

# BudgetIV: Optimal Partial Identification of Causal Effects with Mostly Invalid Instruments

Jordan Penn<sup>1</sup>Gecia Bravo-Hermsdorff<sup>3</sup>Lee M. Gunderson<sup>2</sup>Ricardo Silva<sup>2</sup>David S. Watson<sup>1</sup><sup>1</sup>King's College London<sup>2</sup>University College London<sup>3</sup>University of Edinburgh

## Abstract

Instrumental variables (IVs) are widely used to estimate causal effects in the presence of unobserved confounding between an exposure  $X$  and outcome  $Y$ . An IV must affect  $Y$  *exclusively* through  $X$  and be *unconfounded* with  $Y$ . We present a framework for relaxing these assumptions with tuneable and interpretable "budget constraints". Our algorithm returns a feasible set of causal effects that can be identified exactly given perfect knowledge of observable covariance statistics. This feasible set might contain disconnected sets of possible solutions for the causal effect. We discuss conditions under which this set is *sharp*, i.e., contains all and only effects consistent with the background assumptions and the joint distribution of observable variables. Our method applies to a wide class of semiparametric models, and we demonstrate how its ability to select specific subsets of instruments confers an advantage over convex relaxations in both linear and nonlinear settings. We adapt our algorithm to form confidence sets that are asymptotically valid under a common statistical assumption from the Mendelian randomization literature.

Fisher, 1935). However, due to financial, ethical or physical constraints, confounding factors affecting  $X$  and  $Y$  often cannot be held fixed or observed and adjusted for (Rubin, 1974; Holland, 1986). In these situations, instrumental variables (IVs) are commonly used to infer causal effects.

IVs are a set of pre-treatment covariates  $Z$  that are (A1) associated with  $X$ ; (A2) unconfounded with  $Y$ ; and (A3) whose causal effect on  $Y$  is exclusively mediated through  $X$  (Angrist et al., 1996; Heckman and Vytlačil, 1999) (see Sect. 2.1 for formal definition). Due to the presence of latent confounders, one cannot test whether (A2) and (A3) hold for any particular pre-treatment covariate. However, with multiple candidate instruments, a simple fact may be exploited: if different sets of covariates predict different values of the causal parameter  $\theta$  when treated as IVs, then at most one of the sets contains valid IVs.

Kang et al. (2016) exploit this fact by showing the following in linear models with a scalar treatment: if more than 50% of the pre-treatment covariates are correctly assumed to be valid IVs, then  $\theta$  is point-identifiable (i.e., given an infinite number of observations from the model, the value of  $\theta$  can be uniquely determined). A number of majority rule based approaches to inference with invalid IVs extend this work (Bowden et al., 2016a; Hartwig et al., 2017; Bucur et al., 2020; Hartford et al., 2021).

However, in settings where (a) 50% or fewer candidate IVs are valid, or (b) more than one causal parameter is needed (e.g., because the treatment is multidimensional), we show that  $\theta$  can be at best *partially* identified. For example, even with oracle access to the joint distribution  $P(X, Y, Z_1, Z_2)$ , it may be undecidable whether  $Z_1$  is a valid IV and  $Z_2$  is not or vice versa. Though some values of  $\theta$  may be excluded, the parameter will not converge to a single value even in the limit of infinite data.

Rather than assuming that all or a majority of  $Z$  are valid IVs, our approach is to set a minimum propor-

## 1 INTRODUCTION

The causal effect of a "treatment" (or "exposure")  $X$  on an "outcome"  $Y$  captures the change in the distribution of  $Y$  when we intervene on  $X$  (Pearl, 2009). In a randomized control trial,  $X$  is controlled explicitly so that the causal effect can be identified (Neyman, 1923;

Proceedings of the 28<sup>th</sup> International Conference on Artificial Intelligence and Statistics (AISTATS) 2025, Mai Khao, Thailand. PMLR: Volume 258. Copyright 2025 by the author(s).

tion which are. In fact, we consider a more general setting (c), in which “degrees” of IV validity are set by user-input thresholds. It is not assumed *a priori* which thresholds apply to which candidate IVs. For each threshold, however, a user-input “budget” constrains how many candidates are allowed to violate the IV assumptions up to that degree. Bucur et al. (2020) and Xue et al. (2023) propose point estimators for  $\theta$  under budget-like constraints, restraining fewer than 50% of candidate IVs to be valid. As we will see in Sect. 2, however,  $\theta$  is only partially identified in this setting. The allocation of budgets to  $Z$  are also partially identified in general, which can lead to disconnected sets of solutions for  $\theta$ .

Our discussion applies when IVs are invalid due to violation of (A2), (A3) or both. We show that under the assumption of homogeneous causal effects (Holland, 1986), which implies an additive separation of the  $X$ - and  $Z$ -signals into  $Y$  (Newey and Powell, 2003), a latent vector statistic summarizes bias in the estimated  $\theta$  due to violation of both assumptions. This statistic has a simple interpretation: the residual covariance between pre-treatment covariates  $Z$  and the outcome  $Y$  not explained by the causal effect of the treatment  $X$ . Vancak and Sjölander (2023) make a similar proposal to quantify the violation of IV assumptions in a single sensitivity parameter. However, their approach is limited to a single scalar  $Z$ , so does not utilize the notion of IV candidates and whether they give consistent estimates of  $\theta$ .

We provide an algorithm, `budgetIV`, that provably finds *sharp* partial identification sets for causal effect parameters under fixed budget constraints. We show that the set can be inferred from relevant summary statistics, which reduce to the familiar covariance parameters  $\text{Cov}(X, Z)$  and  $\text{Cov}(Y, Z)$  for linear models. In the special case of scalar treatments, we provide a polytime procedure for computing sharp solution sets. However, under setting (b), we prove that partial identification is NP-hard in the number of instruments, budget thresholds and causal effect parameters. Finally, we provide a method for handling finite-sample uncertainty. In Sect. 4, we show how `budgetIV` can be adapted to form (asymptotically) valid confidence sets under the so-called “no measurement error” assumption from the Mendelian randomization literature (Bowden et al., 2016b).

The remainder of this paper is structured as follows. We formalize our problem in Sect. 2. We introduce the `budgetIV` algorithm in Sect. 3 and propose related inference procedures in Sect. 4. Experimental results are presented in Sect. 5. Following a literature review in Sect. 6, we conclude with a brief discussion in Sect. 7. Proofs, pseudocode, and experimental details are given

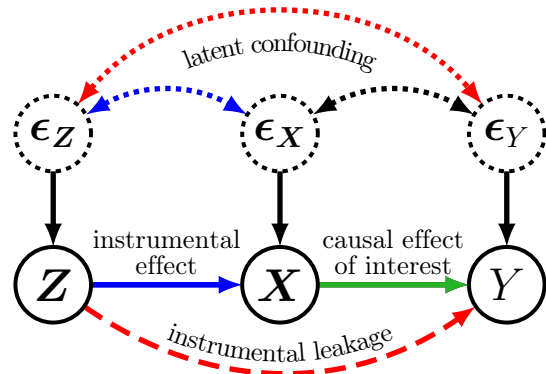


Figure 1: **Acyclic directed mixed graph for our problem setup.** Solid circles represent observable variables and dashed circles latent variables. Bidirected arrows are interpreted as any mutual dependence between noise residuals. The dotted black arrow indicates the unobserved confounding between  $X$  and  $Y$ . The relevance assumption (A1) requires at least one of the blue arrows. The red arrows contribute to violations of the exogeneity conditions (A2) and (A3). The green arrow indicates the causal effect of interest.

in the appendix.

## 2 PROBLEM SETUP

Our observable variables include a set of candidate instruments  $Z \in \Omega_Z \subseteq \mathbb{R}^{d_Z}$ ; treatments  $X \in \Omega_X \subseteq \mathbb{R}^{d_X}$ ; and a univariate outcome  $Y \in \Omega_Y \subseteq \mathbb{R}$ . We assume that the ground truth structural equation model (SEM) between these variables takes the following form:

$$Z := f_z(\epsilon_z), \quad (1)$$

$$X := f_x(Z, \epsilon_x), \quad (2)$$

$$Y := \theta^* \cdot \Phi(X) + g_y(Z, \epsilon_y), \quad (3)$$

where  $\epsilon_z$ ,  $\epsilon_x$ , and  $\epsilon_y$  are noise residuals that capture any and all effects from latent variables, including unobserved confounding. The additive separability of  $X$  and  $Z$  in  $Y$  is guaranteed by assuming that the treatment effect is homogeneous, i.e.,  $P(y | \text{do}(x)) - P(y | \text{do}(x_0))$  is independent of  $Z$  and unobserved confounding between  $X$  and  $Y$ .

The function  $\Phi : \Omega_X \mapsto \Omega_\Phi \subseteq \mathbb{R}^{d_\Phi}$  may provide a basis-expansion of a nonlinear treatment effect or a representation of a high-dimensional treatment  $X$ . Though some authors have studied the underspecified regime in which  $d_\Phi > d_Z$  (Pfister and Peters, 2022; Ailer et al., 2023), we restrict attention to the more common case where  $d_\Phi \leq d_Z$ . This guarantees the identification result of Thm. 1, first shown in linear models by Koopmans (1949).

The vector  $\theta^* \in \mathbb{R}^{d_\Phi}$  is the causal parameter of interest. Together with  $\Phi$ , it defines the average treatment effect:

$$\text{ATE}(\mathbf{x}; \mathbf{x}_0) = \theta^* \cdot (\Phi(\mathbf{x}) - \Phi(\mathbf{x}_0)),$$

which represents the expected change in  $Y$  if we were to replace the intervention  $do(\mathbf{X} = \mathbf{x}_0)$  with  $do(\mathbf{X} = \mathbf{x})$ .

## 2.1 A Sensitivity Parameter for IV Violation

We call  $\mathbf{Z}$  a set of valid IVs if they satisfy the following properties:

- (A1) *Association*:  $\mathbf{Z} \not\perp \mathbf{X}$ ,
- (A2) *Unconfoundedness*:  $\mathbf{Z} \perp \epsilon_{\mathbf{y}}$ ,
- (A3) *Exclusion*:  $\mathbf{Z} \perp Y \mid \{\mathbf{X}, \epsilon_{\mathbf{y}}\}$ .

In several well-studied regimes (Koopmans, 1949; Imbens and Angrist, 1994; Newey and Powell, 2003; Heckman and Vytlacil, 2005; Saengkyongam et al., 2022), the assumptions above allow for point identification of the target parameter  $\theta^*$ . For our SEM (Eqs. (1) to (3)), we prove point identification under a stronger version of (A1) and an assumption that combines aspects of (A2) and (A3), despite being strictly weaker than either. To show this, we introduce the following parameters:

$$\begin{aligned} \gamma_{\mathbf{g}} &:= \text{Cov}(g_{\mathbf{y}}(\mathbf{Z}, \epsilon_{\mathbf{y}}), \mathbf{Z}), \\ \beta_{\Phi} &:= \text{Cov}(\Phi(\mathbf{X}), \mathbf{Z}), \\ \beta_{\mathbf{y}} &:= \text{Cov}(Y, \mathbf{Z}). \end{aligned}$$

By taking the covariance between each term in Eq. (3) with  $\mathbf{Z}$ , we see that target parameter  $\theta^*$  is related to  $\gamma_{\mathbf{g}}$  by  $\gamma_{\mathbf{g}} = \beta_{\mathbf{y}} - \theta^* \cdot \beta_{\Phi}$ . Thus,  $\gamma_{\mathbf{g}}$  is the residual covariance between  $Y$  and  $\mathbf{Z}$  not explained by the ground truth causal effect of  $\mathbf{X}$  on  $Y$ . Notice that if both (A2) and (A3) are satisfied, then  $\gamma_{\mathbf{g}} = \mathbf{0}$ .

Given that  $d_{\Phi} \leq d_{\mathbf{Z}}$ , the following constraints are sufficient for identification:

- (B1\*) *Association (strong)*:  $\text{rank}(\beta_{\Phi}) = d_{\Phi}$   
and  $(\beta_{\Phi})_i \neq 0$  ( $\forall Z_i \in \mathbf{Z}$ );
- (B2\*) *Exogeneity (strong)*:  $\gamma_{\mathbf{g}} = \mathbf{0}$ ,

where  $(\beta_{\Phi})_i := \text{Cov}(\Phi(\mathbf{X}), Z_i)$  and ‘‘exogeneity’’ refers to the conjunction of (A2) and (A3).

**Theorem 1 (Identifiability).** Assume Eqs. (1) to (3) and claims (B1\*), (B2\*) hold for some  $d_{\Phi} \leq d_{\mathbf{Z}}$ . Assume the existence of a ground truth joint distribution  $P(\mathbf{X}, Y, \mathbf{Z})$  with finite covariance parameters  $\beta_{\Phi}^*, \beta_{\mathbf{y}}^*$ . Then the causal parameter  $\theta^*$  can be identified exactly as the unique solution to  $\beta_{\mathbf{y}}^* - \theta^* \cdot \beta_{\Phi}^* = \mathbf{0}$ .

We can relax (B1\*) and (B2\*) further, in particular using  $\gamma_{\mathbf{g}}$  to model violations of exogeneity:

- (B1) *Association (relaxed)*:  $\beta_{\Phi} \neq \mathbf{0}$ ;
- (B2) *Exogeneity (relaxed)*:  $\gamma_{\mathbf{g}} \in \Gamma$ .

Such a relaxation may still allow for partial identification, which we will see in the following section. In Sect. 2.3 we introduce budget constraints as tuneable and interpretable choices for  $\Gamma$ .

## 2.2 Formalizing Optimal Partial Identification

We say that the causal parameter  $\theta^*$  is *partially identified* when more than one (but not all) of its possible values are consistent with the data and our structural assumptions. In this section, we define a certain notion of *optimality* for general partial identification problems, and establish the ingredients for an optimal solution in our setting.

The following definitions allow us to state and prove these results. Let  $\mathcal{M}$  be a class of SEMs and  $m^* \in \mathcal{M}$  the ground truth model. Each  $m \in \mathcal{M}$  implies a joint distribution  $P_m(\mathbf{V})$  over the observables  $\mathbf{V} \in \Omega_{\mathbf{V}}$ . (In our setting  $\mathbf{V} = \{\mathbf{X}, Y, \mathbf{Z}\}$  and  $\Omega_{\mathbf{V}} = \Omega_{\mathbf{X}} \times \Omega_Y \times \Omega_{\mathbf{Z}}$ .)

A constraint  $c : \mathcal{M} \mapsto \{0, 1\}$  is a logical formula that either does or does not hold for any given model. For instance,  $c$  may bound the range of some parameter(s) in  $m$  or impose conditional independence on certain variables in  $\mathbf{V}$ . Let  $\mathcal{C}$  be a set of such constraints, with  $\mathcal{C}^* := \{c \in \mathcal{C} : c(m^*) = 1\}$  denoting ground truth. (In our setting,  $\mathcal{C}$  includes the relaxed association and exogeneity assumptions (B1) and (B2).)

An *observable* statistic  $s$  is a functional of the joint distribution,  $s : \{P_m(\mathbf{V}) : m \in \mathcal{M}\} \mapsto \Omega_s$ , with ground truth value  $s^* := s(P_{m^*}(\mathbf{V}))$ . Examples include (conditional) moments or correlations between variables. (In our setting,  $s$  comprises the cross-covariance parameters  $\beta_{\mathbf{y}}, \beta_{\Phi}$ .) The target parameter  $q^* := q(m^*)$  is the ground truth for some *latent* statistic  $q : \mathcal{M} \mapsto \Omega_q$ , which cannot be determined by  $P_m(\mathbf{V})$  alone. (In our setting,  $q^*$  is the causal parameter  $\theta^*$ .)

Given a constraint  $c \in \mathcal{C}$  and statistic  $s$ , the ‘‘plausible’’ values of  $q^*$  form a *solution set*  $\mathcal{T}(c, s) \subseteq \Omega_q$ . Such sets are the image of a *solution map* for  $q$ ,  $\mathcal{T} : \mathcal{C} \times \Omega_s \mapsto \mathcal{P}(\Omega_q)$ , where  $\mathcal{P}$  denotes the power set.

We define an *optimal* solution map in terms of soundness, completeness, and minimality criteria. We defer discussion of computational complexity and finite sample inference to Sect. 3 and 4, respectively.

**Definition 1 (Soundness).** A solution map  $\mathcal{T}$  is *sound* if, for any ground truth model  $m^* \in \mathcal{M}$ , given statistic  $s^* := s(P_{m^*}(\mathbf{V}))$  and constraint  $c^* \in \mathcal{C}^*$ , we have  $q^* \in \mathcal{T}(c^*, s^*)$ .

This condition ensures that our solution map cannot

exclude the target  $q^*$  when provided with ground truth inputs.

**Definition 2 (Completeness).** A solution map  $\mathcal{T}$  is *complete* if, for any ground truth model  $m^* \in \mathcal{M}$ , given  $s^*$  and any  $c^* \in \mathcal{C}^*$ , the following holds. For all  $q \in \mathcal{T}(c^*, s^*)$ , there is at least one model  $m_q \in \mathcal{M}$  for which  $q = q(m_q)$  and  $s(P_{m_q}(\mathbf{V})) = s^*$ .

This condition ensures that no sound map can exclude more values of  $q \in \Omega_q \setminus \{q^*\}$ . Solution sets that are both sound and complete are said to be *sharp*, as originally defined by Manski (1990, 2003).

For the next definition, we introduce a partial order on observable statistics with respect to their information content. We say that  $s'$  is at least as informative as  $s$  if there exists a deterministic function  $f : \Omega_{s'} \mapsto \Omega_s$  such that, for any  $m \in \mathcal{M}$ ,  $s(P_m(\mathbf{V})) = f(s'(P_m(\mathbf{V})))$ . In this case, we write  $s' \preceq s$ .

**Definition 3 (Minimality).** A sharp solution map  $\mathcal{T} : \mathcal{C} \times \Omega_s \mapsto \Omega_q$  is *minimal* if, for any other sharp solution map  $\mathcal{T}' : \mathcal{C} \times \Omega_{s'} \mapsto \Omega_q$  that takes a different input statistic  $s'$ , we have  $s' \preceq s$ .

Minimality ensures that no sharp map could be constructed using strictly less information. With Defs. 1 to 3 in place, we can show the existence of an optimal solution map for our problem.

**Theorem 2 (Optimal solution map).** Let the model class  $\mathcal{M}$  consist of all SEMs consistent with Eqs. (1) to (3) for some given  $\Phi$  and  $d_\Phi \leq d_Z$ . Assume the existence of a ground truth joint distribution  $P(\mathbf{X}, Y, \mathbf{Z})$  with finite covariance parameters  $\beta_\Phi^*, \beta_y^*$ . Define the affine function  $h(\theta) := \beta_y - \theta \cdot \beta_\Phi$ , for some  $\beta_y, \beta_\Phi$ . Then the solution map  $\mathcal{T}$  defined by:

$$\mathcal{T}(c, s = (\beta_\Phi, \beta_y)) = \{\theta \in \mathbb{R}^{d_\Phi} : h(\theta) \in \Gamma\}, \quad (4)$$

is sound, complete, and minimal with respect to background constraint set  $\mathcal{C} = \{\mathbb{I}[\gamma_g \in \Gamma] : \Gamma \subseteq \mathbb{R}^{d_Z}\}$ , where  $\mathbb{I}[\cdot]$  denotes the indicator function.

In Appx. B.1 we examine how extra assumptions can result in smaller solution sets. In particular, we show that this is the case for models with categorical (Balke and Pearl, 1997) or bounded (Manski, 1990) outcomes.

### 2.3 Budget Background Constraints

We introduce sensible choices for the search space  $\Gamma$ . For any  $\gamma_g \in \mathbb{R}^{d_Z}$  and any  $\alpha \in [0, 1)$ , if  $\gamma'_g := \alpha\gamma_g$ , the degree of IV violation attributed to each candidate and the combined violation is strictly weaker for  $\gamma'_g$  than for  $\gamma_g$ . Below, we formalize the notion that if  $\gamma_g$  is plausible, then so too is  $\gamma'_g$ .

**Definition 4 (Star domain).** A set of points  $A \subseteq \mathbb{R}^d$  is a *star domain* if there exists a point  $\mathbf{a} \in A$  such

that the line between  $\mathbf{a}$  and any other point  $\mathbf{a}' \in A$  is contained within  $A$ . Equivalently, a star domain is a space  $A \subseteq \mathbb{R}^d$  with a nonempty convex kernel:

$$\text{ck}(A) := \{\mathbf{a} \in A \mid \forall \mathbf{a}' \in A, \eta \in [0, 1] : \eta \mathbf{a} + (1 - \eta) \mathbf{a}' \in A\}.$$

**Principle 1 (Starfish principle).** Relax a structural assumption by adding a latent sensitivity parameter,  $\delta$ , whose direction is related to the mechanisms of violation and whose magnitude increases with the degree of violation by these mechanisms. Bound the support of this parameter within a star domain for the FeasIble Search space—i.e., a starFISH—whose convex kernel includes the point  $\delta = \mathbf{0}$ .

This principle encapsulates current literature restricting the number<sup>1</sup> (Kang et al., 2016; Hartwig et al., 2017; Silva and Shimizu, 2017; Hartford et al., 2021) or total effect (Ramsahai, 2012; Conley et al., 2012; Silva and Evans, 2016; Watson et al., 2024; Jiang and Kocaoglu, 2024) of invalid instruments. Convex restrictions lead to convex identified sets when the sensitivity parameter is an affine function of the causal effect parameters. If this affine function intersects the convex kernel, the intersection is connected (by definition). In the general case, the identified set for the quantities of interest may be disconnected (see Fig. 2). However, each disjoint subset might also provide a distinct mechanistic interpretation for the plausible violations.

Here, we introduce budget constraints, which are based on two user-specified components: thresholds,  $\tau := (\tau_1, \tau_2, \dots, \tau_K)$ , which describe degrees of IV invalidity; and budgets,  $\mathbf{b} := (b_1, b_2, \dots, b_K)$ , which determine the minimum number of IVs assumed to be at least as valid as the corresponding threshold. The thresholds must be nonnegative and increasing:  $0 \leq \tau_1 < \tau_2 < \dots < \tau_K < \infty$ . Integer budgets are increasing and strictly positive:  $0 < b_1 < b_2 < \dots < b_K \leq d_Z$ . Thresholds and budgets with the same index  $i \in [K]$  form pairs. The number of thresholds–budget pairs  $K$  is chosen by the user but we require  $K \leq d_Z$ .

We define the  $d_Z$ -dimensional latent statistic  $U(\gamma_g)$  through the following relationship:

$$U_i = \ell \iff \tau_{\ell-1} \leq |\gamma_{g_i}| \leq \tau_\ell.$$

Let  $U_{i1}, U_{i2}, \dots, U_{iK}$  represent the encoding of  $U_i$  such that  $U_{i\ell} = \mathbb{I}[U_i \leq \ell]$  for all  $\ell \in [K]$ , where  $\mathbb{I}[\cdot]$  is the indicator function. We define the budget constraint restriction  $\Gamma(\tau, \mathbf{b})$  as the set of vectors  $\gamma_g$  satisfying, for all  $\ell \in [K]$ ,  $\sum_{i=1}^{d_Z} U_{i\ell} \geq b_\ell$  for some choice of budgets and thresholds. We denote the set of  $U$  with encodings

<sup>1</sup>Notice that the  $L_0$  norm constraint  $\|h(\theta)\|_0 \leq b$  is a star domain with a convex kernel including the point  $\theta = \mathbf{0}$ .

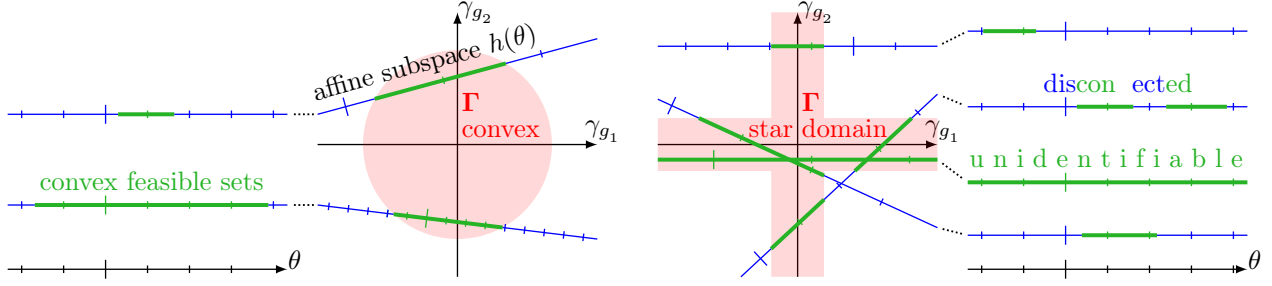


Figure 2: **The topology of a feasible set depends on the shape of the background constraints.**

Plots of constrained search spaces  $\Gamma$  (shaded) and lines  $h(\theta)$  corresponding to  $d_Z = 2, d_\Phi = 1$ . The intersection between a line and shaded region determines a feasible set of  $\theta$ . The constraints are a subspace of  $\mathbb{R}^{d_Z}$  while  $h(\theta)$  are  $d_\Phi$ -dimensional affine subspaces. (Left) Convex  $\Gamma$  entails convex feasible sets. (Right) Budget constraints  $\Gamma(\tau, \mathbf{b})$  form a star domain. They are non-convex in general and are unbounded for  $b_K < d_Z$  (i.e., some  $\gamma_{g_i}$  may be unconstrained). This can lead to disconnected or even unidentifiable causal effect. In Appx. B.2 we show that unidentifiability occurs only under violation of (B1\*) and can be tested in polytime.

satisfying the inequality above by  $\Sigma_{\mathbf{b}}$ . In particular, if equality holds, we say  $\mathbf{U} \in \Sigma_{\mathbf{b}}^{(\max)}$  is a maximally relaxed assignment. These definitions will be useful in the next section.

Since  $\mathbf{U}(\gamma_g)$  depends on  $\theta^*$  it is only partially identifiable in general. We define the set of plausible budget assignments:

$$\{\mathbf{U}(h(\theta)) : \theta \in \mathbb{R}^{d_Z}, h(\theta) \in \Gamma(\mathbf{b}, \tau)\}.$$

For instance, if  $K = 1$  and  $\tau_1 = 0$ , we assume at least  $b_1$  candidate IVs are valid (having  $\gamma_{g_i} = 0$ ), but allow  $d_Z - b_1$  to be unrestricted. This extends the  $L_0$ -norm constraint from Kang et al. (2016) to violation of (A2) and/or (A3).

### 3 THE ALGORITHM

The objective of `budgetIV` is twofold: (O1) discover plausible budget assignments; and (O2) partially identify the ATE. We simplify (O1), which for general  $d_Z > d_\Phi > 1$  allows the solution set for  $\theta$  to be calculated using subset selection over  $\Sigma_{\mathbf{b}}^{(\max)}$ . We conjecture that subset selection cannot be avoided, which means finding any  $\theta$  in the solution set is NP-hard (w.r.t.  $d_Z$  and  $d_\Phi$ ), while finding the entire solution set is #P-hard. As it stands, `budgetIV` should not be run for  $d_Z, d_\Phi > 10$ . When  $d_\Phi = 1$ , however, we show the solution set can be found in polynomial time using our algorithm `budgetIV_scalar`. Pseudocode for both algorithms is given in Appx. C; executable code is available online.<sup>2</sup>

**Simplifying (O1) for  $d_\Phi > 1$ .** Since  $\mathbf{U}(\gamma)$  is uniquely defined for each  $\gamma \in \mathbb{R}^{d_Z}$ , the feasible region  $\Gamma(\tau, \mathbf{b})$  decomposes into subsets  $\Gamma_U$  for which

$\gamma \in \Gamma_U \implies \mathbf{U}(\gamma) = \mathbf{U}$ . However, in Appx. E, we show that  $\Gamma(\tau, \mathbf{b})$  can also be thought of as a union of overlapping cuboids. Indeed, in Fig. 2 we see a feasible region made up of two overlapping rectangles. Each cuboid  $\tilde{\Gamma}_{\tilde{\mathbf{U}}}$  is indexed by a maximal budget assignment  $\tilde{\mathbf{U}} \in \Sigma_{\mathbf{b}}^{(\max)}$ .

#### 3.1 budgetIV with Oracle $\beta$ Parameters

We summarize the `budgetIV` algorithm below.

- (1) For each  $\tilde{\mathbf{U}} \in \Sigma^{(\max)}$ , test for the intersection between  $h(\theta) = \beta_y - \theta \cdot \beta_\Phi$  and  $\tilde{\Gamma}_{\tilde{\mathbf{U}}}$ .
- (2) If there is an intersection, solve the following linear program to find the bounds on ATE( $\mathbf{x}; \mathbf{x}_0$ ) for each  $\mathbf{x} \in \Omega_{\mathbf{X}}$  of interest and baseline treatment  $\mathbf{x}_0$ :

$$\begin{aligned} & \min/\max \theta \cdot (\Phi(\mathbf{x}) - \Phi(\mathbf{x}_0)). \\ & \theta: h(\theta) \in \tilde{\Gamma}_{\tilde{\mathbf{U}}} \end{aligned}$$

- (3) Let  $\theta^{-/+}(\mathbf{x})$  denote the argmin/argmax solution (respectively) to the optimization problem in step (2) at point  $\mathbf{x}$ . The function  $U_i(h(\theta))$  returns the value of the latent variable  $U_i$  at the point  $h(\theta)$ . Compute  $\tilde{\mathbf{U}}$ , defined by its components:

$$\tilde{U}_i := \max_{\mathbf{x} \in \Omega_{\mathbf{X}}} \max_{\diamond \in \{-, +\}} U_i(\theta^\diamond(\mathbf{x})).$$

- (4) Return all unique  $\tilde{\mathbf{U}}$  (O2) along with the corresponding ATE bounds (O1).

**Polytime `budgetIV_scalar`.** The algorithm above describes an ILP that requires a linear search over  $|\Sigma_{\mathbf{b}}^{(\max)}|$ , which itself is bounded above by  $d_Z^{d_Z}$ , corresponding to the case where  $K = d_Z$ . Though standard solvers are highly optimized, this task can quickly become intractable with many candidate instruments.

<sup>2</sup><https://github.com/jpenn2023/budgetIVr>.

When  $d_\Phi = 1$ , however, exact partial identification of the ATE can be computed in  $\tilde{O}(d_Z K)$  time. This is done by noticing that the value of  $\mathbf{U}$  associated with  $h(\theta)$  can only change as a function of  $\theta$  at the  $2d_Z K$  points:

$$\Theta_{ik}^\pm := \frac{\pm\tau_k - (\beta_y)_i}{(\beta_\Phi)_i},$$

where we ignore any  $i$  for which  $(\beta_\Phi)_i = 0$ . Similarly, the full set of plausible budget assignments can be calculated in polynomial time. See Appx. C for further details.

### 3.2 The $L_0$ -norm Constraint

Kang et al. (2016) show that point identification is possible in the linear setting with univariate exposure  $X$  and some invalid instruments—provided at least half are valid. We show that this result cannot be generalized to multidimensional exposures. However, we place tight bounds on the cardinality of the feasible set with respect to the minimum number of valid IVs  $b$ . We use the following shorthand:

$$\mathcal{T}_{L_0}(b) := \mathcal{T}(c = \mathbb{I}[\gamma_g \in \Gamma(0, b)], s = \{\beta_\Phi, \beta_y\}).$$

**Theorem 3 (*t*-point identification).** Assume Eqs. 1, 2, 3 and claims (B1\*), (B2) hold for some  $\Gamma(0, b)$ . Then, for all  $b < d_Z - d_\Phi + 1$ , the cardinality of the feasible set  $t := |\mathcal{T}_{L_0}(b)| \in \mathbb{N}$  is bounded above by:

$$t \leq \frac{d_Z!}{(d_\Phi - 1)!(d_Z - d_\Phi + 1)!} \left\lfloor \frac{d_Z - d_\Phi + 1}{d_Z - b - d_\Phi + 1} \right\rfloor.$$

This bound is tight in the sense that equality holds for some values of  $\beta_\Phi, \beta_y$ . We extend these results to incorporate violation of (B1\*) in the proof (Appx. A.3).

**Corollary 3.1 (No point identification for  $d_\Phi > 1$ ).** There is no value of  $b$  for which  $t \leq 1$  is guaranteed for all  $|\mathcal{T}_{L_0}|$  when  $d_\Phi > 1$ .

**Corollary 3.2 (*t*-point identification for  $d_\Phi = 1$ ).** In the case of  $d_\Phi = 1$ , we have:

$$|\mathcal{T}_{L_0}(b)| =: t \leq \left\lfloor \frac{d_Z}{d_Z - b} \right\rfloor \leq d_Z.$$

This reduces to point identifiability when  $b < d_Z/2$ .

## 4 INFERENCE

In this section we consider finite-sample uncertainty. We seek confidence sets with (asymptotically) valid coverage over the solution set  $\{\theta \in \mathbb{R}^{d_Z} : h(\theta) \in \Gamma(\tau, \mathbf{b})\}$ .

Since  $\mathcal{T}$  is a deterministic map, it is possible to use a confidence set  $(\mathbf{B}_\Phi, \mathbf{B}_y)$  over the covariance parameters

$\beta_\Phi, \beta_y$ . The confidence set over the causal parameter is then all  $\theta \in \mathbb{R}^{d_\Phi}$  satisfying:

$$\forall (\beta_\Phi, \beta_y) \in (\mathbf{B}_\Phi, \mathbf{B}_y) : (\beta_y - \theta \cdot \beta_\Phi) \in \Gamma(\tau, \mathbf{b}).$$

Provided the estimators  $\hat{\beta}_\Phi, \hat{\beta}_y$  have finite variance, an asymptotically valid confidence set can be constructed by modeling these estimators as multivariate normal (see Appx. A.4). However, we choose to make the following simplifications taken from the applied IV literature.

### 4.1 Coverage with Summary Statistics under the NOME Assumption

Candidate IVs are often selected from a pool of covariates, with inclusion based solely on marginal association with the exposure. For instance, Mendelian randomization (MR) is a popular approach in genetic epidemiology whereby genetic variants  $Z$  are used as IVs to determine the causal effect of phenotype(s)  $X$  on a health outcome  $Y$ . The  $Z \rightarrow X$  link is usually established by a genome-wide association study (GWAS). Empirically, the chosen  $Z$  tend to be less strongly associated with  $Y$  than with  $X$ . This can occur for various reasons: strong latent confounding between  $X$  and  $Y$ ; a weak causal effect of  $X$  on  $Y$ ; low variation in  $Y$  (e.g., rare diseases); and/or smaller sample sizes for evaluating  $P(Z, Y)$  than  $P(Z, X)$  (Pierce and Burgess, 2013).

In such cases, finite-sample error in  $\hat{\theta}$  is mostly explained by finite-sample error in  $\hat{\beta}_y$ . This has led to the introduction of a *no measurement error* (NOME) assumption in MR studies (Bowden et al., 2016b, 2018), under which one assumes finite sample error is only due to error in  $\hat{\beta}_y$ . For complex choices of  $\Phi$ , it may also be the case that  $p$ -values for testing  $(\hat{\beta}_\Phi)_{ij} = 0$  are lower than those for testing  $(\hat{\beta}_y)_i = 0$ , which would justify the NOME assumption. We have adapted `budgetIV` to construct confidence sets for  $\theta$  with (asymptotically) valid probabilities to cover the entire solution set  $\{\theta \in \mathbb{R}^{d_\Phi} : h(\theta) \in \Gamma(\tau, \mathbf{b})\}$  under the NOME assumption.

In particular, we use a Bonferroni adjustment to construct a box-shaped confidence set over  $\beta_y$ . For a target coverage  $(1 - \alpha) \times 100\%$ , we take a union over  $(1 - \alpha/d_Z) \times 100\%$  confidence intervals corresponding to each  $(\beta_y)_j$ . Let  $(\delta\beta_y)_i$  denote the half width of the confidence interval over  $(\beta_y)_i$ . While box-shaped confidence sets are conservative, they neatly superimpose with the  $\tau$ -thresholds:

$$\begin{aligned} \left| (\hat{\beta}_y)_i \pm (\delta\beta_y)_i - (\theta \cdot \hat{\beta}_\Phi)_i \right| &\leq \tau \\ \iff \left| (\hat{\beta}_y)_i - (\theta \cdot \hat{\beta}_\Phi)_i \right| &\leq \tau + (\delta\beta_y)_i. \end{aligned}$$

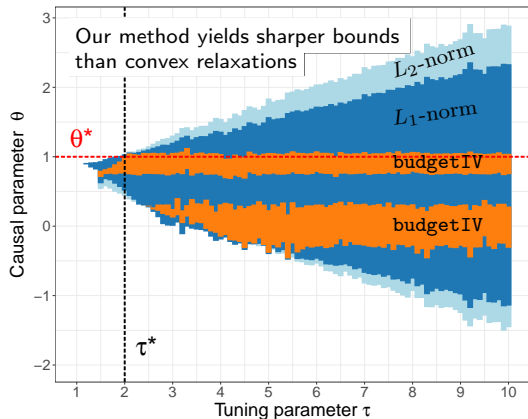


Figure 3: **Our method yields sharper bounds than convex relaxations in linear models.** Bounds on  $\theta$  for a series of linear models with scalar exposure  $X$  and  $\gamma_g^* = (-2, -0.4)$ . Plug-in estimators are used throughout. Orange bounds come from `budgetIV` with  $\tau = (0.6, \tau)$ ; dark blue from the  $L_1$ -norm constraint  $\|\gamma_g\|_1 \leq \tau + 0.6$ ; and light blue from the  $L_2$ -norm constraint  $\|\gamma_g\|_2 \leq \sqrt{\tau^2 + 0.6^2}$ . We vary  $\tau$  linearly from 0 to 10—each bound is an experiment.

We use this relationship to, we add slack  $(\pm\delta\beta_y)_i$  to the corresponding face of each  $\tilde{\Gamma}_{\tilde{U}}$  for  $\tilde{U} \in \Sigma^{(\max)}$  to form  $\hat{\Gamma}_{\tilde{U}}$ , and optimize over these domains to construct corresponding confidence intervals for the ATE( $\mathbf{x}; \mathbf{x}_0$ ) (see Sect. 3.1). Our approach has the following guarantee.

**Theorem 4 (Coverage).** Fix the target level  $\alpha \in (0, 1)$ . Let  $\hat{\Gamma}_\alpha := \bigcup_{\tilde{U} \in \Sigma^{(\max)}} \hat{\Gamma}_{\tilde{U}}$  be formed from the  $(1 - \alpha/d_Z) \times 100\%$  confidence intervals for each component of  $\hat{\beta}_y$  as described above, where the intervals are estimated from a dataset of  $n$  samples drawn iid from  $P(\mathbf{X}, Y, \mathbf{Z})$ . Using the shorthand:

$$\hat{\mathcal{T}}_\alpha := \mathcal{T}\left(c = \mathbb{I}[\gamma_g \in \hat{\Gamma}_\alpha], s = (\beta_\Phi^*, \hat{\beta}_y)\right),$$

we have, as  $n \rightarrow \infty$ :

$$P(\theta^* \in \hat{\mathcal{T}}_\alpha) \geq 1 - \alpha.$$

## 5 EXPERIMENTS

The full simulation studies are detailed in Appx. D, along with additional results.

**Sharper than convex relaxations.** In Fig. 3 we study the effect of varying the thresholds  $\tau$  on the feasible set. We consider a linear Gaussian model with  $d_X = 1$  and  $d_Z = 2$  where (A2) is violated through correlation between  $\mathbf{Z}$  and  $\epsilon_y$ , while (A3) is satisfied. We fix the ground truth parameters  $\theta^* = 1$ ,  $\gamma_g^* =$

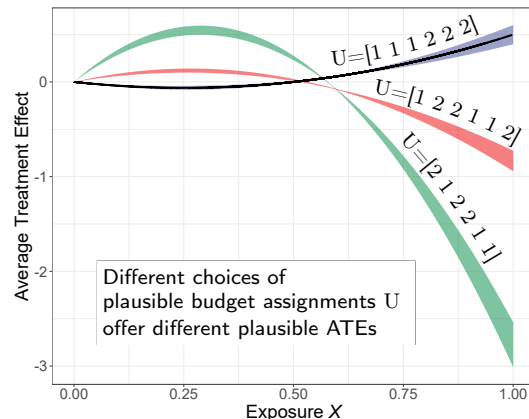


Figure 4: **Budget constraints provide information about the structure of the problem.** Feasible values of the ATE relative to a baseline of  $x_0 = 0$  as exposure  $X$  varies in a nonlinear SEM with  $d_Z = 6$ . The true ATE is given by the solid black curve. Each colored region corresponds to a unique intersection of  $\gamma_g$  and the star domain  $\Gamma$ . The union of such intersections at each value of  $X$  produces a disconnected feasible set.

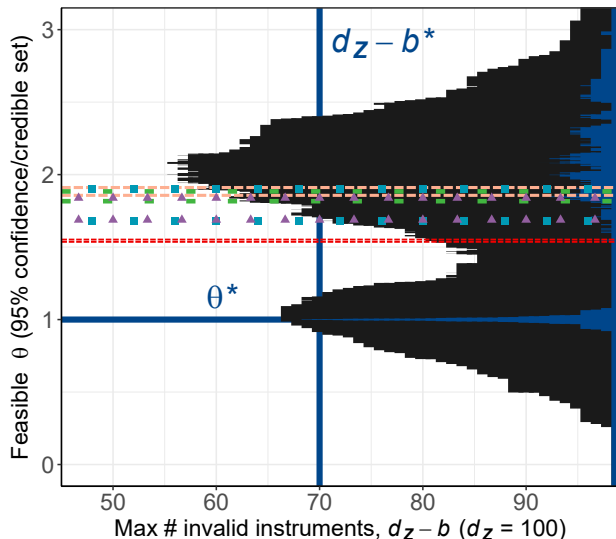
$(-2, 0.4)$  and  $\beta_\Phi^* = (2, -4)$ . The remaining parameters are randomized while  $\Gamma(\tau, b = 1)$  is varied.

The results illustrate that if the degree of validity for some instruments can be bounded, `budgetIV` may be insensitive to weak constraints on other instruments. By contrast, when used for partial identification, the convex relaxations suffer from these bottlenecks of limited background knowledge, and the ATE bounds grow linearly as a result.

**Budget constraints highlight possible mechanisms.** Fig. 4 shows bounds on the ATE for feasible values of  $U$  with a quadratic ground truth  $\Phi^*(X)$  under the presence of (A2) and (A3) violations. We study a collection of violations in Appx. D.

Feasible sets returned by `budgetIV` represent a union of convex bounds, each of which corresponds to a unique causal hypothesis that cannot be determined by the data or budget constraints alone (see Fig. 4). This drastically reduces the search space of possible causal mechanisms for the IV candidates. If one or more of these solutions can be ruled out by expert knowledge—e.g., if a monotonicity assumption is justified (Angrist et al., 1996)—then the method can be rerun with added constraints, pruning the search space still further. In this way, `budgetIV` can help practitioners evaluate causal systems in a dynamic and principled manner, aiding in hypothesis generation and experimental design.

In Fig. 4 the ground truth corresponds to the only



Method	Confidence/Credible Set	$\ni \theta^*$
MR-Egger	[1.680, 1.898]	✗
MR-Median	[1.817, 1.884]	✗
MBE	[1.857, 1.911]	✗
IVW	[1.688, 1.839]	✗
MASSIVE	[1.536, 1.552]	✗
budgetIV	$[0.914, 1.160] \cup [1.705, 2.403]$	✓
Oracle	[0.997, 1.003]	✓

Figure 5: **budgetIV captures the true causal effect when most candidates are invalid IVs.** Results from a simulation study with  $d_Z = 100$  candidate IVs, 70 of which violate (A3), benchmarking (Black) 95% coverage of the feasible set according `budgetIV_scalar`, where the budget constraints  $\Gamma(\tau = 0, b)$  are varied along the  $x$ -axis. (Blue) The optimal solution set relative to the constraint  $\Gamma(\tau = 0.001, b)$  (for visibility) captures the true causal effect if and only if the choice of  $b$  doesn't exclude the ground truth  $\gamma_g^*$ . (Others) Confidence intervals for benchmark methods produce do not include  $\theta^*$ .

plausible ATE with positive convexity. With additional expert knowledge, one might infer that  $Z_1, Z_2$ , and  $Z_3$  are the true set of valid instruments.

**Consistent inference under relaxed assumptions.** Fig. 5 shows a benchmarking experiment with  $d_Z = 100$  candidate IVs, 30 of which are valid and the remaining 70 violate (A3). The data is simulated via a linear with a scalar exposure  $X$  and multivariate Gaussian  $(\epsilon_x, \epsilon_y, \mathbf{Z} := \epsilon_z)$ . Linearity, (A3) violation, scalar  $X$  and Gaussian exogenous variables reflect the modeling assumptions in the benchmarking methods. We apply `budgetIV_scalar` for computational efficiency.

We apply a series of budget constraints  $\Gamma(\tau = 0, b)$

where  $b$  takes each value from 1 to 100 ( $b = 0$  corresponds to no constraint). We calculate  $\delta\beta_y$  from a 95% confidence set using 100,000 samples, which is a typical GWAS sample size.

The resulting confidence sets decompose into three kinds: ( $d_Z - b < 56$ ) empty, falsifying the budget constraints; ( $56 \leq d_Z - b < 67$ ) nonempty but excluding  $\theta^*$ ; and ( $d_Z - b \geq 67$ ) containing  $\theta^*$ . Whenever  $b^* < b$  (so that  $\gamma_g^* \in \Gamma(0, b)$ ) our confidence set captures  $\theta^*$ . Similarly, the solution set, calculated using  $\beta_y^*$  and  $\beta_\Phi^*$ , is shown to contain  $\theta^*$  if and only if  $b < b^*$ .

The methods we benchmark against give biased results. Inverse variance weighting (IVW) is a classical method that relies on all candidates being valid IVs. MR-median (Bowden et al., 2016a) assumes median  $(\gamma_{g_i}^*)_{i \in [d_Z]} = 0$ . MR-Egger assumes candidates are invalid through independent mechanisms (e.g., not affecting  $Y$  via shared mediators or being confounded with  $Y$  via shared confounders). These assumptions do not hold in our simulation study so we expect the results to be biased. Intriguingly, the stated assumptions for the Bayesian approach MASSIVE (Bucur et al., 2020) and the mode based estimate MBE (Hartwig et al., 2017) hold in our experiment.

The MASSIVE estimator assumes an  $L_0$ -norm constraint  $\Gamma(\tau = 0, b)$  with a tuneable value of  $b$ , which we set to  $b^*$ . We have shown  $\theta^*$  is only partially identified under such constraints. MASSIVE applies Bayesian model averaging over plausible sets of valid IVs. The resulting posterior distribution sits between the two disjoint confidence intervals returned by `budgetIV_scalar`, which themselves correspond to different sets of valid IVs.

MBE assumes mode  $(\gamma_{g_i}^*)_{i \in [d_Z]} = 0$ . Fig. 5 shows that this assumption holds because the optimal solution set peaks at  $\theta = \theta^* := 1$ . They estimate the corresponding modal causal effect  $\theta$  using the summary statistics  $\hat{\beta}_y$ ,  $\hat{\beta}_\Phi$ , standard errors and the bandwidth selection rule of Bickel and Levina (2008). Performing statistical inference on an estimator of the mode is not straightforward (Genovese et al., 2015). Hartwig et al. (2017) base their confidence intervals a normal approximation bootstrap. However, mode estimators generally do not converge to a normal distribution, which may explain the underestimated uncertainty. Indeed Fig. 7 in Appx. D.3 shows an expanded grid of experiments under which the estimates from MBE are highly variable.

## 6 RELATED WORK

There is a substantial literature on partial identification for IV models with restricted outcome domains. Early work includes seminal papers by Manski (1990)



and Balke and Pearl (1997). Later work considered bounded violations of the IV assumptions (Ramsahai, 2012; Silva and Evans, 2016; Jiang and Kocaoglu, 2024). Others have proposed generic methods for bounding counterfactual probabilities in discrete settings—either via polynomial programs (Duarte et al., 2023; Michael C. Sachs and Gabriel, 2023) or Markov chain Monte Carlo (Zhang et al., 2022)—with applications to IV models.

In the continuous setting, partial identification for pseudo-IV models has typically been formulated with respect to some convex relaxation of either (A2) or (A3), in both linear (Conley et al., 2012; Watson et al., 2024) and nonlinear SEMs (Newey and Powell, 2003; Gunsilius, 2019). In recent years, several authors have described more generic solutions based on stochastic gradient descent (Kilbertus et al., 2020; Hu et al., 2021; Padh et al., 2023). Unlike the linear programming approach of `budgetIV`, these methods are not guaranteed to converge on global optima.

In the MR literature, various methods are designed to try to handle linkage disequilibrium and/or pleiotropy. One strategy is to include a large number of candidate instruments and assume that biases will tend to cancel out in the limit (Kolesár et al., 2015; Bowden et al., 2015). Others take a feature selection approach in which  $Z$ 's may be rejected on the basis of statistical tests (Chu et al., 2001; Kang et al., 2020) or  $L_1$ -penalized regression (Kang et al., 2016; Guo et al., 2018; Windmeijer et al., 2019; Xue et al., 2023). Alternatively, instruments may be pooled into a single feature (Kuang et al., 2020). An especially popular choice is the modal validity assumption (Hartwig et al., 2017; Hartford et al., 2021), which asserts that the most common causal effect estimate is consistent. The goal in these works is point identification, which may be unrealistic if underlying assumptions fail.

Bayesian methods for causal inference in IV models are well-established. Priors can be used to encode uncertainty with regard to latent parameters that track either (A2) (Shaplund et al., 2019) or (A3) violations (Bucur et al., 2020; Gkatzionis et al., 2021).

In recent work, Vancak and Sjölander (2023) define a single sensitivity parameter for violation of (A2) and/or (A3). We extend this approach by using a latent statistic of a nonparametric function. Star-domain restrictions for partial identification have also been proposed (Molinari, 2008), though not as a guiding principle.

For a summary of the constraints and affordances of various relaxed IV methods, see Appx. F, Table 1.

## 7 DISCUSSION

The `budgetIV` optimization problem is NP-hard for any  $d_{\Phi} > 1$ . Therefore, with multivariate  $\Phi$ , the method becomes impractical for large  $d_Z$ . Approximate solutions that rely on grid search may be preferable in such cases. Many MR studies are built on just one or a handful of genetic variants, so a cap on  $d_Z$  may not be overly restrictive in such settings.

The additive separability of causal effects in pseudo-IV models has been assumed by various authors (Newey and Powell, 2003; Saengkyongam et al., 2022; Christiansen et al., 2021). We have considered the special case of a homogeneous treatment effect, though our method can be further generalized if  $\Phi$  is promoted to known functions of  $\mathbf{X}$ ,  $\mathbf{Z}$  and/or  $\epsilon_Z$ . Extending our approach to compute feasible sets of *conditional* average treatment effects may result in more informative outputs (Cai et al., 2007; Hartford et al., 2021; Levis et al., 2023).

We rely on the assumption  $d_{\Phi} \leq d_Z$  and assume our choice of  $\Phi$  is sufficient to describe the ground truth causal effect of  $\mathbf{X}$  on  $Y$  exactly. Future work could investigate error arising from misspecification of  $\Phi$  or incorporate thresholds on this kind of error into the background assumptions. In particular, basis expansions that have been truncated to satisfy the bound on  $d_{\Phi}$  may be of interest.

While we have presented our method as inferring the ATE from the distribution of instrument validity, the correspondence can be thought of as bidirectional. One promising direction for future work is to invert such methods for the purpose of instrument discovery (Silva and Shimizu, 2017). Another direction could be to use background knowledge about the functional form of the ATE (when  $d_{\Phi} > 1$ ) to restrict the solution set.

There are several other extensions to `budgetIV` that could be of interest. Sharper partial identification may be achieved in models with restricted outcome domains. The finite-sample properties of our method might also be improved. We used a Bonferroni correction to construct the confidence set over  $\beta_y$ . Since this is conservative, a future direction may involve adaptively selecting the confidence thresholds for each  $(\beta_y)_i$  to minimize the width of the bounds on the ATE. On the other hand, coverage under finite sample uncertainty without the NOME assumption remains an open problem.

## References

- Ailer, E., Hartford, J., and Kilbertus, N. (2023). Sequential underspecified instrument selection for cause-effect estimation. In *Proceedings of the 40th International Conference on Machine Learning*, pages

408–420.

- Angrist, J. D., Imbens, G. W., and Rubin, D. B. (1996). Identification of causal effects using instrumental variables. *J. Am. Stat. Assoc.*, 91(434):444–455.
- Angrist, J. D. and Pischke, J. (2009). *Mostly Harmless Econometrics*. Princeton University Press.
- Balke, A. and Pearl, J. (1997). Bounds on treatment effects from studies with imperfect compliance. *J. Am. Stat. Assoc.*, 92(439):1171–1176.
- Bickel, P. J. and Levina, E. (2008). Regularized estimation of large covariance matrices. *Ann. Stat.*, 36(1):199 – 227.
- Bowden, J., Davey Smith, G., and Burgess, S. (2015). Mendelian randomization with invalid instruments: effect estimation and bias detection through Egger regression. *Int. J. Epidemiol.*, 44(2):512–525.
- Bowden, J., Davey Smith, G., Haycock, P. C., and Burgess, S. (2016a). Consistent estimation in mendelian randomization with some invalid instruments using a weighted median estimator. *Genet. Epidemiol.*, 40(4):304–314.
- Bowden, J., Del Greco M, F., Minelli, C., Davey Smith, G., Sheehan, N. A., and Thompson, J. R. (2016b). Assessing the suitability of summary data for two-sample Mendelian randomization analyses using MR-Egger regression: the role of the I<sup>2</sup> statistic. *Int. J. Epidemiol.*, 45(6):1961–1974.
- Bowden, J., Del Greco M, F., Minelli, C., Zhao, Q., Lawlor, D. A., Sheehan, N. A., Thompson, J., and Davey Smith, G. (2018). Improving the accuracy of two-sample summary-data mendelian randomization: moving beyond the NOME assumption. *Int. J. Epidemiol.*, 48(3):728–742.
- Bowden, R. J. and Turkington, D. A. (1984). *Instrumental variables*. Cambridge University Press, Cambridge.
- Bucur, I. G., Claassen, T., and Heskes, T. (2020). MASSIVE: Tractable and robust Bayesian learning of many-dimensional instrumental variable models. In *Proceedings of the 36th Conference on Uncertainty in Artificial Intelligence*, pages 1049–1058.
- Cai, Z., Kuroki, M., and Sato, T. (2007). Non-parametric bounds on treatment effects with non-compliance by covariate adjustment. *Stat. Med.*, 26(16):3188–3204.
- Christiansen, R., Pfister, N., Jakobsen, M., Gnecco, N., and Peters, J. (2021). A causal framework for distribution generalization. *IEEE transactions on pattern analysis and machine intelligence*.
- Chu, T., Scheines, R., and Spirtes, P. L. (2001). Semi-instrumental variables: A test for instrument admissibility. In *Proceedings of the 17th Conference on Uncertainty in Artificial Intelligence*, pages 83–90.
- Conley, T. G., Hansen, C. B., and Rossi, P. E. (2012). Plausibly exogenous. *Rev. Econ. Stat.*, 94(1):260–272.
- Duarte, G., Finkelstein, N., Knox, D., Mummolo, J., and Shpitser, I. (2023). An automated approach to causal inference in discrete settings. *J. Am. Stat. Assoc.*, pages 1–16.
- Fisher, R. A. (1935). *The Design of Experiments*. Oliver and Boyd.
- Genovese, C. R., Perone-Pacifico, M., Verdinelli, I., and Wasserman, L. (2015). Non-parametric inference for density modes. *Journal of the Royal Statistical Society Series B: Statistical Methodology*, 78(1):99–126.
- Gkatzionis, A., Burgess, S., Conti, D. V., and Newcombe, P. J. (2021). Bayesian variable selection with a pleiotropic loss function in Mendelian randomization. *Stat. Med.*
- Gunsilius, F. (2019). A path-sampling method to partially identify causal effects in instrumental variable models. *arXiv:1910.09502*.
- Guo, Z., Kang, H., Tony Cai, T., and Small, D. S. (2018). Confidence intervals for causal effects with invalid instruments by using two-stage hard thresholding with voting. *J. R. Stat. Soc. Series B Stat. Methodol.*, 80(4):793–815.
- Hartford, J. S., Veitch, V., Sridhar, D., and Leyton-Brown, K. (2021). Valid causal inference with (some) invalid instruments. In *Proceedings of the 38th International Conference on Machine Learning*, pages 4096–4106.
- Hartwig, F. P., Davey Smith, G., and Bowden, J. (2017). Robust inference in summary data Mendelian randomization via the zero modal pleiotropy assumption. *Int. J. Epidemiol.*, 46(6):1985–1998.
- Heckman, J. J. and Vytlacil, E. (2005). Structural equations, treatment effects, and econometric policy evaluation. *Econometrica*, 73(3):669–738.
- Heckman, J. J. and Vytlacil, E. J. (1999). Local instrumental variables and latent variable models for identifying and bounding treatment effects. *Proc. Natl. Acad. Sci.*, 96(8):4730–4734.
- Holland, P. W. (1986). Statistics and causal inference. *J. Am. Stat. Assoc.*, 81(396):945–960.
- Hu, Y., Wu, Y., Zhang, L., and Wu, X. (2021). A generative adversarial framework for bounding confounded causal effects. In *Proceedings of the 35th AAAI Conference on Artificial Intelligence*.

- Imbens, G. W. and Angrist, J. D. (1994). Identification and estimation of local average treatment effects. *Econometrica*, 62(2):467–475.
- Jiang, Z. and Kocaoglu, M. (2024). Conditional common entropy for instrumental variable testing and partial identification. In *Proceedings of the 41st International Conference on Machine Learning*, volume 235 of *Proceedings of Machine Learning Research*, pages 21824–21843. PMLR.
- Kang, H., Lee, Y., Cai, T. T., and Small, D. S. (2020). Two robust tools for inference about causal effects with invalid instruments. *Biometrics*.
- Kang, H., Zhang, A., Cai, T. T., and Small, D. S. (2016). *J. Am. Stat. Assoc.*, 111(513):132–144.
- Kilbertus, N., Kusner, M. J., and Silva, R. (2020). A class of algorithms for general instrumental variable models. In *Advances in Neural Information Processing Systems*.
- Kolesár, M., Chetty, R., Friedman, J., Glaeser, E., and Imbens, G. W. (2015). Identification and inference with many invalid instruments. *J. Bus. Econ. Stat.*, 33(4):474–484.
- Koopmans, T. C. (1949). Identification problems in economic model construction. *Econometrica*, 17(2):125–144.
- Kuang, Z., Sala, F., Sohoni, N., Wu, S., Córdova-Palomera, A., Dunmmon, J., Priest, J., and Re, C. (2020). Ivy: Instrumental variable synthesis for causal inference. In *Proceedings of the 23rd International Conference on Artificial Intelligence and Statistics*, page 398–410.
- Levis, A. W., Bonvini, M., Zeng, Z., Keele, L., and Kennedy, E. H. (2023). Covariate-assisted bounds on causal effects with instrumental variables. *arXiv preprint*, 2301.12106.
- Malloy, M. L., Tripathy, A., and Nowak, R. D. (2021). Optimal confidence regions for the multinomial parameter. *arXiv preprint*, 2002.01044.
- Manski, C. F. (1990). Nonparametric bounds on treatment effects. *Am. Econ. Rev.*, 80(2):319–323.
- Manski, C. F. (2003). *Partial identification of probability distributions*. Springer, New York.
- Michael C. Sachs, Gustav Jonzon, A. S. and Gabriel, E. E. (2023). A general method for deriving tight symbolic bounds on causal effects. *J. Comput. Graph. Stat.*, 32(2):567–576.
- Molinari, F. (2008). Partial identification of probability distributions with misclassified data. *Journal of Econometrics*, 144(1):81–117.
- Nevo, A. and Rosen, A. M. (2012). Identification with imperfect instruments. 94(3):659–671.
- Newey, W. and Powell, J. (2003). Instrumental variable estimation of nonparametric models. *Econometrica*, 71(5):1565–1578.
- Neyman, J. (1923). On the Application of Probability Theory to Agricultural Experiments. Essay on Principles. Section 9. *Statistical Science*, 5(4):465 – 472.
- Padh, K., Zeitler, J., Watson, D., Kusner, M., Silva, R., and Kilbertus, N. (2023). Stochastic causal programming for bounding treatment effects. In *Proceedings of the 2nd Conference on Causal Learning and Reasoning*, pages 142–176.
- Pearl, J. (2009). *Causality*. Cambridge University Press, New York.
- Pfister, N. and Peters, J. (2022). Identifiability of sparse causal effects using instrumental variables. In *Proceedings of the 38th Conference on Uncertainty in Artificial Intelligence*, pages 1613–1622.
- Pierce, B. L. and Burgess, S. (2013). Efficient design for mendelian randomization studies: Subsample and 2-sample instrumental variable estimators. *American Journal of Epidemiology*, 178(7):1177–1184.
- Ramsahai, R. R. (2012). Causal bounds and observable constraints for non-deterministic models. *J. Mach. Learn. Res.*, 13(29):829–848.
- Rubin, D. (1974). Estimating causal effects of treatments in randomized and nonrandomized studies. *J. Educ. Psychol.*, 66.
- Saengkyongam, S., Henckel, L., Pfister, N., and Peters, J. (2022). Exploiting independent instruments: Identification and distribution generalization. In *Proceedings of the 39th International Conference on Machine Learning*, pages 18935–18958.
- Shapland, C. Y., Thompson, J. R., and Sheehan, N. A. (2019). A Bayesian approach to Mendelian randomisation with dependent instruments. *Stat. Med.*, 38(6):985–1001.
- Silva, R. and Evans, R. (2016). Causal inference through a witness protection program. *J. Mach. Learn. Res.*, 17(56):1–53.
- Silva, R. and Shimizu, S. (2017). Learning instrumental variables with structural and non-Gaussianity assumptions. *J. Mach. Learn. Res.*, 18(120).
- Vaart, A. W. v. d. (1998). *Asymptotic Statistics*. Cambridge Series in Statistical and Probabilistic Mathematics. Cambridge University Press.
- Vancak, V. and Sjölander, A. (2023). Sensitivity analysis of g-estimators to invalid instrumental variables. *Statistics in Medicine*, 42(23):4257–4281.
- Watson, D. S., Penn, J., Gunderson, L. M., Bravo-Hermsdorff, G., Mastouri, A., and Silva, R. (2024).

Bounding causal effects with leaky instruments. In *Proceedings of the 40th Conference on Uncertainty in Artificial Intelligence*.

Windmeijer, F., Farbmacher, H., Davies, N., and Davey Smith, G. (2019). On the use of the lasso for instrumental variables estimation with some invalid instruments. *J. Am. Stat. Assoc.*, 114(527):1339–1350.

Wright, P. G. (1928). *The Tariff on Animal and Vegetable Oils*. Macmillan, New York.

Xue, H., Shen, X., and Pan, W. (2023). Causal inference in transcriptome-wide association studies with invalid instruments and GWAS summary data. *Journal of the American Statistical Association*, 118(543):1525–1537.

Zhang, J., Tian, J., and Bareinboim, E. (2022). Partial counterfactual identification from observational and experimental data. In *Proceedings of the 39th International Conference on Machine Learning*, volume 162, pages 26548–26558.

## A PROOFS

### A.1 Proof of Thm. 1

This is a special case of a standard result from the classical IV literature, which goes back to [Koopmans \(1949\)](#):  $\theta$  has a unique solution under the assumptions (B1\*) and (B2\*) because  $\beta_y - \theta \cdot \beta_\Phi = \mathbf{0}$  is a complete system of simultaneous equations. With finite samples, it is common to estimate  $\theta$  solving this equation via two-stage least squares (2SLS) or (for  $d_\Phi = 1$ ) inverse variance weighting (IVW). The 2SLS estimator is given by:

$$\theta_{(2\text{SLS})} := (\beta_\Phi \cdot \Sigma_{ZZ}^{-1} \cdot \beta_\Phi^\top)^{-1} \cdot \beta_\Phi \cdot \Sigma_{ZZ}^{-1} \cdot \beta_y,$$

where  $\Sigma_{ZZ}$  denotes the  $d_Z \times d_Z$  covariance matrix between the instruments. For details, see [Wright \(1928\)](#); [Bowden and Turkington \(1984\)](#); [Angrist and Pischke \(2009\)](#).

### A.2 Proof of Thm. 2

Let the structural Eqs. (1) to (3) hold for some ground truth functions  $f_z^*$ ,  $f_x^*$ ,  $\Phi^*$ ,  $g_y^*$ , and causal parameter  $\theta^* \in \mathbb{R}^{d_\Phi}$ .

**Soundness** We prove soundness of the feasible map  $\mathcal{T}$  for which:

$$\mathcal{T}(c = \mathbb{I}[\gamma_g \in \Gamma], s = (\beta_\Phi, \beta_y)) = \{\theta \in \mathbb{R}^{d_\Phi} : h(\theta) \in \Gamma\},$$

where  $h(\theta) = \beta_y - \theta \cdot \beta_\Phi$ . For ease of notation, we use the shorthand  $\mathcal{T}(\Gamma, \beta_\Phi, \beta_y)$  hereinafter.

It follows immediately from Eq. (3) and the left-linearity of the covariance operator that:

$$\begin{aligned} \gamma_g^* &:= \text{Cov}(g_y(\mathbf{Z}, \epsilon_y), \mathbf{Z}) = \text{Cov}(Y, \mathbf{Z}) - \text{Cov}(\theta^* \cdot \Phi(\mathbf{X}), \mathbf{Z}) \\ &= \beta_y^* - \theta^* \cdot \beta_\Phi^*. \end{aligned}$$

Therefore,  $\theta^* \in \mathcal{T}(\Gamma, \beta_\Phi^*, \beta_y^*)$  whenever  $\gamma_g^* = (\beta_y^* - \theta^* \cdot \beta_\Phi^*) \in \Gamma$  as required.

**Completeness** To prove completeness of  $\mathcal{T}$  we have to show that any  $\theta \in \mathcal{T}(\Gamma, \beta_\Phi, \beta_y)$  cannot be rejected by any statistic of the observed joint distribution  $P(\mathbf{Z}, \mathbf{X}, Y)$ . Notice that in specifying Eqs. (1) to (3), we have not made any a priori assumptions about the joint distribution  $P_\epsilon(\epsilon_z, \epsilon_x, \epsilon_y)$  or the function classes to which  $f_z^*$ ,  $f_x^*$ ,  $g_y^*$  belong. In Appx. B.1 we extend this proof to find conditions under which  $\mathcal{T}$  remains complete when further structural assumptions are made.

Consider any  $\theta^\dagger \in \mathbb{R}^{d_\Phi}$ . Then the following holds:

(Z) There exists at least one function  $f_z^\dagger$  with the following property. Given any  $z \in \Omega_Z$ , either: (i)  $P(\mathbf{x}, y, z) = 0$  for all  $\mathbf{x} \in \Omega_X$  and  $y \in \Omega_Y$ ; or (ii) there exists at least one value  $\epsilon_z^\dagger$  that solves the equation:

$$f_z^\dagger(\epsilon_z^\dagger) = z.$$

Note that the ground truth  $f_z^*$  is one valid choice of  $f_z^\dagger$ .

We can, therefore, define a function  $\epsilon_z^\dagger(z)$  that satisfies the above equation for all  $z$  for which case (i) is false. We have not demanded either  $f_z^\dagger$  or  $\epsilon_z^\dagger(z)$  to be unique.

(X) Likewise, there exists at least one function  $f_x^\dagger$  with the following property. Given any  $\mathbf{x} \in \Omega_X$ ,  $z \in \Omega_Z$ , either: (i)  $P(\mathbf{x}, y, z) = 0$  for all  $y \in \Omega_Y$ , or (ii) there exists at least one value  $\epsilon_x^\dagger$  that solves the equation:

$$f_x^\dagger(z, \epsilon_x^\dagger) = \mathbf{x},$$

and  $f_x^*$  is one valid choice of  $f_x^\dagger$ . We can, therefore, define at least one function  $\epsilon_x^\dagger(z, \mathbf{x})$  that satisfies the above equation for all  $\mathbf{x}$  and  $z$  for which case (i) is false.

(Y) Likewise, for any  $\mathbf{x} \in \Omega_X$ ,  $y \in \Omega_Y$ ,  $z \in \Omega_Z$ , either (i)  $P(\mathbf{x}, y, z) = 0$ , or (ii) for any  $\theta^\dagger \in \mathbb{R}^{d_\Phi}$  there exists at least function  $g_y^\dagger$  and one value  $\epsilon_y^\dagger$  that solves the equation:

$$g_y^\dagger(z, \epsilon_y^\dagger) = y - \theta^\dagger \cdot \Phi(\mathbf{x}).$$

For instance we could have  $g_y(z, \epsilon_y) = \epsilon_y$  and  $\epsilon_y^\dagger := y - \theta^\dagger \cdot \Phi(\mathbf{x})$  is one such choice. We can therefore define a function  $\epsilon_y^\dagger(z, \mathbf{x}, y; \theta^\dagger)$  for each  $\theta^\dagger \in \mathbb{R}^{d_\Phi}$  that solves the above equation for any  $\mathbf{x}, y$  and  $z$ .

The full joint distribution that generates  $P(\mathbf{X}, Y, \mathbf{Z})$  from the structural equations can be factorized as follows:

$$P(\mathbf{X}, Y, \mathbf{Z}, \epsilon_x, \epsilon_y) = P(\epsilon_x, \epsilon_y \mid \mathbf{X}, Y, \mathbf{Z}) P(\mathbf{X}, Y, \mathbf{Z}).$$

We can therefore define (for any  $\theta^\dagger$ ) at least one joint distribution consistent with the structural assumptions,  $\theta^\dagger$ , and the observed joint distribution  $P(\mathbf{X}, Y, \mathbf{Z})$ :

$$P^\dagger(\mathbf{X}, Y, \mathbf{Z}, \epsilon_x, \epsilon_y) := D^\dagger(\epsilon_x \mid \mathbf{X}, \mathbf{Z}) \delta(\epsilon_y - \epsilon_y^\dagger(\mathbf{X}, Y, \mathbf{Z}; \theta)) P(\mathbf{X}, Y, \mathbf{Z}),$$

where:

$$D^\dagger(\epsilon_x \mid \mathbf{x}, \mathbf{z}) = \begin{cases} \delta(\epsilon_x - \epsilon_x^\dagger(\mathbf{x}, \mathbf{z})) & \exists y \in \Omega_Y : P(\mathbf{x}, y, \mathbf{z}) \neq 0 \\ 0 & \text{otherwise,} \end{cases}$$

and  $\delta$ 's are Dirac delta measures over the domains of their arguments.

Therefore, we cannot reject any  $\theta \in \mathbb{R}^{d_\Phi}$  based on the structural equations and the ground truth  $P(\mathbf{X}, Y, \mathbf{Z})$  alone, since there exists some joint distribution consistent with both.

Assume the validity of a background constraint  $\Gamma$ , so that  $\text{Cov}(\mathbf{Z}, g_y) \equiv (\beta_y^* - \theta^* \cdot \beta_\Phi^*) \in \Gamma$ . We can only use  $\Gamma$  to exclude  $\theta$  for which  $(\beta_y^* - \theta \cdot \beta_\Phi^*) \notin \Gamma$ . Therefore, it follows immediately that the feasible map  $\mathcal{T}$  for which  $\mathcal{T}(\Gamma, \beta_\Phi^*, \beta_y^*) = \{\theta : (\beta_y^* - \theta \cdot \beta_\Phi^*) \in \Gamma\}$  is complete.

**Minimality** We have proven that any and all  $\theta \in \mathcal{T}(\Gamma, \beta_\Phi, \beta_y)$  satisfy  $(\beta_y - \theta \cdot \beta_\Phi) \in \Gamma$ . Since  $\Gamma, \beta_\Phi$  and  $\beta_y$  can be varied arbitrarily and independently of each other, it is clear the  $d_{\mathbf{Z}d_\Phi} + d_{\mathbf{Z}}$  independent, real parameters required to specify  $\beta_\Phi$  and  $\beta_y$  are needed to specify the feasible map  $\mathcal{T}$ . This completes the proof of optimality.

### A.3 Proof of Thm. 3

We prove a slight generalization of the theorem that accounts for violation of (B1\*). As we will see, the the tight bound depends on a quantity  $B$  that equals 0 when (B1\*) is satisfied.

#### We can account for violation of (B1\*) by defining a reduced problem

We extend the theorem as stated by allowing for (B1\*) violation. By simple linear algebra, the left and right null spaces and accounting for (B1\*) violation can be identified polynomial time.

If  $\beta_\Phi$  has a nonempty left null space, then each point  $h(\theta) \in h := h[\theta]$  (i.e., the affine space that is the image of the affine map  $h(\theta)$ ) corresponds to a continuum of possible  $\theta \in \Theta \subseteq \mathbb{R}^{d_\Phi}$ , where  $\Theta$  is unbounded. This means any nonempty feasible set  $\mathcal{T}_{L_0}$  will also be unbounded, and depending on the choice of  $\Phi$ , the ATE may be vacuous. We choose to ignore this kind of violation of (B1\*), which corresponds to  $\beta_\Phi$  not being full rank.

We decompose the set of covariance-irrelevant candidate instruments  $I = \{i \in [d_{\mathbf{Z}}] : (\beta_\Phi)_i = 0\}$  into the following subsets:

$$\begin{aligned} I_{=0} &:= \{i \in [d_{\mathbf{Z}}] : (\beta_\Phi)_i = 0 \wedge (\beta_y)_i = 0\}, \\ I_{\neq 0} &:= \{i \in [d_{\mathbf{Z}}] : (\beta_\Phi)_i = 0 \wedge (\beta_y)_i \neq 0\}, \end{aligned}$$

where the former corresponds to irrelevant candidate instruments that are uncorrelated with the outcome while the latter corresponds to those correlated with the outcome. Any  $\gamma_g \in h := h[\theta]$  will satisfy  $(\gamma_g)_{i \in I_{=0}} = 0$  and  $(\gamma_g)_{i \in I_{\neq 0}} \neq 0$ . Therefore, we can count whether these components are always or never 0 irrespective of  $\theta$ .

This motivates the following definitions:

$$\begin{aligned} D_{\mathbf{Z}} &:= d_{\mathbf{Z}} - |I_{=0}| - |I_{\neq 0}|, \\ B &:= b - |I_{\neq 0}|, \end{aligned}$$

through which we define the reduced problem of finding  $\theta$  for which:

$$\|H(\theta)\|_0 \leq B,$$

where:

$$H(\boldsymbol{\theta}) = \sum_{j \in [d_{\mathbf{Z}}] \setminus I} ((\beta_y)_i - \boldsymbol{\theta} \cdot (\boldsymbol{\beta}_{\Phi})_i) \mathbf{e}_i.$$

The solutions for  $\boldsymbol{\theta}$  are exactly the same solutions as those to  $\|h(\boldsymbol{\theta})\|_0 \leq b$ . Notice that  $H := H[\boldsymbol{\theta}]$  is a  $d_{\Phi}$ -dimensional affine subspace of  $\mathbb{R}^{D_{\mathbf{Z}}}$  (we can ignore  $i \in I$ ). Defining  $J = [d_{\mathbf{Z}}] \setminus I$ , we see that none of the basis vectors  $\mathbf{e}_j$  for  $j \in J$  are orthogonal to  $H$ .

Notice that  $D_{\mathbf{Z}} = d_{\mathbf{Z}}$ ,  $B = b$  and  $H(\boldsymbol{\theta}) = h(\boldsymbol{\theta})$  iff (B1\*) holds.

**We begin with  $d_{\Phi} = 1$**

If  $d_{\Phi} = 1$  then  $H$  represents a line embedded in the  $D_{\mathbf{Z}}$ -dimensional Euclidean space. Consider a point  $\boldsymbol{\gamma} := H(\boldsymbol{\theta}) \in H$  for some  $\boldsymbol{\theta} \in \mathbb{R}$ . This point satisfies  $\|\boldsymbol{\gamma}\|_0 \leq B$  iff there are at least  $D_{\mathbf{Z}} - B$  many values of  $j \in J$  for which  $\gamma_j = 0$ .

Since there are no  $j \in J$  for which  $H$  is orthogonal to  $J$ , we know there is exactly one solution  $\theta_{(j)}$  to the equation  $(\beta_y)_j - \theta(\beta_{\Phi})_j = 0$ . Thus the equivalence relation  $\sim$  defined by:

$$\forall i, j \in J : i \sim j \iff \theta_{(i)} = \theta_{(j)},$$

forms a partition over  $[J]$ . Each equivalence class  $\langle j \rangle$  of this partition must have cardinality at least  $D_{\mathbf{Z}} - B$  for  $\theta_{(j)}$  to be in the feasible set. Since  $|J| = D_{\mathbf{Z}}$ , the number of unique  $\theta_{(j)}$  solving the constraint,  $n$ , is bounded above by:

$$\left\lfloor \frac{D_{\mathbf{Z}}}{D_{\mathbf{Z}} - B} \right\rfloor. \quad (5)$$

This proves Cor. 3.2, including the special case of point identification when  $B > D_{\mathbf{Z}}/2$  first shown by Kang et al. (2016).

**For  $d_{\Phi} > 1$  we enumerate the plausible unique points satisfying the budget constraints**

**Consider demanding a single instrument is valid** In general,  $H$  corresponds to a  $d_{\Phi}$ -dimensional affine space, not orthogonal to  $\mathbf{e}_j$  for any  $j \in J$ . For some choice of  $j \in J$ ,  $H$  will intersect the  $(D_{\mathbf{Z}} - 1)$ -dimensional affine space  $\{\boldsymbol{\gamma} \in \mathbb{R}^{D_{\mathbf{Z}}} : \gamma_j = 0\}$  to form  $H_{\setminus j} := \{\boldsymbol{\gamma}' \in H : \gamma'_j = 0\}$ .  $H_{\setminus j}$  is necessarily a  $(d_{\Phi} - 1)$ -dimensional affine space by the rule for the dimension of the intersection of two affine spaces and the fact that  $\mathbf{e}_j \not\perp H$  (where  $\perp$  means orthogonal). Notice that  $\mathbf{e}_j \perp H_{\setminus j}$ .

It is possible that  $H_{\setminus j} = H_{\setminus k}$  for some  $k \in J \setminus j$ , which motivates introducing another equivalence relation  $\overset{d_{\Phi}}{\sim}$ :

$$\forall j, k \in J : j \overset{d_{\Phi}}{\sim} k \iff H_{\setminus j} = H_{\setminus k},$$

which induces a partition over  $J$ . The equivalence classes  $\langle j \rangle_{d_{\Phi}}$  represent a single  $(d_{\Phi} - 1)$ -dimensional affine space that consists of points  $\boldsymbol{\gamma}$  that simultaneously satisfy the constraints  $\gamma_k = 0$  for all  $k$  in the equivalence class. In other terms:

$$H_{\setminus j} = H_{\langle j \rangle_{d_{\Phi}}} := \{\boldsymbol{\gamma} \in H : \forall k \in \langle j \rangle_{d_{\Phi}} (\gamma_k = 0)\}.$$

The number of unique equivalence classes with respect to  $\overset{d_{\Phi}}{\sim}$ ,  $n_{d_{\Phi}}$ , is anywhere between 1 and  $D_{\mathbf{Z}}$ . However, if any of the equivalence classes has cardinality at least  $D_{\mathbf{Z}} - B$ , then there is at least one  $(d_{\Phi} - 1)$ -dimensional solution to the  $L_0$ -norm constraint.

**What if we demand more instruments are valid?** We could equally ask how many sets of constraints  $(\gamma_q)_{q \in Q} = 0$  for  $q \in Q \subseteq J$  lead to unique  $(d_{\Phi} - 2)$ -dimensional affine spaces  $H_{\setminus Q} := \{\boldsymbol{\gamma} \in H : \forall k \in Q (\gamma_k = 0)\}$ . The  $Q$  such that  $H_{\setminus Q}$  is a  $(d_{\Phi} - 2)$ -dimensional affine space is the union of at least two equivalence classes  $\langle j \rangle_{d_{\Phi}}$  and  $\langle k \rangle_{d_{\Phi}} \neq \langle j \rangle_{d_{\Phi}}$ . Therefore the number of unique  $(d_{\Phi} - 2)$ -dimensional affine spaces  $H_{\setminus Q}$ , denoted  $n_{d_{\Phi}-1}$ , is anywhere between 1 and  $\binom{n_{d_{\Phi}}}{2}$ .

Likewise, the number of sets  $Q$  generating unique 1-dimensional affine spaces (lines)  $H_{\setminus Q}$ , denoted  $n_1$ , is anywhere between 1 and  $\binom{d_{\Phi}}{d_{\Phi}-1}$ . This implies the upper bound:

$$n_1 \leq \binom{D_{\mathbf{Z}}}{d_{\Phi}-1} = \frac{D_{\mathbf{Z}}!}{(d_{\Phi}-1)!(D_{\mathbf{Z}}-d_{\Phi}+1)!},$$

where this upper bound corresponds to all such  $Q \subset J$  having length  $d_{\Phi}-1$ , which is the minimum length  $Q$  can take.

**Bounding the number of unique points for a given line** Consider some  $Q$  for which  $H_{\setminus Q} := \{\gamma \in H : \forall k \in Q(\gamma_k = 0)\}$  is a line. Then define  $Q' := \{i \in J : \forall \gamma \in H_{\setminus Q}(\gamma_k = 0)\}$ , which expresses every constraint satisfied along the line. We ask how many unique points along  $H_{\setminus Q}$  can satisfy the  $L_0$ -norm constraint  $\|H_{\setminus Q}\|_0 \leq b$ . Along the entire line, there are  $|Q'|$  many components for which  $\gamma_k = 0$ . For a point to satisfy the constraint, we must have at least  $D_{\mathbf{Z}} - |Q'|$  many  $\gamma_k = 0$  for  $k \in J \setminus Q'$ . Therefore, by the same arguments that lead to Eq. (5), the number of unique points solving the constraint along  $H_Q$ , denoted  $p_Q$ , is bounded by:

$$p_Q \leq \left\lfloor \frac{D_{\mathbf{Z}} - |Q'|}{D_{\mathbf{Z}} - |Q'| - B} \right\rfloor.$$

**Putting the pieces together** As discussed earlier, the minimum length of  $Q'$  is  $d_{\Phi}-1$ . This value maximizes  $p_Q$ . The maximum number of unique lines  $n_1$  also increases as the length  $|Q'|$  for each line decreases. Therefore, the upper bound:

$$n \leq \frac{D_{\mathbf{Z}}!}{(d_{\Phi}-1)!(D_{\mathbf{Z}}-d_{\Phi}+1)!} \left\lfloor \frac{D_{\mathbf{Z}} - d_{\Phi} + 1}{D_{\mathbf{Z}} - d_{\Phi} + 1 - B} \right\rfloor,$$

which corresponds to all unique lines  $H_Q$  having  $Q' = d_{\Phi}-1$ , is valid and tight.

#### A.4 Proof of Thm. 4

We begin with the assumption that the estimated statistic  $\hat{\beta}_{\Phi}$  is exactly equal to the observable ground truth  $\beta_{\Phi}^*$ . Furthermore, since elliptical confidence sets require estimating  $\text{Cov}\left((\hat{\beta}_y)_i, (\hat{\beta}_y)_j\right)$  for pairs  $i \neq j$ —estimates for these are not openly available with GWAS summary statistics—we construct box constraints using the Bonferroni correction. Our  $(1-\alpha) \times 100\%$  confidence set over  $\beta_{\mathbf{y}}$ , consists of all  $\beta_{\mathbf{y}}$  for which each component  $(\beta_y)_i$  is in the corresponding  $(1-\alpha/d_{\mathbf{Z}}) \times 100\%$  confidence interval, defined below.

Calculate the estimator  $\hat{\beta}_{\mathbf{y}} := \widehat{\text{Cov}}(Y, \mathbf{Z})$  using  $N$  i.i.d. samples from  $P(\mathbf{X}, Y, \mathbf{Z})$ . It is well known from the central limit theorem (Vaart, 1998) that if this estimator has finite marginal standard errors  $\text{Se}(\hat{\beta}_y)_i = \sqrt{\text{Var}(\hat{\beta}_y)_i}$ , then the following convergence in distribution holds as our choice of  $N$  approaches  $\infty$ :

$$\frac{\sqrt{N} \left( (\hat{\beta}_y)_i - (\beta_y^*)_i \right)}{\text{Se}(\hat{\beta}_y)_i} \xrightarrow{d} \mathcal{N}(0, 1),$$

where  $\beta_{\mathbf{y}}^*$  denotes the ground truth. A similar statement can be made about  $\beta_{\Phi}$  under the assumption that  $\text{Se}(\beta_{\Phi})_{ij}$  are finite for all  $i \in [d_{\Phi}], j \in [d_{\mathbf{Z}}]$ . However, by the NOME assumption, we choose to neglect any variation in  $\hat{\beta}_{\Phi}$ .

We model  $\text{Se}(\hat{\beta}_y)_i$  with the plug in estimator  $\widehat{\text{Se}}(\hat{\beta}_y)_i$  and construct the required confidence intervals using quantiles of  $\mathcal{N}(0, 1)$ . The confidence intervals are symmetric about  $(\hat{\beta}_y)_i$  and have width  $2(\delta\beta_y)_i$  and we denote the resulting confidence set  $\mathbf{B}_{\mathbf{y}}^{\alpha}$ .

Since  $\mathcal{T}(C = \{\gamma_{\mathbf{g}} \in \Gamma\}, s = \{\beta_{\Phi}^*, \hat{\beta}_{\mathbf{y}}\})$  is a deterministic functional of  $\hat{\beta}_{\mathbf{y}}$ , we can compute a confidence set over  $\mathcal{T}$  with asymptotically valid coverage explicitly:

$$\mathcal{T}^{\alpha}(c = \mathbb{I}(\gamma_{\mathbf{g}} \in \Gamma), s = \{\beta_{\Phi}^*, \mathbf{B}_{\mathbf{y}}^{\alpha}\}) = \{\theta \in \mathbb{R}^{d_{\Phi}} : \exists \beta_{\mathbf{y}} \in \mathbf{B}_{\mathbf{y}}^{\alpha}((\beta_{\mathbf{y}} - \theta \cdot \beta_{\Phi}^*) \in \Gamma)\}.$$



Confidence sets over a functional that are constructed in this way appear, for example, in Duarte et al. (2023) and Malloy et al. (2021).

We can rewrite  $\mathcal{T}^\alpha(c, s)$  defined above as the following set:

$$\left\{ \boldsymbol{\theta} \in \mathbb{R}^{d_\Phi} : \exists \boldsymbol{\beta}_y \in \mathbf{B}_y^\alpha, \tilde{U} \in \Sigma_b^{(\max)} \left( (\boldsymbol{\beta}_y - \boldsymbol{\theta} \cdot \boldsymbol{\beta}_\Phi^*) \in \tilde{\Gamma}_{\tilde{U}} \right) \right\}.$$

For any particular  $\tilde{U} \in \Sigma_b$ , define:

$$\tau_i^{\tilde{U}} = \tau_\ell \iff \tilde{U}_i = \ell.$$

Then  $\boldsymbol{\theta} \in \mathcal{T}^\alpha(c, s)$ , iff there exists a  $\tilde{U} \in \Sigma_b$  for which:

$$\left| (\hat{\beta}_y)_i \pm (\delta\beta_y)_i - \boldsymbol{\theta} \cdot \hat{\beta}_x \right| \leq \tau_i^{\tilde{U}} \iff \left| (\hat{\beta}_y)_i - \boldsymbol{\theta} \cdot \hat{\beta}_x \right| \leq \tau_i^{\tilde{U}} + (\delta\beta_y)_i.$$

We can therefore define  $\hat{\Gamma}_\alpha := \bigcup_{\tilde{U} \in \Sigma_b} \hat{\Gamma}_{\tilde{U}}^\alpha$ , where:

$$\hat{\Gamma}_{\tilde{U}}^\alpha := \left\{ \boldsymbol{\gamma} \in \mathbb{R}^{d_Z} : \forall i \in [d_Z] \left( |\gamma_i| \leq \tau_i^{\tilde{U}} + (\delta\beta_y)_i \right) \right\},$$

and the confidence set can be rewritten as:

$$\mathcal{T}^\alpha \left( c = \mathbb{I}[\boldsymbol{\gamma}_g \in \Gamma], s = \{\boldsymbol{\beta}_\Phi^*, \hat{\mathbf{B}}_y^\alpha\} \right) = \mathcal{T} \left( c' = \mathbb{I}[\boldsymbol{\gamma}_g \in \hat{\Gamma}_\alpha], s = \{\boldsymbol{\beta}_\Phi^*, \hat{\boldsymbol{\beta}}_y\} \right).$$

Because  $\mathcal{T}^\alpha$  includes all  $\boldsymbol{\theta} \in \mathcal{T}(c, s = \{\boldsymbol{\beta}_\Phi^*, \boldsymbol{\beta}_y^*\})$  with probability at least  $(1 - \alpha) \times 100\%$ , we have the following guarantee:

$$P \left( \boldsymbol{\theta}^* \in \mathcal{T} \left( c', s = \{\boldsymbol{\beta}_\Phi^*, \hat{\boldsymbol{\beta}}_y\} \right) \right) \geq 1 - \alpha,$$

provided  $\boldsymbol{\gamma}_g^* \in \Gamma$  (i.e., for all  $c \in \mathcal{C}^*$ ).

## B ADDITIONAL THEOREMS

### B.1 A Condition for the Optimality of the Feasible Map under Stronger Structural Assumptions

Here we show that with stricter structural assumptions, sharper feasible maps can be obtained. In particular, we focus on the case in which the model class  $\mathcal{M}'$  consists of all SCMs satisfying Eqs. (1) to (3) for some predetermined restricted function classes  $\mathcal{F}_z \ni f_z$ ,  $\mathcal{F}_x \ni f_x$  and  $\mathcal{G}_y \ni g_y$ .

We show that restrictions on  $\mathcal{F}_z$  and  $\mathcal{F}_x$  alone do not affect the sharpness of  $\mathcal{T}$  and we construct necessary and sufficient conditions for restrictions over  $\mathcal{G}_y$  to lead to the existence of sharper feasible maps.

The fact that  $\mathcal{G}_y$  is the important restriction mirrors results in the (valid) IV literature, in which nonparametric bounds can be put on treatment effects based solely on outcomes  $Y$  being categorical (Balke and Pearl, 1997) or bounded continuous (Manski, 1990). In fact, the following result shows that if  $\Omega_Y$  is any proper subset of the real line (e.g., positive but not absolutely bounded), then a sharper feasible map exists.

**Theorem 5** (Necessary and sufficient condition for optimality under stricter assumptions). Let  $\mathcal{M}'$  be all SCMs consistent with Eqs. (1) to (3) under restricted function classes as described above. The feasible map  $\mathcal{T}$  described in Thm. 2 remains sound against the constraint set  $\mathcal{C} := \{\mathbb{I}(\boldsymbol{\gamma}_g \in \Gamma) : \Gamma \in \mathbb{R}^{d_Z}\}$  under any such  $\mathcal{M}'$ . Furthermore, the map remains complete if and only if for all  $\mathbf{z} \in \Omega_Z$ , the section of  $g_y$  at that  $\mathbf{z}$  is onto the full real line  $\mathbb{R}$ .

**Corollary 5.1.** (Restricted outcome domain) In particular,  $\mathcal{T}$  is not complete if  $\Omega_Y$  is a bounded subset of  $\mathbb{R}$ . This includes if  $\Omega_Y$  is bounded and categorical, e.g.,  $\Omega = \{0, 1\}$ .

For any  $\mathbf{z} \in \Omega_Z$ , the "section" of  $g_y$  at  $\mathbf{Z} = \mathbf{z}$  is defined as the function  $g_{y|\mathbf{z}}$  defined over the domain  $\Omega_{\boldsymbol{\epsilon}_y}$  for which  $g_{y|\mathbf{z}}(\boldsymbol{\epsilon}_y) = g_y(\mathbf{z}, \boldsymbol{\epsilon}_y)$ . The necessary and sufficient condition for completeness is that the image of each  $g_{y|\mathbf{z}}$ , denoted  $g_{y|\mathbf{z}}[\boldsymbol{\epsilon}_y]$ , is the full real line  $\mathbb{R}$ . Before proving the theorem, we introduce an intuitive example when this image restriction does not hold.

### B.1.1 An Example in which the Map is Incomplete

Suppose  $\mathcal{M}'$  requires that at some point  $\mathbf{z} \in \Omega_{\mathbf{Z}}$ , the section:

$$g_{y|\mathbf{z}}(\epsilon_y) := g_y(\mathbf{Z} = \mathbf{z}, \epsilon_y) = \frac{a}{1 + e^{-\epsilon_y}},$$

for some constant  $a$ . Intuitively, at this value  $\mathbf{z}$ , the proportion of  $\text{Var}(Y)$  due to confounding is restricted: there is a value of  $\mathbf{z}$  for which we know the confounding between  $\mathbf{X}$  and  $Y$  is restricted.

Then, for any  $\mathbf{x} \in \Omega_{\mathbf{X}}$  and  $y \in \Omega_Y$  for which  $P(\mathbf{X} = \mathbf{x}, Y = y \mid \mathbf{Z} = \mathbf{z}) \neq 0$ , we may write:

$$y = \boldsymbol{\theta} \cdot \boldsymbol{\Phi}(\mathbf{x}) + \frac{a}{1 + e^{\epsilon_y}}. \quad (6)$$

Assuming  $\boldsymbol{\Phi}(\mathbf{x}) \neq 0$ , the above equation imposes a restriction on  $\boldsymbol{\theta}$ . To see this, consider the case whereby  $d_{\Phi} = 1$  and notice that:

$$\frac{1}{1 + e^{-\epsilon_y}} \in [0, 1].$$

Then, the scalar  $\theta$  is bounded by Eq. (6):

$$\theta \in \left[ \frac{y}{\Phi(\mathbf{x})}, \frac{y}{\Phi(\mathbf{x})} + \frac{a}{\Phi(\mathbf{x})} \right].$$

Therefore, with only a single value of  $\mathbf{z}$ , a strong restriction on  $\theta$  has been imposed through the structural equations and  $P(\mathbf{X}, Y \mid \mathbf{Z} = \mathbf{z})$  alone.

### B.1.2 Proof of Thm. 5

This proof mirrors that of the original optimality result Thm. 2 in Appx. A.2. We explicitly describe the effect of the restrictions  $f_{\mathbf{z}} \in \mathcal{F}_{\mathbf{z}}$ ,  $f_{\mathbf{x}} \in \mathcal{F}_{\mathbf{x}}$  and  $g_y \in \mathcal{G}_y$  at each stage.

Let the structural Eqs. (1) to (3) hold for some ground truth functions  $f_{\mathbf{z}} \in \mathcal{F}_{\mathbf{z}}$ ,  $f_{\mathbf{x}} \in \mathcal{F}_{\mathbf{x}}$ ,  $\boldsymbol{\Phi}^*$ ,  $g_y \in \mathcal{G}_y$  and causal parameter  $\boldsymbol{\theta}^* \in \mathbb{R}^{d_{\Phi}}$ .

**Soundness** The proof of soundness is not affected by the functional restrictions.

It follows immediately from Eq. (3) and the left-linearity of the covariance operator that:

$$\begin{aligned} \boldsymbol{\gamma}_{\mathbf{g}}^* &:= \text{Cov}(g_y(\mathbf{Z}, \boldsymbol{\epsilon}_y), \mathbf{Z}) = \text{Cov}(Y, \mathbf{Z}) - \text{Cov}(\boldsymbol{\theta}^* \cdot \boldsymbol{\Phi}(\mathbf{X}), \mathbf{Z}) \\ &= \boldsymbol{\beta}_{\mathbf{y}}^* - \boldsymbol{\theta}^* \cdot \boldsymbol{\beta}_{\boldsymbol{\Phi}}^*. \end{aligned}$$

Therefore,  $\boldsymbol{\theta}^* \in \mathcal{T}(\boldsymbol{\Gamma}, \boldsymbol{\beta}_{\boldsymbol{\Phi}}^*, \boldsymbol{\beta}_{\mathbf{y}}^*)$  whenever  $\boldsymbol{\gamma}_{\mathbf{g}}^* = (\boldsymbol{\beta}_{\mathbf{y}}^* - \boldsymbol{\theta}^* \cdot \boldsymbol{\beta}_{\boldsymbol{\Phi}}^*) \in \boldsymbol{\Gamma}$  as required.

**Completeness** Consider any  $\boldsymbol{\theta}^\dagger \in \mathbb{R}^{d_{\Phi}}$ . Then the following holds:

(Z) There exists at least one function  $f_{\mathbf{z}}^\dagger$  with the following property. Given any  $\mathbf{z} \in \Omega_{\mathbf{Z}}$ , either: (i)  $P(\mathbf{x}, y, \mathbf{z}) = 0$  for all  $\mathbf{x} \in \Omega_{\mathbf{X}}$  and  $y \in \Omega_Y$ ; or (ii) there exists at least one value  $\boldsymbol{\epsilon}_{\mathbf{z}}^\dagger$  that solves the equation:

$$f_{\mathbf{z}}^\dagger(\boldsymbol{\epsilon}_{\mathbf{z}}^\dagger) = \mathbf{z}.$$

We know this must hold because the ground truth  $f_{\mathbf{z}}^* \in \mathcal{F}_{\mathbf{z}}$  is a valid choice of  $f_{\mathbf{z}}^\dagger$ .

We can, therefore, define a function  $\boldsymbol{\epsilon}_{\mathbf{z}}^\dagger(\mathbf{z})$  that satisfies the above equation for all  $\mathbf{z}$  for which case (i) is false. We have not demanded either  $f_{\mathbf{z}}^\dagger$  or  $\boldsymbol{\epsilon}_{\mathbf{z}}^\dagger(\mathbf{z})$  to be unique.

(X) Likewise, there exists at least one function  $f_{\mathbf{x}}^\dagger$  with the following property. Given any  $\mathbf{x} \in \Omega_{\mathbf{X}}$ ,  $\mathbf{z} \in \Omega_{\mathbf{Z}}$ , either: (i)  $P(\mathbf{x}, y, \mathbf{z}) = 0$  for all  $y \in \Omega_Y$ , or (ii) there exists at least one value  $\boldsymbol{\epsilon}_{\mathbf{x}}^\dagger$  that solves the equation:

$$f_{\mathbf{x}}^\dagger(\mathbf{z}, \boldsymbol{\epsilon}_{\mathbf{x}}^\dagger) = \mathbf{x},$$

and, again,  $f_{\mathbf{x}}^* \in \mathcal{F}_{\mathbf{x}}$  is a valid choice of  $f_{\mathbf{x}}^\dagger$ . We can, therefore, define at least one function  $\boldsymbol{\epsilon}_{\mathbf{x}}^\dagger(\mathbf{z}, \mathbf{x})$  that satisfies the above equation for all  $\mathbf{x}$  and  $\mathbf{z}$  for which case (i) is false.

(Y) In the proof for Thm. 2, we showed that for any  $\mathbf{x} \in \Omega_X$ ,  $y \in \Omega_Y$ ,  $z \in \Omega_Z$ , either (i)  $P(\mathbf{x}, y, z) = 0$ , or (ii) for any  $\boldsymbol{\theta}^\dagger \in \mathbb{R}^{d_\Phi}$  there exists at least one value  $\epsilon_y^\dagger$  that solves the equation:

$$g_y^\dagger(z, \epsilon_y^\dagger) = y - \boldsymbol{\theta}^\dagger \cdot \Phi(\mathbf{x}).$$

However, for some choices of  $\mathcal{G}_y$  this is not true. The statement can only be true if, whenever  $\Phi(\mathbf{x}) \neq \mathbf{0}$ , we have  $g_{y|z}[\epsilon_y] = \mathbb{R}$ . Otherwise, there is a choice of  $\boldsymbol{\theta}^\dagger$  and a value  $g \in \mathbb{R} \setminus g_{y|z}[\epsilon_y]$  such that  $y - \boldsymbol{\theta}^\dagger \cdot \Phi(\mathbf{x}) = g$ .

We define the set of potentially problematic samples:

$$\Omega_{\mathfrak{G}} = \{(\mathbf{x}, y, z) \in \Omega_X \times \Omega_Y \times \Omega_Z : P(\mathbf{X}, Y, \mathbf{Z}) \neq 0 \wedge \Phi(\mathbf{x}) \neq \mathbf{0}\}, \quad (7)$$

and then the values that are "missing" in at least one section of  $g_y$ :

$$\mathfrak{G} = \bigcup_{(\mathbf{x}, y, z) \in \Omega_{\mathfrak{G}}} \{g \in \mathbb{R} \setminus g_{y|z}[\epsilon_y]\}. \quad (8)$$

For each  $g \in \mathfrak{G}$  there are a set of disallowed values for  $\boldsymbol{\theta}$ :

$$\Theta_g := \{\boldsymbol{\theta} \in \mathbb{R}^{d_\Phi} : \exists(\mathbf{x}, y, z) \in \Omega_{\mathfrak{G}} (\boldsymbol{\theta} \cdot \Phi(\mathbf{x})) = y - g\},$$

from which we define  $\Theta_{\mathfrak{G}} := \bigcup_{g \in \mathfrak{G}} \Theta_g$ .

We can therefore construct the sharper feasible set than  $\mathcal{T} = \{\boldsymbol{\theta} : h(\boldsymbol{\theta}) \in \Gamma\}$  by excluding every  $\boldsymbol{\theta}$  which is disallowed:

$$\mathcal{T}'(c = \mathbb{I}[\gamma_g \in \Gamma], s' = P(\mathbf{X}, Y, \mathbf{Z})) = \mathcal{T} \setminus \Theta_{\mathfrak{G}},$$

with  $\mathcal{T}' = \mathcal{T}$  iff for all  $z \in \Omega_z$ ,  $g_{y|z}[\epsilon_y] = \mathbb{R}$  (equivalently,  $\mathfrak{G} = \emptyset$ ).

We know the corresponding feasible map  $\mathcal{T}'$  is sound because it only excludes  $\boldsymbol{\theta} \in \mathbb{R}^{d_\Phi}$  if either (i) the sound feasible map  $\mathcal{T}$  excludes  $\boldsymbol{\theta}$ , or (ii) there does not exist a  $g_y \in \mathcal{G}_y$  (and thus an  $m \in \mathcal{M}'$ ) which is consistent with  $P(\mathbf{x}, y, z)$ . We have therefore proven the necessity of the assumption that the image of each section of  $g_y$  is the full real line for sharpness of  $\mathcal{T}$ .

Sufficiency is simple to prove by following step (Y) in the proof of Thm. 2 in Appx. A.2 to construct a full joint distribution consistent with each  $P(\mathbf{X}, Y, \mathbf{Z})$  and the SCM:

$$P^\dagger(\mathbf{X}, Y, \mathbf{Z}, \epsilon_x, \epsilon_y) := D^\dagger(\epsilon_x | \mathbf{X}, \mathbf{Z}) \delta(\epsilon_y - \epsilon_y^\dagger(\mathbf{X}, Y, \mathbf{Z}; \boldsymbol{\theta})) P(\mathbf{X}, Y, \mathbf{Z}).$$

For definitions of the  $\delta$ 's and  $D$ , visit this part of the appendix.

This concludes the proof of Thm. 5.

**Minimality** We note that  $\mathcal{T}$  is not necessarily minimal for general  $\mathcal{F}_z$ ,  $\mathcal{F}_x$  and  $\mathcal{G}_y$ . If  $\mathcal{T}$  is sharp, it may not be minimal because underlying symmetries in these function classes may constrain which  $\beta_y, \beta_\Phi$  can arise, which allow one to describe the solution set  $\{\boldsymbol{\theta} : (\beta_y^* - \boldsymbol{\theta} \cdot \beta_\Phi^*) \in \Gamma\}$  with fewer parameters. On the other hand,  $\beta_y, \beta_\Phi$  may be insufficient for specifying  $\mathcal{T}'$  when  $g_{y|z}[\epsilon_y] \subset \mathbb{R}$ .

## B.2 Polytime Testability and Necessary Condition for Unidentifiability

We say the ground truth causal effect  $\boldsymbol{\theta}^*$  is *unidentifiable* when the feasible set  $\mathcal{T}(\Gamma, \beta_\Phi^*, \beta_y^*) = \mathbb{R}^{d_\Phi}$ . In plain English, this means nothing can be learned about  $\boldsymbol{\theta}^*$  from the data under the current assumptions.

In this section we show that budget constraints enable an efficient test for unidentifiable causal effects even when  $d_\Phi > 1$ . We also show unidentifiability can only occur under budget constraints if assumption (B1\*) is violated (recall the definition from Sect. 2.1).

**Theorem 6 (Unidentifiability).** Assume Eqs. 1, 2, 3 and budget constraints according to some  $\Gamma(\boldsymbol{\tau}, \mathbf{b})$ . Then unidentifiability of  $\boldsymbol{\theta}^*$  can be decided in  $\mathcal{O}(K d_Z d_\Phi)$  time. Moreover, unidentifiability never occurs under (B1\*).

### B.2.1 Proof of Thm. 6

We prove necessary and sufficient conditions for  $\mathcal{T}(\Gamma(\boldsymbol{\tau}, \mathbf{b}), \boldsymbol{\beta}_\Phi, \boldsymbol{\beta}_y) = \mathbb{R}^{d_\Phi}$ .

Defining the function  $\mathbf{V} : \mathbb{R}^{d_Z} \mapsto [K + 1]^{d_Z}$  through its components:

$$V_i(\boldsymbol{\gamma}_g) = \begin{cases} 1 & |(\boldsymbol{\gamma}_g)_i| \leq \tau_1 \\ l \in \{2, 3, \dots, K\} & \tau_{l-1} \leq |(\boldsymbol{\gamma}_g)_i| \leq \tau_l \\ K + 1 & |(\boldsymbol{\gamma}_g)_i| > \tau_K, \end{cases}$$

we have  $\mathbf{V}(\boldsymbol{\gamma}_g^*) = \mathbf{U}^*$  is the ground truth for the latent variable  $\mathbf{U}$ . We can represent  $\mathbf{V}$  by its one-hot encoding:

$$V_{il}(\boldsymbol{\gamma}_g) = \begin{cases} 1 & V_i(\boldsymbol{\gamma}_g) \leq l \\ 0 & \text{otherwise.} \end{cases}$$

Then budget background search space is explicitly written as:

$$\Gamma(\boldsymbol{\tau}, \mathbf{b}) := \left\{ \boldsymbol{\gamma}_g \in \mathbb{R}^{d_Z} : \forall l \in [K] \left( \sum_{i=1}^{d_Z} V_{il}(\boldsymbol{\gamma}_g) \geq b_l \right) \right\}.$$

Unidentifiability occurs iff  $h(\boldsymbol{\theta}) = \boldsymbol{\beta}_y - \boldsymbol{\theta} \cdot \boldsymbol{\beta}_\Phi$  is contained within  $\Gamma(\boldsymbol{\tau}, \mathbf{b})$  for any  $\boldsymbol{\theta} \in \mathbb{R}^{d_\Phi}$ . Equivalently, unidentifiability occurs iff the affine space  $h := h[\boldsymbol{\theta}]$  formed by the image of  $h(\boldsymbol{\theta})$  is contained within  $\Gamma(\boldsymbol{\tau}, \mathbf{b})$ .

Consider the standard orthonormal basis  $\{\mathbf{e}_i\}_{i=1}^{d_Z}$  for which  $\boldsymbol{\gamma}_g = \sum_{i=1}^{d_Z} (\boldsymbol{\gamma}_g)_i \mathbf{e}_i$ . For any direction  $\mathbf{e}_i$  we have that either (i)  $(\boldsymbol{\beta}_y)_i - (\boldsymbol{\theta} \cdot \boldsymbol{\beta}_\Phi)_i = (\boldsymbol{\beta}_y)_i$  is constant and  $h$  is orthogonal to  $\mathbf{e}_i$ , or (ii)  $(\boldsymbol{\beta}_y)_i - (\boldsymbol{\theta} \cdot \boldsymbol{\beta}_\Phi)_i$  is unbounded. If (i) holds, it must be the case that  $(\boldsymbol{\beta}_\Phi)_i = \mathbf{0}$  and thus (B1\*) is violated, and we say candidate instrument  $Z_i$  is covariance-irrelevant to  $\Phi(\mathbf{X})$ .

We can define the index set for all covariance-irrelevant candidate instruments:

$$I := \{i \in [d_\Phi] : (\boldsymbol{\beta}_\Phi)_i = \mathbf{0}\}.$$

Then,  $h$  is unbounded in all directions  $[d_Z] \setminus I$  but fixed at a single value for each direction in  $I$ . Since  $h$  is affine it is also unbounded in the direction  $\sum_{i \in [d_Z] \setminus I} \mathbf{e}_i$ . Moreover, since points in  $h$  are parameterized by  $\boldsymbol{\theta}$ , there exists values of  $\boldsymbol{\theta}$  for which  $V_i(h(\boldsymbol{\theta})) = K + 1$  for all  $i \in [d_Z] \setminus I$ .

At such a value for  $\boldsymbol{\theta}$ , we have  $\mathbf{V}(h(\boldsymbol{\theta})) = \mathbf{V}(\mathbf{P})$ , where  $\mathbf{P} \in \mathbb{R}^{d_Z}$  is defined by:

$$P_i = \begin{cases} (\boldsymbol{\beta}_y)_i & i \in I \\ \tau_K + 1 & \text{otherwise.} \end{cases}$$

Therefore, the necessary and sufficient condition for unidentifiability is:

$$\sum_{i=1}^{d_Z} V_{il}(\mathbf{P}) \geq b_l,$$

where  $I$  and  $\mathbf{P}$  can be computed and the condition evaluated by straightforward linear algebra in  $\mathcal{O}(Kd_Zd_\Phi)$  time.

## C ALGORITHMS

Alg. 1 depicts a combinatorial search algorithm that finds the exact feasible set of  $\boldsymbol{\theta}$  and the resulting feasible set of ATE( $\mathbf{x}; \mathbf{x}_0$ ) subject to budget background constraints. The ATE is recovered using a grid search and program over the linear weights  $\boldsymbol{\theta}$  with plug-in values of  $\mathbf{x} \in \Omega_{\mathbf{X}}$ .

Alg. 2 depicts a polynomial time algorithm for the case of  $d_\Phi = 1$ . This method utilises the single solution to the equation  $A_i - B_i\theta = \tau$  for  $B_i = 0$ , and our ability to order such solutions along  $\mathbb{R}$ .

**Algorithm 1** budgetIV general case.

---

```

1: for  $\tilde{U} \in \Sigma_b^{(\max)}$  do                                ▷ Iterate through search space (combinatorial)
2:   if  $h(\theta) \in \tilde{\Gamma}_{\tilde{U}}$  for some  $\theta$  then                ▷ Linear CSP with convex constraints
3:      $\check{U}(\tilde{U}) \leftarrow \mathbf{0}$ 
4:     for  $\mathbf{x} \in \Omega_{\mathbf{X}}$  do                                    ▷ Grid evaluation (for  $d_{\mathbf{X}} \ll d_{\Phi}$ )
5:       Calculate  $\text{ATE}_{\tilde{U}}^{+/-}(\mathbf{x}; \mathbf{x}_0)$                 ▷ Linear program with convex constraints
6:        $\theta_{\tilde{U}}^{+/-}(\mathbf{x}; \mathbf{x}_0)$                             ▷ Arguments for above LP
7:       for  $i \in [d_{\mathbf{Z}}]$  do
8:          $\check{U}_i(\tilde{U}) \leftarrow \max\{\check{U}_i(\tilde{U}), U_i(\theta^+), U_i(\theta^-)\}$ 
9:       end for
10:    end for
11:  end if
12: end for
13: return  $\left\{ \left( \text{ATE}_{\tilde{U}}^{+/-}, \check{U}(\tilde{U}) \right) : \tilde{U} \in \Sigma_b^{(\max)} \right\}$ 

```

---

**Algorithm 2** Polytime budgetIV with  $d_{\Phi} = 1$ .

---

```

1: Let  $c'(\theta)$  be the indicator function for  $h(\theta) \in \Gamma(\tau, \mathbf{b})$  and  $\mathbf{U}(\theta)$  be the unique value of  $\mathbf{U}$  for which
    $h(\theta) \in \Gamma_{\mathbf{U}}(\tau, \mathbf{b})$  provided  $c'(\theta)$ .
2: for  $(i, j) \in [d_{\mathbf{Z}}] \times [K]$  do
3:   if  $(\beta_{\Phi})_i = 0 \wedge (\beta_y)_i \leq \tau_j$  then
4:      $b_j \leftarrow b_j - 1$ 
5:   else if  $(\beta_{\Phi})_i \neq 0$  then
6:      $\theta_{ij}^{\pm} \leftarrow \frac{(\beta_y)_i}{(\beta_{\Phi})_i} \mp \frac{\tau_j}{(\beta_{\Phi})_i}$ 
7:   end if
8: end for
9: if  $\forall i \in [K](b_i \leq 0)$  then
10:  return unidentifiable
11: end if
12:  $\Theta \leftarrow \text{sort}(\theta_{ij}^{\pm})$ 
13:  $\Theta' \leftarrow \left( \Theta_1 - 1, \frac{\Theta_2 + \Theta_1}{2}, \frac{\Theta_3 + \Theta_2}{2}, \dots, \frac{\Theta_{|\Theta|} + \Theta_{|\Theta|-1}}{2}, \Theta_{|\Theta|} + 1 \right)$ 
14: intervals  $\leftarrow \{ \langle [\Theta_i, \Theta_{i+1}], \mathbf{U}(\Theta'_i) \rangle : c'(\Theta'_i) \}$ 
15: points  $\leftarrow \{ \langle \Theta_i, \mathbf{U}(\Theta_i) \rangle : c'(\Theta_i) \wedge \neg c'(\Theta'_i) \neg c'(\Theta'_{i+1}) \}$ 
16: if intervals  $\cup$  points  $== \emptyset$  then
17:  return infeasible
18: else
19:  return intervals  $\cup$  points
20: end if

```

---

## D EXPERIMENTS

### D.1 Linear Simulation Study

Fig. 3 is the result of a simple experiment to understand the differences between the nonconvex budget background constraints and any convex relaxation thereof. We consider violation of (A2) in a linear model where association between  $X$  and  $\mathbf{Z}$  is wholly due to unobserved confounding:

$$\begin{aligned}\mathbf{Z} &:= \boldsymbol{\epsilon}_{\mathbf{Z}}, \\ X &:= \epsilon_x, \\ Y &:= \theta^* X + \epsilon_y.\end{aligned}$$

We consider  $d_X = 1$  and  $d_{\mathbf{Z}} = 2$ , take the exogenous variables to have a joint distribution  $(\boldsymbol{\epsilon}_{\mathbf{Z}}, \epsilon_x, \epsilon_y) \sim \mathcal{N}(\mathbf{0}, \boldsymbol{\Sigma}_{\boldsymbol{\epsilon}\boldsymbol{\epsilon}})$ , where the covariance is given by the terms:

$$\boldsymbol{\Sigma}_{\boldsymbol{\epsilon}\boldsymbol{\epsilon}} := \begin{pmatrix} \eta_{Z_1}^2 & \rho_{Z_1 Z_2} \eta_{Z_1} \eta_{Z_2} & \rho_{Z_1 \epsilon_x} \eta_{Z_1} \eta_{\epsilon_x} & \rho_{Z_1 \epsilon_y} \eta_{Z_1} \eta_{\epsilon_y} \\ & \eta_{Z_2}^2 & \rho_{Z_2 \epsilon_x} \eta_{Z_2} \eta_{\epsilon_x} & \rho_{Z_2 \epsilon_y} \eta_{Z_2} \eta_{\epsilon_y} \\ & & \eta_{\epsilon_x}^2 & \rho_{\epsilon_x \epsilon_y} \eta_{\epsilon_x} \eta_{\epsilon_y} \\ & & & \eta_{\epsilon_y}^2 \end{pmatrix}.$$

We fix the following throughout the study:

$$\begin{aligned}\theta^* &:= 1, \\ \boldsymbol{\beta}_{\mathbf{x}}^* &:= \text{Cov}(X, \mathbf{Z}) = (2, -4), \\ \boldsymbol{\gamma}_{\mathbf{g}}^* &:= \text{Cov}(\epsilon_y, \mathbf{Z}) = (-2, 0.4),\end{aligned}$$

which, in turn, imply  $\boldsymbol{\beta}_{\mathbf{y}}^* = (0, 4.4)$ .

#### D.1.1 Sweep Through Increasingly Uncertain Background Constraints

We consider three kinds of background constraints: (a) budget constraints, (b) an  $L_2$ -norm relaxation (a similar setting to [Watson et al. \(2024\)](#)), and (c) an  $L_1$ -norm relaxation (which is motivated by the relaxation in [Kang et al. \(2016\)](#), under the setting where point identification is not guaranteed):

- (a)  $\boldsymbol{\gamma}_{\mathbf{g}} \in \mathbf{\Gamma}(\boldsymbol{\tau} = (\tau, 0.6), \mathbf{b} = (1, 2))$ ,
- (b)  $\|\boldsymbol{\gamma}_{\mathbf{g}}\|_2 \leq \tau$ ,
- (c)  $\|\boldsymbol{\gamma}_{\mathbf{g}}\|_1 \leq \tau$ .

We adjust  $\tau$  from 0 to 10 linearly across a grid of 101 simulation studies.

#### D.1.2 Randomized Parameters for the Covariance Matrix

We select random values for  $\boldsymbol{\Sigma}_{\boldsymbol{\epsilon}\boldsymbol{\epsilon}}$  between simulations. The data generating process is constrained by two requirements: (i) the matrix  $\boldsymbol{\Sigma}_{\boldsymbol{\epsilon}\boldsymbol{\epsilon}}$  must be positive definite, and (ii) the fixed values of  $\boldsymbol{g}^*$  and  $\boldsymbol{\beta}_{\mathbf{x}}^*$ , which demand:

$$(\boldsymbol{\gamma}_{\mathbf{g}}^*)_1 = \rho_{Z_1 \epsilon_y} \eta_{Z_1} \eta_{\epsilon_y}, \quad (9)$$

$$(\boldsymbol{\gamma}_{\mathbf{g}}^*)_2 = \rho_{Z_2 \epsilon_y} \eta_{Z_2} \eta_{\epsilon_y}, \quad (10)$$

$$(\boldsymbol{\beta}_{\mathbf{x}}^*)_1 = \rho_{Z_1 \epsilon_x} \eta_{Z_1} \eta_{\epsilon_x}, \quad (11)$$

$$(\boldsymbol{\beta}_{\mathbf{x}}^*)_2 = \rho_{Z_2 \epsilon_x} \eta_{Z_2} \eta_{\epsilon_x}. \quad (12)$$

We choose to generate  $\boldsymbol{\Sigma}_{\boldsymbol{\epsilon}\boldsymbol{\epsilon}}$  by a rejection method. We sample the marginal variances according to:

$$\eta_{Z_1}, \eta_{Z_2}, \eta_{\epsilon_x}, \eta_{\epsilon_y} \stackrel{\text{i.i.d.}}{\sim} \text{Exp}(1),$$

and calculate the would-be  $\rho$ 's from Eqs. (9) to (12). This enforces condition (ii). We then test whether the resultant  $\boldsymbol{\Sigma}_{\boldsymbol{\epsilon}\boldsymbol{\epsilon}}$  is positive definite. If it is, (i) is satisfied and we perform the experiment; otherwise we resample

the  $\eta$ 's from Exp(1). Fig. 3 shows the resulting feasible sets with plug-in  $\hat{\beta}_\Phi$  and  $\hat{\beta}_y$ , calculated from a dataset of  $N = 10000$  samples.

The feasible sets over  $\theta$  in Fig. 3 were found using the efficient `budgetIV` algorithm for  $d_\Phi = 1$ , an implementation of which is provided in the supplement.

## D.2 Nonlinear Simulation Study

In this experiment, we show another advantage of the budget background constraint approach over convex relaxations. Different values of the decision variables  $\mathbf{U}$  lead to different functions  $\text{ATE}(\mathbf{x}; \mathbf{x}_0)$ .

### D.2.1 Endogenous Equations

The candidate instruments are in the domain  $\Omega_{\mathbf{Z}} = \{0, 1\}^{d_{\mathbf{Z}}}$ , where  $d_{\mathbf{Z}} = 6$ ; the exposure is a scalar in  $\Omega_{\mathbf{X}} = [0, 1]$ ; and the outcome is a scalar in the full real line  $\Omega_Y = \mathbb{R}$ .

The ground truth structural equations are:

$$\begin{aligned} \mathbf{Z} &:= \boldsymbol{\epsilon}_{\mathbf{z}} \\ X &:= f_x(\mathbf{Z}, \boldsymbol{\epsilon}_x), \\ Y &:= \boldsymbol{\theta}^* \cdot \Phi^*(\mathbf{X}) + g_y(\mathbf{Z}, \boldsymbol{\epsilon}_y), \end{aligned}$$

where the functions take the form:

$$\begin{aligned} f_x(\mathbf{Z}, \boldsymbol{\epsilon}_x) &= \text{Expit}(\boldsymbol{\epsilon}_x - \mathbf{m} \cdot \mathbf{Z}), \\ g_y(\mathbf{Z}, \boldsymbol{\epsilon}_y) &= \lambda^{(A3)} \mathbf{Z} \cdot \boldsymbol{\Lambda} \cdot \mathbf{Z} + \lambda^{(A2)} \boldsymbol{\epsilon}_y, \\ \Phi^*(X) &= s_\Phi ((X - 1/4)^2 - 1/16), \end{aligned}$$

and  $\boldsymbol{\theta}^* = 1$ .

The real vector  $\mathbf{m} \in \mathbb{R}^{d_{\mathbf{Z}}}$  and the binary matrix  $\boldsymbol{\Lambda} \in \{0, 1\}^{d_{\mathbf{Z}} \times d_{\mathbf{Z}}}$  are sampled randomly for each experiment. The components of  $\mathbf{m}$  are i.i.d normal and the components of  $\boldsymbol{\Lambda}$  are independent Bernoulli trials according to:

$$\begin{aligned} m_1, \dots, m_{d_{\mathbf{Z}}} &\stackrel{\text{i.i.d}}{\sim} \mathcal{N}(1, 4), \\ \Lambda_{ij} &\sim \begin{cases} \mathcal{N}(0.9, 0.9) & i = j > 3, \\ \mathcal{N}(0.3, 0.3) & i \neq j \text{ and } i, j > 3, \\ 0 & \text{otherwise.} \end{cases} \end{aligned}$$

Notice that candidate instruments  $Z_1, Z_2, Z_3$  all do not violate (A3).

The remaining parameters  $s_\Phi, \lambda^{(A2)}, \lambda^{(A3)} \geq 0$  are restricted to be positive but are otherwise varied across simulation settings (see Appx. D.2.3).

### D.2.2 Generating the Exogenous Variables

We generate  $(\mathbf{Z}, \boldsymbol{\epsilon}_x, \boldsymbol{\epsilon}_y)$  to have a block diagonal covariance structure where so that  $Z_1, Z_2$  and  $Z_3$  are valid instruments.

This is achieved by a Markov chain approach over binary  $Z$  with:

$$\begin{aligned} Z_1 &\sim \text{Ber}(0.05), \\ Z_4 &\sim \text{Ber}(0.05), \\ Z_{j+1} \mid Z_j = 1 &\sim \text{Ber}(0.9), \\ Z_{j+1} \mid Z_j = 0 &\sim \text{Ber}(0.05), \end{aligned}$$

where  $j \in \{2, 3, 5, 6\}$ . Then the noise residuals are given by:

$$\begin{aligned} \boldsymbol{\epsilon}_x &= \boldsymbol{\gamma}_{\boldsymbol{\epsilon}_x} \cdot \mathbf{Z} + \boldsymbol{\epsilon}'_x, \\ \boldsymbol{\epsilon}_y &= \boldsymbol{\gamma}_{\boldsymbol{\epsilon}_y} \cdot \mathbf{Z} + \boldsymbol{\epsilon}'_y, \end{aligned}$$

where:

$$\begin{aligned} (\gamma_{\epsilon_x})_1, \dots, (\gamma_{\epsilon_x})_6 &\stackrel{\text{i.i.d.}}{\sim} \mathcal{N}(0, 1), \\ (\gamma_{\epsilon_y})_1, \dots, (\gamma_{\epsilon_y})_3 &= 0, \\ (\gamma_{\epsilon_x})_4, \dots, (\gamma_{\epsilon_x})_6 &\stackrel{\text{i.i.d.}}{\sim} \mathcal{N}(0, 1), \\ \epsilon' &\sim \mathcal{N}(0, 1). \end{aligned}$$

### D.2.3 Varying the Simulation Setting

The remaining parameters,  $(s_\Phi, \lambda^{(A2)}, \lambda^{(A3)})$ , are set to control the variance  $\text{Var}(Y)$ , signal to noise ratio  $\text{SNR}_Y$ , and the ratio of effects, defined by:

$$R := \frac{\lambda^{(A3)} \|\text{Cov}(\mathbf{Z} \cdot \mathbf{\Lambda} \cdot \mathbf{Z}, \tilde{\mathbf{Z}})\|_2}{\lambda^{(A2)} \|\text{Cov}(\epsilon_y, \tilde{\mathbf{Z}})\|_2},$$

where  $\tilde{\mathbf{Z}} := (Z_4, \dots, Z_6)$ .

Using the following shorthand:

$$\begin{aligned} P &:= \varphi(X) := \sqrt{1 + (X - 0.5)^2}, \\ \gamma_i &:= \mathbf{z}_i^\top \cdot \mathbf{\Lambda} \cdot \mathbf{z}_i \\ r &:= \frac{\|\text{Cov}(\boldsymbol{\gamma}, \tilde{\mathbf{Z}})\|_2}{\|\text{Cov}(\epsilon_y, \tilde{\mathbf{Z}})\|_2}, \\ U &:= u(\mathbf{Z}, \epsilon_y) := \frac{r}{R} \boldsymbol{\gamma}(\mathbf{Z}) + \epsilon_y, \end{aligned}$$

the remaining free parameters  $s_\Phi, \lambda^{(A2)}, \lambda^{(A3)} \geq 0$  are expressed as:

$$\begin{aligned} s_\Phi^2 &= \frac{1}{\text{Var } P} \frac{\text{Var } Y}{1 + \frac{1}{\text{SNR}_Y}}, \\ \lambda^{(A2)} &= \frac{s_\Phi \text{Cov}(P, U)}{\text{Var } U} \left( -1 + \frac{1}{2} \sqrt{1 + \frac{\text{Var } U}{s_\Phi^2 (\text{Cov}(P, U))^2} \frac{\text{Var } Y}{1 + \text{SNR}_Y}} \right), \\ \lambda^{(A3)} &= \frac{R}{r} \lambda^{(A2)}. \end{aligned}$$

We choose to fix  $\text{Var}(Y) := 10$  (this is simply a choice of the scale with which we choose to measure  $Y$ ). We produce a three-by-three grid of experiments with the values of  $\text{SNR}_Y = 1/10, 1/5, 1$  and  $R = 1/2, 1, 2$ .

The remaining parameters are fixed to approximately fit these constraints by a post-hoc method using plug-in empirical estimates:  $\widehat{\text{Var}}(P)$ ,  $\widehat{\text{Var}}(U)$ ,  $\widehat{\text{Cov}}(P, U)$  and a plug-in ratio estimate:

$$\hat{r} := \frac{\|\widehat{\text{Cov}}(\mathbf{Z} \cdot \mathbf{\Lambda} \cdot \mathbf{Z}, \tilde{\mathbf{Z}})\|_2}{\|\widehat{\text{Cov}}(\epsilon_y, \tilde{\mathbf{Z}})\|_2},$$

where a dataset of  $5 \times 10^5$  samples is used to construct the estimates.

### D.2.4 Results from Main Text

With the same  $N = 5 \times 10^5$  sample dataset, we run an implementation of `budgetIV`, included in the supplement, with feasible search space  $\mathbf{\Gamma}(\tau = 0, b = 3)$ . Fig. 6 shows the results for the full simulation grid.



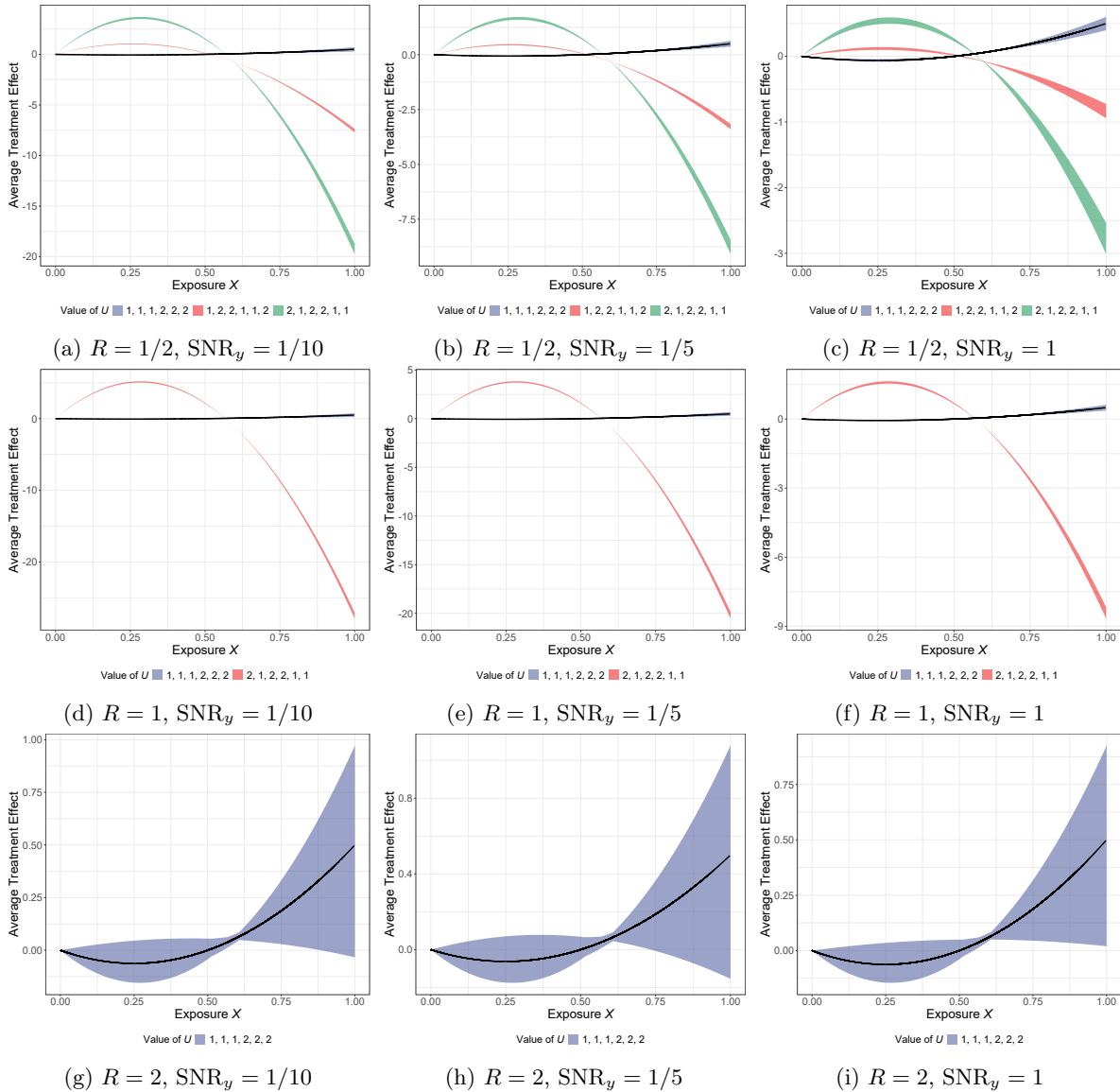


Figure 6: **Simulation grid for the experiment discussed in Appx. D.2.** Some settings (a, b) result in a feasible set in which the ATE is qualitatively very different for each plausible  $\tilde{U}$ . In other settings (g–i)  $U^*$  is identified exactly.  $\Sigma_b^{(\max)} = \binom{6}{3} = 20$ , so budgetIV significantly reduced the space of plausible  $\tilde{U}$  in this experiment.

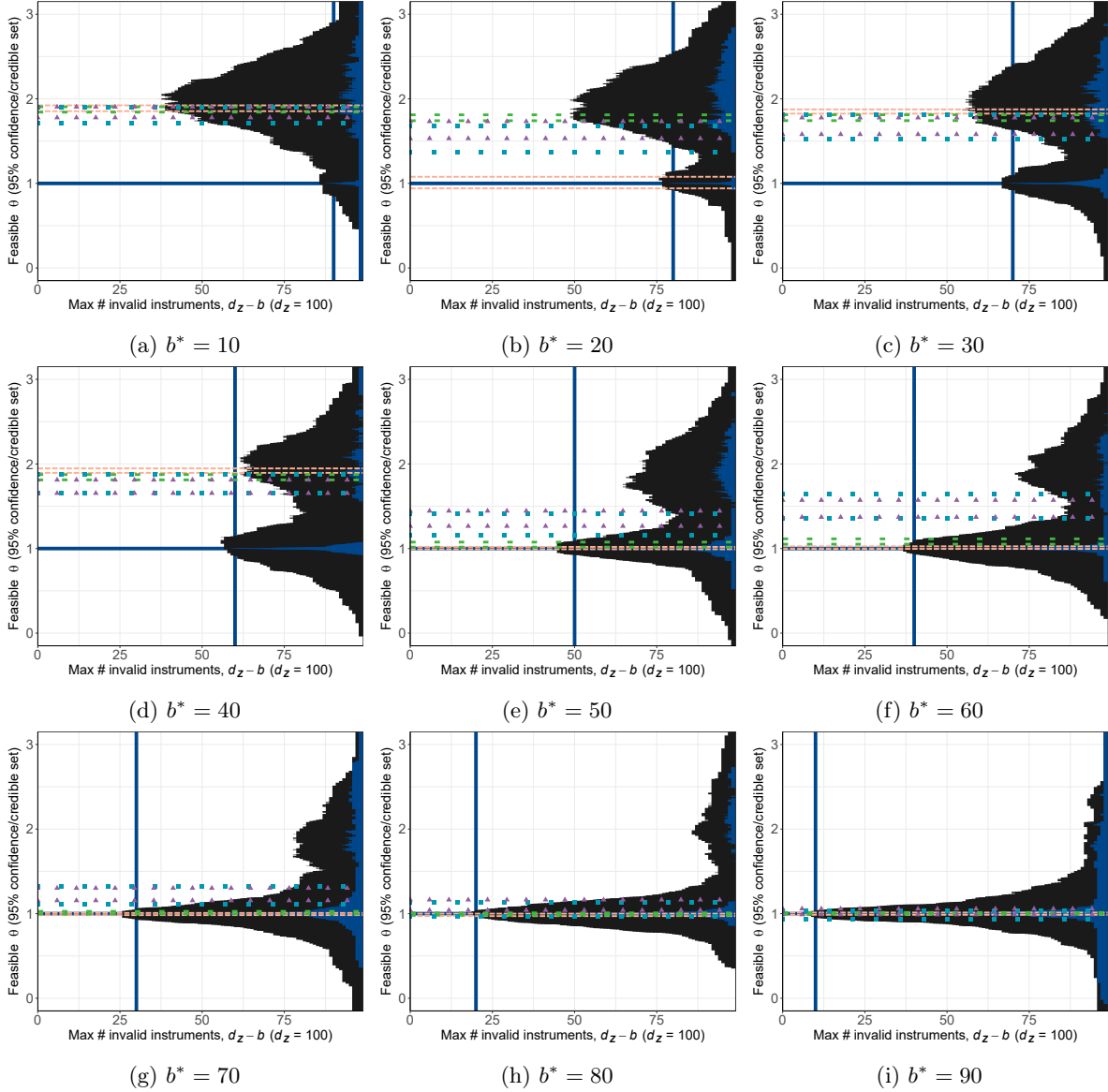


Figure 7: **Simulation grid for the experiment discussed in Appx. D.3.** Each subfigure corresponds to a different ground truth number of valid IVs  $b^*$ . Finite sample confidence sets using `budgetIV_scalar` and **Oracle** are shown for the constraints  $\Gamma(\tau = 0, b)$  where  $b$  is varied along the horizontal axis. Confidence intervals for the benchmarking methods **MR-Egger**, **MR-Median**, **MBE** and **IVW**, under which  $b$  is not an adjustable parameter, are labeled for each experiment. Sub-figures (a) through (e) show variability in the confidence intervals of **MBE** when valid IVs are a minority—despite the required modal assumption of holding for each ground truth model. The approach **MASSIVE**, included in Fig. 5 in the main text, takes significantly more computational resources than the the methods included in this simulation grid.

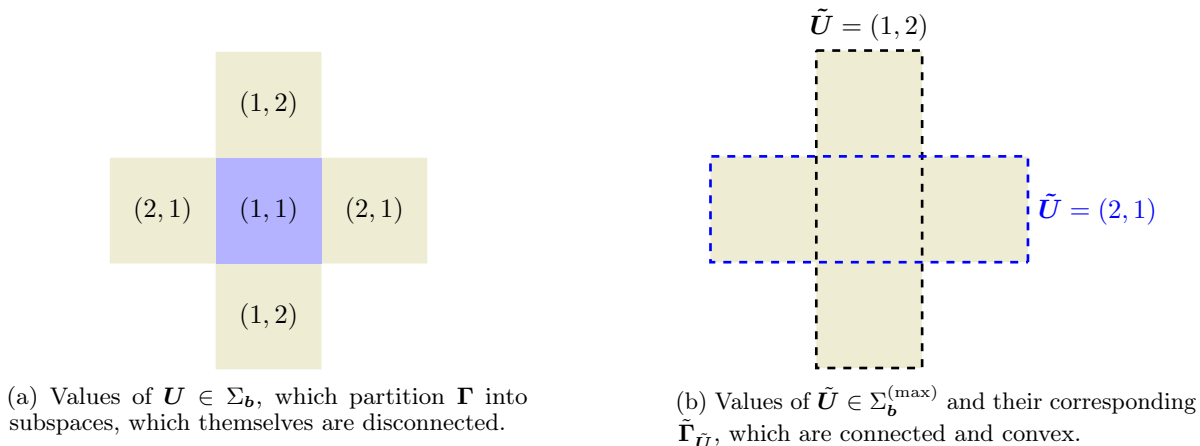


Figure 8: A comparison between representing the search space  $\Gamma(\boldsymbol{\tau} = (\tau_1, \tau_2), \mathbf{b} = (1, 2))$  as a disjoint union over  $\Gamma_{\mathbf{U}}$  or a union over  $\tilde{\Gamma}_{\tilde{\mathbf{U}}}$ .

### D.3 Linear Many Candidate IV Study

**Simulation setting** We performed a grid of simulations using the following linear SEM with (A3) violation among a varying proportion of the 100 candidate IVs:

$$\begin{aligned} \mathbf{Z} &:= \boldsymbol{\epsilon}_{\mathbf{z}}, \\ X &:= \boldsymbol{\delta} \cdot \mathbf{Z} + \epsilon_x, \\ Y &:= \theta^* X + \boldsymbol{\gamma} \cdot \mathbf{Z} + \epsilon_y, \\ b^* &:= \|\boldsymbol{\gamma}_g^*\|_0 \in \{10, 20, 30, 40, 50, 60, 70, 80, 90\}, \\ \theta^* &:= 1. \end{aligned}$$

We set  $(\boldsymbol{\epsilon}_{\mathbf{z}}, \epsilon_x, \epsilon_y)$  to a multivariate normal centered at  $\mathbf{0}$  with no correlation except  $\rho := \text{Corr}(\epsilon_x, \epsilon_y)$  drawn uniformly from  $[-1, 1]$ . The marginal standard deviations are drawn i.i.d. from  $\text{Exp}(1)$ . The first  $b^*$ -many entries of  $\boldsymbol{\gamma}$  are set to 0 and the remaining entries, as well as the entries of  $\boldsymbol{\delta}$ , are samples i.i.d. from the uniform distribution  $\mathcal{U}[1, 2]$ .

A two-sample approach was applied with  $N_x = 1 \times 10^6$  samples used to sample data to generate the summary statistics  $\hat{\boldsymbol{\beta}}_{\mathbf{x}} := \widehat{\text{Cov}}(\mathbf{Z}, X)$  and  $N_y = 1 \times 10^5$  samples used to generate  $\hat{\boldsymbol{\beta}}_{\mathbf{y}} := \widehat{\text{Cov}}(\mathbf{Z}, Y)$ . This was to model the typical occurrence that  $\hat{\boldsymbol{\beta}}_{\mathbf{x}}/\text{SE}(\hat{\boldsymbol{\beta}}_{\mathbf{x}}) \gg \hat{\boldsymbol{\beta}}_{\mathbf{y}}/\text{SE}(\hat{\boldsymbol{\beta}}_{\mathbf{y}})$  (SE is the empirical standard error) because candidate instruments are selected if they have a strong association with  $X$  (see Bowden et al. (2016b)).

This simulation setting reflects common modeling assumptions that are applied in the original experiments for each of the benchmark methods.

**Results** Results for the case  $b^* = 30$  were interpreted in Sect. 5, including the caption to Fig. 5. Further interpretation, including to the stability of the MBE estimator is given in the caption to Fig. 7.

## E SEARCHING THROUGH THE BUDGET CONSTRAINTS

The number of unique decision variables  $|\Sigma_{\mathbf{b}}|$  is given by the number of assignments  $\mathbf{U} \in [K+1]^{d_{\mathbf{z}}}$  for which at least  $b_1$  components of  $\mathbf{U}$  are equal to 1, at least  $b_2$  components are no greater than 2 and so on until exactly  $b_{K+1} := d_{\mathbf{z}}$  components.

Suppose exactly  $d_1$  components are equal to 1, exactly  $d_2 - d_1$  are equal to 2 and so on. This number of combinations with this property is equal to the multinomial coefficient  $\binom{d_{\mathbf{z}}}{d_1, (d_2 - d_1), \dots, (d_{K+1} - d_K)}$ .

Given that we require  $d_i \geq b_i$  for all  $i \in [K + 1]$ , we can write the following:

$$|\Sigma_{\mathbf{b}}| = \sum_{d_1=b_1}^{d_2} \cdots \sum_{d_K=b_K}^{d_{\mathbf{Z}}} \frac{d_{\mathbf{Z}}!}{\prod_{\ell=1}^{K+1} (d_{\ell} - d_{\ell-1})!}.$$

By comparison,  $|\Sigma_{\mathbf{b}}^{(\max)}|$  count  $\tilde{\mathbf{U}} \in [K + 1]^{d_{\mathbf{Z}}}$  for which exactly  $b_1$  components equal 1,  $b_2 - b_1$  components equal 2 and so on. Therefore, denoting  $b_0 := 0$ , we have:

$$|\Sigma_{\mathbf{b}}^{(\max)}| = \frac{d_{\mathbf{Z}}!}{\prod_{\ell=1}^{K+1} (b_{\ell} - b_{\ell-1})!}.$$

The ratio gap between these sizes is given exactly by:

$$\frac{|\Sigma_{\mathbf{b}}|}{|\Sigma_{\mathbf{b}}^{(\max)}|} = \sum_{d_1=b_1}^{d_2} \cdots \sum_{d_K=b_K}^{d_{\mathbf{Z}}} \frac{\prod_{\ell=1}^{K+1} (b_{\ell} - b_{\ell-1})!}{\prod_{\ell=1}^{K+1} (d_{\ell} - d_{\ell-1})!} > 1.$$

Depending on the scaling of  $d_{\mathbf{Z}}$  and  $\mathbf{b}$ , asymptotic gap can be  $\mathcal{O}(1)$  or as extreme as  $\mathcal{O}(d_{\mathbf{Z}}^K)$ .

Fig. 8 shows that the subsets  $\Gamma_{\mathbf{U}} \subset \Gamma(\boldsymbol{\tau}, \mathbf{b})$  for  $\mathbf{U} \in \Sigma_{\mathbf{b}}$  are disconnected, while  $\tilde{\Gamma}_{\tilde{\mathbf{U}}} \subset \Gamma$  for  $\tilde{\mathbf{U}} \in \Sigma_{\mathbf{b}}^{(\max)}$  are cuboids in general. Therefore, it is clear that testing for intersection between  $h(\boldsymbol{\theta})$  and  $\tilde{\Gamma}_{\tilde{\mathbf{U}}}$  is more efficient than searching through the decision variables  $\mathbf{U}$  directly.

## F SUMMARY OF RELAXED IV METHODS

Paper	$d_X > 1$	$d_Z > 1$	Cont.	Nonlin.	2-sample	Inference	Violation	Feasible $\gamma_g$
Conley et al. (2012)	✓	✓	✓	✗	✗	✓	(A3)	Convex set
Nevo and Rosen (2012)	✓	✓	✓	✗	✗	✓	(A2)	Convex set
Ramsahai (2012)	✗	✗	✗	✓	✓	✗	(A2) $\vee$ (A3)	N/A
Bowden et al. (2015)	✗	✓	✓	✗	✓	✓	(A3) <sup>†</sup>	$d_Z \rightarrow \infty$
Kolesár et al. (2015)	✗	✓	✓	✗	✗	✓	(A3) <sup>†</sup>	$d_Z \rightarrow \infty$
Bowden et al. (2016a)	✗	✓	✓	✗	✗	✓	(A3)	Sparse
Kang et al. (2016)	✗	✓	✓	✗	✗	✗	(A3) <sup>*</sup>	Sparse
Silva and Evans (2016)	✗	✓	✗	✓	✓	✓	(A2) & (A3)	N/A
Hartwig et al. (2017)	✗	✓	✓	✗	✓	✓	(A3)	Mode zero
Guo et al. (2018)	✗	✗	✓	✓	✗	✓	(A3)	Mode zero
Windmeijer et al. (2019)	✗	✓	✓	✗	✗	✓	(A3) <sup>*</sup>	Sparse
Shaplant et al. (2019)	✗	✓	✓	✗	✗	✓	(A3) <sup>*</sup>	$L_0$ -norm
Bucur et al. (2020)	✗	✓	✓	✗	✓	-	(A3) <sup>*</sup>	$L_0$ -norm
Kang et al. (2020)	✗	✓	✓	✗	✗	✓	(A2) & (A3)	Point
Kuang et al. (2020)	✗	✓	✗	✗	✗	✗	(A2) <sup>*</sup> & (A3) <sup>*</sup>	Point
Hartford et al. (2021)	✗	✓	✓	✓	✗	✗	(A3) <sup>*</sup>	Mode zero
Vancak and Sjölander (2023)	✗	✗	✓	✓	✗	✓	(A2) & (A3)	Point
Xue et al. (2023)	✗	✓	✓	✗	✓	-	(A2) & (A3)	$L_0$ -norm
Watson et al. (2024)	✗	✓	✓	✗	✓	✓	(A3)	Convex set
budgetIV	✓	✓	✓	✓	✓	✓	(A2) & (A3)	Star domain

Table 1: **Summary of the constraints and affordances of some notable relaxed IV methods.** We exclude approaches that adhere to classical IV assumptions but provide nonparametric partial identification bounds as out of scope (e.g., Balke and Pearl (1997); Kilbertus et al. (2020); Levis et al. (2023)). We do the same with more generic algorithms not specifically designed for the IV setting (e.g., Hu et al. (2021); Duarte et al. (2023); Padh et al. (2023)). The columns denote, in order, if the method allows for: (1) multidimensional exposures; (2) multiple candidate instruments; (3) continuous data; (4) nonlinear structural equations; (5) summary statistic input; and (6) statistical inference (where - corresponds to point estimators in partially identifiable settings). The final two columns indicate (7) which IV assumptions (if any) may be violated and (8) the assumed geometry of the feasible region. The \* in column (7) indicates that only *some* candidate IVs are allowed to violate starred assumptions (typically fewer than 50%), while † indicates that the mechanism behind a candidate IV’s (A3) violation must be independent to the effect of the candidate on the exposure. The N/A entries in column (8) correspond to fully nonparametric models in which  $\gamma_g$  is undefined.