

# An Automated Approach to Causal Inference in Discrete Settings\*

Guilherme Duarte      Noam Finkelstein      Dean Knox  
gjduarte@upenn.edu      noam@jhu.edu      dcknox@upenn.edu

Jonathan Mummolo      Ilya Shpitser  
jmummolo@princeton.edu      ilyas@cs.jhu.edu

First draft: February 10, 2021  
This draft: September 29, 2021

## Abstract

When causal quantities cannot be point identified, researchers often pursue partial identification to quantify the range of possible values. However, the peculiarities of applied research conditions can make this analytically intractable. We present a general and automated approach to causal inference in discrete settings. We show causal questions with discrete data reduce to polynomial programming problems, and we present an algorithm to automatically bound causal effects using efficient dual relaxation and spatial branch-and-bound techniques. The user declares an estimand, states assumptions, and provides data (however incomplete or mismeasured). The algorithm then searches over admissible data-generating processes and outputs the most precise possible range consistent with available information—i.e., *sharp* bounds—including a point-identified solution if one exists. Because this search can be computationally intensive, our procedure reports and continually refines non-sharp ranges that are guaranteed to contain the truth at all times, even when the algorithm is not run to completion. Moreover, it offers an additional guarantee we refer to as  $\varepsilon$ -sharpness, characterizing the worst-case looseness of the incomplete bounds. Analytically validated simulations show the algorithm accommodates classic obstacles, including confounding, selection, measurement error, noncompliance, and nonresponse.

*Keywords:* causal inference, partial identification, constrained optimization, linear programming, polynomial programming

---

\*Guilherme Jardim Duarte is a Ph.D. student in the Operations, Information and Decisions Department, the Wharton School of the University of Pennsylvania. Noam Finkelstein is a Ph.D. student in the Department of Computer Science, Johns Hopkins University. Dean Knox is an Andrew Carnegie Fellow and an assistant professor in the Operations, Information and Decisions Department, the Wharton School of the University of Pennsylvania. Jonathan Mummolo is an assistant professor of Politics and Public Affairs, Princeton University. Ilya Shpitser is the John C. Malone Assistant Professor in the Department of Computer Science, Whiting School of Engineering at the Johns Hopkins University. Authors listed in alphabetical order. For helpful feedback, we thank Peter Aronow, Justin Grimmer, Kosuke Imai, Luke Keele, Gary King, Christopher Lucas, Fredrik Sävje, Brandon Stewart, Eric Tchetgen Tchetgen, and participants in the Harvard Applied Statistics Workshop, the New York University Data Science Seminar, University of Pennsylvania Causal Inference Seminar, PolMeth 2021, and the Yale Quantitative Research Methods Workshop. We gratefully acknowledge financial support from AI for Business and the Analytics at Wharton Data Science and Business Analytics Fund. This research was made possible in part by a grant from the Carnegie Corporation of New York. The statements made and views expressed are solely the responsibility of the authors.

# Contents

<b>1</b>	<b>Introduction</b>	<b>1</b>
<b>2</b>	<b>Related Literature</b>	<b>3</b>
<b>3</b>	<b>Preliminaries</b>	<b>4</b>
3.1	Canonical DAGs . . . . .	5
3.2	Potential Outcomes . . . . .	6
3.3	Principal Stratification . . . . .	6
<b>4</b>	<b>Formulating the Polynomial Program</b>	<b>9</b>
<b>5</b>	<b>Simplifying the Polynomial Program</b>	<b>12</b>
5.1	Eliminating Blocked Disturbances . . . . .	12
5.2	Exploiting the Nested Markov Parameterization . . . . .	12
5.3	Eliminating Additional Constraints and Parameters . . . . .	14
<b>6</b>	<b>Computing <math>\varepsilon</math>-sharp Bounds in Polynomial Programs</b>	<b>15</b>
<b>7</b>	<b>Statistical Inference</b>	<b>16</b>
<b>8</b>	<b>Simulated Examples</b>	<b>18</b>
8.1	Instrumental Variables . . . . .	18
8.2	Coverage of Confidence Bounds . . . . .	20
8.3	More Complex Bounding Problems . . . . .	21
<b>9</b>	<b>Potential Critiques of the Approach</b>	<b>24</b>
<b>10</b>	<b>Future Work with Automated Bounding</b>	<b>25</b>
<b>A</b>	<b>Examples, Algorithms, and Detailed Discussion</b>	<b>30</b>
A.1	Canonicalization of DAGs . . . . .	30
A.2	Functional Models in the Context of Determinism . . . . .	30
A.3	DAG Parameterization for Non-gearred Graphs . . . . .	31
A.4	Example of Program Simplification . . . . .	32
A.5	Constructing the Polynomial Program . . . . .	35
A.6	Optimizing the Polynomial Program . . . . .	36
<b>B</b>	<b>Proofs</b>	<b>39</b>
<b>C</b>	<b>Uncertainty</b>	<b>41</b>
<b>D</b>	<b>Details of Simulated Models</b>	<b>44</b>
D.1	Noncompliance Simulation . . . . .	47
D.2	Outcome-Based Selection Simulation . . . . .	47
D.3	Measurement Error Simulation . . . . .	48
D.4	Outcome Missingness Simulation . . . . .	49
D.5	Joint Missingness Simulation . . . . .	50

The combination of some data and an aching desire for an answer does not ensure that a reasonable answer can be extracted from a given body of data.

---

Tukey (1986, pp. 74–75)

# 1 Introduction

When causal quantities cannot be point identified, researchers often pursue partial identification to quantify the range of possible answers. These solutions are tailored to specific scenarios (e.g. Lee, 2009; Gabriel et al., 2020; Kennedy et al., 2019; Knox et al., 2020; Li and Pearl, 2021; Sjölander et al., 2014), but the idiosyncrasies of applied research can render prior results unusable if identifying assumptions fail or slightly differing causal structures are encountered. This case-by-case approach to deriving causal bounds presents a major obstacle to scientific progress. To increase the pace of discovery, researchers need a general approach that is robust to context-specific peculiarities.

In this paper, we present an automated approach to causal inference in discrete settings which can be applied to all graphical causal models, as well as all observed quantities and domain assumptions in standard use. With our algorithm, users declare an estimand, state assumptions, and provide available data—however incomplete or mismeasured. The algorithm then outputs *sharp bounds*, the most precise possible answer to the causal query given these inputs, including a point estimate if the solution is identified. This approach can accommodate scenarios involving any classic threat to inference, including but not limited to missing data, selection, measurement error, and noncompliance. Our algorithm also has the desirable property of alerting users when assumptions conflict with observed data, indicating a faulty causal theory. Finally, we develop techniques for drawing statistical inferences about estimated bounds. We demonstrate our method using a host of simulations, validating results wherever existing analytic solutions are available.

Our work advances a rich literature on partial identification in causal inference (Robins, 1989; Manski, 1990; Heckman and Vytlačil, 2001; Zhang and Rubin, 2003; Cai et al., 2008; Swanson et al., 2018; Gabriel et al., 2020; Molinari, 2020), outlined in Section 2, which has sometimes cast the task as a constrained optimization problem that can be solved computationally. In pioneering work, Balke and Pearl (1994, 1997) provided a method for calculating sharp bounds when causal queries can be expressed as linear programming problems. However, a wide range of estimands and empirical obstacles result in causal queries that are *not* reducible to linear programs, and a complete computational solution has remained elusive.

When feasible, sharp-bounding approaches offer a principled and transparent approach to causal inference that makes maximum use of available information while acknowledging its limitations. Claims outside the bounds can be immediately rejected, and claims inside the bounds must be explicitly justified by additional assumptions or data that enable tightening. But several obstacles still preclude widespread use of these techniques.

For one, analytic derivation remains intractable for many problems. Within the subclass of linear problems, Balke and Pearl’s (1994) simplex method offers a highly efficient analytic solution, but one that fails to generalize to the many partially observed settings where nonlinearity arises. Analytic nonlinear solutions remain limited to specific results, painstakingly derived case by case (e.g. Knox et al., 2020; Gabriel et al., 2020; Li and Pearl, 2021). Though general sharp bounds can in theory be obtained by various nonlinear optimization techniques (Geiger and Meek, 1999; Zhang and Bareinboim, 2021), such approaches are often computationally infeasible. This is because without exhaustively exploring a vast model space, analysts can obtain local optima that correspond to potentially invalid bounds—i.e., ranges that may fail to contain the truth.

To address these limitations, we first show in Sections 3 and 4 that essentially all common causal queries involving discrete variables can be reduced to polynomial programs—a well-studied class of optimization tasks that nest linear programming as a special case—building on prior results from Geiger and Meek (1999) and Wolfe et al. (2019).<sup>1</sup> While mature techniques have been developed for such tasks (Belotti et al., 2009; Vigerske and Gleixner, 2018; Gamrath et al., 2020), it is well known that solving polynomial programs to global optimality is in general NP-hard. The difficulty of the problem thus highlights the need for efficient algorithms and bounding techniques that remain valid even when analysts are faced with time constraints. In Sections 5–6, we develop a procedure, based on dual relaxation and spatial branch-and-bound relaxation techniques, that provides valid bounds of arbitrary sharpness, for all causal structures, under virtually any information environments and domain assumptions. We show this procedure is guaranteed to achieve complete sharpness with sufficient computation time; in smaller problems, this can occur in a matter of seconds. However, in cases where the time needed to discover sharp bounds is prohibitive—which can occur even in moderately sized problems with severe information fragmentation—our algorithm is *anytime* (Dean and Boddy, 1988), meaning it can be interrupted to obtain non-sharp bounds that are nonetheless guaranteed to be valid. Our technique also offers an additional guarantee we term “ $\varepsilon$ -sharpness,” indicating the worst-case looseness factor of the relaxed bounds relative to the unknown, completely sharp bounds. In Section 7, we provide two approaches for characterizing uncertainty in the estimated bounds, and we demonstrate our technique in a series of simulations in Section 8. Our simulations, validated against previously derived analytically results where possible, show the flexibility of our approach and the ease with which assumptions can be modularly imposed or relaxed. Moreover, we demonstrate how the algorithm can uncover counterintuitive results: in one case, we show a scenario that appears to be partially identified is in fact point identified, improving over widely used bounds (Manski, 1990) and recovering a recent advance in the literature on nonrandom missingness (Miao et al., 2015).

---

<sup>1</sup>Specifically, our results apply to elementary arithmetic functionals or monotonic transformations thereof—a broad set that essentially includes all causal assumptions, observed quantities, and estimands in standard use. For example, the average treatment effect and the log odds ratio can be sharply bounded with our approach, but non-analytic functionals (which are rarely if ever encountered) cannot. Functionals that do not meet these conditions can be approximated to arbitrary precision, if they have convergent power series.

Our approach offers a complete and computationally feasible approach to causal inference in discrete settings. Given a well-defined causal query, valid assumptions, and data, researchers now have a general and automated process to draw causal inferences that are guaranteed to be valid and, with sufficient computation time, provably optimal. As we discuss in Sections 9–10, our approach’s modular nature also allows analysts to conduct principled robustness tests and sensitivity analyses that can identify the most promising avenues for future research, promote research transparency, and accelerate scientific discovery.

## 2 Related Literature

Researchers have long sought to automate causal identification by recasting causal queries as constrained optimization problems that can be solved computationally. Our work is most closely related to [Balke and Pearl \(1994, 1997\)](#), which showed that certain bounding problems in discrete settings—generally corresponding to causal systems in which outcomes and manipulated variables are fully observed—could be formulated as the minimization and maximization of a linear objective function subject to linear equality and inequality constraints. In these cases, causal bounding problems can be reformulated as linear programming problems, which admit both symbolic solutions and highly efficient numerical solutions. Subsequent studies have proven that in particular settings, the bounds produced by this technique are sharp ([Ramsahai, 2012](#); [Bonet, 2001](#); [Heckman and Vytlačil, 2001](#)), and [Sachs et al. \(2020\)](#) shows this approach produces sharp bounds for any such linear problem. These results were extended by [Geiger and Meek \(1999\)](#), which showed that a much broader class of discrete problems can be formulated in terms of polynomial relations—at least, when analysts have precise information about the kinds of disturbances or confounders that may exist, expressed in terms of latent variable cardinalities. These discrete problems include not only the bounds studied in this paper, but the related problem of determining what constraints on the main variables are implied by a causal graph. In addition to the well-known conditional independence constraints implied by d-separation, these can include generalized equality constraints (or Verma constraints; [Verma and Pearl, 1990](#); [Tian and Pearl, 2002](#)). Beyond these equalities, the main variables are also constrained by generalizations of the instrumental inequalities ([Pearl, 1995](#); [Bonet, 2001](#)).

[Geiger and Meek \(1999\)](#) note that in theory, algorithms for quantifier elimination can provide symbolic solutions for these questions. However, the time complexity of quantifier elimination is doubly exponential, rendering it infeasible for all but the simplest cases. At the core of this issue is that symbolic methods provide a general solution, meaning that they must explore the space of all possible inputs. In contrast, numerical methods such as our approach can often eliminate portions of the space that are irrelevant, accelerating computation.

Even so, computation can be time-consuming; polynomial programming is in general NP-hard. In practice, many optimizers are able to rapidly find reasonably good values but cannot guarantee optimality without exhaustively searching the space of candidates.

This approach poses a challenge for obtaining causal bounds, which represent minimal and maximal values of the estimand under all models that are *admissible*, or consistent with observed data and modeling assumptions. If a local optimizer operates on the original problem (the *primal*), proceeding from the interior and widening bounds as more extreme models are discovered, then failing to reach global optimality will result in *invalid bounds*—ranges narrower than the optimal sharp bounds which do not contain all possible solutions.

In this paper, we detail an approach that resolves this obstacle by allowing analysts to obtain valid bounds in limited time. At a high level, our approach is to reexpress causal inference problems in terms of principal strata (Frangakis and Rubin, 2002). To do so, we first present new results on lossless reductions for latent variables of unknown cardinality. We then show that causal estimands, modeling assumptions, and observed information can all be expressed in terms of polynomial expressions, equalities, and inequalities with no loss of information. We show how these systems can be simplified for computational efficiency, then develop an iterative primal-dual algorithm that searches for admissible models from the interior of the bounds (the primal problem) while simultaneously refining a guaranteed-valid outer envelope for the sharp bounds (the dual problem). Even when exhaustive search is computationally infeasible, suboptimal primal and dual values can still be found and improved over time. We show suboptimal dual points allow analysts to report valid *loose bounds*—those that are wider than the unknown sharp bounds. Our method also utilizes the suboptimal primal points, allowing analysts to assess the worst-case *looseness factor* of the reported valid bounds, relative to the unknown sharp bounds.

### 3 Preliminaries

In this section, we define notation and discuss concepts necessary to derive our key results. We first review how arbitrary directed acyclic graphs (DAGs) can be “canonicalized” without loss of information, resulting in an equivalent form with properties amenable to analysis (Evans, 2018). We then describe how graphs in this form give rise to potential outcomes and principal strata (Frangakis and Rubin, 2002), two key building blocks in our analytic strategy.

Suppose that for each i.i.d. unit  $i \in \{1, \dots, N\}$ , the main variables of interest are contained in  $\mathbf{V}_i = \{V_{i,1}, \dots, V_{i,J}\}$ , indexed by  $j$ . We will suppose that the sample space of each main variable,  $\mathcal{S}(V_{i,j})$ , has finite cardinality. These variables may be either observed or unobserved—as we will show, it is often useful to reason about unobserved elements of  $\mathbf{V}_i$  in the context of missing data and measurement error. We will also consider unobserved causal ancestors of  $\mathbf{V}_i$ , collectively denoted  $\mathbf{U}_i = \{U_{i,1}, \dots, U_{i,K}\}$  and indexed by  $k$ , that represent random disturbances or confounders. Without loss of generality, these disturbances—which have unknown, possibly infinite cardinality—are assumed to subsume all phenomena that are causally relevant to  $\mathbf{V}_i$ .<sup>2</sup> As we show in Section 3.3, this assumption is without

---

<sup>2</sup>We note that traditionally, variables in  $\mathbf{V}_i$  are permitted to be affected by exogenous causal noise not represented in the graph. By incorporating all causally relevant factors into  $\mathbf{U}_i$ , we take each variable in  $\mathbf{V}_i$  to be a deterministic function of its parents in the graph, discussed in more detail below.

consequence, because even a continuous and infinite dimensional  $\mathbf{U}_i$  must still map down to the same finite canonical partitions that we describe there. In addition, we will make use of *counterfactual* random variables, which represent hypothetical versions of random variables in  $\mathbf{V}_i$  that would have occurred had, contrary to fact, treatment variables been exogenously set to a specified value. (A more rigorous definition is given in Section 3.2.) By convention, bold letters denote collections of variables; uppercase and lowercase letters respectively denote random variables and their realizations. We will consider population distributions until discussing inference in Section 7.

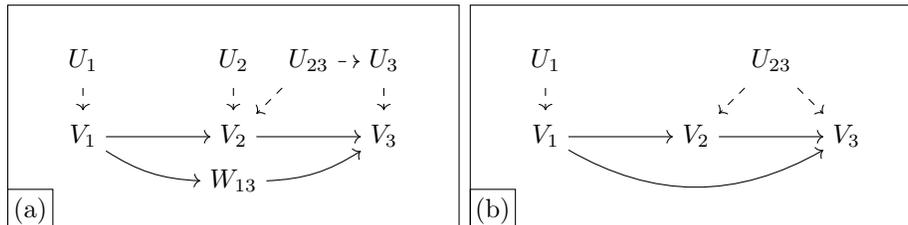
### 3.1 Canonical DAGs

We now discuss how DAGs can be canonicalized, or distilled to minimal form, to clarify which aspects of the structural model can be ignored, greatly simplifying the bounding task. Suppose that causal relationships between all variables in  $\mathbf{V}_i$  and  $\mathbf{U}_i$  are represented by a directed acyclic graph (DAG)  $\mathcal{G}$ . The nonparametric structural equations model (NPSEM) of a DAG states that each main variable  $V_{i,j} \in \mathbf{V}_i$  is a deterministic function of its parents in the graph  $\mathcal{G}$ , denoted  $\mathbf{pa}(V_{i,j})$ . That is, all factors determining  $V_{i,j}$  are contained in  $\mathbf{pa}(V_{i,j})$ , a subset of  $\mathbf{U}_i$  and  $\mathbf{V}_i$ . We denote the function mapping from  $\mathbf{pa}(V_{i,j})$  to  $V_{i,j}$  as  $V_{i,j} = f_j(\mathbf{pa}(V_{i,j}))$ ; we use  $\mathcal{F} = \{f_1, \dots, f_J\}$  to denote the collection of these structural equations, or the *structural causal model* of  $\mathbf{V}$  (Pearl, 2009; Richardson and Robins, 2013). Note that each main variable may be influenced by multiple disturbances, and a single disturbance may influence multiple main variables.

A DAG is said to be in canonical form if (i) all disturbances are exogeneous, i.e. no variable in  $\mathbf{U}_i$  has any parents in  $\mathcal{G}$ ; and (ii) there exists no pair of disturbances,  $U_i$  and  $U'_i$ , such that  $U_i$  influences a subset of the variables influenced by  $U'_i$ . Evans (2018) showed that for any DAG  $\mathcal{G}'$  not in canonical form; there exists a canonical form DAG  $\mathcal{G}$  with an identical distribution over all factual and counterfactual versions of all variables in  $\mathbf{V}_i$ . We can therefore without loss of generality limit our consideration to DAGs in canonical form. An example of a DAG not in canonical form is given in panel Figure 1(a). Panel Figure 1(b) illustrates the canonicalized version of this graph. For convenience, we will refer to the joint distribution over all factual and counterfactual versions of  $\mathbf{V}_i$  as the *full data law*. Moreover, any DAG over  $\mathbf{U}_i$ ,  $\mathbf{V}_i$ , and unobserved ancillary variables  $\mathbf{W}_i$  with unknown cardinality (e.g., confounders or mediators not of direct interest) also has an equivalent canonical DAG with respect to this full data law. A guide for canonicalizing arbitrary DAGs is given in Appendix A.1.

In short, representing the causal graph in canonical form distills the data-generating process (DGP) to its simplest form, eliminating potentially complex networks of disturbances. Removing variables that are irrelevant to the causal goal further simplifies the structure. Without these simplifications, it would be exceedingly difficult, if not intractable, to convert causal problems into polynomial programs that can be readily optimized—the essence of the approach we develop below.

Figure 1: **Canonicalization of a mediation model.** Mediation DAG in non-canonical form (panel a) and canonical form (panel b) that are fully equivalent with respect to their full data law. Unit indices,  $i$ , are suppressed. Canonicalization proceeds as follows: (i) the dependent disturbance  $U_3$  is absorbed into its parent  $U_{23}$ ; (ii) the superfluous  $U_2$  is eliminated as it influences a subset of  $U_{23}$ 's children; and (iii) the irrelevant  $W_{13}$  is absorbed into the  $V_1 \rightarrow V_3$  path as it is neither observed nor of interest. A complete guide to canonicalization is given in Appendix A.1.



## 3.2 Potential Outcomes

The notation of potential outcome functions allows us to compactly express the effects of manipulating a variable's parents or other ancestors. Let  $\mathbf{A} \subset \mathbf{V}$  be a subset of variables that will be intervened upon, fixing them to  $\mathbf{A} = \mathbf{a}$ . When  $\mathbf{A} = \emptyset$ , so that no intervention occurs, then define  $V_{i,j}(\mathbf{a}) = V_{i,j}$ , the natural value. When  $\mathbf{A} \subseteq \mathbf{pa}(V_{i,j})$ , so that only immediate parents are manipulated, then unit  $i$ 's potential outcome function is given by its response function,  $V_{i,j}(\mathbf{a}) = f_j(\mathbf{A} = \mathbf{a}, \mathbf{pa}(V_{i,j}) \setminus \mathbf{A})$ . We will now define more general potential outcome functions by *recursive substitution* (Richardson and Robins, 2013; Shpitser, 2018). For arbitrary interventions on  $\mathbf{A} \subset \mathbf{V}$ , let  $V_{i,j}(\mathbf{a}) = V_{i,j}(\{a_\ell : \mathbf{A}_\ell \in \mathbf{pa}(V_{i,j})\} \cup \{V_{i,j'}(\mathbf{a}) : V_{i,j'} \in \mathbf{pa}(V_{i,j}) \setminus \mathbf{A}\})$ ; here,  $\ell$  is a generic index that sweeps over main variables in the graph. This means that if a parent of  $V_{i,j}$  is directly manipulated, it is set to the corresponding value in  $\mathbf{a}$ . Otherwise, the parent takes on its potential value after intervention on any causally prior variables, or its natural value otherwise. To obtain the parent's potential value, we follow the same definition recursively. When defining potential outcomes, intervention on  $V_{i,j}$  itself is ignored. To illustrate, consider the mediation graph of Figure 1(b). Possible potential outcomes for  $V_{i,3}$  are (i)  $V_{i,3}(\emptyset) = V_{i,3}(V_{i,1}, V_{i,2})$ , the observed distribution; (ii)  $V_{i,3}(v_{i,1}) = V_{i,3}(v_{i,1}, V_{i,2}(v_{i,1}))$ , relating to total effects; and (iii)  $V_{i,3}(v_{i,1}, v_{i,2})$ , relating to controlled effects.

## 3.3 Principal Stratification

Analysts have little information about the disturbances  $\mathbf{U}_i$ , which may take on an infinite number of values. This poses an analytic challenge, as it is difficult to reason about infinite spaces. Here, we review a result that makes the general partial identification problem tractable despite this issue: broadly, when  $\mathbf{V}_i$  are discrete, the model that a DAG encodes can be represented by a finite number of parameters without loss of generality, as long as the reduced space is sufficiently large. We then introduce the *functional parameterization* used for this task, discuss its relationship to principal strata, and review how any marginal

of the full data law can be represented in terms of these parameters.

Finkelstein et al. (2021) show that there are finite state spaces for  $\mathbf{U}_i$  that do not restrict the NPSEM of a DAG for  $\Pr(\mathbf{V}_i = \mathbf{v})$ , i.e. the model over the *factual* main variables. In the following proposition, we extend this result to show that there are finite state spaces for  $\mathbf{U}_i$  that do not restrict the NPSEM of a DAG for the full data law—i.e., the full distribution over all factual and counterfactual versions of the main variables.

**Proposition 1.** *Suppose  $\mathcal{G}$  is a canonical DAG over discrete main variables  $\mathbf{V}_i$  and disturbances  $\mathbf{U}_i$  with infinite cardinality. Then the model over the full data law implied by  $\mathcal{G}$  is unchanged by assuming that the disturbances have finite cardinalities, provided those cardinalities are sufficiently large.*

A proof can be found in Appendix B, along with details on how to obtain a lower bound on non-restrictive cardinalities for the disturbances.

Further, Evans (2018) showed that for a large class of graphs called *geared* graphs, it is possible to develop a functional model that does not alter the causal model of a DAG. In the functional model of a graph, each main variable  $V_i$  is associated with a disturbance  $U_i$  that fully determines how  $V_i$  responds to the values of its remaining parents.<sup>3</sup>

Proposition 1 enables us to develop functional models for graphs that are not geared as well. Finkelstein et al. (2021) presents an algorithm for constructing a concise functional model for non-geared graphs, taking as input a disturbance cardinality that is non-restrictive of the model over factual random variables. By instead substituting the Proposition 1 disturbance cardinality, which may be larger and restricts neither the factual nor the counterfactual random variables, we obtain a functional model that is likewise non-restrictive of the full data law. Intuitively, functional models are closely related to principal stratification (Greenland and Robins, 1986; Frangakis and Rubin, 2002). For example, consider the simple DAG,

$$U_{i,1} \rightarrow V_{i,1} \rightarrow V_{i,2} \leftarrow U_{i,2} \tag{1}$$

in which  $V_{i,1}$  and  $V_{i,2}$  are binary. This relationship is governed by the structural equations  $V_{i,1} = f_1(U_{i,1})$  and  $V_{i,2} = f_2(V_{i,1}, U_{i,2})$ , where the functions  $f_1 : \mathcal{S}(U_{i,1}) \rightarrow \mathcal{S}(V_{i,1})$  and  $f_2 : \mathcal{S}(V_{i,1}) \times \mathcal{S}(U_{i,2}) \rightarrow \mathcal{S}(V_{i,2})$  are deterministic and shared across all units. Thus, the only source of randomness is in the disturbances,  $\mathbf{U}_i = \{U_{i,1}, U_{i,2}\}$ .

Analysts generally do not have direct information about these disturbances. For example,  $U_{i,1}$  could potentially take on any value in  $(-\infty, \infty)$ . However, Proposition 1 states that this variation is irrelevant, because  $V_{i,1}$  can only take on two possible values: 0 and 1. The space of  $U_{i,1}$  can therefore be divided into two *canonical partitions* (Balke and Pearl, 1997)—those that deterministically lead to  $V_{i,1} = 0$  and those that lead to  $V_{i,1} = 1$ —and thus there is no loss of generality in treating  $U_{i,1}$  as if it were binary.<sup>4</sup>

---

<sup>3</sup>Note that if any main variable  $V_i$  has multiple parents in  $\mathbf{U}_i$ , there may be multiple valid functional parameterizations, depending on which disturbance is chosen to determine which main variable. If each main variable has only a single parent in  $\mathbf{U}_i$ , there is only a single functional parameterization.

<sup>4</sup>See Section 8.2 of Pearl (2009) for a related discussion.

The situation with  $V_{i,2}$  is similar but more involved. After the random  $U_{i,2}$  is realized, it induces the *partially applied* response function  $V_{i,2} = f_2(V_{i,1}, U_{i,2} = u_2) = f_2^{(U_{i,2}=u_2)}(V_{i,1})$ , which deterministically governs how  $V_{i,2}$  counterfactually responds to  $V_{i,1}$ . Regardless of how many values the disturbance can take on, this response function must fall into one of only four possible groups, or *principal strata*, each corresponding to a possible function of the form  $f_2^{(U_{i,2}=u_2)} : \mathcal{S}(V_{i,1}) \rightarrow \mathcal{S}(V_{i,2})$  (Angrist et al., 1996). These groups are (i)  $V_{i,2} = 1$  regardless of  $V_{i,1}$ , “always takers”; (ii)  $V_{i,2} = 0$  regardless of  $V_{i,1}$ , “never takers”; (iii)  $V_{i,2} = V_{i,1}$ , “compliers”; and (iv)  $V_{i,2} = 1 - V_{i,1}$ , “defiers”. Thus, from the perspective of  $V_{i,2}$ , any finer-grained variation in  $\mathcal{S}(U_{i,2})$  beyond the canonical partitions is irrelevant. These partitions of  $\mathbf{U}$  are in one-to-one correspondence with principal strata, allowing causal quantities to be expressed in simple algebraic expressions; e.g., the average treatment effect (ATE) in (1) is equal to the proportion of compliers minus the proportion of defiers.<sup>5</sup> By writing down all information in terms of (possibly unknown) strata sizes, we can convert causal inference problems into tractable polynomial programming problems over these variables.

The functional parameterization of this graph has four free parameters: one for the binary  $U_{i,1}$  (or its reduced representation) and three for the quaternary  $U_{i,2}$ .<sup>6</sup> Because the distributions of disturbances are independent in canonical DAGs by virtue of their exogeneity, only their marginal distributions need be parameterized. Each of  $U_{i,1}$  and  $U_{i,2}$  encode full information about how  $V_{i,1}$  and  $V_{i,2}$  respectively respond to their remaining parents. In other words, each setting of  $\mathbf{U}_i$  provides full information not only about each variable  $\mathbf{V}_i$ , but also about each of its potential outcomes. This means that we can represent “cross-world” distributions such as  $\Pr(V_{i,2}(V_{i,1} = 0) = 0, V_{i,2}(V_{i,1} = 1) = 1)$ —the “complier” proportion—in terms of parameters of the marginal distributions of  $U_{i,2}$  alone. As we will see below, this fact will be useful in encoding cross-world type assumptions like monotonicity, as well as for bounding cross-world targets like the natural direct effect or the probability of causation. More generally, any marginal of the full data law may be expressed in terms of the functional parameters.

Finally, consider a more complex example, the mediation DAG of Figure 1(b). The response functions for  $V_{i,1}$  and  $V_{i,2}$  remain unchanged. In contrast,  $V_{i,3}$  is caused by  $\mathbf{pa}(V_{i,3}) = \{V_{i,1}, V_{i,2}\}$  via the structural equation  $V_{i,3} = f_3(V_{i,1}, V_{i,2}, U_{i,23})$ . Substituting in a realization of the disturbance,  $U_{i,23} = u_{i,23}$ , will produce one of sixteen response functions of the form  $f_3^{(U_{i,23}=u_{i,23})} : \mathcal{S}(V_{i,1}) \times \mathcal{S}(V_{i,2}) \rightarrow \mathcal{S}(V_{i,3})$ . More generally, the number of unique response functions grows with (i) the cardinality of the variable, (ii) the number of causal parents it has, and (iii) the parents’ cardinalities. Specifically,  $V_{i,j}$  has  $|\mathcal{S}(V_{i,j})|^{|\mathbf{pa}(V_{i,j})|}$  possible response functions: given a particular input from  $V_{i,j}$ ’s parents, the number of possible outputs for  $V_{i,j}$  is  $|\mathcal{S}(V_{i,j})|$ ; the number of possible inputs from  $V_{i,j}$ ’s parents is

<sup>5</sup>To see this, note that the ATE is given by  $\mathbb{E}[V_{i,2}(V_{i,1} = 1) - V_{i,2}(V_{i,1} = 0)] = \sum_{\text{strata}} \mathbb{E}[V_{i,2}(V_{i,1} = 1) - V_{i,2}(V_{i,1} = 0) \mid \text{strata}] \cdot \Pr(\text{strata}) = 0 \cdot \Pr(\text{always taker}) + 0 \cdot \Pr(\text{never taker}) + 1 \cdot \Pr(\text{complier}) - 1 \cdot \Pr(\text{defier})$ .

<sup>6</sup>These can be thought of as the probabilities of encountering a unit of the “control” type with  $V_{i,1} = 0$  (for  $U_{i,1}$ ) and of encountering units of the “always-taker,” “never-taker,” and “complier” types (for  $U_{i,2}$ ). These parameters determine the probabilities of the remaining types (the “treatment” type for  $U_{i,1}$  and the “defier” type for  $U_{i,2}$ ), as principal strata probabilities must sum to unity.

$|\mathcal{S}(\mathbf{pa}(V_{i,j}))| = \prod_{V_{i,j'} \in \mathbf{pa}(V_{i,j})} |\mathcal{S}(V_{i,j'})|$ , the product of the parents’ cardinalities.

In turn, this determines the minimal cardinality of  $\mathbf{U}$  in a reduced but non-restrictive functional model—roughly speaking, the number of principal strata combinations that exist, if we think of  $\mathbf{U}$  as principal strata. Here, “non-restrictive” means that the simplified model is fully expressive, or that it can represent any possible full data law. For example, to capture the joint response patterns that a unit may have on  $V_{i,2}$  and  $V_{i,3}$ , a reduced version of  $U_{i,23}$  will be capable of expressing any full data law if it has a cardinality of  $|\mathcal{S}(U_{23})| = 4 \times 16$ , because  $V_{i,2}$  has four possible response functions and  $V_{i,3}$  has sixteen.

## 4 Formulating the Polynomial Program

We now turn to the central problem of this paper: sharply bounding causal quantities with missing data. Our approach is to (i) rewrite the causal query into a polynomial expression, (ii) rewrite modeling assumptions and empirical information into polynomial constraints, and (iii) thereby transform the task into a constrained optimization problem that can be solved computationally. *Sharp bounds* are the narrowest range that contain all admissible values for a target quantity, i.e., all values that are consistent with information available to the analyst: structural causal knowledge in the form of a canonical DAG,  $\mathcal{G}$ ; as well as empirical evidence,  $\mathcal{E}$ , and modeling assumptions,  $\mathcal{A}$ , formalized below. We also suppose that the main variables take on values in a known, discrete set,  $\mathcal{S} = \mathcal{S}(\mathbf{V})$ . In this section, we will demonstrate (i) that  $\{\mathcal{G}, \mathcal{E}, \mathcal{A}, \mathcal{S}\}$  restricts the admissible values of the target quantity, and (ii) this range of observationally indistinguishable values can be recovered by polynomial programming.

The causal graph and sample space,  $\mathcal{G}$  and  $\mathcal{S}$ , together imply a set of possible functional models, each fully characterizing the main variables. By Proposition 1, without loss of generality, we can consider a simple functional model in which (i) each counterfactual main variable is a deterministic function of exogeneous, discrete disturbances; (ii) there are a relatively small number of such disturbances; and (iii) disturbances take on a finite number of possible values, corresponding to principal strata of the main variables. When repeatedly sampling units (along with each unit’s random disturbances,  $\mathbf{U}_i$ ), the  $k$ -th disturbance thus follows the categorical distribution with parameters  $\mathcal{P}_{U_k} = \{\Pr(U_{i,k} = u_{i,k}) : u_{i,k}\}$ . By the properties of canonical DAGs, these disturbances are independent. It follows that the parameters  $\mathcal{P}_{\mathbf{U}}$  of the joint disturbance distribution  $\Pr(\mathbf{U}_i = \mathbf{u}_i) = \prod_k \Pr(U_{i,k} = u_{i,k})$  not only fully determine the distribution of each factual main variable—i.e. the potential outcome under no intervention,  $V_{i,j}(\emptyset)$ —they also determine the counterfactual distribution of  $V_{i,j}(\mathbf{a})$  under any intervention  $\mathbf{a}$ , and its joint distribution with other counterfactual variables  $V_{i,j'}(\mathbf{a}')$  under possibly different interventions  $\mathbf{a}'$ . This leads to the following proposition.

**Proposition 2.** *Suppose  $\mathcal{G}$  is a canonical DAG and  $\mathcal{C} = \{V_{i,\ell}(\mathbf{a}_\ell) = v_\ell\}$  is a set of counterfactual statements, indexed by  $\ell$ , that variable  $V_{i,\ell}$  will take on value  $v_\ell$  under manipulation  $\mathbf{a}_\ell$ . Let  $\mathcal{U} \subset \mathcal{S}(\mathbf{U})$  indicate the subset of disturbance realizations that lead deterministically*

to every statement in  $\mathcal{C}$  being satisfied. Then under the structural equation model  $\mathcal{G}$ ,

$$\Pr \left( \bigwedge_{\ell} \mathcal{C}_{\ell} \right) = \sum_{\mathbf{u} \in \mathcal{U}} \prod_{u_k \in \mathbf{u}} \Pr(U_{i,k} = u_k), \quad (2)$$

which is a polynomial equation in  $\mathcal{P}_{\mathbf{U}_i}$ , the parameters of  $\Pr(\mathbf{U} = \mathbf{u})$ .

For example, in the mediation setting of Figure 1(b), Proposition 2 implies that the joint distribution of the factual variables— $V_{i,1}(\emptyset)$ ,  $V_{i,2}(\emptyset)$ , and  $V_{i,3}(\emptyset)$ —is given by

$$\Pr(V_{i,1}(\emptyset) = v_1, V_{i,2}(\emptyset) = v_2, V_{i,3}(\emptyset) = v_3) = \sum_{\{u_1, u_{23}\} \in \mathcal{U}} \Pr(U_1 = u_1) \Pr(U_{23} = u_{23}), \quad (3)$$

where  $\mathcal{U} = \left\{ \{u_1, u_{23}\} : f_1^{(U_1=u_1)}(\emptyset) = v_1, f_2^{(U_{23}=u_{23})}(v_1) = v_2, f_3^{(U_{23}=u_{23})}(v_1, v_2) = v_3 \right\}$  is the set of all disturbances that are consistent with a particular  $\mathbf{V}_i = \{v_1, v_2, v_3\}$ . Alternatively, analysts may be interested in the probability that a randomly drawn unit  $i$  has a positive controlled direct effect when fixing the mediator to  $V_{i,2} = 0$ . This is given by  $\Pr(V_{i,3}(V_{i,1} = 0, V_{i,2} = 0) = 0, V_{i,3}(V_{i,1} = 1, V_{i,2} = 0) = 1)$  and is similarly expressed in terms of the disturbances as  $\sum_{\{u_1, u_{23}\} \in \mathcal{U}'} \Pr(U_{i,1} = u_1) \Pr(U_{i,23} = u_{23})$ , summing over a different subset of the disturbance space,  $\mathcal{U}' = \left\{ \{u_1, u_{23}\} : f_3^{(U_{i,23}=u_{23})}(V_{i,1} = 1, V_{i,2} = 0) = 1, f_3^{(U_{i,23}=u_{23})}(V_{i,1} = 0, V_{i,2} = 0) = 0 \right\}$ .

We now expand this result to include a large class of functionals of marginal probabilities and logical statements about these functionals.

**Corollary 1.** *Suppose  $\mathcal{G}$  is a canonical DAG. Let  $\mathcal{P}_{\mathbf{V}}$  denote the full data law and  $g(\mathcal{P}_{\mathbf{V}})$  denote a functional of  $\mathcal{P}_{\mathbf{V}}$  involving elementary arithmetic operations on constants and marginal probabilities of  $\mathcal{P}_{\mathbf{V}}$ . Then  $g(\mathcal{P}_{\mathbf{V}})$  can be expressed as a polynomial fraction in the parameters of  $\mathcal{P}_{\mathbf{U}}$ ,  $h(\mathcal{P}_{\mathbf{U}})$ , by replacing each marginal probability with its Proposition 2 polynomialization.*

We say functionals of the full data law that fulfill these properties are *polynomial-fractionalizable*, or simply *polynomializable* if the result contains no fractions. The corollary has a number of implications, which we briefly discuss here. First, it demonstrates that a wide array of single-world and cross-world functionals can be expressed as polynomial fractions. These include traditional targets such as the ATE, as well as more complex targets such as the pure direct effect and the probability of causal sufficiency. It also suggests that any non-elementary functional of  $\mathcal{P}_{\mathbf{V}}$  can be approximated to arbitrary precision by a polynomial fraction, provided the functional has a convergent power series.<sup>7</sup>

Next, observe that when (i)  $g(\mathcal{P}_{\mathbf{V}})$  is a polynomial-fractionalizable expression; (ii)  $\star \in \{<, \leq, =, >, \geq, \neq\}$  is a binary comparison operator; and (iii)  $\alpha$  is a constant, then

<sup>7</sup>We note that non-elementary functionals rarely arise in practice, with the exception of target quantities on logarithmic or exponential scales. In such cases, bounds on monotonic transformations of polynomials can be straightforwardly obtained by bounding the underlying polynomial, then applying the transformation. An example of a functional that our approach cannot handle is the non-analytic  $\mathbb{1}(\text{ATE is rational})$ .

statements of the form  $g(\mathcal{P}_{\mathcal{U}}) \star \alpha$  can be equivalently expressed as non-fractional polynomial relations  $h(\mathcal{P}_{\mathcal{U}}) \star 0$ . Finally, by the same token, any polynomial-fractional expression  $h(\mathcal{P}_{\mathcal{U}})$  in the parameters of  $\mathcal{P}_{\mathcal{U}}$  can be reexpressed with (i) a non-fractional polynomial in an expanded parameter space and (ii) a polynomial equation in the same expanded space.<sup>8</sup> We will make extensive use of these properties to convert causal queries to polynomial programs.

In Appendix A, Algorithm 1 provides a step-by-step procedure for obtaining sharp bounds. We begin by transforming a factual or counterfactual target of inference  $\mathcal{T}$  into polynomial form, possibly with the use of additional auxiliary variables to eliminate fractions. To accomplish this task, the procedure utilizes the possibly non-canonical DAG  $\mathcal{G}$  and the possible main-variable outcomes  $\mathcal{S}(\mathbf{V})$  to reexpress  $\mathcal{T}$  in terms of functional parameters that correspond to principal strata proportions. The result is the objective function of the polynomial program. The procedure then polynomializes the sets of constraints on polynomializable functionals resulting from empirical evidence and by modeling assumptions, respectively  $\mathcal{E}$  and  $\mathcal{A}$ . For example, in the binary mediation setting of Figure 1,  $\mathcal{G}$  may be the graph depicted in either panel (a) or (b). If only observational data is available, then  $\mathcal{E}$  consists of eight pieces of evidence, each represented as a statement corresponding to a cell of the factual distribution  $\Pr(V_{i,1}(\emptyset) = v_1, V_{i,2}(\emptyset) = v_2, V_{i,3}(\emptyset) = v_3) = \Pr(V_{i,1} = v_1, V_{i,2} = v_2, V_{i,3} = v_3)$  for observable values in  $\{0, 1\}^3$ . Modeling assumptions include all other information, such as monotonicity or dose-response assumptions; these can be expressed in terms of principal strata. For example, the assumed unit-level monotonicity of the  $V_1 \rightarrow V_2$  relationship (e.g., the “no defiers” assumption of Angrist et al., 1996) can be written as the statement that  $\Pr(V_{i,2}(V_{i,1} = 0) = 1, V_{i,2}(V_{i,1} = 1) = 0) = 0$ . Assumed population-level monotonicity is typically written  $\mathbb{E}[V_{i,2}(V_{i,1} = 1) - V_{i,2}(V_{i,1} = 0)] \geq 0$ , but can equivalently be reformulated in terms of principal strata as  $\Pr(V_{i,2}(V_{i,1} = 1) = 1, V_{i,2}(V_{i,1} = 0) = 0) - \Pr(V_{i,2}(V_{i,1} = 0) = 1, V_{i,2}(V_{i,1} = 1) = 0) \geq 0$ . Finally, the statement that each disturbance  $k$  follows a categorical probability distribution is reexpressed as the polynomial relations  $\Pr(U_k = u_k) \geq 0 : u_k$  and  $\sum_{u_k} \Pr(U_k = u_k) = 1$ .

Algorithm 1 produces an optimization problem with a polynomial objective subject to polynomial constraints. This polynomial programming problem is equivalent to the original causal bounding problem. This leads directly to the following theorem.

**Theorem 1.** *Minimization (maximization) of the polynomial program produced by Algorithm 1 produces sharp lower (upper) bounds on  $\mathcal{T}$  under the sample space  $\mathcal{S}(\mathbf{V})$ , structural equation model  $\mathcal{G}$ , additional modeling assumptions  $\mathcal{A}$ , and empirical evidence  $\mathcal{E}$ .*

Once the causal problem is expressed in polynomial form, a variety of computational solvers can in principle be used to optimize (e.g. IPOPT; Wächter and Biegler, 2006). However, local solvers cannot guarantee valid bounds without exhaustively searching the

---

<sup>8</sup>To see this, let  $s$  be a scalar auxiliary variable and set  $h(\mathcal{P}_{\mathcal{U}}) = s$ , which can be manipulated to obtain a non-fractional polynomial equation, per (ii). The original expression can now be rewritten simply as  $s$ , which is a monomial and hence a polynomial, per (i). Thus, the original polynomial-fractional expression has been reexpressed in terms of (i) a non-fractional polynomial expression and (ii) a non-fractional polynomial equation.

space; when time is limited, these often fail to discover global extrema for the causal estimand, resulting in intervals that may fail to contain the quantity of interest. Moreover, such approaches often become computationally intractable as causal problems grow complex. In the next section, we show how the polynomial program can be simplified to speed computation.

## 5 Simplifying the Polynomial Program

Because solving polynomial programs is in general NP-hard, efficient computation requires us to fully exploit our knowledge of the problem structure. This knowledge allows analysts to reduce the complexity of the program in ways that algebraic presolvers may not necessarily detect. In this section, we discuss several ways to do this.

### 5.1 Eliminating Blocked Disturbances

We begin by using the following observation to limit the number of disturbance distribution parameters involved in the target and constraints.

**Proposition 3.** *Consider the polynomialization of a probability  $\Pr\left(\bigwedge_{\ell} \mathcal{C}_{\ell}\right)$ , where  $\mathcal{C} = \{V_{i,\ell}(\mathbf{a}_{\ell}) = v_{\ell} : \ell\}$ . We say that for intervention  $\mathbf{a}_{\ell}$ , a disturbance  $U_{i,k}$  is blocked from the corresponding counterfactual  $V_{i,\ell}$  if there are no paths from  $U_{i,k}$  to  $V_{i,\ell}$  that do not go through the causally prior members of the intervention  $\mathbf{a}_{\ell}$ . When  $U_{i,k}$  is blocked from  $V_{i,\ell}$  for every  $\ell$ , the corresponding parameters  $\mathcal{P}_{U_k}$  can be eliminated from the polynomialization.*

In other words, Proposition 3 states that each main variable  $V_{i,j}$  is only a function of its ancestors in  $\mathbf{U}$  that affect it through a variable not under intervention. For each marginal probability of an event, the disturbances that do not affect any variable in the event are irrelevant. This allows us to amend the polynomialization of Proposition 2 so that the outer sum ranges only over all possible settings of *relevant* disturbances, reducing the degree of each term in the polynomial. For example, in the mediation graph of Figure 1(b), consider the total effect of the treatment  $V_{i,1}$  on the outcome  $V_{i,3}$ . Here, all probabilities are of the form  $V_{i,3}(v_{i,1} = a_1) = v_3$ .<sup>9</sup> The disturbance  $U_{i,1}$  is therefore blocked from the outcome  $V_{i,3}$ , because the sole path from  $U_{i,1}$  to  $V_{i,3}$  goes through the intervention set  $V_{i,1}$ . This means that whenever  $\Pr(U_1 = u_1)$  appears in the polynomial, it does so in a way that ensures  $\sum_{u_1} \Pr(U_1 = u_1) = 1$  can be factored out and eliminated.

### 5.2 Exploiting the Nested Markov Parameterization

We now consider the common case when the empirical evidence  $\mathcal{E}$  includes *single-world marginal distributions*. This occurs when (i) a factual or counterfactual event  $\bigwedge_{\ell} \{V_{i,\ell}(\mathbf{a}) = v_{\ell}\}$  involves the same intervention,  $\mathbf{a}$ , for every variable of interest; and (ii) the probability

<sup>9</sup>The total effect is given by  $\Pr(V_{i,3}(V_{i,1} = 1) = 1) - \Pr(V_{i,3}(V_{i,1} = 0) = 1)$ , which can equivalently be written  $\Pr(V_{i,3}(V_{i,1} = 1, V_{i,2} = V_{i,2}(V_{i,1} = 1)) = 1) - \Pr(V_{i,3}(V_{i,1} = 0, V_{i,2} = V_{i,2}(V_{i,1} = 0)) = 1)$ .

$\Pr\left(\bigwedge_{\ell}\{V_{i,\ell}(\mathbf{a}) = v_{\ell}\}\right)$  is observed for every event in that state space,  $\{v_{\ell} : \ell\} \in \prod_{\ell} \mathcal{S}(V_{i,\ell})$ . For example, in the binary mediation setting of Figure 1(b), an observational study might obtain information about  $\Pr(V_{i,1}(\emptyset) = v_1, V_{i,2}(\emptyset) = v_2, V_{i,3}(\emptyset) = v_3)$  for every combination of  $v_1$ ,  $v_2$ , and  $v_3$ . Similarly, an experiment that randomly manipulated  $V_{i,1}$  would obtain two such distributions: (i) by observing  $\Pr(V_{i,2}(V_{i,1} = 0) = v_2, V_{i,3}(V_{i,1} = 0) = v_3)$  for all  $v_2$  and  $v_3$ ; and similarly, (ii) by observing  $\Pr(V_{i,2}(V_{i,1} = 1) = v_2, V_{i,3}(V_{i,1} = 1) = v_3)$  for all  $v_2$  and  $v_3$ .

A naïve parameterization might use one parameter for the probability of each principal strata (e.g., four parameters for the proportion of “always takers,” “never takers,” “compliers,” and “defiers”). It is immediately apparent that one naïve parameter is redundant, as its value is already implied by the fact that all distributions must marginalize to unity. Hidden-variable DAG models imply additional equality constraints, each of which can likewise be used to reduce the number of parameters needed to describe the model, thereby reducing the number of polynomial constraints in the program.

These additional equality constraints, which are implied by the structural equations model of  $\mathcal{G}$ , are well understood. In particular,  $\mathcal{G}$  imposes certain *conditional independence* and *generalized equality* constraints (or Verma constraints, Verma and Pearl, 1990; Tian and Pearl, 2002) on these distributions. Each equality constraint can be used to eliminate one parameter of the single-world distribution. The parameterization that takes full advantage of these structural equality constraints to reduce the number of parameters is called the *nested Markov* parameterization (Evans et al., 2019). This parameterization achieves the minimal number of parameters, equal to the dimension of the model of  $\mathcal{G}$ .

Each nested Markov parameter is exactly equal to an identified marginal probability of a single-world event,  $\Pr\left(\bigwedge_{\ell'}\{V_{i,\ell'}(\mathbf{a}') = v_{\ell'}\}\right)$ . Because each of the nested Markov parameters is identified from the initial single-world distribution,  $\Pr\left(\bigwedge_{\ell}\{V_{i,\ell}(\mathbf{a}) = v_{\ell}\}\right)$ , whether indirectly or directly, it can be calculated directly from the empirical evidence  $\mathcal{E}$ . By Proposition 2, this probability remains polynomializable in the parameters  $\mathcal{P}_{\mathcal{U}}$ . This allows us to add one equality constraint to the program per nested Markov parameter, based on its polynomialization, rather than one equality constraint per outcome in the state space. For hidden variable DAGs that imply a large number of equality constraints, this can substantially reduce the number of constraints. Evans et al. (2019) offers a complete guide to obtaining nested Markov formulations of arbitrary single-world distributions.

Each parameter in the nested Markov formulation is the probability of a single-world event that involves fewer main variables and more interventions, compared to any event used by the naïve parameterization. As a result, the corresponding polynomialization will have fewer terms, of lower degree, than the polynomialization of naïve parameters. An example is provided in Appendix A.4. Using this technique, we modify Algorithm 1 by partitioning the empirical evidence into all single-world marginal distributions  $\mathcal{E}_M$  and the remaining evidence  $\mathcal{E}_R$ . The constraints in  $\mathcal{E}_M$  can then be reduced into their nested Markov form before polynomialization. Because the nested Markov parameterization allows for fewer, simpler polynomial constraints in the program, it is important to use it whenever

the empirical evidence permits.

We note that when certain deterministic relationships exist between variables in  $V_i$ , as in the missing-data setting of Figure 5(c–d),<sup>10</sup> these relationships may imply equality constraints not exploited by the nested Markov parameterization. In such cases, it may be possible to further reduce the number of constraints; we do not explore that option here.

### 5.3 Eliminating Additional Constraints and Parameters

Finally, we describe when constraints and parameters can be safely eliminated from a program. We say that parameters  $x$  and  $y$  *co-occur* in a polynomial system if they appear in the same constraint; they *interact* if there exists a sequence of parameters from  $x$  to  $y$  such that every adjacent pair co-occurs.<sup>11</sup> If a constraint’s parameters do not interact with the objective’s parameters, that constraint may be dropped. If a parameter exists only in constraints that have been eliminated, then the parameter has also been eliminated, simplifying the system.

This may be used in conjunction with the structure of  $\mathcal{G}$  to help simplify the program, because different *districts*—components in  $\mathcal{G}$  connected by bidirected arcs (Tian and Pearl, 2002; Richardson, 2003)—do typically do not interact. That is, likelihoods on marginal distributions of  $\mathcal{G}$  have a representation that is decomposable by districts. For example, in Figure 1(b),  $V_{i,1}$  lies in one district; in contrast,  $V_{i,2}$  and  $V_{i,3}$  lie in another district, because they are connected by  $U_{i,23}$ . Because each nested Markov parameter is the probability of a single-world event involving main variables within a single district, its polynomialization will at most involve disturbance parameters in that district. This leads to the following proposition.

**Proposition 4.** *The degree of each polynomial in the nested Markov constraints is bounded from above by the number of latent variables in the corresponding district. Moreover, if two disturbances  $U_{i,k}$  and  $U_{i,k'}$  appear in different districts, their parameters  $\mathcal{P}_{U_k}$  and  $\mathcal{P}_{U_{k'}}$  will not interact in any nested Markov constraint.*

To illustrate, consider the common scenario where an analyst observes the full joint distribution over factual variables,  $\Pr(V_{i,1}(\emptyset) = v_1, \dots, V_{i,J}(\emptyset) = v_J)$ , and seeks to bound a functional relating a treatment  $\mathbf{a}$  to an outcome  $V_{i,j}(\mathbf{a})$  in the same district. As an example, in Figure 1(b), the effect of the mediator  $V_{i,2}$  on the outcome  $V_{i,3}$  is wholly contained within a single district. We can therefore drop all constraints related to nested Markov parameters involving other districts, and thus all disturbance parameters in other districts.

---

<sup>10</sup>In this graph, a latent variable  $Y$  has an observed version  $Y^*$  that deterministically inherits  $Y^* = Y$  when a reporting variable  $R = 1$ , but takes on the missing-value indicator  $Y^* = \text{NA}$  otherwise.

<sup>11</sup>For example, consider the constraints  $x + y = a$ ,  $y + z = b$ . Here,  $x$  and  $y$  co-occur;  $x$  and  $z$  interact.

## 6 Computing $\varepsilon$ -sharp Bounds in Polynomial Programs

We now turn to the practical optimization of the polynomial program defined by Algorithm 1. Theorem 1 states minimization and maximization of this original primal program is equivalent to the initial bounding problem. However, obtaining globally optimal solutions in polynomial programming can be computationally intensive. Worryingly, methods that iteratively improve suboptimal values for the primal problem may fail to produce valid bounds (i.e., bounds containing all possible values of the estimand, including global extrema) without searching the full parameter space,  $\mathcal{P}$ . To address this challenge, we use *dual* methods that construct and iteratively refine an outer envelope around the primal function (i.e. the objective function, or causal quantity of interest). Specifically, we employ a variation of the spatial branch-and-bound method, combined with a piecewise linear envelope, implemented using a variety of optimization frameworks that include SCIP and Couenne (Vigerske and Gleixner, 2018; Gamrath et al., 2020; Belotti et al., 2009). Throughout the optimization process, current suboptimal values for dual minimization and maximization problems are guaranteed to produce valid but loose *outer* bounds; current suboptimal values for the primal problem produce possibly invalid *inner* bounds; and the lower (upper) endpoint of the unknown sharp bounds is guaranteed to lie between the current suboptimal primal and dual minimization (maximization) values. Through simultaneous *primal-dual* optimization, we use these suboptimal inner bounds to precisely quantify worst-case looseness,  $\varepsilon$ , of the suboptimal but valid outer bounds. This allows researchers to assess how more computation may lead to tightened conclusions.

A step-by-step description of our optimization procedure, which we term  $\varepsilon$ -sharp bounding, is given in Algorithm 2 of Appendix A. At a high level, it proceeds as follows. Our procedure takes as inputs the polynomialized objective function  $\mathcal{T}(\mathbf{p})$  and constraint set  $\mathcal{C}(\mathbf{p})$ , obtained from Algorithm 1. It then evaluates a range of models, or points  $\mathbf{p}$  in the model space  $\mathcal{P}$  for which  $\mathcal{C}(\mathbf{p})$  is satisfied. It seeks to identify extreme values of  $\mathcal{T}(\mathbf{p})$  within this subspace. It also accepts two parameters:  $\epsilon^{\text{thresh}}$ , a stopping threshold for the looseness factor stopping, and  $\theta^{\text{thresh}}$ , a stopping threshold for width of the bounds. The algorithm returns two types of information: the bounds for the causal program, and the worst-case looseness factor  $\varepsilon$ .

Primal bounds are denoted  $\underline{P}$  and  $\overline{P}$ , adopting the convention that underlines refer to objects used for minimization and overlines for maximization. These indicate the extreme values of the target estimand in any admissible model—that is, satisfying  $\mathcal{C}(\mathbf{p})$ —that has been located so far. These are initialized at  $+\infty$  and  $-\infty$ , respectively, indicating that no admissible models have been found yet. As optimization proceeds, the primal bounds improve as new, more extreme admissible models are found. We refer to  $[\underline{P}, \overline{P}]$  as the *inner bounds*: the unknown sharp bounds must at least contain these points, which correspond to models that are observationally indistinguishable from the true DGP.

Dual optimization begins by partitioning the parameter space into branches, proceeding separately for the lower and upper bound and respectively producing partitions  $\underline{\mathcal{B}}_b$  and  $\overline{\mathcal{B}}_b$ . At initialization, these consist of a single branch spanning the entire parameter space; each

branch is then recursively divided. The lower and upper parts of the dual envelope, or outer envelope, are denoted  $\underline{\mathcal{D}}$  and  $\overline{\mathcal{D}}$ . These are piecewise linear functions, with pieces corresponding to the branching partitions, that are *relaxations* of the true objective function,  $\mathcal{T}(\mathbf{p})$ , from below and above. These relaxations are made to ensure they will always contain the entire objective function at all points in the parameter space. Within branch  $b$ , the value  $\min\{\underline{\mathcal{D}}_b(\mathbf{p}) : \mathbf{p} \in \overline{\mathcal{B}}_b\}$  indicates the lowest value attained by the lower envelope; thus,  $\underline{\mathcal{T}} = \min_b \{\min\{\underline{\mathcal{D}}_b(\mathbf{p}) : \mathbf{p} \in \overline{\mathcal{B}}_b\}\}$  represents the lowest value attained by the lower envelope anywhere in the parameter space. Conversely,  $\overline{\mathcal{T}} = \max_b \{\max\{\overline{\mathcal{D}}_b(\mathbf{p}) : \mathbf{p} \in \overline{\mathcal{B}}_b\}\}$  represents the highest value of the upper envelope. These extreme points on the dual envelope,  $[\underline{\mathcal{T}}, \overline{\mathcal{T}}]$ , define the dual (outer) bounds. These are the reported causal bounds; whatever the true sharp bounds, they must lie inside the dual bounds, even if the algorithm has not run to completion. We let  $\theta$  equal the bound width, or the difference between the upper and lower dual bounds, and we define the worst-case looseness factor  $\varepsilon$  as the slack (the difference in dual and primal bound widths) divided by the primal bound width.

The algorithm heuristically selects branches in the model space that appear promising, and refines primal and dual bounds in turn. It first searches within the branch for an admissible model; if found, and if the associated causal estimand is more extreme than those previously encountered, it is stored as a new primal bound. Whatever the true nonparametric sharp bounds, they must lie outside the primal bounds because the true bounds must contain the extreme models that define the primal bounds. Then, it divides the branch into sub-branches and refines the dual envelope by tightening the piecewise linear outer-approximation. The algorithm continuously prunes branches of  $\underline{\mathcal{B}}_b$  and  $\overline{\mathcal{B}}_b$  that wholly violate constraints; it also continuously branches and refines the bounds while  $\theta$  and  $\varepsilon$  exceed specified thresholds.

## 7 Statistical Inference

We now turn to statistical inference for the bounds developed above. We say that the results of Algorithm 2 when applied to  $\mathcal{E}$ , the population empirical constraints—i.e., margins of the full data law that are observed without sampling error—are *population bounds*. In practice, the empirical quantities used in these constraints are estimated from finite samples. Our goal in this section is to account for variation in  $\hat{\mathcal{E}}$ , the estimated constraints, that arises over repeated sampling. The results of Algorithm 2 when substituting  $\hat{\mathcal{E}}$  for  $\mathcal{E}$  are referred to as the *estimated bounds*. In this section, we describe how to construct *confidence bounds* that (i) contain the estimated bounds and (ii) contain the population bounds at a rate of at least the confidence level  $\alpha$  over repeated samples.

Recall that each element of empirical evidence  $\mathcal{E}$  is a relation between (i) some population quantity that is an observable functional of the main variables' distribution,  $g(\mathcal{P}_{\mathbf{V}})$ , reexpressed in terms of the disturbance distribution  $\mathcal{P}_{\mathbf{U}}$ ; and (ii) the population value of that observable quantity. In  $\hat{\mathcal{E}}$ , we plug in for (ii) the estimated value of the quantity in finite data. For example, in the mediation graph of Figure 1(b), an analyst with access to a sample of observational data would have

$$\hat{\mathcal{E}} = \left\{ \text{polynomialize} \left( \begin{array}{l} \Pr (V_{i,1}(\varnothing) = v_1, V_{i,2}(\varnothing) = v_2, V_{i,3}(\varnothing) = v_3) \\ = \frac{1}{N} \sum_{i=1}^N \mathbb{1} \{V_{i,1} = v_1, V_{i,2} = v_2, V_{i,3} = v_3\} \end{array} \right) : v_1, v_2, v_3 \right\} \quad (4)$$

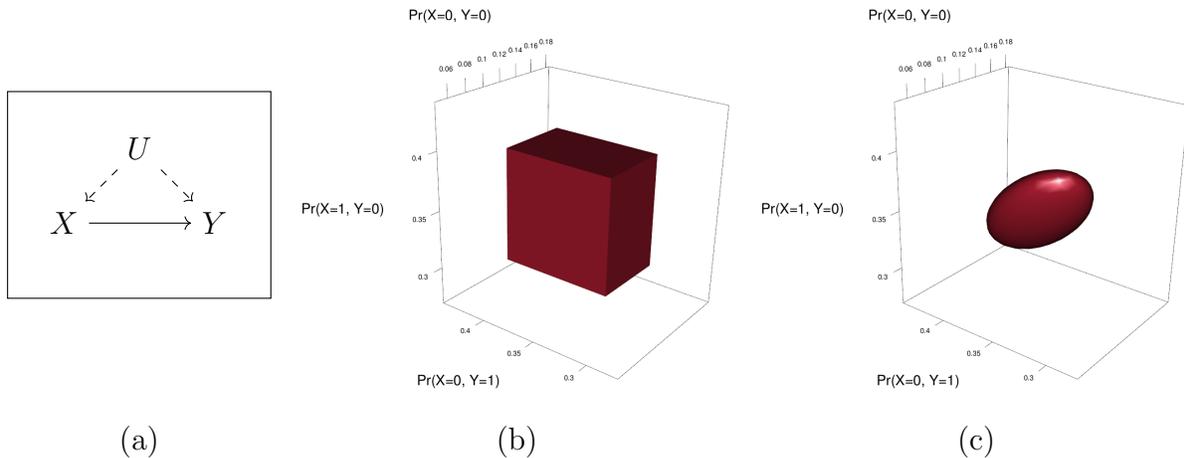
We will refer to the vector of estimated quantities on the right-hand side of  $\hat{\mathcal{E}}$  elements—in the above example, quantities of the form  $\frac{1}{N} \sum_{i=1}^N \mathbb{1} (V_{i,1} = v_1, V_{i,2} = v_2, V_{i,3} = v_3)$ —as  $\hat{\mathbf{E}}$ . We denote the corresponding population quantities as  $\mathbf{E}$ , the right-hand side values in  $\mathcal{E}$ .

To construct confidence bounds we consider the sampling variability of these estimated quantities. We construct regions,  $\text{CR}_\alpha(\hat{\mathbf{E}})$ , containing  $\hat{\mathbf{E}}$  and guaranteed to contain the population quantities  $\mathbf{E}$  with at least probability  $\alpha$  over repeated samples. These regions correspond to population distributions over observed parameters that cannot be rejected at level  $\alpha$ . In Algorithm 2, we then replace the  $\hat{\mathcal{E}}$  constraints with a set of loosened *confidence constraints*  $\text{CR}_\alpha(\hat{\mathcal{E}})$ . In other words, if the population bounds are obtained by optimizing subject to a equality constraint  $\{g_\ell(\mathcal{P}_\mathbf{V}) = \mathbf{E}_\ell\} \in \mathcal{E}$ , and the estimated bounds are obtained with the plug-in version  $\{g_\ell(\mathcal{P}_\mathbf{V}) = \hat{\mathbf{E}}_\ell\} \in \hat{\mathcal{E}}$ , then the confidence bounds will incorporate the interval constraint  $\{g_\ell(\mathcal{P}_\mathbf{V}) \in \text{CR}_\alpha(\hat{\mathbf{E}}_\ell)\} \in \text{CR}_\alpha(\hat{\mathcal{E}})$ .

Because loosening  $\hat{\mathcal{E}}$  to  $\text{CR}_\alpha(\hat{\mathcal{E}})$  can only decrease (increase) the minimum (maximum) value obtained by the polynomial program, confidence bounds always contain the estimated bounds. Similarly, when the confidence region for estimated quantities fully contains their population analogues, then the confidence constraint  $\text{CR}_\alpha(\hat{\mathcal{E}})$  is looser than the population constraint  $\mathcal{E}$ , and resulting confidence bounds also contain the population bounds. However, when the confidence region does not fully contain population quantities due to sampling error, confidence bounds may still contain population bounds. This can occur if the non-covered quantity corresponds to a constraint that is irrelevant to the bounds. Therefore, if the confidence region on the observed quantities has coverage of exactly  $\alpha$ , confidence bounds will contain the population bounds in at least  $\alpha$  of repeated samples.

In discrete settings, the task of obtaining confidence bounds thus reduces to the problem of constructing regions  $\text{CR}_\alpha(\hat{\mathbf{E}})$  for the multinomial proportion, such that  $\Pr(\mathbf{E} \in \text{CR}_\alpha(\hat{\mathbf{E}})) \geq \alpha$ . We focus on two methods for doing so. Drawing on Malloy et al. (2020), we first consider a “Bernoulli-KL” approach that constructs separate confidence regions for each observable atomic event,  $\Pr(\mathbf{V}_i = \mathbf{v})$ , treating it as a “success” in a Bernoulli distribution. The approach rotates through all possible  $\mathbf{v}$  and combines the event-specific regions using a result on the Kullback-Leibler divergence of sampling distributions to the underlying population distribution. The Bernoulli-KL method produces a confidence region for single-world distributions that is guaranteed to have conservative coverage for the multinomial proportion in finite samples. The region can be represented as a system of linear inequality constraints, then incorporated into the polynomial program. Our second approach uses an asymptotically valid confidence region based on the multivariate Gaussian limiting distribution of the Dirichlet (Bienaymé, 1838), which can be represented as a single convex quadratic inequality constraint. Figure 2 visualizes these regions for a simple two-node

Figure 2: **Polynomial confidence regions in a binary graph.** Panel (a) presents a causal graph in which binary  $X$  causes binary  $Y$ , but both are confounded by an unobserved  $U$ .  $N = 1,000$  observations are drawn from this DGP, producing an empirical distribution with proportions  $\frac{1}{N} \sum_{i=1}^N \mathbb{1}(X_i = x, Y_i = y)$ . Panels (b–c) depict confidence regions for  $\Pr(X_i = 0, Y_i = 0)$ ,  $\Pr(X_i = 0, Y_i = 1)$ , and  $\Pr(X_i = 1, Y_i = 0)$ ; the final category,  $\Pr(X_i = 1, Y_i = 1)$  (not depicted), must sum to unity. Panel (b) shows the Bernoulli-KL confidence region, which is conservative in finite samples and can be polynomialized as a set of linear inequalities. Panel (c) shows the Gaussian confidence region, which is asymptotically valid and can be polynomialized as a single convex quadratic inequality.



graph. Simulations reported in Section 8.2 evaluate coverage of the methods for various sample sizes. Appendix C provides details on the implementation of these methods. We also provide a method for polynomializing arbitrary confidence regions, allowing analysts to exploit tighter finite-sample confidence regions.

## 8 Simulated Examples

We now demonstrate our algorithm’s performance via simulations. Several examples correspond to known analytic solutions, offering further validation of our approach. Section 8.1 illustrates how Algorithms 1–2 allow analysts to iteratively state possible assumptions, test their observable implications, and use them to narrow causal bounds under noncompliance. Section 8.2 evaluates our proposals for statistical inference with estimated bounds. Section 8.3 examines several challenges—selection, mismeasurement, and missingness—that pose more complex threats to statistical inference. For clarity of exposition, all simulations use binary variables; our method adapts automatically to categorical variables.

### 8.1 Instrumental Variables

Noncompliance, or deviation between assigned ( $Z_i$ ) and realized ( $X_i$ ) treatment status, is a common obstacle to causal inference in randomized trials. Balke and Pearl (1997) showed that the task of bounding the ATE on an outcome  $Y_i$  in the presence of noncompliance

can be formulated as a linear programming problem, admitting a computational solution to partial identification. However, this approach cannot be extended to bound the local average treatment effect (LATE) among “compliers” that accept the assigned treatment—a principal effect that has received considerable attention—because this estimand corresponds to a nonlinear objective function. Angrist et al. (1996) shows the LATE can be point identified, but only if a number of conditions hold. These conditions include (i) ignorability of  $Z_i$ ; (ii) a non-null effect of  $Z_i$  on  $X_i$ ; (iii) an exclusion restriction, or the absence of a direct effect of  $Z_i$  on  $Y_i$ ; and (iv) monotonicity, or the absence of “defiers” that behave inversely to instructions. In this section, we estimate both the ATE and LATE in settings where assumptions i–ii are satisfied, then probe the implications of assumptions iii–iv. Our results show that while extant methods offer solutions for specific scenarios and estimands, even minor deviations from ideal conditions can render them inapplicable or inaccurate. Below, we show how our algorithm easily accommodates these variations and complications.

Figure 3: **DGPs with noncompliance.** The figure displays three possible causal models corresponding to scenarios in which an encouragement  $Z$  causes treatment  $X$ . Panel (b) corresponds to the true DAG in our simulated dataset, in which the monotonicity assumption is violated (indicated here by the absence of a + symbol) but other key identifying assumptions are satisfied. Panel (a) depicts a DAG assumed by an overcautious analyst that allows for violations of the exclusion restriction. Panel (c) depicts a model assumed by an overconfident analyst in which monotonicity of  $Z \rightarrow X$  is incorrectly invoked.

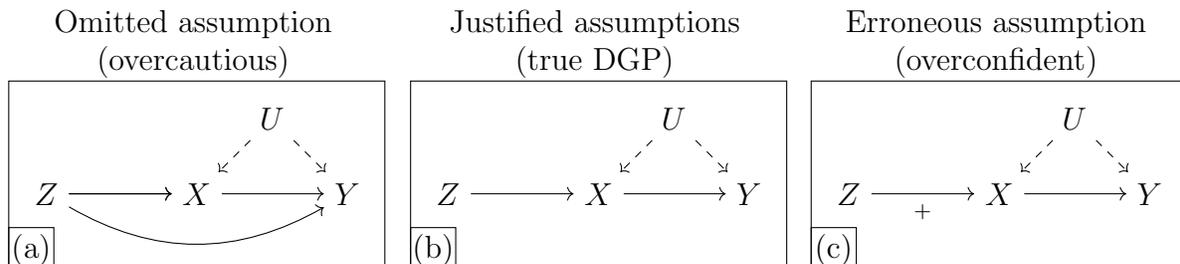


Figure 3 displays three possible DGPs that analysts might assume in a scenario involving noncompliance. We simulate data from the true DGP, shown in panel (b), in which all assumptions in Angrist et al. (1996) are satisfied except monotonicity of  $Z \rightarrow X$ . In this simulation, the true values of the ATE and LATE are  $-0.25$  and  $-0.36$ , respectively. In practice, analysts may proceed with an abundance of caution and make the conservative causal assumptions depicted in panel (a)—a challenging scenario in which a direct effect of the instrument on the outcome cannot be excluded and monotonicity is not assumed. Assuming model (a) and applying our algorithm yields sharp bounds of  $[-0.63, 0.37]$  and  $[-1, 1]$  for the ATE and LATE, respectively. While these bounds are relatively wide—the ATE cannot be signed, and the bounds for LATE are entirely uninformative—the resulting intervals do contain the true estimand values, and they represent the most precise statement possible under assumptions the analyst is willing to defend.

If the analyst was willing to assume the exclusion restriction, per model (b)—perhaps due to domain expertise or an experimental design that ruled out direct effects—our al-

gorithm would bound the ATE at  $[-0.55, -0.15]$ , revealing a negative effect and correctly containing the true value of  $-0.25$ . However, under these circumstances, the bounds on the LATE remain entirely uninformative at  $[-1, 1]$ . This reflects the fact that without strong assumptions, it is difficult to learn about cross-world quantities such as principal effects.

Finally, panel (c) shows a DGP imagined by an overconfident analyst, in which all four identifying assumptions in Angrist et al. (1996) are embraced. Unbeknownst to the analyst, the monotonicity assumption is in fact violated. Helpfully, when asked to estimate bounds, our Algorithm 2 reports that the causal query is *infeasible*. Recall that the true DGP corresponds to model (b), in which defiers are present; because the algorithm fails to locate any DGPs in which the observed information is consistent with the absence of defiers, it provides a clear warning to users that the assumption cannot be defended. However, if the analyst naïvely applied the traditional instrumental variables two-stage least squares estimator, they would not be alerted to this fact. Rather, they would obtain a point estimate of  $-0.74$ , roughly twice the true LATE. Put differently, *the standard IV approach ignores observable implications of underlying assumptions*. In contrast, our algorithm flags faulty theory by identifying infeasible scenarios, forestalling fruitless inquiry.

## 8.2 Coverage of Confidence Bounds

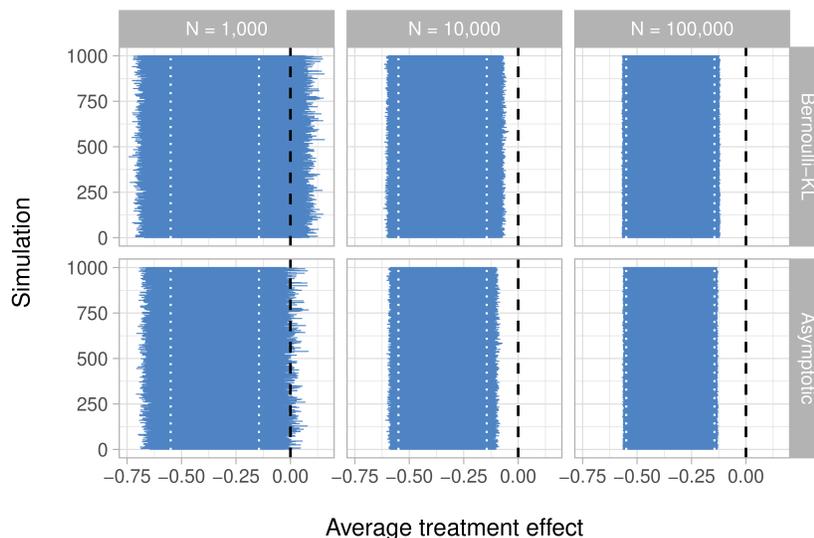
In applied settings, the bounds estimated by our algorithm will be subject to sampling error. We now evaluate the performance of confidence bounds that characterize this uncertainty, constructed according to Section 7, using the instrumental variable model of Figure 3(b). Specifically, we draw samples of  $N = 1,000$ ,  $N = 10,000$ , or  $N = 100,000$  observations from this DGP. For each sample, we then compute estimates of eight quantities:  $\Pr(Z_i = z, X_i = x, Y_i = y)$  for all  $x, y, z \in \{0, 1\}$ . These quantities form the basis of estimated bounds, by the plug-in principle. To quantify uncertainty, we compute 95% confidence regions on the same observed quantities, then convert them to polynomial constraints for inclusion in Algorithm 2. Optimizing subject to these confidence constraints produces confidence bounds, depicted in Figure 4. For each combination of sample size and uncertainty method, we draw 1,000 simulated datasets and run Algorithm 2 once.

Table 1 reports average values of estimated lower (upper) confidence bounds obtained by Algorithm 2 over 1,000 simulated datasets, for varying  $N$ . At all sample sizes, estimated bounds are centered on population bounds. Figure 2 shows confidence bounds obtained across methods and sample sizes. The Bernoulli-KL method produces wider confidence intervals at all  $N$ ; at  $N = 1,000$ , it is generally unable to reject zero, whereas the asymptotic method does so occasionally. Differences in interval width persist but shrink rapidly as sample size grows and both methods collapse on population bounds. As discussed in Section 7, we find more conservative coverage for confidence bounds on the ATE (100% coverage of population bounds), compared to coverage of the underlying confidence regions on the observed quantities (95% joint coverage of observed population quantities for the asymptotic method).

Table 1: **Bias of estimated bounds.** Average lower (upper) estimated bounds simulated datasets of varying size. Average estimated bounds correspond closely to population bounds.

Quantity	$N = 1,000$	$N = 10,000$	$N = 100,000$	Population
Lower bound	-0.549	-0.551	-0.551	-0.550
Upper bound	-0.144	-0.146	-0.146	-0.146

Figure 4: **Coverage of confidence bounds.** Each of 1,000 simulations is depicted with a horizontal line. For each simulation, a horizontal error bar represents a 95% confidence bound obtained per Section 7. All confidence bounds fully contain the population bounds, indicating 100% coverage. The upper (lower) row of panels reflect confidence bounds obtained with the Bernoulli-KL (asymptotic) method. Columns of panels report confidence bounds obtained using samples of various sizes. Vertical dotted white lines show true population lower and upper bounds, which contain the true ATE of  $-0.25$ ; vertical dashed black lines indicate zero.



### 8.3 More Complex Bounding Problems

We now examine four hypothetical DGPs, shown in Figure 5, featuring various threats to inference. Throughout, we target the ATE of  $X$  on  $Y$ . Panel (a) illustrates outcome-based selection: we observe unit  $i$  only if  $S_i = 1$ , where  $S_i$  may be affected by  $Y_i$ . Selection severity,  $\Pr(S_i = 0)$ , is known, but no information about  $\Pr(X_i = x, Y_i = y | S_i = 0)$  is available.  $X_i$  and  $Y_i$  are also confounded by unobserved  $U_i$ . Bounding in this setting is a nonlinear program, with an analytic solution recently derived in Gabriel et al. (2020). Panel (b) illustrates measurement error: an unobserved confounder  $U_i$  jointly causes  $Y_i$  and its proxy  $Y_i^*$ , but only treatment and the proxy outcome are observed. Bounding in this setting is a linear problem. A number of results for linear measurement error were recently presented in Finkelstein et al. (2020); here, we examine the monotonic errors case, where  $Y_i^*(Y_i = 1) \geq Y_i^*(Y_i = 0)$ . Panel (c) depicts missingness in outcomes, i.e. nonresponse

or attrition. Here,  $X_i$  affects both the partially observed  $Y_i$  and response indicator  $R_i$ ; if  $R_i = 1$ , then  $Y_i^* = Y_i$ , but if  $R_i = 0$ , then  $Y_i^*$  takes on the missing value indicator NA. Nonresponse on  $Y_i$  is differentially affected by both  $X_i$  and the value of  $Y_i$  itself (i.e. “missingness not at random,” MNAR); [Manski \(1990\)](#) provides analytic bounds. Finally, panel (d) depicts joint missingness in both treatment and outcome—sometimes a challenge in longitudinal studies with dropout—with MNAR on  $Y_i$ .

Figure 6(a–c) illustrates how Algorithm 2 recovers sharp bounds. Each panel shows progress in time, converging on known analytic results depicted at the right of each plot. Primal bounds (blue) widen over time as more extreme, observationally equivalent models are found. Dual bounds (red) narrow as the outer envelope is tightened. When a region cannot possibly produce a more extreme value than a previously discovered primal point, it is eliminated from consideration. Optimization proceeds by simultaneously searching for more extreme primal points and narrowing the dual envelope. Analysts can terminate the process at any time, reporting guaranteed-valid dual bounds along with their worst-case suboptimality factor,  $\varepsilon$ —or await complete sharpness,  $\varepsilon = 0$ .

Figure 5: **Various threats to inference.** Panels depict (a) outcome-based selection, (b) measurement error, (c) nonresponse and (d) joint missingness. In each graph,  $X$  and  $Y$  are treatment and outcome, respectively. Dotted red regions represent observed information. In (a), the box around  $S$  indicates selection: other variables are only observed conditional on  $S = 1$ . In (b),  $Y^*$  represents a mismeasured version of the unobserved true  $Y$ . In (c),  $R_Y$  indicates reporting, so that  $Y^* = Y$  if  $R = 1$  and is missing otherwise. In (d), both treatment and outcome can be missing; and missingness on  $X$  can affect missingness on  $Y$ .

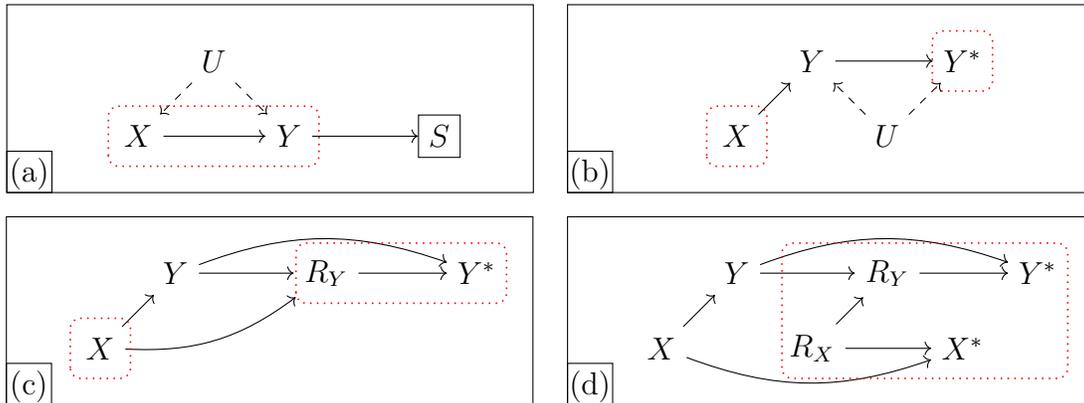
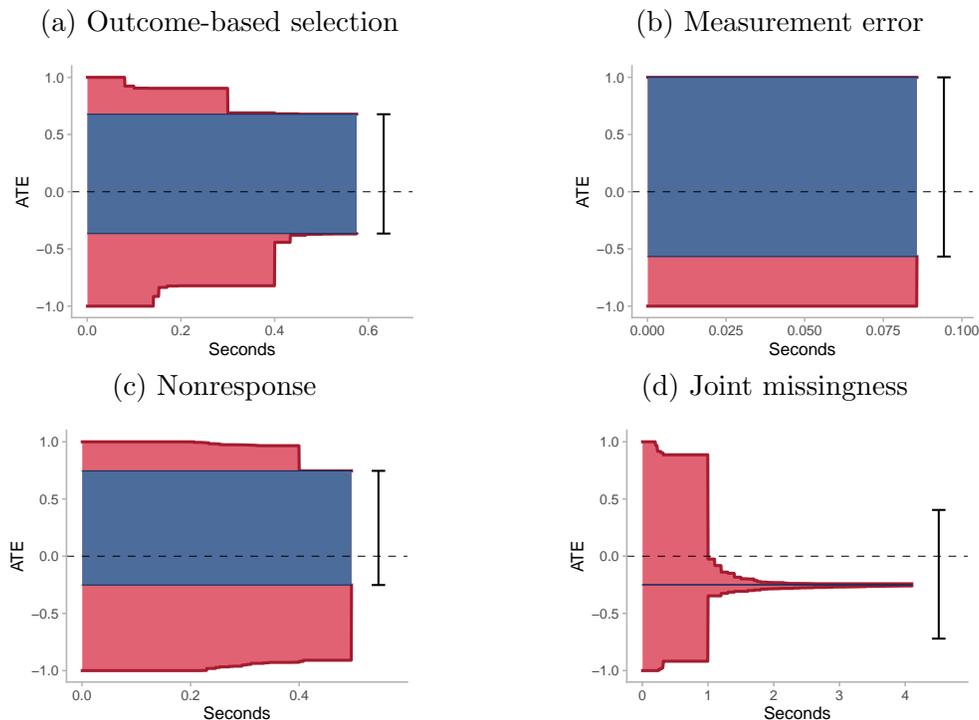


Figure 6: **Computation of ATE bounds.** Progress of Algorithm 2 in simulated Figure 5(a–d) DGPs. Black error bars are known analytic bounds,  $y$ -axes are ATE values, and  $x$ -axes are runtimes of Algorithm 2. Prior analytic bounds are sharp for settings (a–c). In setting (d), Algorithm 2 achieves point identification, but [Manski \(1990\)](#) bounds do not. Red regions are dual bounds, which always contain sharp bounds and the unknown true causal effect; these can only narrow over time, converging on optimality. Blue regions are primal bounds, which can only widen over time as more extreme models are found. Optimization stops when primal and dual bounds meet, indicating bounds are sharp.



In Figure 6(a–c), the algorithm converges on known analytic results. Ultimately, in the selection simulation (a), Algorithm 2 achieves bounds of  $[-0.37, 0.68]$ , correctly recovering Gabriel et al.’s (2020) bounds; in (b), measurement error bounds are  $[-0.57, 1.00]$ , matching Finkelstein et al. (2020); and in (c), outcome missingness bounds are  $[-0.25, 0.75]$ , equaling Manski (1990) bounds. Somewhat counterintuitively, Figure 6(d) shows dual bounds collapsing to a point, correctly point-identifying the ATE at  $-0.25$  despite severe missingness. This surprising result turns out to be a special case of an approach using “shadow variables” recently developed by Miao et al. (2015).<sup>12</sup> This example illustrates that the algorithm is general enough to recover results even when they are not widely known in a particular model; note that the commonly used approach of Manski (1990) produces far looser bounds of  $[-0.72, 0.40]$ , failing to exploit causal structure given in Figure 5(d). This result suggests our approach enables an empirical investigation of complex models where general identification results are not yet available. Situations where bounds converge suggest models where point identification via an explicit functional may be possible, potentially enabling new identification theory.

## 9 Potential Critiques of the Approach

Below, we briefly discuss several potential critiques of our method.

**“The user must know the true causal model.”**

Our algorithm requires users specify a causal graph and assumptions, but in many applications, the true DGP is unknown. This is precisely the obstacle that motivates our approach, which allows for valid inferences in the *absence* of complete information. Rather than assert a faulty “complete” model, the user need only input what they know or believe. The algorithm then outputs the most precise possible solution given that information; key assumptions can be relaxed further using easily incorporated sensitivity analyses, as needed. We note the difficulty of declaring a causal theory, even a partial one, is universal: any attempt to draw causal inferences from data—even in experimental settings—is premised (often implicitly) on underlying causal theory. Making assumptions explicit is not a trade-off relative to other methods, but a boon for research transparency.

**“The bounds may be too wide to be informative.”**

Yes.

When a point-identified solution exists, our algorithm will discover it. As Section 8.3 shows, this can occur in surprising scenarios and may help reveal new identification theory. However, when point-identification is impossible, our approach produces sharp bounds. These bounds may be insufficient for an analyst to achieve a goal such as discerning the sign of a causal effect. This is simply a fact about the limitations of the research design—as we prove, it is impossible to narrow the bounds further without additional information. Again, there is no tradeoff: incorrect point estimates based on faulty assumptions are also

---

<sup>12</sup>Specifically, it can be shown the ATE is identified for the Figure 5(d) graph only among faithful distributions where  $X \rightarrow Y$  is non-null—i.e. almost everywhere in the model space.

uninformative. When sharp bounds incorporating all defensible assumptions are wide, it means progress will require collecting more data or justifying additional assumptions.

**“What about continuous variables?”**

Our approach applies to discrete data, but analysis of continuous variables can often still proceed with some adjustments. Discrete approximations often suffice in applied work. (Indeed, “all data as observed are discrete,” [Rubin, 1981](#), p. 133). When continuous treatments (e.g. birth date, vehicle speed) often affect discrete outcomes (school admittance, police stops) only when exceeding a threshold, discretization is lossless. Moreover, when analyzing discrete treatments and continuous outcomes, much of our theory generalizes to estimands involving expectations of the outcome. Future work may study our method’s applicability to bounded continuous variables with smooth effects.

**“The bounds will take too long to compute.”**

Computation time for sharp bounds may sometimes be prohibitive, but our approach is likely still faster than manual derivation. Notably, the algorithm recovers several recently published analytic results in mere seconds ([Gabriel et al., 2020](#); [Miao et al., 2015](#); [Knox et al., 2020](#)). Second, when computation time is long, our algorithm’s “anytime” guarantee ensures premature termination will still produce valid bounds and report a worst-case looseness factor for the resulting non-sharp bounds.

## 10 Future Work with Automated Bounding

Causal inference is a central goal of science, and several established techniques can estimate causal quantities under ideal conditions. But in many applications, these conditions are simply not satisfied, and developing new analytic solutions is often intractable. For knowledge accumulation to proceed in the messy world of applied statistics, a general solution is needed. We present a tool to automatically produce sharp bounds on causal quantities in settings involving discrete data. Our approach involves a reduction of all such causal queries to polynomial programming problems, enables efficient search over observationally indistinguishable DGPs, and produces sharp bounds on arbitrary causal estimands. This approach is sufficiently general to accommodate a range of classic inferential obstacles.

Beyond providing a general tool for causal inference, our approach aligns closely with recent calls to improve research transparency by requiring the explicit declaration of estimands, identifying assumptions, and theory ([Miguel et al., 2014](#); [Lundberg et al., 2021](#)). With a common understanding of goals and premises, scholars can have meaningful debates over the credibility of research. When aspects of a theory are contested, our approach allows for a fully modular exploration of how assumptions shape empirical conclusions. Scholars can learn whether a particular assumption is empirically consequential, and if so, craft a targeted line of inquiry to probe its validity. Our approach can also act as a safeguard for analysts, flagging assumptions as infeasible when they conflict with observed information. This means hopeless research projects can be abandoned before wasting effort or disseminating untruths.

Future work should seek to reduce computation time for sharp bounds, especially when incorporating point-identified subquantities or additional semi-parametric modeling approaches. Causal inference scholars may also use this method as an exploratory tool to aid in the discovery of new identification theory. These lines of inquiry now represent the major open questions in discrete causal inference.

## References

- Angrist, J. D., G. W. Imbens, and D. B. Rubin (1996). Identification of causal effects using instrumental variables. *Journal of the American Statistical Association* 91(434), 444–455.
- Balke, A. and J. Pearl (1994). Counterfactual probabilities: Computational methods, bounds and applications. In *Uncertainty Proceedings 1994*, pp. 46–54. Elsevier.
- Balke, A. and J. Pearl (1997). Bounds on treatment effects from studies with imperfect compliance. *Journal of the American Statistical Association* 92(439), 1171–1176.
- Belotti, P., J. Lee, L. Liberti, F. Margot, and A. Wächter (2009). Branching and bounds tightening techniques for non-convex MINLP. *Optimization Methods and Software* 24(4–5), 597–634.
- Bienaymé, I. J. (1838). *Mémoire sur la probabilité des résultats moyens des observations: démonstration directe de la règle de Laplace*. Imprimerie Royale.
- Bonet, B. (2001). Instrumentality tests revisited. In J. S. Breese and D. Koller (Eds.), *UAI '01: Proceedings of the 17th Conference in Uncertainty in Artificial Intelligence*, pp. 48–55. Morgan Kaufmann.
- Cai, Z., M. Kuroki, J. Pearl, and J. Tian (2008). Bounds on direct effects in the presence of confounded intermediate variables. *Biometrics* 64(3), 695–701.
- Carathéodory, C. (1907, March). Über den variabilitätsbereich der koeffizienten von potenzreihen, die gegebene werte nicht annehmen. *Mathematische Annalen* 64(1), 95–115.
- Dean, T. L. and M. Boddy (1988). An analysis of time-dependent planning. pp. 49–54. American Association for Artificial Intelligence.
- Evans, R. (2018). Margins of discrete bayesian networks. *Annals of Statistics* 46(6A), 2623–2656.
- Evans, R. J. (2016). Graphs for margins of bayesian networks. *Scandinavian Journal of Statistics* 43(3), 625–648.
- Evans, R. J., T. S. Richardson, et al. (2019). Smooth, identifiable supermodels of discrete dag models with latent variables. *Bernoulli* 25(2), 848–876.
- Finkelstein, N., R. Adams, S. Saria, and I. Shpitser (2020). Partial identifiability in discrete data with measurement error. *arXiv preprint arXiv:2012.12449*.

- Finkelstein, N., E. Wolfe, and I. Shpitser (2021). Non-restrictive cardinalities and functional models for discrete latent variable dags. *Working Paper*.
- Frangakis, C. E. and D. B. Rubin (2002). Principal stratification in causal inference. *Biometrics* 58(1), 21–29.
- Gabriel, E. E., M. C. Sachs, and A. Sjölander (2020). Causal bounds for outcome-dependent sampling in observational studies. *Journal of the American Statistical Association*. DOI: 10.1080/01621459.2020.1832502.
- Gamrath, G., D. Anderson, K. Bestuzheva, W.-K. Chen, L. Eifler, M. Gasse, P. Gemander, A. Gleixner, L. Gottwald, K. Halbig, et al. (2020). The scip optimization suite 7.0.
- Geiger, D. and C. Meek (1999, August). Quantifier elimination for statistical problems. In *Proceedings of Fifteenth Conference on Uncertainty in Artificial Intelligence, Stockholm, Sweden* (Proceedings of Fifteenth Conference on Uncertainty in Artificial Intelligence, Stockholm, Sweden ed.), pp. 226–235.
- Greenland, S. and J. Robins (1986). Identifiability, exchangeability, and epidemiological confounding. *International Journal of Epidemiology* 15, 413–419.
- Heckman, J. and E. Vytlacil (2001). *Instrumental variables, selection models, and tight bounds on the average treatment effect*, pp. 1–15. Physica.
- Kennedy, E. H., S. Harris, and L. J. Keele (2019). Survivor-complier effects in the presence of selection on treatment, with application to a study of prompt icu admission. *Journal of the American Statistical Association* 114(525), 93–104.
- Knox, D., W. Lowe, and J. Mummolo (2020). Administrative records mask racially biased policing. *American Political Science Review* 114, 619–637.
- Kuchibhotla, A. K., S. Balakrishnan, and L. Wasserman (2021). The hulc: Confidence regions from convex hulls.
- Lee, D. (2009). Training, wages, and sample selection: Estimating sharp bounds on treatment effects. *The Review of Economic Studies* 76(3), 1071–1102.
- Li, A. and J. Pearl (2021). Bounds on causal effects and application to high dimensional data. arXiv preprint arXiv:2106.12121.
- Lundberg, I., R. Johnson, and B. M. Stewart (2021). What is your estimand? defining the target quantity connects statistical evidence to theory. *American Sociological Review* 86(3), 532–565.
- Malloy, M. L., A. Tripathy, and R. D. Nowak (2020). Optimal confidence regions for the multinomial parameter. *arXiv preprint arXiv:2002.01044*.
- Manski, C. (1990). Nonparametric bounds on treatment effects. *The American Economic Review* 80(2), 319–323.

- Miao, W., L. Liu, E. T. Tchetgen, and Z. Geng (2015). Identification, doubly robust estimation, and semiparametric efficiency theory of nonignorable missing data with a shadow variable. *arXiv preprint arXiv:1509.02556*.
- Miguel, E., C. Camerer, K. Casey, J. Cohen, K. Esterling, A. Gerber, R. Glennerster, D. Green, M. Humphreys, G. Imbens, and D. Laitin (2014). Promoting transparency in social science research. *Science* 343(6166), 30–31.
- Molinari, F. (2020). Microeconometrics with partial identification. [arXiv:2004.11751](https://arxiv.org/abs/2004.11751).
- Ottmann, T., S. Schuierer, and S. Soundaralakshmi (1995). Enumerating extreme points in higher dimensions. In *Annual Symposium on Theoretical Aspects of Computer Science*, pp. 562–570. Springer.
- Pearl, J. (1995). On the testability of causal models with latent and instrumental variables. *Uncertainty in Artificial Intelligence II. San Francisco, CA: Morgan Kaufmann Publishers*.
- Pearl, J. (2009). *Causality*. New York: Cambridge University Press.
- Ramsahai, R. R. (2012). Causal bounds and observable constraints for non-deterministic models. *Journal of Machine Learning Research* 13(3), 829–848.
- Richardson, T. (2003). Markov properties for acyclic directed mixed graphs. *Scandinavian Journal of Statistics* 30(1), 145–157.
- Richardson, T. S., R. J. Evans, J. M. Robins, and I. Shpitser (2017). Nested Markov properties for acyclic directed mixed graphs. Working paper.
- Richardson, T. S. and J. M. Robins (2013). Single world intervention graphs (swigs) : A unification of the counterfactual and graphical approaches to causality. *Working Paper, Center for Stat. & Soc. Sci., U. Washington* 128(30).
- Robins, J. (1989). The analysis of randomized and non-randomized aids treatment trials using a new approach to causal inference in longitudinal studies. *Health service research methodology: a focus on AIDS*, 113—159.
- Rubin, D. B. (1981). The Bayesian Bootstrap. *The Annals of Statistics* 9(1), 130 – 134.
- Sachs, M., E. Gabriel, and A. Sjölander (2020). Symbolic computation of tight causal bounds.
- Shpitser, I. (2018). Identification in graphical causal models. In M. Maathuis, M. Drton, S. Lauritzen, and M. Wainwright (Eds.), *Handbook of Graphical Models*. CRC Press.
- Sjölander, A., W. Lee, H. Källberg, and Y. Pawitan (2014). Bounds on causal interactions for binary outcomes. *Biometrics* 70(3), 500–505.
- Swanson, S. A., M. A. Hernán, M. Miller, J. M. Robins, and T. S. Richardson (2018). Partial identification of the average treatment effect using instrumental variables: Review of methods for binary instruments, treatments, and outcomes. *Journal of the American Statistical Association* 113(522), 933–947. DOI: 10.1080/01621459.2018.1434530.

- Tian, J. and J. Pearl (2002). On the testable implications of causal models with hidden variables. In *UAI '02: Proceedings of the 18th Conference in Uncertainty in Artificial Intelligence*, pp. 519–527.
- Tukey, J. (1986). Sunset salvo. *The American Statistician* 40(1), 72–76.
- Verma, T. and J. Pearl (1990). Equivalence and synthesis of causal models. In P. P. Bonissone, M. Henrion, L. N. Kanal, and J. F. Lemmer (Eds.), *Proc. of the Conf. on Uncertainty in Artificial Intelligence*, pp. 255–268. Morgan Kaufmann.
- Vigerske, S. and A. Gleixner (2018). Scip: Global optimization of mixed-integer nonlinear programs in a branch-and-cut framework. *Optimization Methods and Software* 33(3), 563–593.
- Wächter, A. and L. T. Biegler (2006). On the implementation of an interior-point filter line-search algorithm for large-scale nonlinear programming. *Mathematical programming* 106(1), 25–57.
- Wolfe, E., R. W. Spekkens, and T. Fritz (2019). The inflation technique for causal inference with latent variables. *Journal of Causal Inference* 7(2).
- Zhang, J. and E. Bareinboim (2021, Feb). Non-parametric methods for partial identification of causal effects. Technical Report R-72, Causal Artificial Intelligence Lab, Columbia University.
- Zhang, J. L. and D. B. Rubin (2003). Estimation of causal effects via principal stratification when some outcomes are truncated by “death”. *Journal of Educational and Behavioral Statistics* 28(4), 353–368.

# A Examples, Algorithms, and Detailed Discussion

## A.1 Canonicalization of DAGs

In this appendix, we summarize the process for obtaining a canonical hidden variable DAG, presented as Definition 4.6 in [Evans \(2016\)](#). Theorem 4.13 in [Evans \(2016\)](#) shows that the marginal model of any hidden variable DAG is the same as that of its canonical hidden variable DAG, and Proposition 7.4 of the same work shows that the same holds for the model for post-intervention distributions, when interventions are restricted to the main variables.

Given a hidden variable DAG  $\mathcal{G}$ , the canonical form of the DAG is constructed by the following procedure:

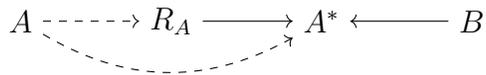
1. Add an edge  $X_j \rightarrow X_{j'}$  for any pair of variables  $X_j, X_{j'}$  such that there is a path from  $X_j$  to  $X_{j'}$  along which all variables between  $X_j$  and  $X_{j'}$  are hidden.  $X_j$  and  $X_{j'}$  can each be hidden or observed.
2. Remove incoming edges to hidden variables.
3. Remove hidden variables whose children are a subset of the children of another hidden variable.

By construction, all latent variables in the canonical DAG will be exogenous.

## A.2 Functional Models in the Context of Determinism

The general approach for obtaining functional models for discrete hidden variable DAGs ([Evans, 2018](#); [Finkelstein et al., 2021](#)) does not take account of the kind of determinism introduced into the model by missingness indicators, and as such may yield a functional model that is *over-parameterized*. Due to the complexity of polynomial programming, it is beneficial to avoid excess parameters where possible. We now briefly explore this issue.

Figure 7: **A graph with determinism.**



Consider the scenario depicted in Figure 7. In this graph,  $A^*$  is a proxy for the unobserved variable  $A$ , which is observed with missingness as indicated by  $R_A$ . When  $R_A = 0$ , then  $A^*$  is deterministically equal to a special value indicating missingness (usually denoted with the special value such as “?” or “NA”). In addition,  $A^*$  is affected by  $B$ . This scenario might arise if  $A$  is measured with missingness *and* measurement error, and the nature of the error is affected by  $B$ . Of note,  $A^*$  is not a fully deterministic function of  $A$  and  $R_A$ , and cannot simply be removed from the functional parameterization, as in traditional missingness without measurement error. However, we can use the fact that it is a *partially* deterministic function of  $R_A$  to reduce the number of parameters needed in the functional model for this graph.

In general, the functional model for this graph would allocate one value of  $\epsilon_{A^*}$ —the exogenous noise that determines  $A^*$  in terms of its parents—for every combination of possible responses of  $A^*$  to its parents. Suppose  $A^*$  takes values in  $\{0, 1, ?\}$ , and  $A$ ,  $R_A$  and  $B$  take values in  $\{0, 1\}$ . This would correspond to  $3^8 = 6561$  possible values of  $\epsilon_{A^*}$ . However, any such value that maps  $R_A = 0$  to  $A^* \in \{0, 1\}$  or  $R_A = 1$  to  $A^* = ?$  is ruled out by the deterministic relationship. As a result,  $\epsilon_{A^*}$  need only specify the response of  $A^*$  in  $\{0, 1\}$  to  $A$  and  $B$  when  $R_A = 1$ . This yields only  $2^4 = 16$  possible values for  $\epsilon_{A^*}$ . This example demonstrates that incorporating known deterministic relationships can yield a non-restrictive functional parameterization with fewer parameters.

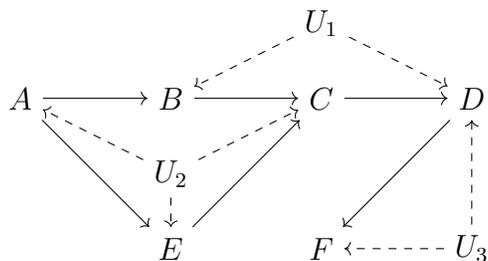
### A.3 DAG Parameterization for Non-g geared Graphs

Most graphs we encounter in practice are *geared* (Evans, 2018), which means they have no *non-trivial bi-directed cycles* Finkelstein et al. (2021). When graphs are not geared, and the target estimand as well as all empirical evidence involves only single world probabilities,

it is possible to improve the complexity of the system. Under these circumstances, it is preferable to obtain non-restrictive bounds on the cardinalities of latent variables according to [Finkelstein et al. \(2021\)](#). All single world probabilities can be expressed in terms of the usual DAG parameters according to the g-formula, and therefore all functionals of such probabilities described in [Corollary 1](#) can be polynomialized as well. If the target or any of the empirical evidence involve cross-world probabilities, we must revert to the functional model approach.

## A.4 Example of Program Simplification

Figure 8: **A graph with conditional independence and Verma constraints.**



Consider the graph presented in [Figure 8](#). We will use this graph to illustrate a number of points raised in the main body of the paper. Suppose we are interested in the ATE of  $E$  on  $C$ . First, we will explicitly construct the functional model of this graph, then use it to generate a simple polynomial program that bounds a causal target. Next, we will employ several of the strategies described in [Section 5](#) to simplify the program, demonstrating the importance of these strategies in obtaining tractable program formulations. Finally, we will observe that a broader class of partial identification problems than previously recognized can be formulated as linear programs.

Suppose all observed variables in the graph above are binary. In constructing a functional model, we first note that  $U_2$  is responsible for determining the values of  $A$ ,  $C$  and  $E$  in response to their parents.  $A$  has no parents,  $E$  has one parent, and  $C$  has two parents. Therefore  $U_2$  takes values in a state space of size  $2^1 \times 2^2 \times 2^4 = 128$ . Next, we suppose

$U_1$  is responsible for determining the value of  $B$  in response to  $A$ , and therefore has size  $2^2 = 4$ .  $U_3$  is left to determine the value of  $F$  in response to  $D$ , and of  $D$  in response to  $U_1$  and  $C$ . It therefore takes values in space of size  $2^8 \times 2^2 = 1024$ .<sup>13</sup>

To construct the polynomial program, we begin with the non-negativity and linear marginalization constraints on the parameters of the distributions of the disturbances (for simplicity, we abstain from eliminating one parameter per distribution using the sum-to-unity constraint):

$$\begin{aligned} \Pr(U_i = u) &\geq 0 && \forall i, \forall u \in \Omega_{U_i} \\ \sum_{u \in \Omega_{U_i}} \Pr(U_i = u) &= 1 && \forall i. \end{aligned}$$

We then add constraints encoding the empirical evidence  $\mathcal{E}$ . For simplicity, we assume that we observe the full joint distribution  $\Pr(A = a, B = b, C = c, D = d, E = e, F = f)$ , which is a vector of size  $2^6 = 64$ , corresponding to 64 equality constraints in the program. There are 3 disturbance variables in this graph, including  $\epsilon_E$ , leading to polynomials in these equality constraints with terms of degree 3. Given the cardinalities of the disturbances, there are  $2^4 \times 2^7 \times 2^{10} = 2,097,152$  possible combinations of disturbance assignments. By a simple exchangeability argument, the same number of possible combinations lead to each outcome in the state space. As there are  $2^6$  outcomes, each of the 64 polynomial equality constraints for  $\mathcal{E}$  will have  $\frac{2^{21}}{2^6} = 2^{15}$  terms, again each of degree 3. This is a very large program.

$$\Pr(A = a, B = b, C = c, D = d, E = e, F = f) = \sum_{u \in \Omega_U} \prod_i \Pr(U_i = u) \mathbb{1}(u \implies a, b, c, d, e, f) \tag{1}$$

We now consider the strategies described in Section 5. First, observe that there are

---

<sup>13</sup>It is also possible to construct a functional model by first taking  $U_3$  to be responsible for determining  $F$  in response to  $D$ , and then  $U_1$  to be responsible for determining  $B$  in response to  $A$  and  $D$  in response to  $C$  and  $U_3$ . By a simple symmetry argument, the two functional models yield the same number of parameters.

only 31 nested Markov parameters for this graph, corresponding to 31 polynomial equality constraints encoding  $\mathcal{E}$ : a substantial savings over the 64 parameters of the naïve parameterization. This reduced parameterization is possible because it encodes standard conditional independences, such as  $F \perp A \mid D$ . In addition, it encodes Verma constraints, which emerge either (i) from independences in post-intervention distributions or (ii) from the irrelevance of an intervention to a particular distribution. In this case,  $A \perp \{D, F\} \mid do(C)$ . As discussed in the main text, each equality constraint can be used to reduce the number of parameters needed in a non-restrictive reduction that can express every possible distribution in the model.

Recall that each nested Markov parameter corresponds to the identified probability of a single world event, where the event is specified in terms of variables in a single district, and the intervention is on all parents of the district relevant to those variables. For example, in this case, one of the nested Markov parameters is  $\Pr(b = 1, f = 1 \mid d = 1, do(a = 1, c = 1))$ . We can now make use of Proposition 3 to reason that each of these polynomial constraints must involve only disturbances from a single district. Therefore in the equations corresponding to nested Markov parameters for the district corresponding to  $U_2$ , parameters of the distributions of  $U_1$  and  $U_3$  will all sum out, and we will be left with equations that are linear in the parameters of  $U_2$ . Likewise, in equations corresponding to nested Markov parameters for the district containing descendants of  $U_1$  and  $U_3$ , parameters for the distribution of  $U_2$  will factor out, and we will be left with a quadratic equation.

Finally, we can make use of Proposition 4 to note that constraints involving nested Markov parameters corresponding to the  $\{U_1, U_3\}$  district can be dropped from the program. This is because they only involve parameters for the distributions of  $U_1$  and  $U_3$ , which do not appear in any constraint involving parameters for the distribution of  $U_2$ . The target, by contrast, involves only parameters for the distribution of  $U_2$ .

As a result of taking the three steps described in Section 5, we have taken this problem from a *polynomial* program involving 1156 parameters to a *linear* program involving only  $2^7 = 128$  parameters and fewer constraints. This example also motivates the following

corollary, which expands the class of partial identification problems that can be formulated as linear programs relative to known results (Balke and Pearl, 1997; Finkelstein et al., 2020; Wolfe et al., 2019).

**Corollary 2.** *Suppose  $\mathcal{G}$  is a hidden variable DAG with observed variables  $\mathbf{V}$ ,  $\mathcal{C} = \{V_\ell(\mathbf{a}_\ell) = v_\ell \mid \ell \in \mathcal{L}\}$  is a set of counterfactual statements, and  $\Pr(\mathcal{C})$  is the target of interest. Further suppose that the full joint distribution  $\Pr(\mathbf{V} = \mathbf{v})$  is observed. Then  $\Pr(\mathcal{C})$  can be sharply bounded given the observed data by optimizing a linear program if all  $\{V_\ell \mid \ell \in \mathcal{L}\}$  are in the same single-latent-variable district.*

*Proof.* Because the common district of  $\mathcal{C}$  contains only a single latent variable, by Proposition 3 the objective will be linear in the parameters of the distribution of that latent variable. By Proposition 4, the constraints will not involve parameters corresponding to other districts. By Algorithm 1, no single term in a constraint will involve multiple parameters for the same latent distribution, meaning that all constraints involving only parameters corresponding to a single-variable district will be linear. The non-negativity and sum-to-unity constraints on the parameters of the latent-variable distribution are also linear. It follows that the objective and all constraints are linear.  $\square$

## A.5 Constructing the Polynomial Program

Algorithm 1 constructs a polynomial program to sharply bound any factual or counterfactual target of inference,  $\mathcal{T}$ , that is a polynomial fraction or monotonic transformation thereof. In addition to  $\mathcal{T}$ , the algorithm takes as input a possibly non-canonical DAG  $\mathcal{G}$ ; empirical evidence  $\mathcal{E}$ , modeling assumptions  $\mathcal{A}$ , and sample space of possible outcomes for the main variables,  $\mathcal{S}(\mathbf{V})$ . It produces an optimization problem with a polynomial objective subject to polynomial constraints. This polynomial programming problem is equivalent to the original causal bounding problem.

---

**Algorithm 1** Constructing a Polynomial Program

---

**Input:** graph  $\mathcal{G}$ , evidence  $\mathcal{E}$ , assumptions  $\mathcal{A}$ , sample space  $\mathcal{S}(\mathbf{V})$ , target  $\mathcal{T}$

**Output:** polynomial program in parameters  $\mathcal{P}_{\mathcal{U}}$  or  $\mathcal{P}_{\mathcal{U}} \cup s$

*Initialization*

- 1: initialize empty constraint set  $\mathcal{C} \leftarrow \emptyset$
- 2:  $\mathcal{G} \leftarrow$  canonicalize  $\mathcal{G}$
- 3:  $\mathcal{P}_{\mathcal{U}} \leftarrow$  parameters of functional model for  $\mathcal{G}$

*Polynomialize objective function*

- 4:  $\mathcal{T} \leftarrow$  polynomial-fractionalize( $\mathcal{T}$ )
- 5: **if**  $\mathcal{T}$  contains fractions **then**
- 6:     polynomialize( $\mathcal{T} = s$ ) and append to  $\mathcal{C}$
- 7:      $\mathcal{T} \leftarrow s$
- 8: **end if**

*Polynomialize constraints*

- 9: **for**  $(g(\mathcal{P}_{\mathbf{V}}) \star \alpha) \in (\mathcal{E} \cup \mathcal{A})$  **do**
- 10:     polynomialize( $g(\mathcal{P}_{\mathbf{V}}) \star \alpha$ ) and append to  $\mathcal{C}$
- 11: **end for**
- 12: **for**  $U_{i,k} \in U_i$  **do**
- 13:     append ( $\mathcal{P}_{U_k}$  is a distribution) to  $\mathcal{C}$
- 14: **end for**

*Optimize*

- 15: **return** optimize  $\mathcal{T}$  subject to  $\mathcal{C}$
- 

## A.6 Optimizing the Polynomial Program

Algorithm 2 provides a step-by-step description of the  $\varepsilon$ -sharp bounding procedure. For ease of reference, we duplicate the Section 6 discussion of the algorithm’s various components here.

Algorithm 2 takes as inputs the polynomialized objective function  $\mathcal{T}(\mathbf{p})$  and constraint set  $\mathcal{C}(\mathbf{p})$ , obtained from Algorithm 1. It then evaluates a range of models, or points  $\mathbf{p}$  in the model space  $\mathcal{P}$  for which  $\mathcal{C}(\mathbf{p})$  is satisfied. It seeks to identify extreme values of  $\mathcal{T}(\mathbf{p})$  within this subspace. It also accepts two parameters:  $\epsilon^{\text{thresh}}$ , a stopping threshold for the looseness factor stopping, and  $\theta^{\text{thresh}}$ , a stopping threshold for width of the bounds. The algorithm returns two types of information: the upper and lower bounds for the causal program, and the worst-case looseness factor  $\varepsilon$ .

Primal bounds are denoted  $\underline{P}$  and  $\overline{P}$ , adopting the convention that underlines refer to objects used for minimization and overlines for maximization. These indicate the extreme

values of the target estimand in any admissible model—that is, satisfying  $\mathcal{C}(\mathbf{p})$ —that has been located so far. These are initialized at  $+\infty$  and  $-\infty$ , respectively, indicating that no admissible models have been found yet. As optimization proceeds, the primal bounds improve as new, more extreme admissible models are found. We refer to  $[\underline{P}, \overline{P}]$  as the *inner bounds*: the unknown sharp bounds must at least contain these points, which correspond to models that are observationally indistinguishable from the true DGP.

Dual optimization begins by partitioning the parameter space into branches, proceeding separately for the lower and upper bound and respectively producing partitions  $\underline{\mathcal{B}}_b$  and  $\overline{\mathcal{B}}_b$ . At initialization, these consist of a single branch spanning the entire parameter space; each branch is then recursively divided. The lower and upper parts of the dual envelope, or outer envelope, are denoted  $\underline{\mathcal{D}}$  and  $\overline{\mathcal{D}}$ . These are piecewise linear functions, with pieces corresponding to the branching partitions, that are *relaxations* of the true objective function,  $\mathcal{T}(\mathbf{p})$ , from below and above. These relaxations are made to ensure they will always contain the entire objective function at all points in the parameter space. Within branch  $b$ , the value  $\min\{\underline{\mathcal{D}}_b(\mathbf{p}) : \mathbf{p} \in \overline{\mathcal{B}}_b\}$  indicates the lowest value attained by the lower envelope; thus,  $\underline{\mathcal{I}} = \min_b \{\min\{\underline{\mathcal{D}}_b(\mathbf{p}) : \mathbf{p} \in \overline{\mathcal{B}}_b\}\}$  represents the lowest value attained by the lower envelope anywhere in the parameter space. Conversely,  $\overline{\mathcal{T}} = \max_b \{\max\{\overline{\mathcal{D}}_b(\mathbf{p}) : \mathbf{p} \in \overline{\mathcal{B}}_b\}\}$  represents the highest value of the upper envelope. These extreme points on the dual envelope,  $[\underline{\mathcal{I}}, \overline{\mathcal{T}}]$ , define the dual (outer) bounds. These are the reported causal bounds; whatever the true sharp bounds, they must lie inside the dual bounds, even if the algorithm has not run to completion. We let  $\theta$  equal the bound width, or the difference between the upper and lower dual bounds, and we define the worst-case looseness factor  $\varepsilon$  as the slack (the difference in dual and primal bound widths) divided by the primal bound width.

The algorithm heuristically selects branches in the model space that appear promising, and refines primal and dual bounds in turn. It first searches within the branch for an admissible model; if found, and if the associated causal estimand is more extreme than those previously encountered, it is stored as a new primal bound. Whatever the true nonparametric sharp bounds, they must lie outside the primal bounds because the true

bounds must contain the extreme models that define the primal bounds. Then, it divides the branch into sub-branches and refines the dual envelope by tightening the piecewise linear outer-approximation. The algorithm continuously prunes branches of  $\underline{\mathcal{B}}_b$  and  $\overline{\mathcal{B}}_b$  that are inconsistent with specified constraints; it also continuously branches and refines the bounds while  $\theta$  and  $\varepsilon$  exceed specified thresholds.

---

**Algorithm 2** Computing  $\varepsilon$ -sharp Bounds

---

**Input:** target  $\mathcal{T}(\mathbf{p})$  and constraint relations  $\mathcal{C}(\mathbf{p})$  in parameters  $\mathbf{p}$ ,  
stopping thresholds  $\varepsilon^{\text{thresh}}$  and  $\theta^{\text{thresh}}$   
**Output:** lower bound  $\underline{\mathcal{T}}$ , upper bound  $\overline{\mathcal{T}}$ , maximum looseness factor  $\varepsilon$

*Initialization*

- 1: branches of parameter space: indexed partitions  $\underline{\mathcal{B}} \leftarrow \{[0, 1]^{\#\{\mathbf{p}\}}\}$ ,  $\overline{\mathcal{B}} \leftarrow \{[0, 1]^{\#\{\mathbf{p}\}}\}$
- 2: dual (outer) bounds: indexed families of functions  $\underline{\mathcal{D}} \leftarrow \{\mathbf{p} \mapsto -\infty\}$ ,  $\overline{\mathcal{D}} \leftarrow \{\mathbf{p} \mapsto +\infty\}$
- 3: primal (inner) bounds:  $\underline{P} = +\infty$  and  $\overline{P} = -\infty$
- 4: bounds width:  $\theta = +\infty$
- 5: bounds looseness factor  $\varepsilon = +\infty$

*Spatial branch and bound*

- 6: **while**  $\varepsilon > \varepsilon^{\text{thresh}}$  **and**  $\theta > \theta^{\text{thresh}}$  **do**
- 7:     **for** extremum in min, max **do**

*Select direction*

- 8:     **if** extremum is min **then**
- 9:         set  $\star \leftarrow \_$  and  $\blackstar \leftarrow \leq$
- 10:     **else if** extremum is max **then**
- 11:         set  $\star \leftarrow \_$  and  $\blackstar \leftarrow \geq$
- 12:     **end if**

*Primal refinement*

- 13:     continue search for local extremum of  $\mathcal{T}(\mathbf{p})$  s.t.  $\mathcal{C}(\mathbf{p})$  is satisfied
- 14:     **if** feasible point is found **and**  $\mathcal{T}(\mathbf{p}) \blackstar P^*$  **then**
- 15:         update primal bound  $P^* \leftarrow \mathcal{T}(\mathbf{p})$
- 16:     **end if**

*Dual refinement*

- 17:     select outermost branch  $b = \arg \text{extremum}_{b'} \{ \text{extremum} \{ \mathcal{D}_{b'}^*(\mathbf{p}) : \mathbf{p} \in \mathcal{B}_{b'}^* \} \}$
- 18:     pop  $\mathcal{B}_b^*$  from  $\mathcal{B}^*$  and subpartition it, pop  $\mathcal{D}_b^*(\mathbf{p})$  from  $\mathcal{D}^*$
- 19:     **for each** subpartition  $\mathcal{B}_{b'}^*$  **in**  $\mathcal{B}_b^*$  **do**
- 20:         push new branch  $\mathcal{B}_{b'}^*$  into  $\mathcal{B}^*$
- 21:         find linear function  $\mathcal{D}_{b'}^*$  s.t.  $\mathcal{D}_{b'}^*(\mathbf{p}) \blackstar \mathcal{T}(\mathbf{p})$  for all  $\mathbf{p} \in \mathcal{B}_{b'}^*$
- 22:         push linear programming relaxation  $\mathcal{D}_{b'}^*$  into  $\underline{\mathcal{D}}$
- 23:     **end for**

*Prune branches that cannot widen bounds*

- 24:     **for each**  $b$  **in**  $1, \dots, |\mathcal{B}^*|$  **do**
- 25:         **if**  $P^* \blackstar \text{extremum} \{ \mathcal{D}_b^*(\mathbf{p}) : \mathbf{p} \in \mathcal{B}_b^* \}$  **or**  $\mathcal{C}(\mathbf{p}) = \text{False}$  for all  $\mathbf{p} \in \mathcal{B}_b^*$  **then**
- 26:             pop  $\mathcal{B}_b^*$  from  $\mathcal{B}^*$ , pop  $\mathcal{D}_b^*$  from  $\mathcal{D}^*$
- 27:         **end if**
- 28:     **end for**

- 29:     **end for**

*Check progress*

- 30:      $\underline{\mathcal{T}} \leftarrow \min_b \{ \min \{ \underline{\mathcal{D}}_b(\mathbf{p}) : \mathbf{p} \in \underline{\mathcal{B}}_b \} \}$ ,  $\overline{\mathcal{T}} \leftarrow \max_b \{ \max \{ \overline{\mathcal{D}}_b(\mathbf{p}) : \mathbf{p} \in \overline{\mathcal{B}}_b \} \}$
- 31:      $\theta \leftarrow \overline{\mathcal{T}} - \underline{\mathcal{T}}$

32:  $\varepsilon \leftarrow \theta / (\bar{P} - \underline{P}) - 1$   
 33: **end while**  
 34: **return**  $\underline{\mathcal{I}}, \bar{\mathcal{T}}, \varepsilon$

---

## B Proofs

### Proof of Proposition 1.

*Proof.* We adapt the proof of Finkelstein et al. (2021) to account for counterfactuals as follows. First, we define *one-step-ahead* counterfactuals,  $V_{i,j}(\mathbf{pa}(V_{i,j}) = \mathbf{a})$ , to be those where all main parents of a variable are subject to intervention  $\mathbf{pa}(V_{i,j}) = \mathbf{a}$ . Next, we note that all other counterfactuals and factials in the full data law are deterministic functions of one-step-ahead variables, after fixing  $\mathbf{U}_i$ . Therefore it is sufficient to reason about only one-step-ahead variables; intervention on other variables is irrelevant to the full data law.

Because the likelihoods of multi-district graphs factorize as the likelihoods of the districts after intervention on their parents (Richardson et al., 2017), we can consider single-district graphs without loss of generality. In multi-district graphs, the bound obtained below can be applied within each district.

Each main variable  $V_{i,j}$  has  $|\mathcal{S}(\mathbf{pa}(V_{i,j}))|$  one-step-ahead counterfactuals, corresponding to possible manipulations of its parents. Each one-step-ahead counterfactual  $V_{i,j}(\mathbf{pa}(V_{i,j}) = \mathbf{a})$  has a cardinality equal to those of the corresponding main variable  $|\mathcal{S}(V_{i,j})|$ . Therefore, the collection of a single variable's one-step-ahead counterfactuals  $\{V_{i,j}(\mathbf{pa}(V_{i,j}) = \mathbf{a}), V_{i,j}(\mathbf{pa}(V_{i,j}) = \mathbf{a}'), \dots\}$ , can take on  $|\mathcal{S}(V_{i,j})|^{|\mathcal{S}(\mathbf{pa}(V_{i,j}))|}$  possible values, and there are  $d \equiv \prod_{V_{i,j} \in \mathbf{V}_i} |\mathcal{S}(V_{i,j})|^{|\mathcal{S}(\mathbf{pa}(V_{i,j}))|}$  values that the full collection of all one-step-ahead variables can take. Any model over this full collection must be a subset of the  $d - 1$  simplex. We let  $\mathbf{V}(\mathbf{pa}(\mathbf{V}))$  denote the collection of one-step-ahead variables.

Suppose the disturbances  $\mathbf{U}_i$  are enumerated as  $\{U_{i,1}, \dots, U_{i,K}\}$ . We will now show that each  $U_{i,k}$  can be assumed to be discrete without altering the model for  $\mathbf{V}(\mathbf{pa}(\mathbf{V}))$  and therefore the full data law. First, for each value  $u_k$  in the domain of  $U_{i,k}$ , we define the distribution  $P_{u_k}(\mathbf{V}(\mathbf{pa}(\mathbf{V}))) = \int_{\mathbf{u}_{\setminus k}} P(\mathbf{V}(\mathbf{pa}(\mathbf{V})) \mid \mathbf{u}_{\setminus k}, u_k) P(\mathbf{u}_{\setminus k})$ , where  $\mathbf{u}_{\setminus k}$  denotes

all disturbances other than  $u_k$ . This fixes  $U_{i,k}$  at the value  $u_k$ , modifying the distribution over  $\mathbf{V}(\mathbf{pa}(\mathbf{V}))$ .

We now make two observations. First, the model for  $\mathbf{V}(\mathbf{pa}(\mathbf{V}))$  contains  $P_{u_k}$  for any  $u_k$ , because  $U_{i,k}$  is not restricted by the model and is therefore permitted to have a point-mass distribution at  $u_k$ . Second, the expected value of  $P_{u_k}$  with respect to  $U_{i,k}$  recovers the original marginal distribution  $P(\mathbf{V}(\mathbf{pa}(\mathbf{V})))$ , which is therefore in the convex hull of the set of distributions  $\mathcal{S}(P_{u_k}) \equiv \{P_{u_k} \mid u_k \in \mathcal{S}(U_{i,k})\}$ .

Carathéodory's Theorem (1907) states that for any point  $P$  in the convex hull of a set  $\mathcal{S}$  in a space of dimension  $d-1$ , there exists a set of  $d-1$  points  $\{P_{u_{k_1}}, \dots, P_{u_{k_{d-1}}}\}$  and weights  $\{w_1, \dots, w_{d-1}\}$  such that  $P = \sum_{\ell=1}^{d-1} w_\ell P_{u_{i_\ell}}$ . It then follows directly that any distribution in the marginal model over  $\mathbf{V}(\mathbf{pa}(\mathbf{V}))$  when latent variables have unrestricted cardinality is also in the marginal model over  $\mathbf{V}(\mathbf{pa}(\mathbf{V}))$  when latent variables have cardinality restricted to  $\prod_{V_{i,j} \in \mathbf{V}_i} |\mathcal{S}(V_{i,j})|^{|\mathcal{S}(\mathbf{pa}(V_{i,j}))|} - 1$  or higher.  $\square$

## Proof of Proposition 2

*Proof.* Using the approach developed in Evans (2018) and generalized to arbitrary graphs in Finkelstein et al. (2021), we can obtain a functional model that is non-restrictive of the causal model of  $\mathcal{G}$  over observed variables. In such a model, each  $V_{i,\ell}(\mathbf{a}_\ell)$  is determined by values of the disturbances  $\mathbf{U}_i$ . By assumption,  $\mathcal{G}$  is in canonical form, rendering all disturbances marginally independent. The proposition then follows from standard probability calculus.  $\square$

## Proof of Proposition 3

*Proof.* Under the conditions specified, no element in  $\mathcal{C}$  involves a function of  $U_{i,k}$ . It follows that whether the disturbances lead to  $\mathcal{C}$  is not a function of the value of  $U_{i,k}$ . As a result, a sum over all parameters of the distribution of  $U_{i,k}$  can be factored out of the product in Equation 2. By the definition of probability distributions, this sum will be equal to 1, rendering the parameters irrelevant to the polynomial.  $\square$

## Proof of Proposition 4

*Proof.* Each of the nested Markov parameters corresponds to the probability that random variables in a single district take certain values after an intervention on parents of the district. It follows from Proposition 3 that no disturbances outside the district corresponding to the nested Markov parameter will appear in the polynomialization of that parameter. From this, it then follows that no disturbances in different districts will interact in constraints corresponding to nested Markov parameters. By Proposition 2, the degree of a polynomialization of the probability of the event is at most the number of relevant disturbances.  $\square$

## C Uncertainty

In this appendix, we provide details on our approach to quantifying the uncertainty of bounds based on estimated empirical inputs,  $\hat{\mathbf{E}} = [\hat{E}_\ell]$ . Recall that the *estimated bounds* are obtained from a polynomial program using equality constraints of the form  $\text{polynomialize}(g_\ell(\mathcal{P}_V) = \hat{E}_\ell)$ , which is equivalent to  $\text{polynomial-fractionalize}(g_\ell(\mathcal{P}_V)) = \hat{E}_\ell$ . Here,  $\hat{E}_\ell$  is the noisily estimated empirical quantity and  $\text{polynomial-fractionalize}(g_\ell(\mathcal{P}_V))$  is the reexpression of that same quantity in terms of principal strata sizes. At a high level, we will proceed by constructing confidence regions  $\text{CR}_\alpha(\hat{\mathbf{E}})$  such that  $\Pr(\mathbf{E} \in \text{CR}_\alpha(\hat{\mathbf{E}})) \geq \alpha$ . To obtain *confidence bounds*, we then replace empirical equality constraints with a looser version that accounts for sampling variation, of the form  $\text{polynomial-fractionalize}(g_\ell(\mathcal{P}_V)) \in \text{CR}_\alpha(\hat{\mathbf{E}})$ .

Observe that because the main variables are discrete,  $\hat{\mathbf{E}}$  is a realization of a multinomial proportion. In what follows, we will assume that empirical evidence arises from a single multinomial distribution, such as a single-world marginal distribution; if multiple independent sets of empirical evidence about differing quantities are available, the procedure generalizes straightforwardly by repeating the procedure within each set and combining the results appropriately.

Based on this idea, we examine two methods for constructing  $\text{CR}_\alpha(\hat{\mathbf{E}})$ . Drawing on Malloy et al. (2020), we first consider a ‘‘Bernoulli-KL’’ approach that constructs separate confidence regions for each observable atomic event,  $\Pr(\mathbf{V}_i = \mathbf{v})$ , treating it as a ‘‘success’’ in a Bernoulli distribution. The approach rotates through all possible  $\mathbf{v}$  and combines the event-specific regions using a result on the Kullback-Leibler divergence of sampling distributions to the underlying population distribution. The Bernoulli-KL method produces a confidence region for single-world distributions that is guaranteed to have conservative coverage for the multinomial proportion in finite samples.

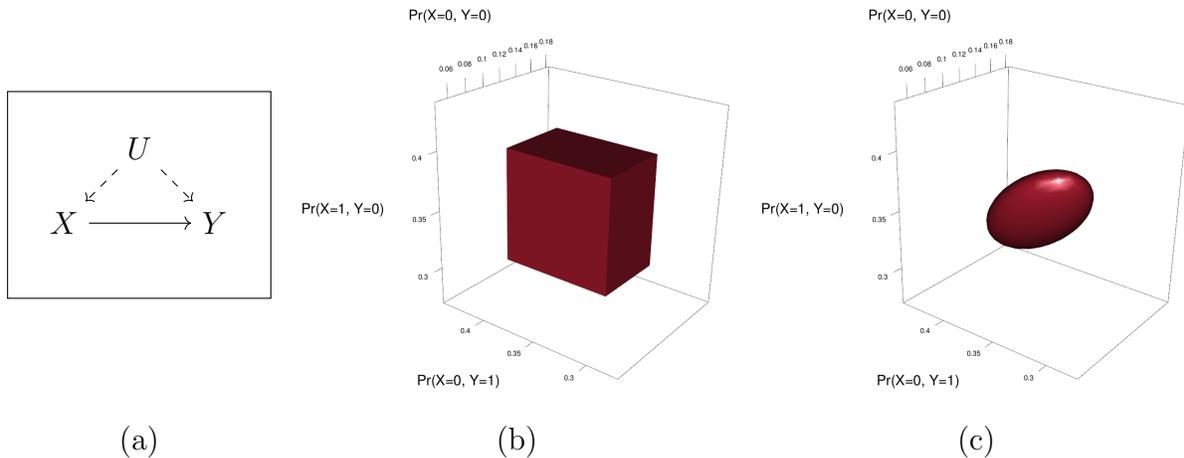
Let  $k \in \{1, \dots, K\}$  index possible atomic events, and denote the probability of the  $k$ -th event as  $p_k = \Pr(\mathbf{V}_i = \mathbf{v}_k)$ . Empirical frequencies are denoted  $\hat{p}_k$ . For the Bernoulli-KL method, we will develop a confidence region of the form  $\text{CR}_\alpha(\hat{\mathbf{E}}) = \bigcap_{k=1}^K [\underline{p}_k, \bar{p}_k]$ , noting that each  $p_k$  can be polynomialized. A visualization of the resulting region is given in Figure 2(b).

We now describe how  $\underline{p}_k$  and  $\bar{p}_k$  can be calculated to ensure that  $\Pr(\mathbf{E} \in \text{CR}_\alpha(\hat{\mathbf{E}})) \geq \alpha$ . At a high level, we will do so by analyzing each of the  $K$  observable events as a Bernoulli distribution. Taking each  $\hat{p}_k$  estimate as given, we identify regions of the unknown  $p_k$  from which the observed  $\hat{p}_k$  diverge substantially. Equation 11 of Malloy et al. (2020) provides bounds on the sampling probability of observing  $\text{KL}([1 - \hat{p}_k, \hat{p}_k], [1 - p_k, p_k])$  in excess of some threshold, where  $\text{KL}([1 - \hat{p}_k, \hat{p}_k], [1 - p_k, p_k]) = \hat{p}_k \log \frac{\hat{p}_k}{p_k} + (1 - \hat{p}_k) \log \frac{1 - \hat{p}_k}{1 - p_k}$ .

In turn, these bounds imply regions of  $p_k$  that can be conservatively rejected. Let  $\underline{p}_k$  be given by the solution to  $\text{KL}([1 - \hat{p}_k, \hat{p}_k], [1 - \underline{p}_k, \underline{p}_k]) = \frac{1}{N} \log \frac{2K}{1 - \alpha}$  subject to  $\underline{p}_k \in [0, \hat{p}_k]$ . Similarly, let  $\bar{p}_k$  be given by  $\text{KL}([1 - \hat{p}_k, \hat{p}_k], [1 - \bar{p}_k, \bar{p}_k]) = \frac{1}{N} \log \frac{2K}{1 - \alpha}$  subject to  $\bar{p}_k \in [\hat{p}_k, 1]$ . It can be seen from Malloy et al. (2020) that when constructing  $\underline{p}_k$  and  $\bar{p}_k$  in this way,  $\Pr\left(\bigcap_{k=1}^K p_k \in [\underline{p}_k, \bar{p}_k]\right) \geq \alpha$  over repeated samples.

Our second approach uses an asymptotic confidence region based on the multivariate Gaussian limiting distribution of the multinomial proportion,  $\mathcal{N}(\mathbf{p}, \text{diag}(\mathbf{p}) - \mathbf{p}\mathbf{p}^\top)$  (Bienaymé, 1838). Because the multinomial proportion must sum to unity, this distribution is degenerate, and it is often more convenient to work with its first  $K - 1$  elements,  $\mathbf{p}_{\setminus K}$ . We con-

Figure 9: **Polynomial confidence regions in a binary graph.** Panel (a) presents a causal graph in which binary  $X$  causes binary  $Y$ , but both are confounded by an unobserved  $U$ .  $N = 1,000$  observations are drawn from this DGP, producing an empirical distribution with proportions  $\frac{1}{N} \sum_{i=1}^N \mathbb{1}(X_i = x, Y_i = y)$ . Panels (b–c) depict confidence regions for  $\Pr(X_i = 0, Y_i = 0)$ ,  $\Pr(X_i = 0, Y_i = 1)$ , and  $\Pr(X_i = 1, Y_i = 0)$ ; the final category,  $\Pr(X_i = 1, Y_i = 1)$  (not depicted), must sum to unity. Panel (b) shows the Bernoulli-KL confidence region, which is conservative in finite samples and can be polynomialized as a set of linear inequalities. Panel (c) shows the Gaussian confidence region, which is asymptotically valid and can be polynomialized as a single convex quadratic inequality.



construct the asymptotic confidence region as  $(\hat{\mathbf{p}}_{\setminus K} - \mathbf{p}_{\setminus K})^\top \left( \text{diag}(\hat{\mathbf{p}}_{\setminus K}) - \hat{\mathbf{p}}_{\setminus K} \hat{\mathbf{p}}_{\setminus K}^\top \right)^{-1} (\hat{\mathbf{p}}_{\setminus K} - \mathbf{p}_{\setminus K}) \leq z$ , where  $z$  is an appropriate critical value of the  $\chi^2$  distribution. A visualization of the resulting region is given in Figure 2(c). As before, each element in  $\mathbf{p}$  is polynomializable, leading to a single confidence constraint that can be straightforwardly incorporated into the optimization routine.

For ease of reference, we duplicate Figure 2 in Figure 9, below. This figure depicts these regions visually for a simple two-node graph, shown in Figure 9(a). The resulting Bernoulli-KL and Gaussian confidence regions are depicted in Figure 9(b–c).

Finally, we describe how arbitrary confidence regions, such as the optimal level-set regions of Malloy et al. (2020) or the exact finite-sample regions of Kuchibhotla et al. (2021), can be polynomialized. At a high level, the proposed method uses a circumscribing polytope, adding faces along the region’s principal axes until the desired tightness is achieved.

One possible approach to doing so is to enumerate candidate  $\mathbf{p}$  along a fine grid, assess each candidate for membership in the confidence region, and compute the convex hull of

the non-rejected points. This procedure produces a system of linear inequalities describing the hull facets. However, it is infeasible for even moderately sized problems, as the time complexity of hull construction can grow exponentially in the dimension of the space,  $K$  (Ottmann et al., 1995). Our approach builds on this basic intuition of circumscribing a complex confidence region with a larger, more tractable polytope. We compute the principal components of the non-rejected points, then identify the two extreme non-rejected points along each axis. Each principal axis is the normal vector for two boundary planes, and each extreme point along that axis defines an boundary plane offset. By repeating this procedure along each principal axis, we obtain a circumscribing confidence region, a parallelepiped that contains the KL confidence region. The gap between the two confidence regions can be rapidly approximated by using number of grid points that lie in the inscribing region but not the original confidence region. By slicing the simplex along additional directions, such as convex combinations of principal axes, this gap can be tightened to arbitrary precision. The resulting polytope defines a system of linear inequalities that can then be incorporated into the polynomial program.

## D Details of Simulated Models

In this section, we detail all models presented in Section 8. For simplicity, all main variables in these models are binary. Simulation parameters are described in terms of principal strata. Principal strata can take one of three forms, depending on the number of parents of the relevant variable. Below, we provide compact notation for referring to these principal strata. Subsequent sections report strata probabilities for each simulation, including joint distributions over strata for multiple variables where confounding exists.

1. **Variables with no parents, which have two strata.** Consider a hypothetical variable  $X_i$  with no parents, as in Figure 5(a). We use  $x_0$  to denote units with  $X_i(\emptyset) = 0$  and  $x_1$  to denote  $X_i(\emptyset) = 1$ .
2. **Variables with a single parent, which have four strata.** Consider a hypotheti-

cal variable  $Y_i$  influenced by parent  $X_i$ , also depicted in Figure 5(a). For compactness, we adopt the convention that counterfactual manipulations of parent variables are presented in the form  $y_{Y_i(X_i=0), Y_i(X_i=1)}$ . For example, (i) we use  $y_{00}$  to denote “never takers” with example,  $Y_i(X_i = 0) = 0$  and  $Y_i(X_i = 1) = 0$ . Similarly, (ii)  $y_{01}$  denotes “compliers” with  $Y_i(X_i = 0) = 0$  and  $Y_i(X_i = 1) = 1$ , (iii)  $y_{10}$  denotes “defiers” with  $Y_i(X_i = 0) = 1$  and  $Y_i(X_i = 1) = 0$ , and  $y_{11}$  denotes “always takers” with  $Y_i(X_i = 0) = 1$  and  $Y_i(X_i = 1) = 1$ .

3. **Variables with two parents, which have sixteen strata.** Consider a hypothetical variable  $Y_i$  influenced by parents  $Z_i$  and  $X_i$ , as in Figure 3(a). Extending the convention described above, we denote these in compact forms ranging from  $y_{0000}$  to  $y_{1111}$ . Specific definitions are provided in Table 2.

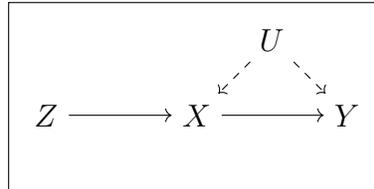
Table 2: **Principal strata for a variable  $Y_i$  with two parents,  $Z_i$  and  $X_i$ .** Each row corresponds to a strata, with compact names given in the first column. For each strata, counterfactual values of  $Y_i$  are given in subsequent columns.

	$Y_i(Z_i = 0, X_i = 0)$	$Y_i(Z_i = 0, X_i = 1)$	$Y_i(Z_i = 1, X_i = 0)$	$Y_i(Z_i = 1, X_i = 1)$
$y_{0000}$	0	0	0	0
$y_{1000}$	1	0	0	0
$y_{0100}$	0	1	0	0
$y_{1100}$	1	1	0	0
$y_{0010}$	0	0	1	0
$y_{1010}$	1	0	1	0
$y_{0110}$	0	1	1	0
$y_{1110}$	1	1	1	0
$y_{0001}$	0	0	0	1
$y_{1001}$	1	0	0	1
$y_{0101}$	0	1	0	1
$y_{1101}$	1	1	0	1
$y_{0011}$	0	0	1	1
$y_{1011}$	1	0	1	1
$y_{0111}$	0	1	1	1
$y_{1111}$	1	1	1	1

## D.1 Noncompliance Simulation

In this section, we describe the DGP for our noncompliance simulation analyzed in Section 8.1. The DGP follows the model of Figure 3(b), reproduced below for ease of reference. Simulation parameters are reported in terms of the joint distribution over principal strata.

Figure 10: **DGP with noncompliance.**



Strata for  $Z$ :

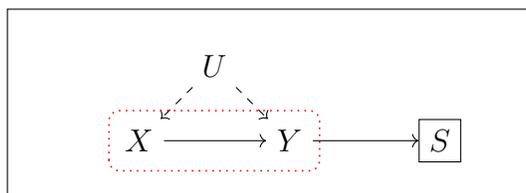
$z_0$	0.649335
$z_1$	0.350665

Strata for  $X$  and  $Y$ :

	$y_{00}$	$y_{10}$	$y_{01}$	$y_{11}$
$x_{00}$	0.000757	0.013034	0.006125	0.002606
$x_{10}$	0.004541	0.074105	0.034526	0.014387
$x_{01}$	0.026040	0.418847	0.195419	0.082264
$x_{11}$	0.004534	0.073950	0.034123	0.014742

## D.2 Outcome-Based Selection Simulation

In this section, we describe the DGP for our outcome-based selection simulation, analyzed in Section 8.3 and Figure 6(a). The DGP follows the model of Figure 5(a), reproduced below for ease of reference. Simulation parameters are reported in terms of the joint distribution over principal strata.



Strata for  $X$  and  $Y$

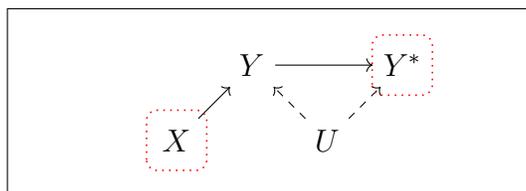
	$y_{00}$	$y_{10}$	$y_{01}$	$y_{11}$
$x_0$	0.124855	0	0.249647	0.124847
$x_1$	0.125375	0	0.249851	0.125425

Strata for  $S$

$S_{10}$	0.50052
$S_{01}$	0.49948

### D.3 Measurement Error Simulation

In this section, we describe the DGP for our measurement error simulation, analyzed in Section 8.3 and Figure 6(b). The DGP follows the model of Figure 5(b), reproduced below for ease of reference. Simulation parameters are reported in terms of the joint distribution over principal strata.



Strata for  $X$

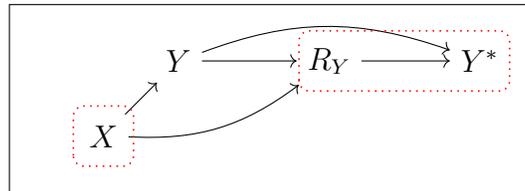
$x_0$	0.499442
$x_1$	0.500558

Strata for  $Y$  and  $Y^*$

	$Y_{00}^*$	$Y_{10}^*$	$Y_{01}^*$	$Y_{11}^*$
$y_{00}$	0	0.167269	0	0
$y_{10}$	0	0	0	0
$y_{01}$	0	0.165838	0.500388	0
$y_{11}$	0	0.166505	0	0

## D.4 Outcome Missingness Simulation

In this section, we describe the DGP for our outcome missingness simulation, analyzed in Section 8.3 and Figure 6(c). The DGP follows the model of Figure 5(c), reproduced below for ease of reference. Simulation parameters are reported in terms of the joint distribution over principal strata.



Strata for  $X$

$x_0$	0.499159
$x_1$	0.500841

Strata for  $Y$

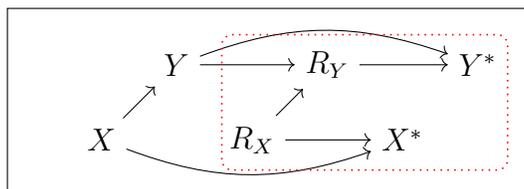
$y_{00}$	0.166371
$y_{10}$	0
$y_{01}$	0.666851
$y_{11}$	0.166778

Strata for  $R$

$r_{0000}$	0
$r_{1000}$	0
$r_{0100}$	0.250368
$r_{1100}$	0.249910
$r_{0010}$	0
$r_{1010}$	0
$r_{0110}$	0
$r_{1110}$	0
$r_{0001}$	0
$r_{1001}$	0
$r_{0101}$	0.250154
$r_{1101}$	0
$r_{0011}$	0
$r_{1011}$	0
$r_{0111}$	0
$r_{1111}$	0.249568

## D.5 Joint Missingness Simulation

In this section, we describe the DGP for our joint missingness simulation, analyzed in Section 8.3 and Figure 6(d). The DGP follows the model of Figure 5(d), reproduced below for ease of reference. Simulation parameters are reported in terms of the joint distribution over principal strata.



Strata for  $X$

$x_0$	0.43464
$x_1$	0.56536

Strata for  $Y$

$y_{00}$	0.485336
$y_{10}$	0.253616
$y_{01}$	0.003768
$y_{11}$	0.257279

Strata for  $R_x$

$r_{x,0}$	0.470201
$r_{x,1}$	0.529798

Strata for  $R_y$

$r_{y,0000}$	0
$r_{y,1000}$	0
$r_{y,0100}$	0.162045
$r_{y,1100}$	0
$r_{y,0110}$	0.177470
$r_{y,0001}$	0.107010
$r_{y,1001}$	0.120311
$r_{y,0101}$	0.255778
$r_{y,1101}$	0.081733
$r_{y,0011}$	0
$r_{y,1011}$	0
$r_{y,0111}$	0.095652
$r_{y,1111}$	0