

Inference for Treatment Effects Conditional on Generalized Principal Strata using Instrumental Variables *

Yuehao Bai

Department of Economics
University of Southern California

yuehao.bai@usc.edu

Shunzhuang Huang

Booth School of Business
University of Chicago

shunzhuang.huang@chicagobooth.edu

Sarah Moon

Department of Economics
Massachusetts Institute of Technology

sarahmn@mit.edu

Andres Santos

Department of Economics
University of California–Los Angeles

andres@econ.ucla.edu

Azeem M. Shaikh

Department of Economics
University of Chicago

amshaikh@uchicago.edu

Edward J. Vytlačil

Department of Economics
Yale University

edward.vytlacil@yale.edu

November 11, 2024

Abstract

In a setting with a multi-valued outcome, treatment and instrument, this paper considers the problem of inference for a general class of treatment effect parameters. The class of parameters considered are those that can be expressed as the expectation of a function of the response type conditional on a generalized principal stratum. Here, the response type simply refers to the vector of potential outcomes and potential treatments, and a generalized principal stratum is a set of possible values for the response type. In addition to instrument exogeneity, the main substantive restriction imposed rules out certain values for the response types in the sense that they are assumed to occur with probability zero. It is shown through a series of examples that this framework includes a wide variety of parameters and assumptions that have been considered in the previous literature. A key result in our analysis is a characterization of the identified set for such parameters under these assumptions in terms of existence of a non-negative solution to linear systems of equations with a special structure. We propose methods for inference exploiting this special structure and recent results in [Fang et al. \(2023\)](#).

KEYWORDS: Multi-valued Treatments, Multi-arm Randomized Controlled Trial, Imperfect Compliance, Principal Strata, Treatment Effects, Endogeneity, Instrumental Variables, Partial Identification.

JEL classification codes: C31, C35, C36

*We thank Jiaying Gu, Sukjin Han, Nathan Kallus, Sokbae Lee, Derek Neal, Kirill Ponomarev, Vitor Possebom, Bernard Salanié and Panos Toulis for helpful comments. Shaikh acknowledges support from NSF Grant SES-2419008; Vytlačil acknowledges financial support from the Tobin Center for Economic Policy at Yale.

1 Introduction

In a setting of a multi-valued outcome, treatment and instrument, we propose a general approach for inference for a broad class of treatment effect parameters. The class of parameters considered are those that can be expressed as the expectation of a function of the response type conditional on a generalized principal stratum. Here, the response type simply refers to the vector of potential outcomes and potential treatments, and a generalized principal stratum is a set of possible values for the response type. In this way, our framework accommodates many parameters that have been considered previously in the literature, including both average and distributional treatment effect parameters, such as the probability of being strictly helped and the probability of not being hurt by the treatment. It further accommodates versions of these parameters conditional on sets of possible values of potential treatments, such as the local average treatment effect in [Imbens and Angrist \(1994\)](#) and average effects conditional on principal strata in [Frangakis and Rubin \(2002\)](#), or conditional on sets of possible values of potential outcomes, such as the parameters considered in [Heckman et al. \(1997\)](#) and [Heckman and Smith \(1998\)](#).

In addition to instrument exogeneity, the main substantive restriction imposed in our analysis is that certain values for the response types occur with probability zero. We show through a series of examples that our framework accommodates a wide variety of assumptions that have been considered in the previous literature, including the following: (i) restrictions in the analysis of randomized controlled trials (RCTs) with non-compliance considered, e.g., in [Imbens and Angrist \(1994\)](#); [Cheng and Small \(2006\)](#); (ii) generalizations of these restrictions considered in [Bai et al. \(2024\)](#); (iii) revealed preference-type restrictions on response types considered in [Kirkeboen et al. \(2016\)](#); [Kline and Walters \(2016\)](#); [Heckman and Pinto \(2018\)](#); and finally, (iv) restrictions on the ordering of potential treatments or on the ordering of potential outcomes considered in [Manski \(1997\)](#); [Manski and Pepper \(1998\)](#); [Machado et al. \(2019\)](#). We note that such restrictions are typically insufficient for identification of such parameters.

A key result in our analysis is a characterization of the identified set for such parameters under these assumptions in terms of existence of a non-negative solution to linear systems of equations with a special structure. Using this result, we exploit results in [Fang et al. \(2023\)](#) to develop a test for the null hypothesis that a pre-specified value for the parameter of interest lies in the identified set. Importantly, the resulting test remains well behaved even in “high-dimensional” settings, meaning it is uniformly consistent in level over a large class of distributions satisfying weak assumptions and permitting the support of the outcome, treatment and instrument or the support of the response types to be large relative to the sample size. Through test inversion, we may also construct confidence regions for the identified set that are uniformly consistent in level.

For a very limited number of special cases of our general framework, other methods for inference have been proposed in the previous literature; see, e.g., [Cheng and Small \(2006\)](#), [Bhattacharya et al. \(2008, 2012\)](#), and [Machado et al. \(2019\)](#). These approaches all rely upon closed-form expressions for the identified set for the parameter of interest, which must be computed on a case-by-case basis. Moreover, as explained in [Remark 4.1](#), such an approach is only computationally feasible when the support of the response types is very small. In contrast, our approach does not rely on such expressions and remains computationally

feasible even in high-dimensional settings. Finally, we note that a by-product of our analysis is that the bounds for the parameter of interest in [Cheng and Small \(2006\)](#) are not sharp; see [Remark 4.2](#) for further discussion. Further, as explained in [Remark 4.4](#), the approach to inference developed in [Cheng and Small \(2006\)](#) is invalid in that the resulting confidence regions may fail to cover the parameter of interest with (approximately) the desired coverage probability even in large samples.

The remainder of the paper is organized as follows. [Section 2](#) introduces our formal setup and notation. In [Section 3](#), we provide several examples of parameters and restrictions previously considered in the literature that are nested by our framework. We present our inference method in [Section 4](#). Proofs of all results can be found in the Appendix.

2 Setup and Notation

Denote by $Y \in \mathcal{Y}$ a multi-valued outcome of interest, by $D \in \mathcal{D}$ a multi-valued endogenous regressor (i.e., the treatment), and by $Z \in \mathcal{Z}$ a multi-valued instrumental variable. To rule out degenerate cases, we assume throughout that $2 \leq |\mathcal{Y}| < \infty$, $2 \leq |\mathcal{D}| < \infty$, and $2 \leq |\mathcal{Z}| < \infty$. For later use, we also define $\mathcal{M} \equiv \mathcal{Y} \times \mathcal{D} \times \mathcal{Z}$. Further denote by $Y(d) \in \mathcal{Y}$ the potential outcome if $D = d \in \mathcal{D}$ and by $D(z) \in \mathcal{D}$ the potential treatment if $Z = z \in \mathcal{Z}$. As usual, we assume that

$$Y = \sum_{d \in \mathcal{D}} Y(d) \mathbb{1}\{D = d\} \quad \text{and} \quad D = \sum_{z \in \mathcal{Z}} D(z) \mathbb{1}\{Z = z\} . \quad (1)$$

Let P denote the distribution of (Y, D, Z) and Q denote the distribution of $((Y(d) : d \in \mathcal{D}), (D(z) : z \in \mathcal{Z}), Z)$. For T defined by [\(1\)](#), we have

$$(Y, D, Z) = T((Y(d) : d \in \mathcal{D}), (D(z) : z \in \mathcal{Z}), Z) ,$$

and therefore $P = QT^{-1}$. In our discussion below, it will be convenient to define $R \equiv (R_y, R_d)$, where $R_y \equiv (Y(d) : d \in \mathcal{D})$ and $R_d \equiv (D(z) : z \in \mathcal{Z})$. Following [Heckman and Pinto \(2018\)](#), we will refer to R_d as the “treatment response type.” By analogy with this terminology, we will also refer to R_y as the “outcome response type” and to R as simply the “response type.”

Below we will require that $Q \in \mathbf{Q}$, where \mathbf{Q} is a class of distributions satisfying assumptions that we will specify. Different choices of \mathbf{Q} represent different assumptions that we impose on the distribution of potential outcomes and potential treatments. In this sense, \mathbf{Q} may be viewed as a model for potential outcomes and potential treatments.

Given P and a model \mathbf{Q} , we define the set of $Q \in \mathbf{Q}$ that can rationalize P as

$$\mathbf{Q}_0(P, \mathbf{Q}) = \{Q \in \mathbf{Q} : P = QT^{-1}\} .$$

We say \mathbf{Q} is consistent with P if and only if $\mathbf{Q}_0(P, \mathbf{Q}) \neq \emptyset$. For every model \mathbf{Q} considered in this paper, every $Q \in \mathbf{Q}$ is assumed to satisfy:

Assumption 2.1 (Instrument Exogeneity). $R \perp\!\!\!\perp Z$ under Q .

Our remaining restrictions on \mathbf{Q} will be formulated in terms of restrictions on possible values of R . These restrictions will be expressed in terms of $\mathcal{R} \subseteq \mathcal{Y}^{|\mathcal{D}|} \times \mathcal{D}^{|\mathcal{Z}|}$. In Section 3, we provide several different choices of \mathcal{R} that have been previously considered in the literature. For a given choice of \mathcal{R} , we will impose the following assumption on every $Q \in \mathbf{Q}$:

Assumption 2.2 (Response Type Restrictions). $Q\{R \in \mathcal{R}\} = 1$.

Following Frangakis and Rubin (2002), sets of the form $\{R_d = r_d\}$ for r_d a possible value of R_d are referred to as “principal strata.” By analogy with this terminology, we will refer to sets of the form $\{R \in \mathcal{R}'\}$ for $\mathcal{R}' \subseteq \mathcal{R}$ as “generalized principal strata.”

Finally, we define our parameters of interest. Using the notation above, the parameters we will consider can be written as

$$\theta(Q) \equiv E_Q[g(R) \mid R \in \mathcal{R}'] \tag{2}$$

for different choices of $g : \mathcal{R} \rightarrow \mathbb{R}$ and $\mathcal{R}' \subseteq \mathcal{R}$. Note that for $\theta(Q)$ in (2) to be well defined we require Q to satisfy $Q\{R \in \mathcal{R}'\} > 0$. A wide variety of parameters can be accommodated in this way. In Section 3, we will describe some specific parameters that have been considered previously in the literature corresponding to different choices of g and \mathcal{R}' , but we note that natural choices of g correspond to the effect of one treatment versus another, i.e., $g(R) = Y(d) - Y(d')$ for $d \in \mathcal{D}$ and $d' \in \mathcal{D}$, and the probability that one treatment leads to a larger outcome than another treatment, i.e., $g(R) = \mathbb{1}\{Y(d) > Y(d')\}$ for $d \in \mathcal{D}$ and $d' \in \mathcal{D}$. We further note that when $\mathcal{R}' = \mathcal{R}$, $\theta(Q)$ defined in (2) simplifies to $E_Q[g(R)]$. For fixed P and \mathbf{Q} , note that the identified set for $\theta(Q)$ under P relative to \mathbf{Q} is given by

$$\Theta_0(P, \mathbf{Q}) \equiv \{\theta(Q) : Q \in \mathbf{Q}_0(P, \mathbf{Q}) \text{ and } Q\{R \in \mathcal{R}'\} > 0\} . \tag{3}$$

$\Theta_0(P, \mathbf{Q})$ is nonempty whenever there exists at least one $Q \in \mathbf{Q}_0(P, \mathbf{Q})$ with $Q\{R \in \mathcal{R}'\} > 0$. By construction, this set is “sharp” in the sense that for any value in the set there exists $Q \in \mathbf{Q}_0(P, \mathbf{Q})$ with $Q\{R \in \mathcal{R}'\} > 0$ for which $\theta(Q)$ equals the prescribed value.

3 Examples

In this section, we show how to accommodate several examples from the previous literature in our framework. Our discussion focuses in particular on Assumption 2.1 and the specification of \mathcal{R} in Assumption 2.2, but, where the cited literature has emphasized specific parameters of interest, we additionally describe how those parameters of interest can be expressed as (2) for suitable choices of g and \mathcal{R}' .

Example 3.1 (*RCT with one-sided non-compliance*). Consider a multi-arm randomized controlled trial (RCT) with noncompliance, where $\mathcal{D} = \mathcal{Z} = \{0, \dots, |\mathcal{D}|-1\}$, $Z = d$ denotes random assignment to treatment d , and $D(d) = d$ denotes that the subject would comply with assignment if assigned to treatment d . In

this example, Q satisfies Assumption 2.1 because Z is randomly assigned. Suppose non-compliance to the assignment is one-sided in the sense that one can always take the control $d = 0$, but, for any other treatment $d \neq 0$, one can only take that treatment if assigned to that treatment. This restriction can be formulated in terms of Assumption 2.2 with $\mathcal{R} = \{(y(0), \dots, y(|\mathcal{D}| - 1), d(0), \dots, d(|\mathcal{D}| - 1)) : d(j) \in \{0, j\} \text{ for all } j \in \mathcal{D}\}$. ■

Example 3.2. Cheng and Small (2006) study the special case of Example 3.1 in which $\mathcal{D} = \mathcal{Z} = \{0, 1, 2\}$. As explained in Example 3.1, Q satisfies Assumption 2.1 because Z is randomly assigned. Their “Monotonicity I” assumption corresponds exactly to Assumption 2.2 with \mathcal{R} defined as in Example 3.1. They further consider imposing the restriction that $Q\{D(1) = 1 \mid D(2) = 2\} = 1$, i.e., subjects who would comply with assignment to treatment 2 would also comply with assignment to treatment 1. They argue that such an assumption is plausible in contexts where the “cost” of compliance with treatment 1 is lower than treatment 2. The combination of these two restrictions can be formulated in terms of Assumption 2.2 with $\mathcal{R} = \{(y(0), y(1), y(2), d(0), d(1), d(2)) : (d(0), d(1), d(2)) \in \{(0, 0, 0), (0, 1, 0), (0, 1, 2)\}\}$. In their application, Cheng and Small (2006) are particularly interested in the following parameters: (i) $E_Q[Y(j) - Y(0) \mid (D(0), D(1), D(2)) = (0, 1, 2)]$ for $j \in \{1, 2\}$; and (ii) $Q\{(D(0), D(1), D(2)) = r_d\}$ for different possible values of r_d . Each of these parameters can be written as (2) for appropriate choices of g and \mathcal{R}' . For example, the parameter in (i) equals $E_Q[g(R) \mid R \in \mathcal{R}']$ for $g(R) = Y(j) - Y(0)$ and $\mathcal{R}' = \{(y(0), y(1), y(2), d(0), d(1), d(2)) \in \mathcal{R} : (d(0), d(1), d(2)) = (0, 1, 2)\}$. ■

Example 3.3 (*Encouragement design*). Consider a multi-arm RCT with noncompliance, where $\mathcal{D} = \mathcal{Z} = \{0, \dots, |\mathcal{D}| - 1\}$, $Z = d$ denotes random assignment to treatment d , and $D(d) = d$ denotes that the subject would comply with assignment if assigned to treatment d . More generally, not necessarily in the context of an RCT, one can interpret $Z = d$ as random encouragement to treatment d and interpret $D(d) = d$ as the subject would take treatment d if encouraged to do so. In this example, Q satisfies Assumption 2.1 because Z is randomly assigned. Bai et al. (2024) generalize the “no-defier” restriction of Imbens and Angrist (1994) as follows:

$$Q\{D(d) \neq d, D(d') = d \text{ for some } d' \neq d\} = 0, \quad (4)$$

i.e., there is zero probability that a subject would not take treatment d if assigned to (encouraged to take) d but would take d if assigned (encouraged) to some other treatment $d' \neq d$. This restriction can be formulated in terms of Assumption 2.2 with $\mathcal{R} = \{(y(0), \dots, y(|\mathcal{D}| - 1), d(0), \dots, d(|\mathcal{D}| - 1)) : d(j) \neq j \Rightarrow d(k) \neq j \forall j, k \in \mathcal{D}\}$. ■

Example 3.4 (*RCT with close substitute*). Kline and Walters (2016) consider an RCT with a “close substitute” to study the effects of preschooling on educational outcomes. In their setting, $D \in \mathcal{D} = \{0, 1, 2\}$, where $D = 0$ denotes home care (no preschool), $D = 2$ denotes a preschool program called Head Start, and $D = 1$ denotes preschools other than Head Start, i.e., the close substitute. Let $Z \in \mathcal{Z} = \{0, 1\}$ denote an indicator variable for an offer to attend Head Start. Assumption 2.1 holds because Z is randomly assigned. Kline and Walters (2016) impose the restriction that

$$Q\{D(1) = 2 \mid D(0) \neq D(1)\} = 1. \quad (5)$$

The condition in (5) states that if a family’s schooling choice changes upon receiving a Head Start offer, then they must choose Head Start when receiving the offer. In other words, it cannot be the case that upon receiving a Head Start offer, a family switches from no preschool to preschools other than Head Start, or the other way around. The restriction can be formulated in terms of Assumption 2.2 with

$$\mathcal{R} = \{(y(0), y(1), y(2), d(0), d(1)) : (d(0), d(1)) \in \{(0, 0), (0, 2), (1, 1), (1, 2), (2, 2)\}\} .$$

Kline and Walters (2016) show that the Wald estimand identifies a weighted combination of what they call “sub-LATEs”:

$$\frac{E[Y | Z = 1] - E[Y | Z = 0]}{P\{D = 2 | Z = 1\} - P\{D = 2 | Z = 0\}} = S_c \text{SubLATE}_{12} + (1 - S_c) \text{SubLATE}_{02} ,$$

where

$$\text{SubLATE}_{12} = E_Q[Y(2) - Y(1) | D(1) = 2, D(0) = 1]$$

$$\text{SubLATE}_{02} = E_Q[Y(2) - Y(0) | D(1) = 2, D(0) = 0]$$

and $S_c = Q\{D(1) = 2, D(0) = 1 | D(1) = 2, D(0) \neq 2\}$. Kline and Walters (2016) establish conditions under which optimal policy depends upon these “sub-LATEs.” Here, S_c and the “sub-LATEs” can all be written as (2) for appropriate choices of g and \mathcal{R}' . ■

Example 3.5. Kirkeboen et al. (2016) study the effects of fields of study on earnings. In their setting, $\mathcal{D} = \{0, 1, 2\}$ represent three fields of study, ordered by their (soft) admission cutoffs from the lowest to the highest. The instrument is $Z \in \{0, 1, 2\}$, with $Z = 1$ when the student crosses the (soft) admission cutoff for field 1, $Z = 2$ when the student crosses the (soft) admission cutoff for field 2, and $Z = 0$ otherwise. The authors assume that Z is exogenous in the sense that Q satisfies Assumption 2.1 and impose the following monotonicity conditions:

$$Q\{D(1) = 1 | D(0) = 1\} = 1 , \tag{6}$$

$$Q\{D(2) = 2 | D(0) = 2\} = 1 . \tag{7}$$

The conditions in (6)–(7) require that crossing the cutoff for field 1 or 2 weakly encourages them towards that field. They further impose the following “irrelevance” conditions:

$$Q\{\mathbb{1}\{D(1) = 2\} = \mathbb{1}\{D(0) = 2\} | D(0) \neq 1, D(1) \neq 1\} = 1 , \tag{8}$$

$$Q\{\mathbb{1}\{D(2) = 1\} = \mathbb{1}\{D(0) = 1\} | D(0) \neq 2, D(2) \neq 2\} = 1 . \tag{9}$$

The condition in (8) states that if crossing the cutoff for field 1 does not cause the student to switch to field 1, then it does not cause them to switch to or away from field 2. A similar interpretation applies to (9). The restrictions of (6)–(9) can be formulated in terms of Assumption 2.2 with $\mathcal{R} = \{(y(0), y(1), y(2), d(0), d(1), d(2)) : (d(0), d(1), d(2)) \in \{(0, 0, 0), (0, 0, 2), (0, 1, 0), (0, 1, 2), (1, 1, 1), (1, 1, 2), (2, 1, 2), (2, 2, 2)\}\}$. ■

Example 3.6 (*Restrictions from WARP*). Heckman and Pinto (2018) consider a setting in which there is a voucher Z that subsidizes in different ways that we specify below the purchase of three different cars, that we denote by A , B and C . They further assume that the voucher is randomly assigned, so that Assumption 2.1 holds. The treatment D corresponds to the purchase of the different cars; let $D = 1$ correspond to the purchase of car A , $D = 2$, correspond to the purchase of car B , and by $D = 3$ correspond to the purchase of car C . In this setting, Heckman and Pinto (2018) consider a series of examples in which they use the Weak Axiom of Revealed Preference (WARP) to restrict treatment response types, each of which can be formulated in terms of Assumption 2.2 for appropriate choice of \mathcal{R} .

- (i) In their leading example, $Z = 0$ corresponds to no voucher, $Z = 1$ corresponds to a voucher that subsidizes the consumption of A , and $Z = 2$ corresponds to a voucher that subsidizes the consumption of either B or C . WARP generates the restriction in Table III of Heckman and Pinto (2018), which can be formulated in terms of Assumption 2.2 with $\mathcal{R} = \{(y(0), y(1), y(2), d(0), d(1), d(2)) : (d(0), d(1), d(2)) \in \{(0, 0, 0), (0, 0, 1), (0, 0, 2), (1, 0, 1), (1, 1, 1), (2, 0, 2), (2, 2, 2)\}\}$.
- (ii) In second example, $Z = 0$ corresponds to no voucher, $Z = 1$ corresponds to a voucher that subsidizes the consumption of B , and $Z = 2$ corresponds to a voucher that subsidizes the consumption of B or C . WARP generates the restriction in Table V of Heckman and Pinto (2018), which can be formulated in terms of Assumption 2.2 with $\mathcal{R} = \{(y(0), y(1), y(2), d(0), d(1), d(2)) : (d(0), d(1), d(2)) \in \{(0, 0, 0), (0, 0, 2), (0, 1, 1), (0, 1, 2), (1, 1, 1), (2, 2, 2), (2, 1, 2)\}\}$.
- (iii) Finally, in a third example, $Z = 0$ corresponds to a voucher that subsidizes the consumption of C , $Z = 1$ corresponds to a voucher that subsidizes the consumption of B , and $Z = 2$ corresponds to a voucher that subsidizes the consumption of B or C . WARP generates the restriction in Table VI of Heckman and Pinto (2018), which can be formulated in terms of Assumption 2.2 with $\mathcal{R} = \{(y(0), y(1), y(2), d(0), d(1), d(2)) : (d(0), d(1), d(2)) \in \{(0, 0, 0), (0, 1, 1), (1, 1, 1), (2, 0, 2), (2, 1, 1), (2, 1, 2), (2, 2, 2)\}\}$.

Heckman and Pinto (2018) provide several additional examples that also fit within our framework. ■

Example 3.7 (*Ordered monotonicity with known direction*). Consider a setting where both the treatment and the instrument admit a natural ordering and write $\mathcal{D} = \{0, \dots, |\mathcal{D}| - 1\}$ and $\mathcal{Z} = \{0, \dots, |\mathcal{Z}| - 1\}$. Suppose further that Z is randomly assigned, so Assumption 2.1 holds. In many applications, it is reasonable to assume that $Q\{D(j) \geq D(k)\} = 1$ for any pair of values for the instrument $j, k \in \mathcal{Z}$ with $j \geq k$. For example, the treatment may be number of condoms purchased, the instrument may be level of subsidy for condoms, and the restriction is that one would purchase at least as many condoms with a greater subsidy as with a lower subsidy. This restriction can be formulated in terms of Assumption 2.2 with $\mathcal{R} = \{(y(0), \dots, y(|\mathcal{D}| - 1), d(0), \dots, d(|\mathcal{D}| - 1)) : d(|\mathcal{Z}| - 1) \geq \dots \geq d(0)\}$. ■

Example 3.8 (*Monotone treatment response*). Consider a setting where the treatment admits a natural ordering and write $\mathcal{D} = \{0, \dots, |\mathcal{D}| - 1\}$. Suppose further that Z is randomly assigned, so Assumption 2.1 holds. Following Manski (1997), in many applications it is reasonable to assume that $Q\{Y(j) \geq Y(k)\} = 1$ for any pair of treatments $j, k \in \mathcal{D}$ with $j \geq k$. For example, in Manski and Pepper (1998), the treatment

is years of schooling and the outcome is $\log(\text{wage})$, and the restriction is that higher years of schooling leads to (weakly) higher $\log(\text{wage})$. This restriction can be formulated in terms of Assumption 2.2 with $\mathcal{R} = \{(y(0), \dots, y(|\mathcal{D}| - 1), d(0), \dots, d(|\mathcal{Z}| - 1)) : y(|\mathcal{D}| - 1) \geq \dots \geq y(0)\}$. ■

Remark 3.1. Angrist and Imbens (1995) consider the monotonicity assumption of Example 3.7 but with the direction of the monotonicity not imposed *a priori*, i.e., they consider the restriction that, for any pair of values for the instrument $j, k \in \mathcal{Z}$, either $Q\{D(j) \geq D(k)\} = 1$ or $Q\{D(j) \leq D(k)\} = 1$. Such a restriction is also analyzed in Vytlacil (2006) and in Heckman and Pinto (2018), where the latter paper terms it “ordered monotonicity.” This restriction is equivalent to imposing $Q\{R \in \mathcal{R}_\pi\} = 1$ for some $\pi \in \Pi(|\mathcal{Z}|)$, where $\Pi(|\mathcal{Z}|)$ is the set of all permutations of $\{0, \dots, |\mathcal{Z}| - 1\}$ and $\mathcal{R}_\pi = \{(y(0), \dots, y(|\mathcal{D}| - 1), d(0), \dots, d(|\mathcal{D}| - 1)) : d(\pi(|\mathcal{Z}| - 1)) \geq \dots \geq d(\pi(0))\}$. A similar generalization can be applied to Example 3.8. In particular, Machado et al. (2019) study such a generalization with $|\mathcal{D}| = |\mathcal{Z}| = 2$. Such examples fall outside of the framework we consider in this paper. ■

Example 3.9 (*Harmless treatment*). Consider a setting where Z is randomly assigned, so Assumption 2.1 holds, and there is a baseline treatment, i.e., the control, corresponding to $0 \in \mathcal{D}$. In many applications, it is reasonable to assume the remaining treatments are all harmless relative to the control. For example, in Angrist et al. (2009), the outcome is academic performance, the control is no treatment, and the non-control treatments are (i) providing students with academic peer-advising service; (ii) financial incentives for good academic performance; and (iii) both (i) and (ii). It is therefore natural to assume that $Q\{Y(d) \geq Y(0)\} = 1$ for all $d \in \mathcal{D}$. This restriction can be formulated in terms of Assumption 2.2 with $\mathcal{R} = \{(y(0), \dots, y(|\mathcal{D}| - 1), d(0), \dots, d(|\mathcal{Z}| - 1)) : y(d) \geq y(0) \text{ for all } d \in \mathcal{D}\}$. ■

Remark 3.2. In many of the examples discussed above, Y may be an ordinal outcome. In such instances, average treatment effects may not be interpretable; researchers may therefore consider instead the following types of parameters: (i) $Q\{Y_j > Y_k \mid R \in \mathcal{R}'\}$, the conditional probability of benefit of treatment j versus treatment k ; (ii) $Q\{Y_j \geq Y_k \mid R \in \mathcal{R}'\}$, the conditional probability of no harm of treatment j versus treatment k ; or (iii) $Q\{Y_j > Y_k \mid R \in \mathcal{R}'\} - Q\{Y_k > Y_j \mid R \in \mathcal{R}'\}$, the conditional relative treatment effect of treatment j versus k . Each of these parameters can be written as (2) for appropriate choices of g . For previous analysis of such parameters when $\mathcal{R}' = \mathcal{R}$ and $|\mathcal{D}| = 2$ or 3 , see Lu et al. (2018), Huang et al. (2019) and Gabriel et al. (2024). ■

4 Inference

Let $(Y_i, D_i, Z_i), i = 1, \dots, n \sim P \in \mathbf{P}$, where \mathbf{P} is a “large” class of possible distributions for P that we will specify below. In this section, we consider the problem of testing

$$H_0 : P \in \mathbf{P}_0 \text{ versus } H_1 : P \in \mathbf{P} \setminus \mathbf{P}_0 \quad (10)$$

at level $\alpha \in (0, 1)$, where, for a pre-specified value of θ_0 ,

$$\mathbf{P}_0 = \{P \in \mathbf{P} : \theta_0 \in \Theta_0(P, \mathbf{Q})\} . \quad (11)$$

As explained further below, using such a test, we will be able to construct confidence regions for $\theta(Q)$ through test inversion.

The key insight underlying the construction of our test is the following lemma, which provides a convenient reformulation of \mathbf{P}_0 in terms of existence of a non-negative solution to a (possibly under-determined) system of linear equations in which the “coefficients” are known. The statement of the lemma involves the following quantity:

$$\beta(P) \equiv ((p_{ydz}(P) : (y, d, z) \in \mathcal{M}), 1, 0)' , \quad (12)$$

where $p_{ydz}(P) \equiv P\{Y = y, D = d | Z = z\}$. Using this notation, we have the following result:

Lemma 4.1. *Suppose every $Q \in \mathbf{Q}$ satisfies Assumptions 2.1 and 2.2. Then, for $\beta(P)$ defined in (12) and a matrix A defined in the beginning of Appendix A that depends only on \mathcal{R} in Assumption 2.2, θ_0 in (11), and the quantities g and \mathcal{R}' in the definition of $\theta(Q)$ in (2),*

$$\mathbf{P}_0 \subseteq \text{cl}(\mathbf{P}_0) \subseteq \{P \in \mathbf{P} : Ax = \beta(P) \text{ for some } x \geq 0\} , \quad (13)$$

where $\text{cl}(\mathbf{P}_0)$ denotes the closure of \mathbf{P}_0 with respect to the usual topology. Furthermore, the final inclusion in (13) is an equality whenever $\mathbf{P}_0 \neq \emptyset$.

Remark 4.1. For the special case in which $Q\{R \in \mathcal{R}'\} > 0$ for all $Q \in \mathbf{Q}$ and is known or identified, in the sense that $\{Q\{R \in \mathcal{R}'\} : Q \in \mathbf{Q}_0(P, \mathbf{Q})\}$ is a singleton for all $P \in \mathbf{P}$ for which $\mathbf{Q}_0(P, \mathbf{Q}) \neq \emptyset$, it is possible to use the characterization of \mathbf{P}_0 in Lemma 4.1 to derive $L(P)$ and $U(P)$ such that $\Theta_0(P, \mathbf{Q}) = [L(P), U(P)]$ for all $P \in \mathbf{P}$ for which $\mathbf{Q}_0(P, \mathbf{Q}) \neq \emptyset$. As in Balke and Pearl (1997), the key idea is to express $L(P)$ and $U(P)$ as the values of suitable linear programs and to use the duals of these programs to obtain expressions for $L(P)$ and $U(P)$ in terms of $p_{ydz}(P)$ through vertex enumeration. Computing $L(P)$ and $U(P)$ in this way rapidly becomes computationally prohibitive as the support of (Y, D, Z) becomes large. We describe this procedure in more detail in Appendix B, and further develop a related procedure when $Q\{R \in \mathcal{R}'\}$ is not identified. Following the approach in Machado et al. (2019), it is possible to use such expressions for inference, but a virtue of the approach we pursue here is that it does not rely on knowledge of these expressions. ■

Remark 4.2. As explained in Example 3.2, Cheng and Small (2006) consider an RCT with one-sided noncompliance and $|\mathcal{D}| = 3$. In this setting, they study a variety of parameters of interest, but are particularly interested in $E_Q[Y(j) - Y(0) | (D(0), D(1), D(2)) = (0, 1, 2)]$ for $j \in \{1, 2\}$, $E_Q[Y(1) - Y(0) | (D(0), D(1), D(2)) = (0, 1, 0)]$ and $E_Q[Y(2) - Y(0) | (D(0), D(1), D(2)) = (0, 0, 2)]$. For each of these parameters, they derive upper and lower bounds, $\tilde{L}(P)$ and $\tilde{U}(P)$, respectively, under their “Monotonicity I” assumption alone and in combination with their “Monotonicity II” assumption. We show in Appendix C that the bounds they derived under “Monotonicity I” assumption alone, in which $Q\{(D(0), D(1), D(2)) = (0, i, j)\}$ for $i \in \{0, 1\}, j \in \{0, 2\}$ is not identified, need not be sharp in the sense that (i) $L(P) \geq \tilde{L}(P)$ and $U(P) \leq \tilde{U}(P)$ for all $P \in \mathbf{P}$ for which $\mathbf{Q}_0(P, \mathbf{Q}) \neq \emptyset$; and (ii) $L(P) > \tilde{L}(P)$ or $U(P) < \tilde{U}(P)$ holds for some $P \in \mathbf{P}$ for which $\mathbf{Q}_0(P, \mathbf{Q}) \neq \emptyset$, where $L(P)$ and $U(P)$ are the sharp bounds derived using our approach in Remark 4.1. ■

In order to test (10), it follows from Lemma 4.1 that we may exploit recently developed tests for the right-hand side of the inclusion in (13). We follow the approach developed in Fang et al. (2023), which has been shown to behave well even when the number of rows and/or columns of A are large relative to the sample size n . For other approaches to the same problem, see Bai et al. (2022).

We begin our description of the proposed test by defining the test statistic. In order to do so, we require some further notation. Denote by \hat{P}_n the empirical distribution of $(Y_i, D_i, Z_i), i = 1, \dots, n$ and define $\hat{\beta}_n \equiv \beta(\hat{P}_n)$. Let $\Omega^e(P)$ be the standard deviation matrix of the limiting distribution of $\sqrt{n}(\hat{\beta}_n - \beta(P))$ and $\Omega^i(P)$ be the standard deviation matrix of the limit in distribution of $\sqrt{n}AA^\dagger(\hat{\beta}_n - \beta(P))$, where A^\dagger is the Moore-Penrose inverse of A ; define $\hat{\Omega}_n^e \equiv \Omega^e(\hat{P}_n)$ and $\hat{\Omega}_n^i \equiv \Omega^i(\hat{P}_n)$. Using this notation, our proposed test rejects for large values of

$$T_n \equiv \max \left\{ \sup_{s \in \hat{\mathcal{V}}_n^e} \sqrt{n} \langle s, (I - AA^\dagger) \hat{\beta}_n \rangle, \sup_{s \in \hat{\mathcal{V}}_n^i} \sqrt{n} \langle A^\dagger s, A^\dagger \hat{\beta}_n \rangle \right\},$$

where $\hat{\mathcal{V}}_n^e \equiv \{s \in \mathbb{R}^{|\mathcal{M}|+2} : \|\hat{\Omega}_n^e s\|_1 \leq 1\}$ and $\hat{\mathcal{V}}_n^i \equiv \{s \in \mathbb{R}^{|\mathcal{M}|+2} : A^\dagger s \leq 0 \text{ and } \|\hat{\Omega}_n^i AA^\dagger s\|_1 \leq 1\}$.

Next, we define the critical value with which we compare T_n . To this end, define, for $\lambda_n \leq 1$, $\hat{\mathbb{U}}_n(s) \equiv \lambda_n \sqrt{n} \langle A^\dagger s, A^\dagger \hat{\beta}_n^r \rangle$, where

$$\hat{\beta}_n^r \in \operatorname{argmin}_b \sup_{s \in \hat{\mathcal{V}}_n^i} |\langle A^\dagger s, A^\dagger (\hat{\beta}_n - b) \rangle| \text{ s.t. } Ax = b \text{ for some } x \geq 0 \text{ and } b = (b'_u, 1, 0)' \text{ for } b_u \in \mathbb{R}^{|\mathcal{M}|}.$$

For $x \in \mathbb{R}$, $\mathcal{V}^e \subseteq \mathbb{R}$, $\mathcal{V}^i \subseteq \mathbb{R}$, $\mathbb{U} : \mathbb{R}^{|\mathcal{M}|+2} \rightarrow \mathbb{R}$, and $P \in \mathbf{P}$, define

$$J_n(x, \mathcal{V}^e, \mathcal{V}^i, \mathbb{U}, P) \equiv P \left\{ \max \left\{ \sup_{s \in \mathcal{V}^e} \langle s, (I - AA^\dagger) \hat{\beta}_n \rangle, \sup_{s \in \mathcal{V}^i} \left(\langle A^\dagger s, A^\dagger (\hat{\beta}_n - \beta(P)) \rangle + \mathbb{U}(s) \right) \right\} \leq x \right\},$$

and for $\alpha \in (0, 1)$, define

$$J_n^{-1}(1 - \alpha, \mathcal{V}^e, \mathcal{V}^i, \mathbb{U}, P) \equiv \inf \{x \in \mathbb{R} : J_n(x, \mathcal{V}^e, \mathcal{V}^i, \mathbb{U}, P) \geq 1 - \alpha\}.$$

In terms of this notation, the critical value is given by $\hat{c}_n(1 - \alpha) \equiv J_n^{-1}(1 - \alpha, \hat{\mathcal{V}}_n^e, \hat{\mathcal{V}}_n^i, \hat{\mathbb{U}}_n, \hat{P}_n)$.

Fang et al. (2023) show that the test $\phi_n(\theta_0) \equiv \mathbb{1}\{T_n > \hat{c}_n(1 - \alpha)\}$ satisfies

$$\limsup_{n \rightarrow \infty} \sup_{P \in \mathbf{P}_0} E_P[\phi_n] \leq \alpha$$

under weak assumptions on \mathbf{P} . We omit a detailed discussion of these assumptions, but emphasize our prior point that the assumptions in Fang et al. (2023) permit that the number of rows and/or columns of A are large relative to the sample size n . This feature is important for the types of applications considered in our paper since these dimensions can grow rapidly as the number of possible values for the outcome, treatment and instrument grow. Indeed, if each of these quantities only take on k values, then the matrix A in the statement of Lemma 4.1 has $k^3 + 2$ rows and $k^k \times k^k$ columns.

Using the test above, we can construct a confidence region for $\theta(Q)$. To this end, define the confidence

region as

$$C_n = \{\theta_0 \in \mathbb{R} : \phi_n(\theta_0) = 0\} .$$

It is straightforward to show that

$$\liminf_{n \rightarrow \infty} \inf_{P \in \mathbf{P}} \inf_{\theta \in \Theta_0(P, \mathbf{Q})} P\{\theta \in C_n\} \geq 1 - \alpha .$$

Remark 4.3. As explained previously, our main substantive restriction, Assumption 2.2, imposes that response types take on certain values with probability zero. Our framework can easily be modified, however, to accommodate the restriction that response types take on certain values with at most (or at least) some pre-specified probability. Our framework thereby permits sensitivity analysis, as in Masten and Poirier (2020). For instance, in Example 3.3, one may relax (4) so that it instead specifies that the left-hand side is at most some pre-specified amount $\epsilon > 0$, and explore how inferences on $\theta(Q)$ change as one varies ϵ . In this way, we can generalize the analysis of Noack (2021) to multi-valued instruments and treatments. To provide another example, consider the setting of Example 3.9. Instead of assuming that each treatment is “harmless,” we may modify the example so that each treatment only causes “harm” with probability at most $\epsilon > 0$. If, for example, the treatment is a surgery, then ϵ may be interpreted as the probability of a surgical complication. ■

Remark 4.4. As explained in Remark 4.2, Cheng and Small (2006) derive for a particular parameter upper and lower bounds, $\tilde{L}(P)$ and $\tilde{U}(P)$, respectively, under different assumptions. They further suggest using the bootstrap for inference on $[\tilde{L}(P), \tilde{U}(P)]$ in the following way. Denote by $\ell(\alpha/2, P)$ the $\alpha/2$ quantile of the distribution of $\tilde{L}(\hat{P}_n)$ under P and by $u(1 - \alpha/2, P)$ the $1 - \alpha/2$ quantile of the distribution of $\tilde{U}(\hat{P}_n)$ under P . In terms of this notation, the proposed confidence region for $[\tilde{L}(P), \tilde{U}(P)]$ is given by $[\ell(\alpha/2, \hat{P}_n), u(1 - \alpha/2, \hat{P}_n)]$. Because the functionals $\tilde{L}(P)$ and $\tilde{U}(P)$ involve the maximum and minimum over functionals of P , they are only directionally differentiable and not fully differentiable when the maximum and minimum are attained at multiple functionals. Hence, it follows from Fang and Santos (2019) that the bootstrap quantiles $\ell(\alpha/2, \hat{P}_n)$ and $u(1 - \alpha/2, \hat{P}_n)$ are not consistent for the population quantiles $\ell(\alpha/2, P)$ and $u(1 - \alpha/2, P)$ that justify the validity of the proposed confidence region. ■

Remark 4.5. It is possible to modify the test described above to test the null hypothesis that the model is correctly specified in the sense that $\mathbf{Q}_0(P, \mathbf{Q}) \neq \emptyset$. In order to do so, let A_0 and $\beta_0(P)$ denote all but the last two rows of A and $\beta(P)$. By arguing as in the proof of Lemma 4.1, it is possible to show that the model is correctly specified if and only if $A_0 x = \beta_0(P)$ for some $x \geq 0$. The test described above can therefore be adapted in a straightforward way to test the null hypothesis that the model is correctly specified. Furthermore, in Appendix D, we show how to use a strategy like the one described in Remark 4.1 to derive a set of analytical inequalities that hold if and only if $\mathbf{Q}_0(P, \mathbf{Q}) \neq \emptyset$. ■

A Proof of Lemma 4.1

We begin by defining A . To this end, let $\mathcal{M} \equiv \mathcal{Y} \times \mathcal{D} \times \mathcal{Z}$ and $\mathcal{N} \equiv \mathcal{Y}^{|\mathcal{D}|} \times \mathcal{D}^{|\mathcal{Z}|}$. We now define a matrix A with $|\mathcal{M}| + 2$ rows and $|\mathcal{N}|$ columns. In order to describe the matrix, index the first $|\mathcal{M}|$ rows of A by $(y, d, z) \in \mathcal{M}$ and the columns of A by $r = ((y(d) : d \in \mathcal{D}), (d(z) \in \mathcal{Z})) \in \mathcal{N}$. The $(y, d, z) \times r$ element of A is given by $\mathbb{1}\{y(d) = y, d(z) = d\}$; the $(|\mathcal{M}| + 1) \times r$ element of A is given by $\mathbb{1}\{r \in \mathcal{R}\}$; finally, the $(|\mathcal{M}| + 2) \times r$ element of A is given by $(g(r) - \theta_0) \mathbb{1}\{r \in \mathcal{R}'\}$.

Suppose $P \in \mathbf{P}_0$. Then, $\theta_0 \in \Theta_0(P, \mathbf{Q})$, so that there exists $Q \in \mathbf{Q}_0(P, \mathbf{Q})$ such that $\theta_0 = \theta(Q)$. Recall that $Q \in \mathbf{Q}_0(P, \mathbf{Q})$ if and only if $Q \in \mathbf{Q}$ and $P = QT^{-1}$. By assumption, $Q \in \mathbf{Q}$, and hence satisfies both Assumptions 2.1 and 2.2. Because Q satisfies Assumption 2.1 and $P = QT^{-1}$,

$$p_{yd|z}(P) \equiv P\{Y = y, D = d|Z = z\} = Q\{Y(d) = y, D(z) = d|Z = z\} = Q\{Y(d) = y, D(z) = d\}. \quad (14)$$

Let q be the column vector with $|\mathcal{N}|$ elements indexed by r , where the r element is $q(r) = Q\{R = r\}$. It follows from (14) and the definition of A that the (y, d, z) row of $Aq = \beta(P)$ holds. Next, because Q satisfies Assumption 2.2,

$$\sum_{r \in \mathcal{R}} q(r) = 1,$$

and it follows from the definition of A that the $|\mathcal{M}| + 1$ row of $Aq = \beta(P)$ holds. Finally, because $\theta_0 = \theta(Q)$, by the definition of the conditional expectation,

$$\theta_0 = \frac{E_Q[g(R) \mathbb{1}\{R \in \mathcal{R}'\}]}{Q\{R \in \mathcal{R}'\}}. \quad (15)$$

Rearranging, we have $E_Q[(g(R) - \theta_0) \mathbb{1}\{R \in \mathcal{R}'\}] = 0$, so that

$$\sum_{r \in \mathcal{R}'} (g(r) - \theta_0) q(r) = 0,$$

and it follows from the definition of A that the $|\mathcal{M}| + 2$ row of $Aq = \beta(P)$ holds. Finally, we show that $\text{cl}(\mathbf{P}_0) = \tilde{\mathbf{P}}_0 \equiv \{P \in \mathbf{P} : Ax = \beta(P) \text{ for some } x \geq 0\}$. Note $\mathbf{P}_0 \subseteq \tilde{\mathbf{P}}_0$ and $\tilde{\mathbf{P}}_0$ is closed as a polyhedron by Theorem 4.9 in [Bertsimas and Tsitsiklis \(1997\)](#), so $\text{cl}(\mathbf{P}_0) \subseteq \tilde{\mathbf{P}}_0$. On the other hand, consider $\tilde{P} \in \tilde{\mathbf{P}}_0$. By definition, there exists $\tilde{Q} \in \mathbf{Q}$ such that $\tilde{P} = \tilde{Q}T^{-1}$ and $E_{\tilde{Q}}[g(R) \mathbb{1}\{R \in \mathcal{R}'\}] = \theta_0 \tilde{Q}\{R \in \mathcal{R}'\}$. By assumption, \mathbf{P}_0 is nonempty, so there exists $P \in \mathbf{P}_0$ such that $P = QT^{-1}$ for some $Q \in \mathbf{Q}$, $Q\{R \in \mathcal{R}'\} > 0$, and $\theta(Q) = \theta_0$. For a sequence $\lambda_n > 0$ and $\lambda_n \downarrow 0$, define $P_n = \lambda_n P + (1 - \lambda_n) \tilde{P}$ and $Q_n = \lambda_n Q + (1 - \lambda_n) \tilde{Q}$. Then, $P_n = Q_n T^{-1}$, $Q_n\{R \in \mathcal{R}'\} > 0$, and $E_{Q_n}[g(R) \mathbb{1}\{R \in \mathcal{R}'\}] = \theta_0 Q_n\{R \in \mathcal{R}'\}$, so $\theta(Q_n) = \theta_0$. Therefore, $P_n \in \mathbf{P}_0$ for each n . Because $P_n \rightarrow \tilde{P}$ in the usual topology as $n \rightarrow \infty$, $\tilde{P} \in \text{cl}(\mathbf{P}_0)$, and the result follows. ■

B Details for Remark 4.1

In this section, we describe a method for deriving analytical expressions for $\Theta_0(P, \mathbf{Q})$ as an interval whose lower and upper endpoints are some functions of P , $L(P)$ and $U(P)$ respectively. If Q satisfies Assumption 2.1 and $P = QT^{-1}$, then $P\{Z = z\} = Q\{Z = z\}$ and (14) holds. Therefore, in what follows, we disregard the distribution of Z under Q and identify Q with the distribution of R , and in turn with the column vector q . Correspondingly, we identify P with $p = \{p_{ydz} : (y, d, z) \in \mathcal{M}\}$. Let A_0 and $\beta_0(P)$ denote the first $|\mathcal{M}| + 1$ rows of A and $\beta(P)$. Because we assume every $Q \in \mathbf{Q}$ satisfies Assumption 2.1,

$$\mathbf{Q}_0(P, \mathbf{Q}) = \{q : A_0q = \beta_0(P), q \geq 0\} . \quad (16)$$

B.1 When $Q\{R \in \mathcal{R}'\}$ is known or identified

Suppose $Q\{R \in \mathcal{R}'\} > 0$ for all $Q \in \mathbf{Q}$ and is known or identified, so that it is a function of P for all P such that $\mathbf{Q}_0(P, \mathbf{Q}) \neq \emptyset$. Denote such a function by $a(P)$. Recall q introduced in the proof of Lemma 4.1, which is indexed by r such that $q(r) = Q\{R = r\}$. For all $Q \in \mathbf{Q}_0(P, \mathbf{Q})$, $\theta(Q) = \frac{1}{a(P)} \sum_{r \in \mathcal{R}'} g(r)q(r)$. Let c be a column vector indexed by r , with the r element given by $g(r)/a(P)$. Here we suppress the dependence of c on P for convenience. With the notation above,

$$\Theta_0(P, \mathbf{Q}) = c' \mathbf{Q}_0(P, \mathbf{Q}) .$$

Before proceeding, we show $\Theta_0(P, \mathbf{Q})$ is a closed interval. Because $\mathbf{Q}_0(P, \mathbf{Q})$ is a polyhedron (in standard form), Corollary 2.5 in Bertsimas and Tsitsiklis (1997) implies $\Theta_0(P, \mathbf{Q}) = c' \mathbf{Q}_0(P, \mathbf{Q})$ is a polyhedron in \mathbb{R} . Furthermore, $\mathbf{Q}_0(P, \mathbf{Q})$ is obviously bounded, so $\Theta_0(P, \mathbf{Q})$ is also bounded. A bounded polyhedron in \mathbb{R} is simply the intersection of bounded closed intervals, so it is itself a bounded closed interval.

Next, note $L(P)$ is the solution to

$$\begin{aligned} \min_q \quad & c'q \\ \text{subject to} \quad & A_0q = \beta_0(P) \\ & q \geq 0 . \end{aligned} \quad (17)$$

Because the feasible set $\mathbf{Q}_0(P, \mathbf{Q})$ in (17) is in addition bounded, the optimal cost has to be finite, so it follows from Theorem 2.8 in Bertsimas and Tsitsiklis (1997) that (17) has an optimal solution. Therefore, by the strong duality theorem (Theorem 4.4, Bertsimas and Tsitsiklis, 1997), the optimal value in (17) is the same as the optimal value of its dual:

$$\begin{aligned} \max_r \quad & \beta_0(P)'r \\ \text{subject to} \quad & A_0'r \leq c . \end{aligned} \quad (18)$$

Note that in (18) the constraints do not involve P . It follows from the resolution theorem (Exercise 4.47,

Bertsimas and Tsitsiklis, 1997) that the feasible set in (18) can be written as

$$\left\{ \sum_{1 \leq j \leq J} \theta_j r_j^{\text{ex}} + \sum_{1 \leq \ell \leq L} \lambda_\ell r_\ell^{\text{ray}} : \theta_j \geq 0, \sum_{1 \leq j \leq J} \theta_j = 1, \lambda_\ell \geq 0 \right\},$$

where each r_j^{ex} is a vertex and each r_ℓ^{ray} spans an extreme ray. Because (18) cannot be unbounded, $b' r_\ell^{\text{ray}} \leq 0$ for each ℓ , so the optimal solution of (18) must have $\lambda_\ell = 0$ for all ℓ . Therefore, it follows from Theorem 2.8 of Bertsimas and Tsitsiklis (1997) that the optimal value for (18) is

$$\max_{1 \leq j \leq J} b' r_j^{\text{ex}}.$$

The conclusion follows because the class of vertices r_j^{ex} are determined by A_0 and c and do not depend on P . $U(P)$ can be obtained in a similar fashion.

B.2 When $Q\{R \in \mathcal{R}'\}$ is partially identified

Suppose now $Q\{R \in \mathcal{R}'\}$ is not point identified. In this case, $\Theta_0(P, \mathbf{Q})$ is not necessarily a closed interval, since the set $\{Q : Q \in \mathbf{Q}_0(P, \mathbf{Q}), Q\{R \in \mathcal{R}'\} > 0\}$ is not necessarily closed, but the latter remains a connected set. Recall (15) that

$$\theta(Q) = \frac{E_Q[g(R)\mathbb{1}\{R \in \mathcal{R}'\}]}{Q\{R \in \mathcal{R}'\}},$$

so $\theta(Q)$ is a continuous function of Q whenever $Q\{R \in \mathcal{R}'\} > 0$. As a result, the identified set $\Theta_0(P, \mathbf{Q}) \equiv \{\theta(Q) : Q \in \mathbf{Q}_0(P, \mathbf{Q}), Q\{R \in \mathcal{R}'\} > 0\}$ as the image of a continuous function on a connected set is also connected. Since $\Theta_0(P, \mathbf{Q}) \subseteq \mathbb{R}$, it is an interval (despite not necessarily closed). The lower endpoint of it, $L(P)$, can be obtained by solving the optimization problem

$$L(P) = \inf_{\substack{Q \in \mathbf{Q}_0(P, \mathbf{Q}) \\ Q\{R \in \mathcal{R}'\} > 0}} \frac{E_Q[g(R)\mathbb{1}\{R \in \mathcal{R}'\}]}{Q\{R \in \mathcal{R}'\}},$$

and similarly for the upper endpoint. The above problem can be written as the following two-step optimization problem which is easier to analyze:

$$\begin{aligned} \inf_{\substack{Q \in \mathbf{Q}_0(P, \mathbf{Q}) \\ Q\{R \in \mathcal{R}'\} > 0}} \frac{E_Q[g(R)\mathbb{1}\{R \in \mathcal{R}'\}]}{Q\{R \in \mathcal{R}'\}} &\iff \inf_{\pi \in \Pi(\mathcal{R}', P, \mathbf{Q}) \setminus \{0\}} \left\{ \inf_{Q \in \mathbf{Q}_0(P, \mathbf{Q}) \cap \Delta(\pi)} \frac{E_Q[g(R)\mathbb{1}\{R \in \mathcal{R}'\}]}{\pi} \right\} \\ &\iff \inf_{\pi \in \Pi(\mathcal{R}', P, \mathbf{Q}) \setminus \{0\}} \frac{1}{\pi} \left\{ \inf_{Q \in \mathbf{Q}_0(P, \mathbf{Q}) \cap \Delta(\pi)} E_Q[g(R)\mathbb{1}\{R \in \mathcal{R}'\}] \right\}, \end{aligned}$$

where for any $\pi \in [0, 1]$, we define

$$\Delta(\pi) \equiv \{Q \in \mathbf{Q} : \pi = Q\{R \in \mathcal{R}'\}\},$$

and

$$\Pi(\mathcal{R}', P, \mathbf{Q}) \equiv \{Q\{R \in \mathcal{R}'\} : Q \in \mathbf{Q}_0(P, \mathbf{Q})\}, \quad (19)$$

as the identified set for $Q\{R \in \mathcal{R}'\}$. Solving for $\Pi(\mathcal{R}', P, \mathbf{Q})$, however, is a special case of Section B.1 by specifying $\mathcal{R}' = \mathcal{R}$ and $g(R) = \mathbb{1}\{R \in \mathcal{R}'\}$ in $\theta(Q) = E_Q[g(R)|R \in \mathcal{R}']$, so that $Q\{R \in \mathcal{R}'\} = 1$ by Assumption 2.2. Therefore, $\Pi(\mathcal{R}', P, \mathbf{Q})$ is a closed interval whose lower (resp. upper) endpoint is the maximum (resp. minimum) of a finite number of linear functions of $(p_{ydz}(P))$.

For the inner minimization problem, $\mathbf{Q}_0(P, \mathbf{Q}) \cap \Delta(\pi)$ is characterized by

$$\mathbf{Q}_0(P, \mathbf{Q}) \cap \Delta(\pi) = \{q : A(\mathcal{R}')q = \beta(P, \pi), q \geq 0\},$$

where $A(\mathcal{R}')$ is a $(|\mathcal{M}| + 2) \times |\mathcal{N}|$ matrix, whose first $|\mathcal{M}| + 1$ rows are the same as A_0 , and the last row is given by $\mathbb{1}\{r \in \mathcal{R}'\}$, and $\beta(P, \pi) \equiv ((p_{ydz}(P) : (y, d, z) \in \mathcal{M}, 1, \pi)'$. Using the same analysis in Section B.1,

$$\inf_{Q \in \mathbf{Q}_0(P, \mathbf{Q}) \cap \Delta(\pi)} E_Q[g(R)\mathbb{1}\{R \in \mathcal{R}'\}]$$

is given by the maximum of a finite number of linear functions of $(p_{ydz}(P)) \cup (\pi)$, denoted as $L(P, \pi)$. As a result, $L(P)$ has the following form:

$$L(P) = \inf_{\pi \in \Pi(\mathcal{R}', P, \mathbf{Q}) \setminus \{0\}} \frac{1}{\pi} L(P, \pi), \quad (20)$$

and $U(P)$ can be obtained in a similar fashion.

C Details for Remark 4.2

The (multi-step) linear programming approach in Appendix B.2 is procedurally similar to the approach to bound the average causal effect within principal strata in Cheng and Small (2006) (hereafter CS). We provide a detailed comparison of our approach with that in CS. In particular, we show that the CS approach may not lead to a sharp bound both by formula and a numerical example, and provide a way to correct their approach as an application of our procedure.

Consider the RCT with one-sided noncompliance setting with $\mathcal{D} = \mathcal{Z} = \{0, 1, 2\}$ and $\mathcal{Y} = \{0, 1\}$, which is the same as the setup in CS with their ‘‘Monotonicity I’’ assumption. Define the strata

$$\mathcal{R}'_{0ij} = \{(y(0), y(1), y(2), d(0), d(1), d(2)) \in \mathcal{R} : (d(0), d(1), d(2)) = (0, i, j)\},$$

for $i \in \{0, 1\}$ and $j \in \{0, 2\}$, and $\pi_{0ij} = Q\{R \in \mathcal{R}'_{0ij}\}$. Then

$$\pi_{0ij} = \sum_{r \in \mathcal{R}} q(r) \mathbb{1}\{r \in \mathcal{R}'_{0ij}\}, \quad (21)$$

where we write $q(r) = Q\{R = r\}$. As in CS, we also write $p_{dz} = P\{D = d|Z = z\}$ for $d \in \mathcal{D}$ and $z \in \mathcal{Z}$. The working example in this section will be $\theta = E_Q[Y(1) - Y(0)|R \in \mathcal{R}'_{012}]$, but similar analysis extends to other parameters of interest.

CS step 1. The first step in CS is to minimize and maximize each π_{0ij} under the following constraints:

$$\begin{aligned}
p_{1|1} &= \pi_{012} + \pi_{010} \\
p_{0|1} &= \pi_{002} + \pi_{000} \\
p_{2|2} &= \pi_{002} + \pi_{012} \\
p_{0|2} &= \pi_{000} + \pi_{010} \\
1 &= \pi_{012} + \pi_{010} + \pi_{002} + \pi_{000} \\
0 &\leq \pi_{012}, \pi_{010}, \pi_{002}, \pi_{000} \leq 1.
\end{aligned} \tag{22}$$

Note that $p_{d|z} = p_{0d|z} + p_{1d|z}$, and write π_{0ij} in terms of $q(r)$ using (21). Then we see that $\{Q : Q \text{ satisfies (22)}\} \subseteq \{Q : Q \in \mathbf{Q}_0(P, \mathbf{Q})\}$. As a result, let $\tilde{\Pi}(\mathcal{R}'_{012}, P, \mathbf{Q})$ be the interval whose lower (resp. upper) endpoint is the solution of minimizing (resp. maximizing) π_{0ij} subject to Q satisfies (22), we have $\tilde{\Pi}(\mathcal{R}'_{012}, P, \mathbf{Q}) \supseteq \Pi(\mathcal{R}'_{012}, P, \mathbf{Q})$, for $\Pi(\mathcal{R}', P, \mathbf{Q})$ defined in (19). Indeed, $\tilde{\Pi}(\mathcal{R}'_{012}, P, \mathbf{Q})$ is given by CS as

$$\tilde{\Pi}(\mathcal{R}'_{012}, P, \mathbf{Q}) = [\max\{0, p_{1|1} - p_{0|2}\}, \min\{p_{1|1}, p_{2|2}\}], \tag{23}$$

while applying the procedure in Remark 4.1 on \mathcal{R}'_{012} , we get

$$\Pi(\mathcal{R}'_{012}, P, \mathbf{Q}) = \left[\max \left\{ \begin{array}{c} 0 \\ p_{01|1} + p_{11|1} - (p_{00|2} + p_{10|2}) \\ p_{10|0} - p_{10|1} - p_{10|2} \\ 1 - p_{00|1} - p_{00|2} - p_{10|0} \end{array} \right\}, \min \left\{ \begin{array}{c} p_{01|1} + p_{11|1} \\ p_{02|2} + p_{12|2} \\ 1 - p_{00|1} - p_{10|2} \\ 1 - p_{10|1} - p_{00|2} \end{array} \right\} \right], \tag{24}$$

from which it's obvious that $\tilde{\Pi}(\mathcal{R}'_{012}, P, \mathbf{Q}) \supseteq \Pi(\mathcal{R}'_{012}, P, \mathbf{Q})$. As we will show later, there exists a P such that $\mathbf{Q}_0(P, \mathbf{Q}) \neq \emptyset$ under which $\tilde{\Pi}(\mathcal{R}'_{012}, P, \mathbf{Q}) \not\supseteq \Pi(\mathcal{R}'_{012}, P, \mathbf{Q})$. One can verify that the same holds for other $\tilde{\Pi}(\mathcal{R}'_{0ij}, P, \mathbf{Q})$ in Section 3.1.1 of CS.

CS step 2. The second step in CS is to bound $E_Q[Y(d)|R \in \mathcal{R}'_{0ij}]$ for $d \in \mathcal{D}$ assuming π_{0ij} is known. Specifically, they minimize and maximize $E_Q[Y(d)|R \in \mathcal{R}'_{0ij}]$ subject to the following constraints:

$$\begin{aligned}
E_P[Y|D = 1, Z = 1] &= \frac{\pi_{012}}{\pi_{012} + \pi_{010}} E_Q[Y(1)|R \in \mathcal{R}'_{012}] + \frac{\pi_{010}}{\pi_{012} + \pi_{010}} E_Q[Y(1)|R \in \mathcal{R}'_{010}] \\
E_P[Y|D = 2, Z = 2] &= \frac{\pi_{012}}{\pi_{012} + \pi_{002}} E_Q[Y(2)|R \in \mathcal{R}'_{012}] + \frac{\pi_{002}}{\pi_{012} + \pi_{002}} E_Q[Y(2)|R \in \mathcal{R}'_{002}] \\
E_P[Y|D = 0, Z = 1] &= \frac{\pi_{002}}{\pi_{002} + \pi_{000}} E_Q[Y(0)|R \in \mathcal{R}'_{002}] + \frac{\pi_{000}}{\pi_{002} + \pi_{000}} E_Q[Y(0)|R \in \mathcal{R}'_{000}] \\
E_P[Y|D = 0, Z = 2] &= \frac{\pi_{010}}{\pi_{010} + \pi_{000}} E_Q[Y(0)|R \in \mathcal{R}'_{010}] + \frac{\pi_{000}}{\pi_{010} + \pi_{000}} E_Q[Y(0)|R \in \mathcal{R}'_{000}] \\
E_P[Y|D = 0, Z = 0] &= \pi_{012} E_Q[Y(0)|R \in \mathcal{R}'_{012}] + \pi_{010} E_Q[Y(0)|R \in \mathcal{R}'_{010}] \\
&\quad + \pi_{002} E_Q[Y(0)|R \in \mathcal{R}'_{002}] + \pi_{000} E_Q[Y(0)|R \in \mathcal{R}'_{000}] \\
0 &\leq E_Q[Y(d)|R \in \mathcal{R}'_{0ij}] \leq 1 \text{ for all } d \in \mathcal{D} \text{ and } i, j.
\end{aligned} \tag{25}$$

The constraints in (25) can be reduced to a form comparable to $Q \in \mathbf{Q}_0(P, \mathbf{Q}) \cap \Delta(\pi_{012})$. Note that

$$E_P[Y|D = d, Z = d] = P\{Y = 1|D = d, Z = z\} = \frac{p_{1d|z}}{p_{d|z}} = \frac{p_{1d|z}}{p_{0d|z} + p_{1d|z}},$$

and $\pi_{012} + \pi_{010} = p_{1|1} = p_{01|1} + p_{11|1}$. Also, write $E_Q[Y(1)|R \in \mathcal{R}'_{012}] = E_Q[Y(1)\mathbb{1}\{R \in \mathcal{R}'_{012}\}]/\pi_{012}$, the first equation in (25) can be written as

$$\begin{aligned} \frac{p_{11|1}}{p_{01|1} + p_{11|1}} &= \frac{E_Q[Y(1)\mathbb{1}\{R \in \mathcal{R}'_{012}\}]}{p_{01|1} + p_{11|1}} + \frac{E_Q[Y(1)\mathbb{1}\{R \in \mathcal{R}'_{010}\}]}{p_{01|1} + p_{11|1}} \\ \iff p_{11|1} &= \sum_{r \in \mathcal{R}} y(1)q(r)\mathbb{1}\{r \in \mathcal{R}'_{012} \cup \mathcal{R}'_{010}\}. \\ \iff p_{11|1} &= \sum_{r \in \mathcal{R}} \mathbb{1}\{y(1) = 1, d(1) = 1\}q(r). \end{aligned}$$

where we write $r = ((y(0), y(1), y(2), d(0), d(1), d(2)))$. Do the same to all equations in (25), use $\sum_{r \in \mathcal{R}} q(r) = 1$ and $\sum_{y,d} p_{yd|z} = 1$ for all $z \in \mathcal{Z}$, we see that $\{Q : Q \text{ satisfies (25)}\} = \{Q \in \mathbf{Q}_0(P, \mathbf{Q}) \cap \Delta(\pi_{012})\}$. Under the constraints in (25), CS derive the following bound for $E_Q[Y(1)|R \in \mathcal{R}'_{012}]$ given π_{012} :

$$\tilde{I}_{1|012}(P, \pi_{012}) = [\tilde{L}_{1|012}(P, \pi_{012}), \tilde{U}_{1|012}(P, \pi_{012})] = \left[\max \left\{ 0, 1 - \frac{1 - p_{1|11}}{\pi_{012}/p_{1|1}} \right\}, \min \left\{ 1, \frac{p_{1|11}}{\pi_{012}/p_{1|1}} \right\} \right], \quad (26)$$

where $p_{1|zd} = P\{Y = 1|D = d, Z = z\}$, and $\tilde{L}_{1|012}(P, \pi_{012})$ (resp. $\tilde{U}_{1|012}(P, \pi_{012})$) is the minimum (resp. maximum) of $E_Q[Y(1)|R \in \mathcal{R}'_{012}]$ subject to (25). Since $\{Q : Q \text{ satisfies (25)}\} = \{Q \in \mathbf{Q}_0(P, \mathbf{Q}) \cap \Delta(\pi_{012})\}$, $\tilde{I}_{1|012}(P, \pi_{012})$ is identical to our results from the inner optimization problem in Remark 4.1 after multiplying π_{012} . In fact, applying our results on $g(R) = Y(1)$ with $R' = \mathcal{R}'_{012}$, the two endpoints from the inner optimization problem are

$$I_{1|012}(P, \pi) = [L_{1|012}(P, \pi), U_{1|012}(P, \pi)] = \left[\max \{0, \pi - p_{01|1}\}, \min \{\pi, p_{11|1}\} \right],$$

and (26) can be written as

$$\begin{aligned} \tilde{I}_{1|012}(P, \pi_{012}) &= \left[\max \left\{ 0, 1 - \frac{1 - p_{1|11}}{\pi_{012}/p_{1|1}} \right\}, \min \left\{ 1, \frac{p_{1|11}}{\pi_{012}/p_{1|1}} \right\} \right] \\ &= \left[\max \left\{ 0, 1 - \frac{1 - p_{11|1}/(p_{01|1} + p_{11|1})}{\pi_{012}/(p_{01|1} + p_{11|1})} \right\}, \min \left\{ 1, \frac{p_{11|1}/(p_{01|1} + p_{11|1})}{\pi_{012}/(p_{01|1} + p_{11|1})} \right\} \right] \\ &= \left[\max \left\{ 0, 1 - \frac{p_{01|1}}{\pi_{012}} \right\}, \min \left\{ 1, \frac{p_{11|1}}{\pi_{012}} \right\} \right] = \pi_{012} I_{1|012}(P, \pi_{012}). \end{aligned}$$

The same result holds for other $E_Q[Y(d)|R \in \mathcal{R}'_{0ij}]$.

CS step 3. The final step in CS is to combine the bounds in (23) and (26) to bound $E_Q[Y(d_1) - Y(d_2)|R \in \mathcal{R}'_{0ij}]$. Specifically, the bound for $E_Q[Y(1) - Y(0)|R \in \mathcal{R}'_{012}]$ is given by $[\tilde{L}(P), \tilde{U}(P)]$ where

$$\tilde{L}(P) = \inf_{\pi_{012} \in \tilde{\Pi}(\mathcal{R}'_{012}, P, \mathbf{Q}) \setminus \{0\}} \left\{ \tilde{L}_{1|012}(P, \pi_{012}) - \tilde{U}_{0|012}(P, \pi_{012}) \right\},$$

$$\tilde{U}(P) = \sup_{\pi_{012} \in \tilde{\Pi}(\mathcal{R}'_{012}, P, \mathbf{Q}) \setminus \{0\}} \left\{ \tilde{U}_{1|012}(P, \pi_{012}) - \tilde{L}_{0|012}(P, \pi_{012}) \right\}.$$

This step uses the crude bound that given π_{012} ,

$$\tilde{L}_{1|012}(P, \pi_{012}) - \tilde{U}_{0|012}(P, \pi_{012}) \leq E_Q[Y(1)|R \in \mathcal{R}'_{012}] - E_Q[Y(0)|R \in \mathcal{R}'_{012}] \leq \tilde{U}_{1|012}(P, \pi_{012}) - \tilde{L}_{0|012}(P, \pi_{012}).$$

It turns out that by specifying $g(R) = Y(1) - Y(0)$ with $R' = \mathcal{R}'_{012}$ in our approach, the two solutions from the inner optimization is indeed $L_{10|012}(P, \pi_{012}) = \pi_{012}(\tilde{L}_{1|012}(P, \pi_{012}) - \tilde{U}_{0|012}(P, \pi_{012}))$ and $U_{10|012}(P, \pi_{012}) = \pi_{012}(\tilde{U}_{1|012}(P, \pi_{012}) - \tilde{L}_{0|012}(P, \pi_{012}))$. It is interesting to study when one gets such a “separable” bound that it suffices to consider bounds for each $E_Q[Y(j)|R \in \mathcal{R}']$ in order to study $E_Q[Y(j_1) - Y(j_2)|R \in \mathcal{R}']$, but it is beyond the scope of the current paper. The final step involves a one-dimensional grid search over $\pi_{012} \in \tilde{\Pi}(\mathcal{R}'_{012}, P, \mathbf{Q})$, which is the same as the outer optimization in (20).

Summary. Putting the three steps together, we see that CS derive the tight bound of $E_Q[Y_{j_1} - Y_{j_2}|R \in \mathcal{R}'_{0ij}]$ given π_{0ij} in their step 2, even if they use a crude bound for the difference of $E_Q[Y_{j_1}|R \in \mathcal{R}'_{0ij}]$ and $E_Q[Y_{j_2}|R \in \mathcal{R}'_{0ij}]$. Nevertheless, their bound for the conditional probability $\pi_{0ij} = Q\{R \in \mathcal{R}'_{0ij}\}$, $\tilde{\Pi}(\mathcal{R}'_{0ij}, P, \mathbf{Q})$, is loose compared to $\Pi(\mathcal{R}'_{0ij}, P, \mathbf{Q})$. As a result, their bounds for $E_Q[Y_{j_1} - Y_{j_2}|R \in \mathcal{R}'_{0ij}]$ are not sharp in general. Moreover, if in their approach we use $\Pi(\mathcal{R}'_{0ij}, P, \mathbf{Q})$ as the bound for π_{0ij} in place of $\tilde{\Pi}(\mathcal{R}'_{0ij}, P, \mathbf{Q})$, we would reach the same bound for the conditional ATE. We illustrate these two points using a numerical example below.

	$p_{00 0}$ 0.764	$p_{10 0}$ 0.236	
$p_{00 1}$ 0.412	$p_{10 1}$ 0.107	$p_{01 1}$ 0.301	$p_{11 1}$ 0.180
$p_{00 2}$ 0.117	$p_{10 2}$ 0.169	$p_{02 2}$ 0.475	$p_{12 2}$ 0.239

Table 1: Distribution P for Appendix C.

$q(000, 000)$ 0.002	$q(000, 010)$ 0.002	$q(000, 002)$ 0.017	$q(000, 012)$ 0.025
$q(001, 000)$ 0.002	$q(001, 010)$ 0.002	$q(001, 002)$ 0.002	$q(001, 012)$ 0.195
$q(010, 000)$ 0.101	$q(010, 010)$ 0.002	$q(010, 002)$ 0.272	$q(010, 012)$ 0.120
$q(011, 000)$ 0.002	$q(011, 010)$ 0.004	$q(011, 002)$ 0.014	$q(011, 012)$ 0.002
$q(100, 000)$ 0.002	$q(100, 010)$ 0.002	$q(100, 002)$ 0.022	$q(100, 012)$ 0.011
$q(101, 000)$ 0.034	$q(101, 010)$ 0.062	$q(101, 002)$ 0.015	$q(101, 012)$ 0.002
$q(110, 000)$ 0.002	$q(110, 010)$ 0.002	$q(110, 002)$ 0.006	$q(110, 012)$ 0.002
$q(111, 000)$ 0.024	$q(111, 010)$ 0.041	$q(111, 002)$ 0.002	$q(111, 012)$ 0.007

Table 2: Distribution Q for Appendix C.

A numerical example. Consider the P distribution in Table 1 and the Q distribution in Table 2, where we write $q(y_0y_1y_2, d_0d_1d_2) = Q\{Y(d) = y_d, D(z) = d_z, (d, z) \in \mathcal{D} \times \mathcal{Z}\}$ and omit any $q(\cdot) = 0$. One can check that $Q \in \mathbf{Q}$ and $P = QT^{-1}$ so that $\mathbf{Q}_0(P, \mathbf{Q}) \neq \emptyset$. Moreover, $Q\{R \in \mathcal{R}\} = 1$ for \mathcal{R} defined in Example 3.1. We are interested in $E_Q[Y(1) - Y(0)|R \in \mathcal{R}'_{012}]$. Firstly $Q\{R \in \mathcal{R}'_{012}\} > 0$. Under this distribution, we can calculate that $\tilde{\Pi}(\mathcal{R}'_{012}, P, \mathbf{Q}) = [0.195, 0.481]$, while $\Pi(\mathcal{R}'_{012}, P, \mathbf{Q}) = [0.235, 0.419]$, as shown in the black and blue lines in Figure 1. The green and red lines are the upper and lower bound of $E_Q[Y(1) - Y(0)|R \in \mathcal{R}'_{012}]$ as a function of π_{012} , i.e., $\tilde{U}_{1|012}(P, \pi_{012}) - \tilde{L}_{0|012}(P, \pi_{012})$ and $\tilde{L}_{1|012}(P, \pi_{012}) - \tilde{U}_{0|012}(P, \pi_{012})$, or equivalently $U_{10|012}(P, \pi_{012})/\pi_{012}$ and $L_{10|012}(P, \pi_{012})/\pi_{012}$. As a result, if we use $\tilde{\Pi}(\mathcal{R}'_{012}, P, \mathbf{Q})$ to do the outer optimization, we get a bound $[-0.219, 0.923]$, which is the minimum of the red line and maximum of the green line *within the two black boundaries*. Using $\Pi(\mathcal{R}'_{012}, P, \mathbf{Q})$ for the outer optimization, however, results in a bound $[-0.219, 0.766]$ as the minimum of the red line and the maximum of the green line *within the two blue boundaries*. It is also interesting to note that the bound for $E[Y(1) - Y(0)|R \in \mathcal{R}'_{012}]$ given π_{012} is not always valid when $\pi_{012} \in \tilde{\Pi}(\mathcal{R}'_{012}, P, \mathbf{Q})$. The green line actually goes below the red line to the right of the plot, though it is valid when $\pi_{012} \in \Pi(\mathcal{R}'_{012}, P, \mathbf{Q})$.

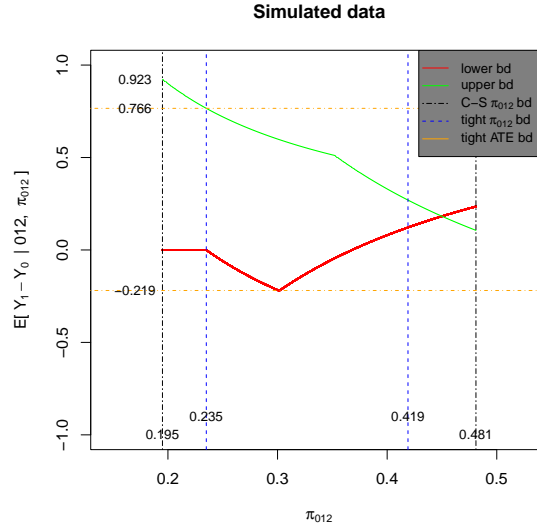


Figure 1: Visualization of $E_Q[Y(1) - Y(0)|R \in \mathcal{R}'_{012}]$ given π_{012} under the simulated distribution in Table 1.

D Details for Remark 4.5

In this section, we discuss a method to obtain sharp testable restrictions of \mathbf{Q} in terms of analytical inequalities. In what follows, we will make heavy use of the fact that a nonempty bounded polyhedron can be represented in two ways. Recall from Definition 2.1 of Bertsimas and Tsitsiklis (1997) that a polyhedron in \mathbf{R}^k is a set $\{x \in \mathbf{R}^k : Ax \leq b\}$, known as the H -representation of a polyhedron. If the polyhedron is nonempty and bounded, then Theorem 2.9 of Bertsimas and Tsitsiklis (1997) implies that it can equivalently be represented as the convex hull of its (finite number of) vertices, known as the V -representation of a polyhedron.

As in Appendix B, we identify Q with the column vector $q = (q(r) : r \in \mathcal{R})$, and P with $p = \{p_{yd|z} : (y, d, z) \in \mathcal{M}\}$. Further let A_1 denote the first $|\mathcal{M}|$ rows of A_0 and let a_0 denote the last row of A_0 . Correspondingly, we note

$$\mathbf{Q} = \{q : a_0'q = 1, q \geq 0\},$$

which is a bounded polyhedron in H -representation. Next, let $\mathbf{P}(\mathbf{Q}) = \{P : P = QT^{-1}, Q \in \mathbf{Q}\}$, and note $\mathbf{P}(\mathbf{Q}) = A_1\mathbf{Q}$. $\mathbf{P}(\mathbf{Q})$ is obviously a bounded polyhedron and is nonempty as long as \mathbf{Q} is nonempty. In that case, Theorem 2.9 of [Bertsimas and Tsitsiklis \(1997\)](#) implies it is the convex hull of its (finite number of) vertices.

The previous discussion leads to the following algorithm for obtaining $\mathbf{P}(\mathbf{Q})$ in terms of inequalities, i.e., its H -representation. For a given polyhedron, we can compute one representation from the other using the `mpt3` package in MATLAB. To obtain the H -representation of $\mathbf{P}(\mathbf{Q})$, we use the following algorithm:

Algorithm D.1.

Step 1: Collect the set of all vertices of \mathbf{Q} , denoted by $V = \{V_i : 1 \leq i \leq n\}$.

Step 2: Compute $A_1V = \{A_1V_i : 1 \leq i \leq n\}$.

Step 3: From Lemma D.1, define $A_1\mathbf{Q} = \text{co}(A_1V)$.

Step 4: Obtain the H representation of $A_1\mathbf{Q}$.

In the algorithm, we have used the following lemma that allows us to define $\mathbf{P}(\mathbf{Q})$ through the vertices of \mathbf{Q} :

Lemma D.1. *Suppose \mathbf{Q} is a nonempty bounded polyhedron and V is the set of vertices of \mathbf{Q} . Then, $A_1\mathbf{Q} = \text{co}(A_1V)$.*

PROOF. We first show $\text{co}(A_1V) \subseteq A_1\mathbf{Q}$. Indeed, each $P \in \text{co}(A_1V)$ could be written as $\sum_{1 \leq i \leq n} \lambda_i A_1V_i = A_1 \sum_{1 \leq i \leq n} \lambda_i V_i$ where $\lambda_i \geq 0$ for all i and $\sum_i \lambda_i = 1$, but $\sum_{1 \leq i \leq n} \lambda_i V_i \in \mathbf{Q}$ since \mathbf{Q} is a polyhedron and hence convex. To show $A_1\mathbf{Q} \subseteq \text{co}(A_1V)$, fix $Q \in \mathbf{Q}$. Because $\mathbf{Q} = \text{co}(V)$, $Q = \sum_{1 \leq i \leq n} \lambda_i V_i$, where $\lambda_i \geq 0$ for all i and $\sum_i \lambda_i = 1$, so $A_1Q = \sum_{1 \leq i \leq n} \lambda_i A_1V_i \in \text{co}(A_1V)$. ■

References

- ANGRIST, J., LANG, D. and OREOPOULOS, P. (2009). Incentives and services for college achievement: Evidence from a randomized trial. *American Economic Journal: Applied Economics*, **1** 136–63.
- ANGRIST, J. D. and IMBENS, G. W. (1995). Two-stage least squares estimation of average causal effects in models with variable treatment intensity. *Journal of the American statistical Association*, **90** 431–442.
- BAI, Y., HUANG, S., MOON, S., SHAIKH, A. and VYTLACIL, E. J. (2024). On the identifying power of generalized monotonicity for average treatment effects. Tech. rep., National Bureau of Economic Research.
- BAI, Y., SANTOS, A. and SHAIKH, A. M. (2022). On testing systems of linear inequalities with known coefficients. Tech. rep.
- BALKE, A. and PEARL, J. (1997). Bounds on treatment effects from studies with imperfect compliance. *Journal of the American Statistical Association*, **92** 1171–1176.
- BERTSIMAS, D. and TSITSIKLIS, J. N. (1997). *Introduction to linear optimization*, vol. 6.
- BHATTACHARYA, J., SHAIKH, A. M. and VYTLACIL, E. (2008). Treatment Effect Bounds under Monotonicity Assumptions: An Application to Swan-Ganz Catheterization. *American Economic Review*, **98** 351–356. URL <https://www.aeaweb.org/articles?id=10.1257/aer.98.2.351>.
- BHATTACHARYA, J., SHAIKH, A. M. and VYTLACIL, E. (2012). Treatment effect bounds: An application to Swan–Ganz catheterization. *Journal of Econometrics*, **168** 223–243. URL <https://www.sciencedirect.com/science/article/pii/S0304407612000024>.
- CHENG, J. and SMALL, D. S. (2006). Bounds on causal effects in three-arm trials with non-compliance. *Journal of the Royal Statistical Society: Series B (Statistical Methodology)*, **68** 815–836.
- FANG, Z. and SANTOS, A. (2019). Inference on directionally differentiable functions. *The Review of Economic Studies*, **86** 377–412.
- FANG, Z., SANTOS, A., SHAIKH, A. M. and TORGOVITSKY, A. (2023). Inference for large-scale linear systems with known coefficients. *Econometrica*, **91** 299–327.
- FRANGAKIS, C. E. and RUBIN, D. B. (2002). Principal stratification in causal inference. *Biometrics*, **58** 21–29.
- GABRIEL, E. E., SACHS, M. C. and JENSEN, A. K. (2024). Sharp symbolic nonparametric bounds for measures of benefit in observational and imperfect randomized studies with ordinal outcomes. *Biometrika* asae020.
- HECKMAN, J. J. and PINTO, R. (2018). Unordered monotonicity. *Econometrica*, **86** 1–35.
- HECKMAN, J. J., SMITH, J. and CLEMENTS, N. (1997). Making the most out of programme evaluations and social experiments: Accounting for heterogeneity in programme impacts. *The Review of Economic Studies*, **64** 487–535.

- HECKMAN, J. J. and SMITH, J. A. (1998). Evaluating the welfare state. *ECONOMETRIC SOCIETY MONOGRAPHS*, **16** 241–318.
- HUANG, E. J., FANG, E. X., HANLEY, D. F. and ROSENBLUM, M. (2019). Constructing a confidence interval for the fraction who benefit from treatment, using randomized trial data. *Biometrics*, **75** 1228–1239.
- IMBENS, G. W. and ANGRIST, J. D. (1994). Identification and estimation of local average treatment effects. *Econometrica*, **62** 467–475.
- KIRKEBOEN, L. J., LEUVEN, E. and MOGSTAD, M. (2016). Field of Study, Earnings, and Self-Selection. *The Quarterly Journal of Economics*, **131** 1057–1111. URL <https://doi.org/10.1093/qje/qjw019>.
- KLINE, P. and WALTERS, C. R. (2016). Evaluating public programs with close substitutes: The case of Head Start. *The Quarterly Journal of Economics*, **131** 1795–1848.
- LU, J., DING, P. and DASGUPTA, T. (2018). Treatment effects on ordinal outcomes: Causal estimands and sharp bounds. *Journal of Educational and Behavioral Statistics*, **43** 540–567.
- MACHADO, C., SHAIKH, A. M. and VYTLACIL, E. J. (2019). Instrumental variables and the sign of the average treatment effect. *Journal of Econometrics*.
- MANSKI, C. F. (1997). Monotone treatment response. *Econometrica: Journal of the Econometric Society* 1311–1334.
- MANSKI, C. F. and PEPPER, J. V. (1998). Monotone instrumental variables with an application to the returns to schooling.
- MASTEN, M. A. and POIRIER, A. (2020). Inference on breakdown frontiers. *Quantitative Economics*, **11** 41–111.
- NOACK, C. (2021). Sensitivity of late estimates to violations of the monotonicity assumption. *arXiv preprint arXiv:2106.06421*.
- VYTLACIL, E. (2006). Ordered discrete-choice selection models and local average treatment effect assumptions: Equivalence, nonequivalence, and representation results. *The Review of Economics and Statistics*, **88** 578–581.