of all theoretical results? [Yes] See Sec. 4.2
580 and Appendix A..
- 581 (b) Did you include complete proofs of all theoretical results? [Yes] See Appendix A.
- 582 3. If you ran experiments...
- 583 (a) Did you include the code, data, and instructions needed to reproduce the main experi-
584 mental results (either in the supplemental material or as a URL)? [Yes] (both).
- 585 (b) Did you specify all the training details (e.g., data splits, hyperparameters, how they
586 were chosen)? [Yes] See Appendix E and G
- 587 (c) Did you report error bars (e.g., with respect to the random seed after running ex-
588 periments multiple times)? [Yes] One standard deviation using 5 random seeds, see
589 Sec. 5.
- 590 (d) Did you include the total amount of compute and the type of resources used (e.g., type
591 of GPUs, internal cluster, or cloud provider)? [Yes] See Appendix E.
- 592 4. If you are using existing assets (e.g., code, data, models) or curating/releasing new assets...
- 593 (a) If your work uses existing assets, did you cite the creators? [Yes] See Sec. 5.2.
- 594 (b) Did you mention the license of the assets? [No] No license is provided on the OHIE
595 website.
- 596 (c) Did you include any new assets either in the supplemental material or as a URL? [Yes]
597 See Appendix D.
- 598 (d) Did you discuss whether and how consent was obtained from people whose data you’re
599 using/curating? [No] We only used simulated and publicly available data.
- 600 (e) Did you discuss whether the data you are using/curating contains personally identifiable
601 information or offensive content? [No] The OHIE data is publicly available, personally
602 identifiable information are censored.
- 603 5. If you used crowdsourcing or conducted research with human subjects...
- 604 (a) Did you include the full text of instructions given to participants and screenshots, if
605 applicable? [No] Not applicable.
- 606 (b) Did you describe any potential participant risks, with links to Institutional Review
607 Board (IRB) approvals, if applicable? [No] Not applicable.
- 608 (c) Did you include the estimated hourly wage paid to participants and the total amount
609 spent on participant compensation? [No] Not applicable.

Estimating individual treatment effects under unobserved confounding using binary instruments

Appendix

Anonymous Author(s)

Affiliation

Address

email

1	Contents	
2	A Proofs	2
3	A.1 Proof of Theorem 1 (multiple robustness property)	2
4	A.2 Proof of Theorem 2 (Convergence rate of MRIV)	3
5	A.3 Proof of Theorem 3 (Convergence rate of the Wald estimator)	4
6	B Theoretical analysis under sparsity assumptions	5
7	C Simulated data	6
8	D Oregon health insurance experiment	8
9	E Details for baseline methods	9
10	E.1 ITE methods for unconfoundedness	9
11	E.2 General IV methods	9
12	E.3 Wald estimator	11
13	F Visualization of predicted ITEs	12
14	G Implementation details and hyperparameter tuning	13
15	H Results for semi-synthetic data	15
16	I Results for cross-fitting	16

17 **A Proofs**

18 We start by deriving an auxiliary Lemma. That is, we derive an explicit expression for the Stage 2
 19 oracle pseudo outcome regression $\mathbb{E}[\hat{Y}_0 | X = x]$ of MRIV.

Lemma 4.

$$\begin{aligned} & \mathbb{E}[\hat{Y}_0 | X = x] \\ &= \frac{\pi(x)}{\hat{\delta}_A(x)\hat{\pi}(x)} (\mu_1^Y(x) - \mu_1^A(x) \hat{\tau}_{\text{init}}(x)) + \frac{(1 - \pi(x))}{\hat{\delta}_A(x)(1 - \hat{\pi}(x))} (\mu_0^A(x) \hat{\tau}_{\text{init}}(x) - \mu_0^Y(x)) \\ & \quad + \frac{\hat{\mu}_0^A(x) \hat{\tau}_{\text{init}}(x) - \hat{\mu}_0^Y(x)}{\hat{\delta}_A(x)} \left(\frac{\pi(x)}{\hat{\pi}(x)} - \frac{1 - \pi(x)}{1 - \hat{\pi}(x)} \right) + \hat{\tau}_{\text{init}}(x) \end{aligned} \quad (1)$$

Proof.

$$\begin{aligned} & \mathbb{E}[\hat{Y}_0 | X = x] \\ &= \pi(x) \mathbb{E} \left[\frac{Y - A \hat{\tau}_{\text{init}}(X) - \hat{\mu}_0^Y(X) + \hat{\mu}_0^A(X) \hat{\tau}_{\text{init}}(X)}{\hat{\delta}_A(X) \hat{\pi}(X)} \middle| X = x, Z = 1 \right] \\ & \quad + (1 - \pi(x)) \mathbb{E} \left[\frac{Y - A \hat{\tau}_{\text{init}}(X) - \hat{\mu}_0^Y(X) + \hat{\mu}_0^A(X) \hat{\tau}_{\text{init}}(X)}{\hat{\delta}_A(X) (1 - \hat{\pi}(X))} \middle| X = x, Z = 0 \right] + \hat{\tau}_{\text{init}}(x) \end{aligned} \quad (2)$$

$$\begin{aligned} &= \frac{\pi(x)}{\hat{\delta}_A(x) \hat{\pi}(x)} (\mu_1^Y(x) - \mu_1^A(x) \hat{\tau}_{\text{init}}(x) - \hat{\mu}_0^Y(x) + \hat{\mu}_0^A(x) \hat{\tau}_{\text{init}}(x)) \\ & \quad + \frac{1 - \pi(x)}{\hat{\delta}_A(x) (1 - \hat{\pi}(x))} (\mu_0^Y(x) - \mu_0^A(x) \hat{\tau}_{\text{init}}(x) - \hat{\mu}_0^Y(x) + \hat{\mu}_0^A(x) \hat{\tau}_{\text{init}}(x)) + \hat{\tau}_{\text{init}}(x) \end{aligned} \quad (4)$$

20 Rearranging the terms yields the desired result. \square

21 **A.1 Proof of Theorem 1 (multiple robustness property)**

22 We use Lemma 4 to show that under each of the three conditions it follows that $\mathbb{E}[\hat{Y}_0 | X = x] = \tau(x)$.

1.

$$\begin{aligned} & \mathbb{E}[\hat{Y}_0 | X = x] \\ &= \frac{\pi(x)}{\delta_A(x) \hat{\pi}(x)} (\mu_1^Y(x) - \mu_1^A(x) \tau(x) + \mu_0^A(x) \tau(x) - \mu_0^Y(x)) \end{aligned} \quad (5)$$

$$+ \frac{(1 - \pi(x))}{\delta_A(x) (1 - \hat{\pi}(x))} (\mu_0^A(x) \tau(x) - \mu_0^Y(x) - \mu_0^A(x) \tau(x) + \mu_0^Y(x)) + \tau(x) \quad (6)$$

$$= \frac{\pi(x)}{\delta_A(x) \hat{\pi}(x)} (\delta_Y(x) - \delta_Y(x)) + \tau(x) = \tau(x). \quad (7)$$

2.

$$\mathbb{E}[\hat{Y}_0 | X = x] = \frac{(\mu_1^Y(x) - \mu_1^A(x) \hat{\tau}_{\text{init}}(x))}{\delta_A(x)} + \frac{(\mu_0^A(x) \hat{\tau}_{\text{init}}(x) - \mu_0^Y(x))}{\delta_A(x)} + \hat{\tau}_{\text{init}}(x) \quad (8)$$

$$= \frac{\delta_Y(x) - \hat{\tau}_{\text{init}}(x) \delta_A(x)}{\delta_A(x)} + \hat{\tau}_{\text{init}}(x) = \tau(x). \quad (9)$$

3.

$$\mathbb{E}[\hat{Y}_0 | X = x] = \frac{(\mu_1^Y(x) - \mu_1^A(x) \tau(x))}{\hat{\delta}_A(x)} + \frac{(\mu_0^A(x) \tau(x) - \mu_0^Y(x))}{\hat{\delta}_A(x)} + \tau(x) \quad (10)$$

$$= \frac{\delta_Y(x)}{\hat{\delta}_A(x)} - \tau(x) \frac{\delta_A(x)}{\hat{\delta}_A(x)} + \tau(x) = \tau(x) \quad (11)$$

23 **A.2 Proof of Theorem 2 (Convergence rate of MRIV)**

24 To prove Theorem 2, we need an additional assumption on the second stage regression estimator $\hat{\mathbb{E}}_n$.
 25 We refer to Kennedy [8] (Theorem 1) for a detailed discussion on this assumption.

26 **Assumption 5** (From Theorem 1 of Kennedy [8]). The following two statements hold:

- 27 1. $\hat{\mathbb{E}}_n[W + c | X = x] = \hat{\mathbb{E}}_n[W | X = x] + c$ for any random W and constant c
 28 2. If $\mathbb{E}[W | X = x] = E[V | X = x]$ then

$$\mathbb{E} \left[\left(\hat{\mathbb{E}}_n[W | X = x] - \mathbb{E}[W | X = x] \right)^2 \right] \asymp \mathbb{E} \left[\left(\hat{\mathbb{E}}_n[V | X = x] - \mathbb{E}[V | X = x] \right)^2 \right]. \quad (12)$$

29 *Proof of Theorem 2.* Using Assumption 5, we can apply Theorem 1 of Kennedy [8] and obtain

$$\mathbb{E} \left[(\hat{\tau}_{\text{init}}(x) - \tau(x))^2 \right] \lesssim \mathcal{R}(x) + \mathbb{E} [\hat{r}(x)^2], \quad (13)$$

30 where $\mathcal{R}(x) = \mathbb{E} \left[(\tilde{\tau}_{MR}(x) - \tau(x))^2 \right]$ is the oracle risk of the second stage regression and $r(x) =$
 31 $\mathbb{E}[\hat{Y}_0 | X = x] - \tau(x)$. We can apply Lemma 4 to obtain

$$\begin{aligned} \hat{r}(x) &= \frac{\pi(x)}{\hat{\delta}_A(x) \hat{\pi}(x)} (\mu_1^Y(x) - \mu_1^A(x) \hat{\tau}_{\text{init}}(x)) + \frac{(1 - \pi(x))}{\hat{\delta}_A(x) (1 - \hat{\pi}(x))} (\mu_0^A(x) \hat{\tau}_{\text{init}}(x) - \mu_0^Y(x)) \\ &\quad + \frac{\hat{\mu}_0^A(x) \hat{\tau}_{\text{init}}(x) - \hat{\mu}_0^Y(x)}{\hat{\delta}_A(x)} \left(\frac{\pi(x)}{\hat{\pi}(x)} - \frac{1 - \pi(x)}{1 - \hat{\pi}(x)} \right) + \hat{\tau}_{\text{init}}(x) - \tau(x) \end{aligned} \quad (14)$$

$$\begin{aligned} &= \left(\frac{\mu_1^Y(x) - \mu_0^Y(x)}{\hat{\delta}_A(x)} \right) \frac{\pi(x)}{\hat{\pi}(x)} + \frac{\mu_0^Y(x) - \hat{\mu}_0^Y(x)}{\hat{\delta}_A(x)} \left(\frac{\pi(x)}{\hat{\pi}(x)} - \frac{1 - \pi(x)}{1 - \hat{\pi}(x)} \right) + (\hat{\tau}_{\text{init}}(x) - \tau(x)) \\ &\quad + \left(\frac{(\mu_0^A(x) - \mu_1^A(x)) \hat{\tau}_{\text{init}}(x)}{\hat{\delta}_A(x)} \right) \frac{\pi(x)}{\hat{\pi}(x)} + \frac{(\hat{\mu}_0^D(x) - \mu_0^D(x)) \hat{\tau}_{\text{init}}(x)}{\hat{\delta}_A(x)} \left(\frac{\pi(x)}{\hat{\pi}(x)} - \frac{1 - \pi(x)}{1 - \hat{\pi}(x)} \right) \end{aligned} \quad (15)$$

$$\begin{aligned} &= \frac{\delta_Y(x) \pi(x)}{\hat{\delta}_A(x) \hat{\pi}(x)} + \frac{(\mu_0^Y(x) - \hat{\mu}_0^Y(x)) (\pi(x) - \hat{\pi}(x))}{\hat{\delta}_A(x) \hat{\pi}(x) (1 - \hat{\pi}(x))} + (\hat{\tau}_{\text{init}}(x) - \tau(x)) \\ &\quad - \frac{\delta_A(x) \pi(x) \hat{\tau}_{\text{init}}(x)}{\hat{\delta}_A(x) \hat{\pi}(x)} + \frac{(\hat{\mu}_0^A(x) - \mu_0^A(x)) \hat{\tau}_{\text{init}}(x) (\pi(x) - \hat{\pi}(x))}{\hat{\delta}_A(x) \hat{\pi}(x) (1 - \hat{\pi}(x))} \end{aligned} \quad (16)$$

$$\begin{aligned} &= \frac{(\pi(x) - \hat{\pi}(x))}{\hat{\delta}_A(x) \hat{\pi}(x) (1 - \hat{\pi}(x))} [(\mu_0^Y(x) - \hat{\mu}_0^Y(x)) + (\hat{\mu}_0^A(x) - \mu_0^A(x)) \hat{\tau}_{\text{init}}(x)] \\ &\quad + (\hat{\tau}_{\text{init}}(x) - \tau(x)) + \frac{\pi(x) \delta_A(x)}{\hat{\pi}(x) \hat{\delta}_A(x)} (\tau(x) - \hat{\tau}_{\text{init}}(x)) \end{aligned} \quad (17)$$

$$\begin{aligned} &= \frac{(\pi(x) - \hat{\pi}(x))}{\hat{\delta}_A(x) \hat{\pi}(x) (1 - \hat{\pi}(x))} [(\mu_0^Y(x) - \hat{\mu}_0^Y(x)) + (\hat{\mu}_0^A(x) - \mu_0^A(x)) \hat{\tau}_{\text{init}}(x)] \\ &\quad + (\tau(x) - \hat{\tau}_{\text{init}}(x)) \left(\delta_A(x) - \hat{\delta}_A(x) \right) \pi(x) + (\tau(x) - \hat{\tau}_{\text{init}}(x)) (\pi(x) - \hat{\pi}(x)) \hat{\delta}_A(x). \end{aligned} \quad (18)$$

32 Applying the inequality $(a + b)^2 \leq 2(a^2 + b^2)$ together with Assumption 4 and the fact that $\pi(x) \leq 1$
 33 yields

$$\begin{aligned} \hat{r}(x)^2 &\leq \frac{4}{\epsilon^4 \rho^2} (\pi(x) - \hat{\pi}(x))^2 \left[(\mu_0^Y(x) - \hat{\mu}_0^Y(x))^2 + (\hat{\mu}_0^A(x) - \mu_0^A(x))^2 K^2 \right] \\ &\quad + 4 (\tau(x) - \hat{\tau}_{\text{init}}(x))^2 \left(\delta_A(x) - \hat{\delta}_A(x) \right)^2 + 4 (\tau(x) - \hat{\tau}_{\text{init}}(x))^2 (\pi(x) - \hat{\pi}(x))^2. \end{aligned} \quad (19)$$

34 By setting $\tilde{K} = \max\{K, 1\}$, we obtain

$$\begin{aligned} \hat{r}(x)^2 &\leq \frac{4\tilde{K}^2}{\epsilon^4 \rho^2} \left((\pi(x) - \hat{\pi}(x))^2 \left[(\mu_0^Y(x) - \hat{\mu}_0^Y(x))^2 + (\hat{\mu}_0^A(x) - \mu_0^A(x))^2 + (\hat{\tau}_{\text{init}}(x) - \tau(x))^2 \right] \right. \\ &\quad \left. + (\tau(x) - \hat{\tau}_{\text{init}}(x))^2 (\delta_A(x) - \hat{\delta}_A(x))^2 \right). \end{aligned} \quad (20)$$

35 Applying expectations on both sides yields

$$\mathbb{E} \left[(\hat{\tau}_{\text{init}}(x) - \tau(x))^2 \right] \quad (21)$$

$$\begin{aligned} &\lesssim \mathcal{R}(x) + \mathbb{E} \left[(\hat{\tau}_{\text{init}}(x) - \tau(x))^2 \right] \left(\mathbb{E} \left[(\hat{\delta}_A(x) - \delta_A(x))^2 \right] + \mathbb{E} \left[(\hat{\pi}(x) - \pi(x))^2 \right] \right) \\ &\quad + \mathbb{E} \left[(\hat{\pi}(x) - \pi(x))^2 \right] \left(\mathbb{E} \left[(\hat{\mu}_0^Y(x) - \mu_0^Y(x))^2 \right] + \mathbb{E} \left[(\hat{\mu}_0^A(x) - \mu_0^A(x))^2 \right] \right), \end{aligned} \quad (22)$$

36 because $(\hat{\pi}(x), \hat{\delta}_A(x)) \perp\!\!\!\perp (\hat{\mu}_0^Y(x), \hat{\mu}_0^A(x), \hat{\tau}_{\text{init}}(x))$ due to sample splitting. The claim follows now
37 by applying Assumption 3. \square

38 A.3 Proof of Theorem 3 (Convergence rate of the Wald estimator)

39 *Proof.* We define $\tilde{C} = \max\{C, 1\}$ and obtain the upper bound

$$(\hat{\tau}_W(x) - \tau(x))^2 \quad (23)$$

$$= \left(\frac{(\hat{\mu}_1^Y(x) - \mu_1^Y(x)) \delta_A(x) + (\mu_0^Y(x) - \hat{\mu}_0^Y(x)) \delta_A(x) + (\delta_A(x) - \hat{\delta}_A(x)) \delta_Y(x)}{\delta_A(x) \hat{\delta}_A(x)} \right)^2 \quad (24)$$

$$\leq \frac{4\tilde{C}^2}{\rho^2 \tilde{\rho}^2} \left[(\hat{\mu}_1^Y(x) - \mu_1^Y(x))^2 + (\hat{\mu}_0^Y(x) - \mu_0^Y(x))^2 + (\delta_A(x) - \hat{\delta}_A(x))^2 \right] \quad (25)$$

$$\begin{aligned} &\leq \frac{8\tilde{C}^2}{\rho^2 \tilde{\rho}^2} \left[(\hat{\mu}_1^Y(x) - \mu_1^Y(x))^2 + (\hat{\mu}_0^Y(x) - \mu_0^Y(x))^2 + (\hat{\mu}_1^A(x) - \mu_1^A(x))^2 \right. \\ &\quad \left. + (\hat{\mu}_0^A(x) - \mu_0^A(x))^2 \right], \end{aligned} \quad (26)$$

40 where we used the inequality $(a + b)^2 \leq 2(a^2 + b^2)$ several times. Taking expectations and applying
41 the smoothness assumptions yields the result. \square

42 B Theoretical analysis under sparsity assumptions

43 In Sec. 4.2, we analyzed MRIV theoretically by imposing smoothness assumptions on the underlying
 44 data generating process. In particular, we derived a multiple robust convergence rate and showed
 45 that MRIV outperforms the Wald estimator if the oracle ITE is smoother than its components. In
 46 this section, we derive similar results by relying on a different set of assumptions. Instead of using
 47 smoothness, we make assumptions on the level of sparsity of the ITE components. This assumption
 48 is often imposed in high-dimensional settings ($n < p$) and is in line with previous literature on
 49 analyzing ITE estimators [4, 8].

50 In the following, we say a function $f(x)$ is k -sparse, if it is linear in $x \in \mathbb{R}^p$ and it only depends
 51 on $k < \min\{n, p\}$ predictors. [22] showed, that in this case the minimax rate of $f(x)$ is given by
 52 $\frac{k \log(p)}{n}$. The linearity assumption can be relaxed to an additive structural assumption, which we omit
 53 here for simplicity. In the following, we replace the smoothness conditions in Assumption 3 with
 54 sparsity conditions.

55 **Assumption 6** (Sparsity). We assume that (1) the nuisance components $\mu_i^Y(\cdot)$ are α -sparse, $\mu_i^A(\cdot)$
 56 and $\delta_A(\cdot)$ are β -sparse, and $\pi(\cdot)$ is δ -sparse; (2) all nuisance components are estimated with their
 57 respective minimax rate of $\frac{k \log(p)}{n}$, where $k \in \{\alpha, \beta, \delta\}$; and (3) the oracle ITE $\tau(\cdot)$ is γ -sparse and
 58 the initial ITE estimator $\hat{\tau}_{\text{init}}$ converges with rate $r_\tau(n)$.

59 We restate now our result from Theorem 3 for MRIV using the sparsity assumption.

60 **Theorem 5** (MRIV upper bound under sparsity). *We consider the same setting as in Theorem 2*
 61 *under the sparsity assumption 6. If the second-stage estimator $\hat{\mathbb{E}}_n$ yields the minimax rate $\frac{\gamma \log(p)}{n}$*
 62 *and satisfies Assumption 5, the oracle risk is upper bounded by*

$$\mathbb{E} \left[(\hat{\tau}_{\text{MRIV}}(x) - \tau(x))^2 \right] \lesssim \frac{\gamma \log(p)}{n} + r_\tau(n) \frac{(\beta + \delta) \log(p)}{n} + \frac{(\alpha + \beta) \delta \log^2(p)}{n^2}.$$

63 *Proof.* Follows immediately from the proof of Theorem 2, i.e., from Eq.(21) by applying Ass- 6. \square

64 Again, we obtain a multiple robust convergence rate for MRIV in the sense that MRIV achieves a fast
 65 rate even if the initial estimator or several nuisance estimators converge slowly. More precisely, for a
 66 fast convergence rate of $\hat{\tau}_{\text{MRIV}}(x)$, it is sufficient if either: (1) $r_\tau(n)$ decreases fast and δ is small;
 67 (2) $r_\tau(n)$ decreases fast and α and β are small; or (3) all α , β , and δ are small.

68 We derive now the corresponding rate for the Wald estimator.

69 **Theorem 6** (Wald oracle upper bound). *Given estimators $\hat{\mu}_i^Y(x)$ and $\hat{\mu}_i^A(x)$. Let $\hat{\delta}_A(x) = \hat{\mu}_1^A(x) -$*
 70 *$\hat{\mu}_0^A(x)$ satisfy Assumption 4. Then, under Assumption 6 the oracle risk of the Wald estimator $\hat{\tau}_W(x)$*
 71 *is bounded by*

$$\mathbb{E} [(\hat{\tau}_W(x) - \tau(x))^2] \lesssim \frac{(\alpha + \beta) \log(p)}{n} \quad (27)$$

72 *Proof.* Follows immediately from the proof of Theorem 3, i.e., from Eq.(23) by applying Ass- 6. \square

73 If $\alpha = \beta = \delta$, we obtain the rates

$$\mathbb{E} \left[(\hat{\tau}_{\text{MRIV}}(x) - \tau(x))^2 \right] \lesssim \frac{\gamma \log(p)}{n} + \frac{\alpha^2 \log^2(p)}{n^2} \quad \text{and} \quad \mathbb{E} [(\hat{\tau}_W(x) - \tau(x))^2] \lesssim \frac{\alpha \log(p)}{n}, \quad (28)$$

74 which means that $\hat{\tau}_{\text{MRIV}}(x)$ outperforms $\hat{\tau}_W(x)$ for $\gamma < \alpha$, i.e., if the oracle ITE is more sparse than
 75 its components.

76 **C Simulated data**

77 In the following, we describe how we simulate synthetic data for the experiments in Sec. 5.1 from the
 78 main paper. As mentioned therein, we simulate the ITE components from Gaussian processes using
 79 the prior induced by the Matern kernel [12]

$$K_{\ell,\nu}(x_i, x_j) = \frac{1}{\Gamma(\nu)2^{\nu-1}} \left(\frac{\sqrt{2\nu}}{\ell} \|x_i - x_j\|_2 \right)^\nu K_\nu \left(\frac{\sqrt{2\nu}}{\ell} \|x_i - x_j\|_2 \right), \quad (29)$$

80 where $\Gamma(\cdot)$ is the Gamma function and $K_\nu(\cdot)$ is the modified Bessel function of second kind. Here, ℓ
 81 is the length scale of the kernel and ν controls the smoothness of the sampled functions.

82 We set $\ell = 1$ and sample functions $\delta_Y \sim \mathcal{GP}(0, K_{\ell,\gamma})$, $\mu_0^Y \sim \mathcal{GP}(0, K_{\ell,\alpha})$, $f_1 \sim \mathcal{GP}(0, K_{\ell,\beta})$,
 83 $f_0 \sim \mathcal{GP}(0, K_{\ell,\beta})$ and $g \sim \mathcal{GP}(0, K_{\ell,\beta})$. Then, we define $\mu_1^Y = \delta_Y + \mu_0^Y$, $\mu_1^A = 0.3 \cdot \sigma \circ f_1 + 0.7$,
 84 $\mu_0^A = 0.3 \cdot \sigma \circ f_0$, $\delta_A = \mu_1^A - \mu_0^A$, $\mu_0^Y = c_0 \delta_A$, and $\pi = \sigma \circ g$. Finally, we set the oracle ITE to

$$\tau = \frac{\mu_1^Y - \mu_0^Y}{\mu_1^A - \mu_0^A} = \frac{\delta_Y}{\delta_A}. \quad (30)$$

85 Note that we can create a setup where the ITE τ is smoother than its components by using a small
 86 α/β ratio. An example is shown in Fig. 1.

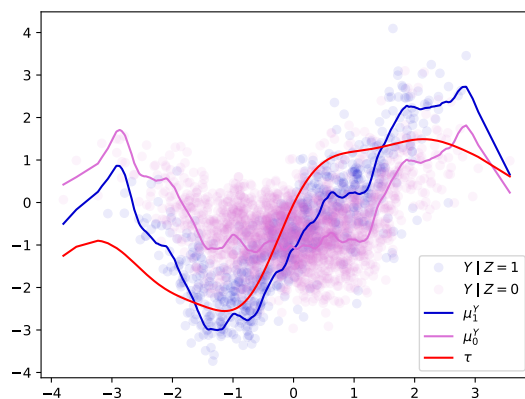


Figure 1: Gaussian process simulation for $\alpha = 1.5$ and $\beta = 50$.

87 In the following, we describe how we generate data the (X, Z, A, Y) using the ITE components
 88 $\mu_i^Y(x)$, $\mu_i^A(x)$, and $\pi(x)$. We begin by sampling n observed confounder $X \sim \mathcal{N}(0, 1)$, unobserved
 89 confounders $U \sim \mathcal{N}(0, 0.2^2)$, and instruments $Z \sim \text{Bernoulli}(\pi(X))$. Then, we obtain treatments
 90 via

$$A = Z \mathbb{1}\{U + \epsilon_A > \alpha_1(X)\} + (1 - Z) \mathbb{1}\{U + \epsilon_A > \alpha_0(X)\} \quad (31)$$

91 with indicator function $\mathbb{1}$, noise $\epsilon_A \sim \mathcal{N}(0, 0.1^2)$, and $\alpha_i(X) = \Phi^{-1}(1 - \mu_i^A(X)) \sqrt{0.1^2 + 0.2^2}$,
 92 where Φ^{-1} denotes the quantile function of the standard normal distribution. Finally, we generate the
 93 outcomes via

$$Y = A \left(\frac{(\mu_1^A(X) - 1)\mu_0^Y(X) - \mu_0^A(X)\mu_1^Y(X) + \mu_1^Y(X)}{\delta_A(X)} \right) \quad (32)$$

$$+ (1 - A) \left(\frac{\mu_1^A(X)\mu_0^Y(X) - \mu_0^A(X)\mu_1^Y(X)}{\delta_A(X)} \right) + \alpha_U U + \epsilon_Y, \quad (33)$$

94 where $\epsilon_Y \sim \mathcal{N}(0, 0.3^2)$ is noise and $\alpha_U > 0$ is a parameter indicating the level of unobserved
 95 confounding. This choice of A and Y in Eq. (31) and Eq. (32), respectively, implies that $\tau(x)$ is
 96 indeed the ITE, i. e., it holds that $\tau(x) = \mathbb{E}[Y(1) - Y(0) | X = x]$.

97 **Lemma 7.** *Let (X, Z, A, Y) be sampled from the the previously described procedure. Then, it holds*
 98 *that*

$$\mu_i^A(x) = \mathbb{E}[A \mid Z = i, X = x] \quad \text{and} \quad \mu_i^Y(x) = \mathbb{E}[Y \mid Z = i, X = x]. \quad (34)$$

99 *Proof.* The first claim follows from

$$\mathbb{E}[A \mid Z = i, X = x] = \mathbb{P}(U + \epsilon_A > \alpha_i(x)) = 1 - \Phi(\Phi^{-1}(1 - \mu_i^A(x))) = \mu_i^A(x), \quad (35)$$

100 because $U + \epsilon_A \sim \mathcal{N}(0, \sqrt{0.1^2 + 0.2^2})$. The second claim follows from

$$\mathbb{E}[Y \mid Z = i, X = x] = \mu_i^A(x) \left(\frac{(\mu_1^A(x) - 1)\mu_0^Y(x) - \mu_0^A(x)\mu_1^Y(x) + \mu_1^Y(x)}{\delta_A(x)} \right) \quad (36)$$

$$+ (1 - \mu_i^A(x)) \left(\frac{\mu_1^A(x)\mu_0^Y(x) - \mu_0^A(x)\mu_1^Y(x)}{\delta_A(x)} \right) \quad (37)$$

$$= \frac{\mu_i^Y(x)\delta_A(x)}{\delta_A(x)} = \mu_i^Y(x). \quad (38)$$

101

□

102 **D Oregon health insurance experiment**

103 The so-called *Oregon health insurance experiment*¹ (OHIE) [6] was an important RCT with non-
 104 compliance. It was intentionally conducted as large-scale effort among public health to assess the
 105 effect of health insurance on several outcomes such as health or economic status. In 2008, a lottery
 106 draw offered low-income, uninsured adults in Oregon participation in a Medicaid program, providing
 107 health insurance. Individuals whose names were drawn could decide to sign up for the program.

108 In our analysis, the lottery assignment is the instrument Z , the decision to sign up for the Medicaid
 109 program is the treatment A , and an overall health score is the outcome Y . The outcome was obtained
 110 after a period of 12 months during in-person interviews. We use the following covariates X : age,
 111 gender, language, the number of emergency visits before the experiment, and the number of people
 112 the individual signed up with. The latter is used to control for peer effects, and it is important to
 113 include this variable in our analysis as it is the only variable influencing the propensity score (see
 114 below). We extract $\sim 10,000$ observations from the OHIE data and plot the histograms of all variables
 115 in Fig. 2. We can clearly observe the presence of non-compliance within the data, because the
 116 ratio of treated / untreated individuals is much lower than the corresponding ratio for the treatment
 117 assignment.

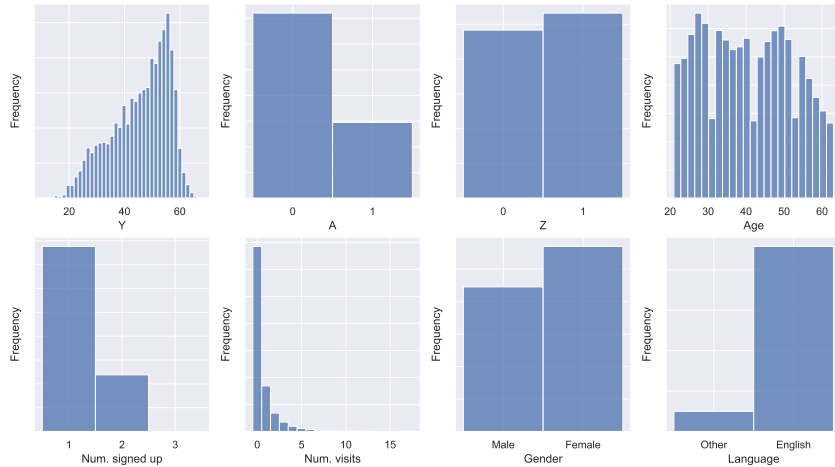


Figure 2: Histograms of each variable in our sample from OHIE.

118 The data collection in the OHIE was done follows: After excluding individuals below the age
 119 of 19, above the age of 64, and individuals with residence outside of Oregon, 74,922 individuals
 120 were considered for the lottery. Among those, 29,834 were selected randomly and were offered
 121 participation in the program. However, the probability of selection depended on the number of
 122 household members on the waiting list: for instance, an individual who signed up with another person
 123 was twice as likely to be selected. From the 74,922 individuals, 57,528 signed up alone, 17,236
 124 signed up with another person, and 158 signed up with two more people on the waiting list. Thus, the
 125 probability of being selected conditional on the number of household members on the waiting list
 126 follows the multivariate version of Wallenius’ noncentral hypergeometric distribution [2].

127 **Propensity score:** We computed the propensity score as follows. To account for the Wallenius’
 128 noncentral hypergeometric distribution, we use the R package *BiasedUrn* to calculate the propensity
 129 score $\pi(x) = \mathbb{P}(Z = 1 | X = x)$. We obtained

$$\pi(x) = \begin{cases} 0.345, & \text{if individual } x \text{ signed up alone,} \\ 0.571, & \text{if individual } x \text{ signed up with one more person,} \\ 0.719, & \text{if individual } x \text{ signed up with two more people.} \end{cases} \quad (39)$$

130 During the training of both MRIV and DRIV, we use the calculated values from Eq. (39) for the
 131 propensity score.

¹Data available here: <https://www.nber.org/programs-projects/projects-and-centers/oregon-health-insurance-experiment>

132 **E Details for baseline methods**

133 In this section, we give a brief overview on the baselines which we used in our experiments. We
 134 implemented: (1) ITE methods for unconfoundedness [8, 13]; (2) general IV methods, i.e., IV
 135 methods developed for IV settings with multiple or continuous instruments and treatments [1, 7, 14,
 136 15, 20, 21]; and (3) two instantiations of the Wald estimator for the binary IV setting [16].

137 **E.1 ITE methods for unconfoundedness**

138 Many ITE methods assume *unconfoundedness*, i.e., that all confounders are observed in the data.
 139 Formally, the unconfoundedness assumption can be expressed in the potential outcomes framework
 140 as

$$Y(1), Y(0) \perp\!\!\!\perp A \mid X. \quad (40)$$

141 Under unconfoundedness, the ITE is identified as

$$\tau(x) = \mu_1(x) - \mu_0(x) \quad \text{with} \quad \mu_i(x) = \mathbb{E}[Y \mid A = i, X = x]. \quad (41)$$

142 Methods that assume unconfoundedness proceed by estimating $\mu_i(x) = \mathbb{E}[Y \mid A = i, X = x]$ from
 143 Eq. (41). However, if unobserved confounders U exist, it follows that

$$\tau(x) = \mathbb{E}[Y \mid A = 1, X = x, U] - \mathbb{E}[Y \mid A = 0, X = x, U] \neq \mu_1(x) - \mu_0(x), \quad (42)$$

144 which means that estimators that assume unconfoundedness are generally biased. Nevertheless, we
 145 include two baselines that assume unconfoundedness into our experiments: TARNet [13] and the
 146 DR-learner [8].

147 **TARNet** [13]: TARNet [13] is a neural network that estimates the ITE components $\mu_i(x)$ from
 148 Eq. 41 by learning a shared representation $\Phi(x)$ and two potential outcome heads $h_i(\Phi(x))$. We train
 149 TARNet by minimizing the loss

$$\mathcal{L}(\theta) = \sum_{i=1}^n L(h_{a_i}(\Phi(x_i, \theta_\Phi), \theta_{h_i}), y_i), \quad (43)$$

150 where $\theta = (\theta_{h_1}, \theta_{h_0}, \theta_\Phi)$ denotes the model parameters and L denotes squared loss if Y is continuous
 151 or binary cross entropy loss if Y is binary.

152 *Note regarding balanced representations:* In [13], the authors propose to add an additional regular-
 153 ization term inspired from domain adaptation literature, which forces TARNet to learn a balanced
 154 representation $\Phi(x)$, i.e., that minimizes the distance the treatment and control group in the feature
 155 space. They showed that this approach leads to minimization of a generalization bound on the ITE
 156 estimation error if the representation is invertible.

157 In our experiments, we refrained from learning balanced representations because minimizing the
 158 regularized loss from [13] does not necessarily result in an invertible representation and thus may
 159 even harm the estimation performance. For a detailed discussion, we refer to [4]. Furthermore,
 160 by leaving out the regularization, we ensure comparability between the different baselines. If
 161 balanced representations are desired, the balanced representation approach could also be extended to
 162 MRIV-Net, as we also build MRIV-Net on learning shared representations.

163 **DR-learner** [8]: The DR-learner [8] is a meta learner that takes arbitrary estimators of the ITE
 164 components μ_i and the propensity score $\pi(x) = \mathbb{P}(A = 1 \mid X = x)$ as input and performs a pseudo
 165 outcome regression by using the pseudo outcome

$$\hat{Y}_0 = \left(\frac{A}{\hat{\pi}(X)} - \frac{1-A}{1-\hat{\pi}(X)} \right) Y + \left(1 - \frac{A}{\hat{\pi}(X)} \right) \hat{\mu}_1(X) - \left(1 - \frac{1-A}{1-\hat{\pi}(X)} \right) \hat{\mu}_0(X). \quad (44)$$

166 In our experiments, we use TARNet as base method to provide initial estimators $\hat{\mu}_i(X)$. We further
 167 learn propensity score estimates $\hat{\pi}(X)$ by adding a separate representation to TARNet as done in
 168 [13].

169 **E.2 General IV methods**

170 **2SLS** [20]: 2SLS [20] is a linear two-stage approach. First, the treatments A are regressed on the
 171 instruments Z and fitted values \hat{A} are obtained. In the second stage, the outcome Y is regressed on \hat{A} .
 172 We implement 2SLS using the scikit-learn package.

173 **KIV** [14]: Kernel IV [14] generalizes 2SLS to nonlinear settings. KIV assumes that the data is
 174 generated by

$$Y = f(A) + U, \quad (45)$$

175 where U is an additive unobserved confounder and f is some unknown (potentially nonlinear)
 176 structural function. KIV then models the structural function via

$$f(a) = \mu^t \psi(a) \quad \text{and} \quad \mathbb{E}[\psi(A) | Z = z] = V\phi(z), \quad (46)$$

177 where ψ and ϕ are feature maps. Here, kernel ridge regressions instead of linear regressions are used
 178 in both stages to estimate μ and V .

179 Following [14] we use the exponential kernel [12] and set the length scale to the median inter-point
 180 distance. KIV does not provide a direct way to incorporate the observed confounders X . Hence, we
 181 augment both the instrument and the treatment with X , which is consistent with previous work [1,
 182 21]. We also use two different samples for each stage as recommended in [14].

183 **DFIV** [21]: DFIV [21] is a similar approach KIV in generalizing 2SLS to nonlinear setting by
 184 assuming Eq. (45) and Eq. (46). However, instead of using kernel methods, DFIV models the features
 185 maps ψ_{θ_A} and ϕ_{θ_Z} as neural networks with parameters θ_A and θ_Z , respectively. DFIV is trained by
 186 iteratively updating the parameters θ_A and θ_Z . The authors also provide a training algorithm that
 187 incorporates observed confounders X , which we implemented for our experiments. During training,
 188 we used two different datasets for each of the two stages as described in in the paper.

189 **DeepIV** [7]: DeepIV [7] also assumes additive unobserved confounding as in Eq. (45), but leverages
 190 the identification result [10]

$$\mathbb{E}[Y | X = x, Z = z] = \int h(a, x) dF(a | x, z), \quad (47)$$

191 where $h(a, x) = f(a, x) + \mathbb{E}[U | X = x]$ is the target counterfactual prediction function. DeepIV
 192 estimates $F(a | x, z)$, i.e., the conditional distribution function of the treatment A given observed
 193 covariates X and instruments Z , by using neural networks. Because we consider only binary
 194 treatments, we simply implement a (tunable) feed-forward neural network with sigmoid activation
 195 function. Then, DeepIV proceeds by learning a second stage neural network to solve the inverse
 196 problem defined by Eq. (47).

197 **DeepGMM** [1]: DeepGMM [1] adopts neural networks for IV estimation inspired by the (optimally
 198 weighted) Generalized Method of Moments. The DeepGMM estimator is defined as the solution of
 199 the following minimax game:

$$\hat{\theta} \in \arg \min_{\theta \in \Theta} \sup_{\tau \in \mathcal{T}} \frac{1}{n} \sum_{i=1}^n f(z_i, \tau)(y_i - g(a_i, \theta)) - \frac{1}{4n} \sum_{i=1}^n f^2(z_i, \tau)(y_i - g(a_i, \tilde{\theta}))^2, \quad (48)$$

200 where $f(z_i, \cdot)$ and $g(a_i, \cdot)$ are parameterized by neural networks. As recommended in [1], we solve
 201 this optimization via adversarial training with the Optimistic Adam optimizer [5], where we set the
 202 parameter $\tilde{\theta}$ to the previous value of θ .

203 **DMLIV** [15]: DMLIV [15] assumes that the data is generated via

$$Y = \tau(X)A + f(X) + U, \quad (49)$$

204 where τ is the ITE f some function of the observed covariates. First, DMLIV estimates the functions
 205 $q(X) = \mathbb{E}[Y | X]$, $h(Z, X) = \mathbb{E}[A | Z, X]$, and $p(X) = \mathbb{E}[A | X]$. Then, the ITE is learned by
 206 minimizing the loss

$$\mathcal{L}(\theta) = \sum_{i=1}^n (y_i - \hat{q}(x_i) - \hat{\tau}(x_i, \theta)(\hat{h}(z_i, x_i) - \hat{p}(x_i)))^2, \quad (50)$$

207 where $\hat{\tau}(X, \cdot)$ is some model for $\tau(X)$. In our experiments, we use (tunable) feed-forward neural
 208 networks for all estimators.

209 **DRIV** [15]: DRIV [15] is a meta learner, originally proposed in combination with DMLIV. It requires
 210 initial estimators for $q(X)$, $p(X)$, $\pi(X) = \mathbb{E}[Z | X = x]$, and $f(X) = \mathbb{E}[A \cdot Z | X = x]$ as well
 211 as an initial ITE estimator $\hat{\tau}_{\text{init}}(X)$ (e.g., from DMLIV). The ITE is then estimated by a pseudo
 212 regression on the following doubly robust pseudo outcome:

$$\hat{Y}_{\text{DR}} = \hat{\tau}_{\text{init}}(X) + \frac{(Y - \hat{q}(X) - \hat{\tau}_{\text{init}}(X)(\hat{h}(X) - \hat{p}(X)))Z - \hat{\pi}(X)}{\hat{f}(X) - \hat{p}(X)\hat{r}(X)}. \quad (51)$$

213 We implement all regressions using (tunable) feed-forward neural networks.

214 **Comparison between DRIV vs. MRIV:** There are two key differences between our paper and [15]:
 215 (i) Our MRIV is multiply robust, while DRIV is only doubly robust. (ii) We derive a multiple robust
 216 convergence rate, while the rate in [15] is not robust with respect to the nuisance rates.

217 **Ad (i):** Both MRIV and DRIV perform a pseudo-outcome regression on the efficient influence
 218 function (EIF) of the ATE. The key difference: DRIV uses the doubly robust parametrization of the
 219 EIF from [11], whereas we use the multiply robust parametrization of the EIF from [17]². Hence,
 220 our MRIV frameworks extends DRIV in a non-trivial way to achieve multiple robustness (rather
 221 than doubly robustness). Thus, our estimator is consistent in the union of *three* different model
 222 specifications rather than *two*.³

223 **Ad (ii):** Here, we compare the convergence rates from DRIV and our MRIV and, thereby, show the
 224 strengths of our MRIV. To this end, let us assume that the pseudo regression function is γ -smooth and
 225 that we use the same second-stage estimator \hat{E}_n with minimax rate $n^{-\frac{2\gamma}{2\gamma+p}}$ for both DRIV and MRIV.
 226 If the nuisance parameters $q(X)$, $p(X)$, $f(X)$, and $\pi(X)$ are α -smooth and further are estimated
 227 with minimax rate $n^{\frac{-2\alpha}{2\alpha+p}}$, Corollary 4 from [15] states that DRIV converges with rate

$$\mathbb{E} \left[(\hat{\tau}_{\text{DRIV}}(x) - \tau(x))^2 \right] \lesssim n^{\frac{-2\gamma}{2\gamma+p}} + n^{\frac{-4\alpha}{2\alpha+p}}.$$

228 In contrast, MRIV assumes estimation of the nuisance parameters $\mu_0^Y(x)$ with rate $n^{\frac{-2\alpha}{2\alpha+p}}$, $\mu_0^A(x)$
 229 and $\delta_A(x)$ with rate $n^{\frac{-2\beta}{2\beta+p}}$, and $\pi(x)$ with rate $n^{\frac{-2\delta}{2\delta+p}}$. If the initial estimator $\hat{\tau}_{\text{init}}(x)$ converges with
 230 rate $r_\tau(n)$, our Theorem 2 yields the rate

$$\mathbb{E} \left[(\hat{\tau}_{\text{MRIV}}(x) - \tau(x))^2 \right] \lesssim n^{\frac{-2\gamma}{2\gamma+p}} + r_\tau(n) \left(n^{\frac{-2\beta}{2\beta+p}} + n^{\frac{-2\delta}{2\delta+p}} \right) + n^{-2\left(\frac{\alpha}{2\alpha+p} + \frac{\delta}{2\delta+p}\right)} + n^{-2\left(\frac{\beta}{2\beta+p} + \frac{\delta}{2\delta+p}\right)}.$$

231 If all nuisance parameters converge with the same minimax rate of $n^{\frac{-2\alpha}{2\alpha+p}}$, the rates of DRIV and
 232 our MRIV coincide. However, different to DRIV, our rate is additionally multiple robust in spirit of
 233 Theorem 1. This presents a crucial strength of our MRIV over DRIV: For example, if δ is small (slow
 234 convergence of $\hat{\pi}(x)$), our MRIV still with fast rate as long as α and β are large (i.e., if the other
 235 nuisance parameters are sufficiently smooth).

236 E.3 Wald estimator

237 Finally, we consider the Wald estimator [16] for the binary IV setting. More precisely, we estimate
 238 the ITE components $\mu_i^Y(x)$ and $\mu_i^A(x)$ separately and plug them into

$$\tau(x) = \frac{\hat{\mu}_1^Y(x) - \hat{\mu}_0^Y(x)}{\hat{\mu}_1^A(x) - \hat{\mu}_0^A(x)}. \quad (52)$$

239 We consider two versions of the Wald estimator:

240 **Linear:** We use linear regressions to estimate the $\mu_i^Y(x)$ and logistic regressions to estimate the
 241 $\mu_i^A(x)$.

242 **BART:** We use Bayesian additive regression trees [3] trees to estimate the $\mu_i^Y(x)$ and random forest
 243 classifier to estimate the $\mu_i^A(x)$.

²For a detailed discussion on multiple robustness and the importance of the EIF parametrization, we refer to [18], Section 4.5.

³On a related note, a similar, important contribution of developing multiply robust method was recently made for the average treatment effect. Here, the estimator of [11] was extended by the estimator of [17] to allow for multi robustness. Yet, this different from our work in that it focuses on the average treatment effect, while we study the individual treatment effect in our paper.

244 **F Visualization of predicted ITEs**

245 We plot the predicted ITEs for the different baselines and MRIV-Net in Fig. 3 (for $n = 3000$). As
 246 expected, the linear methods (2SLS and linear Wald) are not flexible enough to provide accurate
 247 ITE estimates. We also observe that the curve of MRIV-Net without MRIV is quite wiggly, i.e., the
 248 estimator has a relatively large variance. This variance is reduced when the full MRIV-Net is applied.
 249 As a result, curve is much smoother. This is reasonable because MRIV does not estimate the ITE
 250 components individually, but estimates the ITE directly via the Stage 2 pseudo outcome regression.
 251 Overall, this confirms the superiority of our proposed framework.

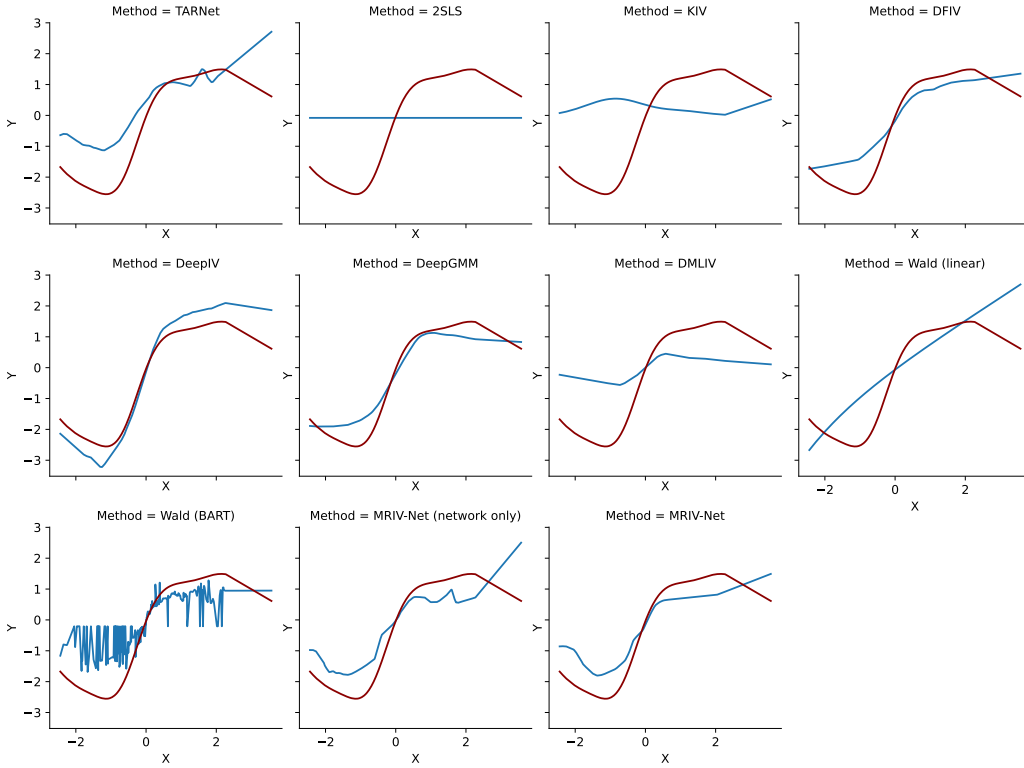


Figure 3: Predicted ITEs (blue) and oracle ITE (red) for different baselines.

252 **G Implementation details and hyperparameter tuning**

253 **Implementation details for deep learning models:** To make the performance of the deep learning
 254 models comparable, we implemented all feed-forward neural networks (including MRIV-Net) as
 255 follows: We use two hidden layers with RELU activation functions. We also incorporated a dropout
 256 layer for each hidden layer. We trained all models with the Adam optimizer [9] using 100 epochs.
 257 Exceptions are only DFIV and DeepGMM, where we used 200 epochs for training, accounting for
 258 slower convergence of the respective (adversarial) training algorithms. For DeepGMM, we further
 259 used Optimistic Adam [5] as in the original paper.

260 **Training times:** We report the approximate times needed to train the deep learning models on
 261 our simulated data with $n = 5000$ in Table 1. For training, we used an AMD Ryzen Pro 7 CPU.
 262 Compared to DMLIV and DRIV, the training of MRIV-Net is faster because only a single neural
 263 network is trained.

Table 1: Training times for deep learning models (in seconds).

TARNet	TARNet + DR	DFIV	DeepIV	DeepGMM	DMLIV	DMLIV + DRIV	MRIV-Net
~10.62	~28.57	~164.98	~30.21	~17.31	~74.98	~91.12	~32.20

264 **Hyperparameter tuning:** We performed hyperparameter tuning for all deep learning models
 265 (including MRIV-Net), KIV, and the BART Wald estimator on all datasets. For all methods except
 266 KIV and DFIV, we split the data into a training set (80%) and a validation set (20%). We then
 267 performed 40 random grid search iterations and chose the set of parameters that minimized the
 268 respective training loss on the validation set. In particular, the tuning procedure was the same for
 269 all baselines, which ensures that the performance gain of MRIV-Net is due to the method itself
 270 and not due to larger flexibility. Exceptions are only KIV and DFIV, for which we implemented
 271 the customized hyperparameter tuning algorithms proposed in [14] and [21] to ensure consistency
 272 with prior literature. For the meta learners (DR-learner, DRIV, and MRIV), we first performed
 273 hyperparameter tuning for the base methods and nuisance models, before tuning the pseudo outcome
 274 regression neural network by using the input from the tuned models. The tuning ranges for the
 275 hyperparameter are shown in Table 2. These include both the hyperparameter rangers shared across
 276 all neural networks and the model-specific hyperparameters. For reproducibility purposes, we publish
 277 the selected hyperparameters in our GitHub project as `.yaml` files.⁴

Table 2: Hyperparameter tuning ranges.

MODEL	HYPERPARAMETER	TUNING RANGE
Feed-forward neural networks (Shared parameter ranges for all deep learning baselines)	Hidden layer size(es)	$p, 5p, 10p, 20p, 30p$ (simulated data) $p, 3p, 5p, 8p, 10p$ (OHIE)
	Learning rate	0.0001, 0.0005, 0.001, 0.005, 0.01
	Batch size	64, 128, 256
	Dropout probability	0, 0.1, 0.2, 0.3
KIV	λ (Ridge penalty first stage)	5, 6, 7, 8, 9, 10, 12
	ξ (Ridge penalty second stage)	5, 6, 7, 8, 9, 10, 12
DFIV	λ_1 (Ridge penalty first stage)	0.0001, 0.001, 0.01, 0.1 (simulated data) 0.01, 0.05, 0.1 (OHIE)
	λ_2 (Ridge penalty second stage)	0.0001, 0.001, 0.01, 0.1 (simulated data) 0.01, 0.05, 0.1 (OHIE)
DeepGMM Wald (BART)	λ_f (learning rate multiplier)	0.5, 1, 1.5, 2, 5
	Number of trees (BART)	20, 30, 40, 50
	Number of trees (Random forest classifier)	20, 30, 40, 50

p = network input size

278 **Hyperparameter robustness checks:** We also investigate the robustness of MRIV-Net with respect
 279 to hyperparameter choice. To to this, we fix the optimal hyperparameter constellation for our simulated
 280 data for $n = 3000$ and perturb the hidden layer sizes, learning rate, dropout probability, and batch size.

⁴Codes are in the supplementary materials. Codes are also available at <https://anonymous.4open.science/r/MRIV-Net-0AC4> (Upon acceptance, we replace the link and point to a public GitHub repository).

281 The results are shown in Fig. 4. We observe that the RMSE only changes marginally when perturbing
282 the different hyperparameters, indicating that our method is to a certain degree robust against
283 hyperparameter misspecification. Furthermore, our results indicate that the performance improvement
284 of MRIV-Net over the baselines observed in our experiments is not due to hyperparameter tuning,
285 but to our method itself.

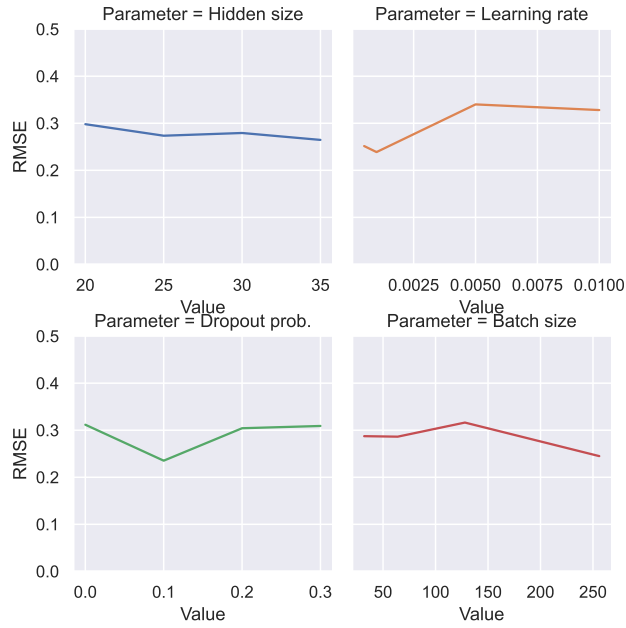


Figure 4: Robustness checks for different hyperparameters of MRIV-Net.

286 **H Results for semi-synthetic data**

287 In the main paper, we evaluated MRIV-Net both on synthetic and real-world data. Here, we provide
 288 additional results by constructing a semi-synthetic dataset on the basis of OHIE. It is common practice
 289 in causal inference literature to use semi-synthetic data for evaluation, because it combines advantages
 290 of both synthetic and real-world data. On the one hand, the real-world data part ensures that the
 291 data distribution is realistic and matches those in practice. On the other hand, the counterfactual
 292 ground-truth is still available, which makes it possible to measure the performance of ITE methods.

293 We construct our semi-synthetic data as follows: First, we extract the covariates $X \in \mathbb{R}^5$ and instru-
 294 ments $Z \in \{0, 1\}$ of our OHIE dataset from Sec. D. Then, we construct the treatment components
 295 $\mu_i^A(x)$ via

$$\mu_1^A(X) = 0.3 \cdot \sigma(X_1) + 0.7 \quad \text{and} \quad \mu_0^A(X) = 0.3 \cdot \sigma(X_1), \quad (53)$$

296 where X_1 is the (standardized) age and $\sigma(\cdot)$ is the sigmoid function. The outcome components are
 297 constructed via

$$\mu_1^Y(X) = 0.5X_1^2 + \sum_{i=2}^5 X_i^2 \quad \text{and} \quad \mu_0^Y(X) = -0.5X_1^2 + \sum_{i=2}^5 X_i^2. \quad (54)$$

298 We then sample treatments A and outcomes Y as in Eq. (31) and Eq. (32). Lemma 7 ensures that
 299 $\mu_i^Y(X) = \mathbb{E}[Y \mid Z = i, X]$ and $\mu_i^A(X) = \mathbb{E}[A \mid Z = i, X]$.

300 Given the above, the oracle ITE becomes

$$\tau(X) = \frac{X_1^2}{0.7}. \quad (55)$$

301 Note that $\tau(X)$ is sparse in the sense that it only depends on age, while the outcome components
 302 depend on all five covariates. Following our theoretical analysis in Sec. B, MRIV-Net should thus
 303 outperform methods that aim at estimating the components directly. This is confirmed in Table 3,
 304 where we show the results for all baselines and MRIV-Net on the semi-synthetic data. Indeed, we
 305 observe that MRIV-Net outperforms all other baselines, confirming both the superiority of our method
 306 as well as our theoretical results under sparsity assumptions from Sec. B.

Table 3: Results for semi-synthetic data.

Method	$n = 3000$	$n = 5000$	$n = 8000$
(1) STANDARD ITE			
TARNet [13]	1.66 ± 0.11	1.58 ± 0.07	1.57 ± 0.11
TARNet + DR [13, 8]	1.31 ± 0.28	1.22 ± 0.37	1.12 ± 0.15
(2) GENERAL IV			
2SLS [19]	1.34 ± 0.06	1.31 ± 0.03	1.32 ± 0.02
KIV [14]	1.97 ± 0.10	1.92 ± 0.05	1.93 ± 0.05
DFIV [21]	1.67 ± 0.44	1.63 ± 0.47	1.45 ± 0.17
DeepIV [7]	1.24 ± 0.26	0.99 ± 0.22	0.84 ± 0.19
DeepGMM [1]	1.39 ± 0.03	1.37 ± 0.16	1.18 ± 0.16
DMLIV [15]	2.12 ± 0.10	2.09 ± 0.09	2.02 ± 0.11
DMLIV + DRIV [15]	1.22 ± 0.10	1.18 ± 0.19	1.00 ± 0.08
(3) WALD ESTIMATOR [16]			
Linear	1.42 ± 0.24	1.28 ± 0.07	1.32 ± 0.07
BART	1.48 ± 0.24	1.29 ± 0.04	1.06 ± 0.13
MRIV-Net (network only)	1.11 ± 0.15	0.84 ± 0.14	0.95 ± 0.21
MRIV-Net (ours)	0.71 ± 0.24	0.75 ± 0.18	0.78 ± 0.26

Reported: RMSE (mean ± standard deviation). Lower = better (best in bold)

307 **I Results for cross-fitting**

308 Here, we repeat our experiments from the main paper but now make use of *cross-fitting*. Recall that,
 309 in Theorem 2, we assume that the nuisance parameter estimation and the pseudo-outcome regression
 310 are performed on three independent samples. We now address this through *cross-fitting*. To this end,
 311 our aim is to show that our proposed MRIV framework is again superior.

312 For MRIV, we proceeded as follows: We split the sample \mathcal{D} into three equally sized samples \mathcal{D}_1 , \mathcal{D}_2 ,
 313 and \mathcal{D}_3 . We then trained $\hat{\tau}_{init}(x)$, $\hat{\mu}_0^Y(x)$, and $\hat{\mu}_0^A(x)$ on \mathcal{D}_1 , $\hat{\delta}_A(x)$ and $\hat{\pi}(x)$ on \mathcal{D}_2 , and performed
 314 the pseudo-outcome regression on \mathcal{D}_3 . Then, we repeated the same training procedure two times, but
 315 performed the pseudo-outcome regression on \mathcal{D}_2 and \mathcal{D}_1 . Finally, we averaged the resulting three
 316 ITE estimators. For DRIV, we implemented the cross-fitting procedure described in [15]. For the
 317 DR-learner, we followed [8].

318 The results are in Table H. Importantly, the results confirm the effectiveness of our proposed MRIV.
 319 Overall, we find that our proposed MRIV outperforms DRIV for the vast majority of base methods
 320 when performing cross-fitting. Furthermore, MRIV-Net is highly competitive even when comparing
 321 it with the cross-fitted estimators. This shows that our heuristic to learn separate representations
 322 instead of performing sample splits works in practice. In sum, the results confirm empirically that our
 323 MRIV is superior.

Table 4: Results for base methods with different meta-learners (i.e., DRIV, and our MRIV) using cross-fitting and results for MRIV-Net without cross-fitting.

Base methods \ Meta-learners	$n = 3000$		$n = 5000$		$n = 8000$	
	DRIV	MRIV (ours)	DRIV	MRIV (ours)	DRIV	MRIV (ours)
(1) STANDARD ITE						
TARNet [13]	0.30 ± 0.02	0.36 ± 0.16	0.18 ± 0.06	0.16 ± 0.03	0.21 ± 0.08	0.13 ± 0.04
TARNet + DR-learner [13, 8]		0.85 ± 0.11		0.66 ± 0.08		0.67 ± 0.12
(2) GENERAL IV						
2SLS [19]	0.42 ± 0.11	0.33 ± 0.09	0.20 ± 0.07	0.23 ± 0.11	0.24 ± 0.10	0.14 ± 0.02
KIV [14]	0.47 ± 0.18	0.45 ± 0.15	0.20 ± 0.06	0.19 ± 0.08	0.22 ± 0.04	0.15 ± 0.03
DFIV [21]	0.35 ± 0.05	0.28 ± 0.09	0.22 ± 0.10	0.18 ± 0.08	0.24 ± 0.12	0.16 ± 0.04
DeepIV [7]	0.38 ± 0.09	0.44 ± 0.16	0.20 ± 0.07	0.19 ± 0.07	0.20 ± 0.08	0.12 ± 0.02
DeepGMM [1]	0.42 ± 0.09	0.42 ± 0.16	0.19 ± 0.04	0.19 ± 0.07	0.22 ± 0.06	0.13 ± 0.02
DMLIV [15]	0.44 ± 0.09	0.46 ± 0.16	0.21 ± 0.04	0.19 ± 0.07	0.21 ± 0.05	0.14 ± 0.02
(3) WALD ESTIMATOR [16]						
Linear	0.47 ± 0.23	0.36 ± 0.12	0.24 ± 0.05	0.20 ± 0.08	0.22 ± 0.05	0.15 ± 0.02
BART	0.43 ± 0.12	0.39 ± 0.12	0.14 ± 0.05	0.13 ± 0.05	0.23 ± 0.08	0.15 ± 0.02
MRIV-Net/w network only (ours)	0.35 ± 0.12	0.26 ± 0.11	0.19 ± 0.13	0.15 ± 0.03	0.18 ± 0.08	0.13 ± 0.03

Reported: RMSE (mean ± standard deviation). Lower = better (best in bold)

324 **References**

- 325 [1] Andrew Bennett, Nathan Kallus, and Tobias Schnabel. “Deep generalized method of moments
326 for instrumental variable analysis”. In: *NeurIPS*. 2019.
- 327 [2] Jean Chesson. “A non-central multivariate hypergeometric distribution arising from biased
328 sampling with application to selective predation”. In: *Journal of Applied Probability* 13.4
329 (1976), pp. 795–797.
- 330 [3] Hugh A. Chipman, Edward I. George, and Robert E. McCulloch. “BART: Bayesian additive
331 regression trees”. In: *The Annals of Applied Statistics* 4.1 (2010), pp. 266–298.
- 332 [4] Alicia Curth and Mihaela van der Schaar. “Nonparametric estimation of heterogeneous treat-
333 ment effects: From theory to learning Algorithms”. In: *AISTATS*. 2021.
- 334 [5] Constantinos Daskalakis et al. “Training GANs with optimism”. In: *ICLR*. 2018.
- 335 [6] Amy Finkelstein et al. “The oregon health insurance experiment: Evidence from the first year”.
336 In: *The Quarterly Journal of Economics* 127.3 (2012), pp. 1057–1106.
- 337 [7] Jason Hartford et al. “Deep IV: A flexible approach for counterfactual prediction”. In: *ICML*.
338 2017.
- 339 [8] Edward H. Kennedy. “Optimal doubly robust estimation of heterogeneous causal effects”. In:
340 *arXiv preprint* (2020).
- 341 [9] Diederik P. Kingma and Jimmy Ba. “Adam: A method for stochastic optimization”. In: *ICLR*.
342 2015.
- 343 [10] Whitney K. Newey and James L. Powell. “Instrumental variable estimation of nonparametric
344 models”. In: *Econometrica* 71.5 (2003), pp. 1565–1578.
- 345 [11] Ryo Okui et al. “Doubly robust instrumental variable regression”. In: *Statistica Sinica* 22.1
346 (2012), pp. 173–205.
- 347 [12] Carl Edward Rasmussen and Christopher K. I. Williams. *Gaussian processes for machine
348 learning*. 3. print. Adaptive computation and machine learning. Cambridge, Mass.: MIT Press,
349 2008.
- 350 [13] Uri Shalit, Fredrik D. Johansson, and David Sontag. “Estimating individual treatment effect:
351 Generalization bounds and algorithms”. In: *ICML*. 2017.
- 352 [14] Rahul Singh, Maneesh Sahani, and Arthur Gretton. “Kernel instrumental variable regression”.
353 In: *NeurIPS*. 2019.
- 354 [15] Vasilis Syrgkanis et al. “Machine learning estimation of heterogeneous treatment effects with
355 instruments”. In: *NeurIPS*. 2019.
- 356 [16] Abraham Wald. “The fitting of straight lines if both variables are subject to error”. In: *Annals
357 of Mathematical Statistics* 11.3 (1940), pp. 284–300.
- 358 [17] Linbo Wang and Eric J. Tchetgen Tchetgen. “Bounded, efficient and multiply robust estimation
359 of average treatment effects using instrumental variables”. In: *Journal of the Royal Statistical
360 Society: Series B* 80.3 (2018), pp. 531–550.
- 361 [18] Yixin Wang and David M. Blei. “The blessings of multiple causes”. In: *Journal of the American
362 Statistical Association* 114.528 (2019), pp. 1574–1596.
- 363 [19] Jeffrey M. Wooldridge. *Introductory Econometrics: A modern approach*. Routledge, 2013.
- 364 [20] Phillip G. Wright. *The tariff on animal and vegetable oils*. New York: Macmillan, 1928.
- 365 [21] Liyuan Xu et al. “Learning deep features in instrumental variable regression”. In: *ICLR*. 2021.
- 366 [22] Yun Yang and Surya T. Tokdar. “Minimax-optimal nonparametric regression in high dimen-
367 sions”. In: *The Annals of Statistics* 43.2 (2015), pp. 652–674.