

---

# Partial Identification of Counterfactual Distributions

---

Anonymous Author(s)

Affiliation

Address

email

## Abstract

1 This paper investigates the problem of bounding counterfactual queries from a  
2 combination of observational data and qualitative assumptions about the underlying  
3 data-generating model. These assumptions are usually represented in the form  
4 of a causal diagram (Pearl, 1995). We show that all counterfactual distributions  
5 (over finite observed variables) in an arbitrary causal diagram could be generated  
6 by a special family of structural causal models (SCMs), compatible with the  
7 same causal diagram, where unobserved (exogenous) variables are discrete, taking  
8 values in a finite domain. This entails a reduction in which the space where the  
9 original, arbitrary SCM lives can be mapped to a dual, more well-behaved space  
10 where the exogenous variables are discrete, and more easily parametrizable. Using  
11 this reduction, we translate the bounding problem in the original space into an  
12 equivalent optimization program in the new space. Solving such programs leads to  
13 optimal bounds over unknown counterfactuals. Finally, we develop effective Monte  
14 Carlo algorithms to approximate these optimal bounds from a finite number of  
15 observational data. Our algorithms are validated extensively on synthetic datasets.

## 16 1 Introduction

17 This paper studies the problem of inferring counterfactual queries from the combination of non-  
18 experimental data (e.g., observational studies) and qualitative assumptions about the data-generating  
19 process. These assumptions are represented in the form of a *causal diagram* [32], which is a  
20 directed acyclic graph where arrows indicate the potential existence of functional relationships among  
21 corresponding variables; some variables are unobserved. This problem arises in diverse fields such  
22 as artificial intelligence, statistics, cognitive science, economics, and the health and social sciences.  
23 For example, when investigating the gender discrimination in college admission, one may ask “what  
24 would the admission outcome be for a female applicant had she been a male?” Such a counterfactual  
25 query contains conflicting information: in the real world the applicant is female, in the hypothetical  
26 world she was not. Therefore, it is not immediately clear how to design effective experimental  
27 procedures for evaluating counterfactuals, let alone how to compute them from observations alone.

28 The problem of identifying counterfactual distributions from the combination of data and a causal  
29 diagram has been studied in the causal inference literature. First, there exist a complete proof system  
30 for reasoning about counterfactual queries [19]. While such a system, in principle, is sufficient in  
31 evaluating any identifiable counterfactual expression, it lacks a proof guideline which determines the  
32 feasibility of such evaluation efficiently. There are algorithms to determine whether a counterfactual  
33 distribution is inferrable from all possible controlled experiments [41]. There exist also algorithms  
34 for identifying path-specific effects from experimental data [1] and observational data [42].

35 In practice, however, the combination of quantitative knowledge and observed data does not always  
36 permit one to point-identify the target counterfactual queries. Partial identification methods concern  
37 with deriving informative bounds over the target counterfactual probability, even when the target

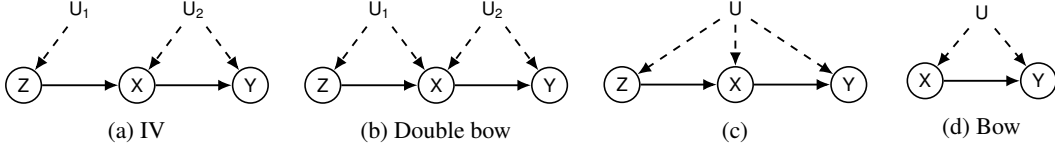


Figure 1: DAGs (a-d) containing a treatment  $X$ , an outcome  $Y$ , an ancestor  $Z$ , and exogenous variables  $U$ ;  $Z$  in (a) is also referred to as an instrumental variable.

38 itself is non-identifiable. Several algorithms have been developed to bound counterfactuals from the  
 39 combination of observational and experimental data [30, 36, 3, 4, 14, 35, 23, 24, 16, 25, 49].

40 In this work, we build on the approach introduced by Balke & Pearl in [3], which involves direct  
 41 discretization of the exogenous domains, also referred to as the principal stratification [17, 34]. Con-  
 42 sider the causal diagram of Fig. 1a where  $X, Y, Z$  are binary variables in  $\{0, 1\}$ ;  $U$  is an unobserved  
 43 variable taking values in an arbitrary continuous domain. [3] showed that domains of  $U$  could be  
 44 discretized into 16 equivalent classes without changing the original counterfactual distributions and  
 45 the graphical structure in Fig. 1a. For instance, despite it being induced by an arbitrary distribution  
 46  $P^*(u)$  over a continuous domain of the exogenous variable  $U$ , the observational distribution  $P(x, y|z)$   
 47 must be reproduced by a generative model of the form  $P(x, y|z) = \sum_u P(x|u, z)P(y|x, u)P(u)$ ,  
 48 where  $P(u)$  is a discrete distribution over a finite exogenous domain  $\{1, \dots, 16\}$ .

49 Using the finite-state representation of unobserved variables, [4] derived tight bounds on treatment  
 50 effects under the condition of noncompliance in Fig. 1a. [11, 21] applied the parsimony of finite-state  
 51 representation in a Bayesian framework, to obtain credible intervals for the posterior distribution of  
 52 causal effects in noncompliance settings. Despite their optimal guarantees, these bounds are only  
 53 applicable to the specific noncompliance setting in Fig. 1a. For the most general cases, a systematic  
 54 procedure for bounding counterfactual queries in arbitrary causal diagrams is still missing.

55 Our goal in this paper is to overcome these challenges. We investigate the expressive power of *discrete*  
 56 *structural causal models* (SCMs) [33] where each unobserved variable is drawn from a discrete  
 57 distribution, takes values in a finite set of states. We show that when inferring about counterfactual  
 58 distributions (over finite observed variables) in an arbitrary causal diagram, one could restrict domains  
 59 of unobserved variables to a finite space without loss of generality. This observation allows us to  
 60 develop novel partial identification algorithms to bound unknown counterfactual probabilities from  
 61 the observational data. More specifically, our contributions are as follows. (1) We introduce a  
 62 special family of discrete SCMs, with finite unobserved domains, and show that it could represent  
 63 all categorical counterfactual distributions in an arbitrary causal diagram. (2) Using this result, we  
 64 translate the original partial identification task into equivalent polynomial programs. Solving such  
 65 programs leads to informative bounds over unknown counterfactual probabilities, which are provably  
 66 optimal. (3) We develop an effective Monte Carlo algorithm to approximate optimal counterfactual  
 67 bounds from a finite number of observational data. Finally, our algorithms are validated extensively  
 68 on synthetic datasets. Given space constraints, all proofs are provided in Appendices A and B.

## 69 1.1 Preliminaries

70 We introduce in this section some basic notations and definitions that will be used throughout the  
 71 paper. We use capital letters to denote variables ( $X$ ), small letters for their values ( $x$ ) and  $\Omega_X$  for  
 72 their domains. For an arbitrary set  $X$ , let  $|X|$  be its cardinality. For convenience, we denote by  $P(x)$   
 73 probabilities  $P(X = x)$ ; for an arbitrary subdomain  $\mathcal{X} \subseteq \Omega_X$ ,  $P(\mathcal{X}) \equiv P(X \in \mathcal{X})$ . Finally, the  
 74 indicator function  $\mathbb{1}_{X=x}$  returns 1 if an event  $X = x$  holds true; otherwise  $\mathbb{1}_{X=x} = 0$ .

75 The basic semantical framework of our analysis rests on *structural causal models* (SCMs) [33,  
 76 Ch. 7]. An SCM  $M$  is a tuple  $\langle V, U, F, P \rangle$  where  $V$  is a set of endogenous variables and  $U$  is  
 77 a set of exogenous variables.  $F$  is a set of functions where each  $f_V \in F$  decides values of an  
 78 endogenous variable  $V \in V$  taking as argument a combination of other variables in the system. That  
 79 is,  $v \leftarrow f_V(pa_V, u_V)$ ,  $Pa_V \subseteq V$ ,  $U_V \subseteq U$ . Exogenous variables  $U \in U$  are mutually independent,  
 80 values of which are drawn from the exogenous distribution  $P(u)$ . Naturally,  $M$  induces a joint  
 81 distribution  $P(v)$  over endogenous variables  $V$ , called the *observational distribution*. Each SCM  
 82 is associated with a causal diagram  $\mathcal{G}$  (e.g., Fig. 1), which is a directed acyclic graph (DAG) where

83 solid nodes represent endogenous variables  $V$ , empty nodes represent exogenous variables  $U$  and  
 84 arrows represent the arguments  $Pa_V, U_V$  of each function  $f_V$ .

85 An intervention on an arbitrary subset  $X \subseteq V$ , denoted by  $\text{do}(x)$ , is an operation where values of  
 86  $X$  are set to constants  $x$ , regardless of how they are ordinarily determined. For an SCM  $M$ , let  
 87  $M_x$  denote a submodel of  $M$  induced by intervention  $\text{do}(x)$ . For any subset  $Y \subseteq V$ , the *potential*  
 88 *response*  $Y_x(u)$  is defined as the solution of  $Y$  in the submodel  $M_x$  given  $U = u$ . Drawing values  
 89 of exogenous variables  $U$  following the probability measure  $P$  induces a *counterfactual variable*  $Y_x$ .  
 90 Specifically, the event  $Y_x = y$  (for short,  $y_x$ ) can be read as “ $Y$  would be  $y$  had  $X$  been  $x$ ”. For any  
 91 subsets  $Y, \dots, Z, X, \dots, W \subseteq V$ , the distribution over counterfactuals  $Y_x, \dots, Z_w$  is defined as:

$$P(y_x, \dots, z_w) = \int_{\Omega_U} \mathbb{1}_{Y_x(u)=y} \wedge \dots \wedge \mathbb{1}_{Z_w(u)=z} dP(u). \quad (1)$$

92 Distributions of the form  $P(y_x)$  is called the *interventional distribution*; when the treatment set  
 93  $X = \emptyset$ ,  $P(y)$  coincides with the *observational distribution*. Throughout this paper, we assume  
 94 that endogenous variables  $V$  are discrete and finite; while exogenous variables  $U$  could take any  
 95 (continuous) value. The counterfactual distribution  $P(y_x, \dots, z_w)$  defined above is thus a categorical  
 96 distribution. For a more detailed survey on SCMs, we refer readers to [33] Ch. 7].

## 97 2 Discretization of Structural Causal Models

98 For a DAG  $\mathcal{G}$  with endogenous  $V$  and exogenous variables  $U$ , let  $P^*$  denote the collection of all  
 99 counterfactual distributions over variables  $V$ . Formally,

$$P^* = \{P(y_x, \dots, z_w) \mid \forall Y, \dots, Z, X, \dots, W \subseteq V\}. \quad (2)$$

100 Let  $\mathcal{M}$  be the family of all the SCMs compatible with the causal diagram  $\mathcal{G}$ , i.e.,  $\mathcal{M} =$   
 101  $\{\forall M \mid \mathcal{G}_M = \mathcal{G}\}$ .<sup>1</sup> Counterfactual distributions in  $\mathcal{G}$  are defined as the collection  $\{P_M^* : \forall M \in \mathcal{M}\}$   
 102 that contains all counterfactual probabilities induced by SCMs  $M$  in the candidate family  $\mathcal{M}$ . In this  
 103 section, we will show that counterfactual distributions in any causal diagram  $\mathcal{G}$  could be generated by  
 104 an alternative family of “generic” SCMs compatible with  $\mathcal{G}$ , which we will define later.

105 **Definition 1** (Counterfactual-Equivalence). For a DAG  $\mathcal{G}$ , let  $\mathcal{M}, \mathcal{N}$  be two sets of SCMs compatible  
 106 with  $\mathcal{G}$ .  $\mathcal{M}$  and  $\mathcal{N}$  are said to be *counterfactually equivalent* (for short, ctf-equivalent) if for any  
 107  $M \in \mathcal{M}$ , there exists an alternative  $N \in \mathcal{N}$  such that  $P_M^* = P_N^*$ , and vice versa.

108 Our analysis rests on a special family of SCMs where values of each exogenous variable are drawn  
 109 from a discrete distribution over a finite set of states.

110 **Definition 2.** An SCM  $M = \langle V, U, F, P \rangle$  is said to be a discrete SCM if

- 111 1. Values of every  $U \in U$  are drawn from a discrete distribution  $P(u)$  over a domain  $\Omega_U$ ; let  
 112  $\theta_u$  denote the probability  $P(U = u)$ , for any  $u \in \Omega_U$ .
- 113 2. Values of every  $V \in V$  are decided by function  $v \leftarrow f_V(pa_V, u_V) \equiv \xi_V^{(pa_V, u_V)}$ , where for  
 114  $\forall pa_V, u_V, \xi_V^{(pa_V, u_V)}$  is a constant in the finite domain  $\Omega_V$ .

115 Given a causal diagram  $\mathcal{G}$ , our goal is to construct a family of discrete SCMs  $\mathcal{N}$  that is counter-  
 116 factually equivalent to the original family of SCMs  $\mathcal{M}$ . Our construction utilizes a special type of  
 117 clustering of nodes in the diagram, called the confounded component [45].

118 **Definition 3.** For an DAG  $\mathcal{G}$ , a subset  $C \subseteq V$  is a c-component if any pair  $X, Y \in C$  is connected  
 119 in  $\mathcal{G}$  by a *bi-directed path* of the form  $V_1 \leftrightarrow V_2 \leftrightarrow \dots \leftrightarrow V_n, n = 1, 2, \dots$ , where (1)  $V_1 = X,$   
 120  $V_n = Y$ ; (2)  $\{V_1, \dots, V_n\} \subseteq V$ ; and (3) each  $V_i \leftrightarrow V_j$  is a sequence  $V_i \leftarrow U_k \rightarrow V_j$  and  $U_k \in U$ .

121 A c-component  $C$  in  $\mathcal{G}$  is maximal if there exists no other c-component that contains  $C$ . We denote  
 122 by  $\mathcal{C}(\mathcal{G})$  the collection of all maximal c-components in  $\mathcal{G}$ . Naturally, c-components in  $\mathcal{C}(\mathcal{G})$  form a  
 123 partition over endogenous variables  $V$ , which, in turn, defines a partition  $\{\cup_{V \in C} U_V \mid \forall C \in \mathcal{C}(\mathcal{G})\}$   
 124 over exogenous variables  $U$ . Therefore, for every  $U \in U$ , there must exist a unique c-component  
 125 in  $\mathcal{C}(\mathcal{G})$ , denoted by  $C_U$ , such that  $U \in \cup_{V \in C_U} U_V$ . For example, exogenous variables  $U_1, U_2$  in  
 126 Fig. 1a corresponds to c-components  $C_{U_1} = \{Z\}$  and  $C_{U_2} = \{X, Y\}$  respectively; while the causal  
 127 diagram of Fig. 1b only has a single c-component  $\{X, Y, Z\}$ .

<sup>1</sup>We will use the subscript  $M$  to represent the restriction to a specific SCM  $M$ . Therefore,  $\mathcal{G}_M$  represents the  
 causal diagram associated with SCM  $M$ ; so does the collection of counterfactuals  $P_M^*$ .

128 **Theorem 1.** For a DAG  $\mathcal{G}$ , consider the following conditions<sup>2</sup>: (1)  $\mathcal{M}$  is the set of all SCMs  
129 compatible with  $\mathcal{G}$ ; (2)  $\mathcal{N}$  is the set of all discrete SCMs compatible with  $\mathcal{G}$  where for every  $U \in \mathbf{U}$ ,  
130 its cardinality  $|\Omega_U| = \prod_{V \in \mathcal{C}_U} |\Omega_{Pa_V} \mapsto \Omega_V|$ , i.e., the number of functions mapping from  $Pa_V$  to  
131  $V$  for every variable  $V$  in the c-component  $\mathcal{C}_U$ . Then,  $\mathcal{M}$  and  $\mathcal{N}$  are counterfactually equivalent.

132 Thm. 1 establishes the expressive power of discrete SCMs in representing counterfactual distributions  
133 in a causal diagram  $\mathcal{G}$ . It implies that the counterfactual distribution  $P(\mathbf{y}_x, \dots, \mathbf{z}_w)$  in any SCM  $M$   
134 could be generated using a generic model as follows, for  $d_U = \prod_{V \in \mathcal{C}_U} |\Omega_{Pa_V} \mapsto \Omega_V|$ ,

$$P(\mathbf{y}_x, \dots, \mathbf{z}_w) = \sum_{U \in \mathbf{U}} \sum_{u=1, \dots, d_U} \mathbb{1}_{\mathbf{Y}_x(\mathbf{u})=\mathbf{y}} \wedge \dots \wedge \mathbb{1}_{\mathbf{Z}_w(\mathbf{u})=\mathbf{z}} \prod_{U \in \mathbf{U}} \theta_u. \quad (3)$$

135 Among above quantities,  $\theta_u$  are parameters of the exogenous distribution  $P(u)$  over a finite domain  
136  $\{1, \dots, d_U\}$ . Counterfactual variables  $\mathbf{Y}_x(\mathbf{u})$  are recursively defined as follows:

$$\mathbf{Y}_x(\mathbf{u}) = \{Y_x(\mathbf{u}) \mid \forall Y \in \mathbf{Y}\}, \text{ where } Y_x(\mathbf{u}) = \begin{cases} \mathbf{x}_Y & \text{if } Y \in \mathbf{X} \\ \xi_Y^{\{\{V_x(\mathbf{u}) \mid V \in Pa_Y\}, u_Y\}} & \text{otherwise} \end{cases} \quad (4)$$

137 where  $\mathbf{x}_Y$  is the value assigned to variable  $Y$  in constants  $\mathbf{x}$ . As an example, consider the causal  
138 diagram  $\mathcal{G}$  described in Fig. 1b where  $X, Y, Z$  are binary variables in  $\{0, 1\}$ . Since  $\mathcal{G}$  has a single c-  
139 component  $\{X, Y, Z\}$ , exogenous variables  $U_1, U_2$  must share the same cardinality  $d$  in the proposed  
140 family of discrete SCMs  $\mathcal{N}$ . It follows from Thm. 1 the counterfactual distribution  $P(z, x_{z'}, y_{x'})$  in  
141 any SCM compatible with  $\mathcal{G}$  could be written as follows:

$$P(z, x_{z'}, y_{x'}) = \sum_{u_1, u_2=1}^d \mathbb{1}_{\xi_Z^{(u_1)}=z} \wedge \mathbb{1}_{\xi_X^{(z', u_1, u_2)}=x} \wedge \mathbb{1}_{\xi_Y^{(x', u_2)}=y} \theta_{u_1} \theta_{u_2}, \quad (5)$$

142 where  $\xi_Z^{(u_1)}, \xi_X^{(z', u_1, u_2)}, \xi_Y^{(x', u_2)}$  are parameters taking values in  $\{0, 1\}$ ;  $\theta_{u_i}, i = 1, 2$ , are probabilities  
143 of the discrete distribution  $P(u_i)$  over the finite domain  $\{1, \dots, d\}$ . The cardinality  $d = |\Omega_Z| \times$   
144  $|\Omega_Z \mapsto \Omega_X| \times |\Omega_X \mapsto \Omega_Y| = 32$ . The total cardinalities of domains for  $U_1, U_2$  are thus  $2d = 64$ .

145 **Comparison with related work** One could naïvely apply the discretization procedure in [3] and  
146 obtain a family of discrete SCMs that are sufficient in representing distributions in an causal diagram.  
147 However, such parametrization is not necessarily complete. To witness, consider again the causal  
148 diagram in Fig. 1b with binary  $X, Y, Z$ . Applying the discretization in [3] leads to a family of discrete  
149 SCMs compatible with a different diagram in Fig. 1c where the cardinality of exogenous variable  
150  $U$  is equal to  $d = 32$  (see Appendix D for details). However, this parametrization fails to capture  
151 some critical constraints over counterfactual distributions since it does not maintain the original  
152 structure of the causal diagram. For instance, counterfactual variables  $Z$  and  $Y_x$  in the original  
153 diagram of Fig. 1b are independent due to independence restrictions [33, Ch. 7.3.2]; while  $Z$  and  
154  $Y_x$  in Fig. 1c are generally correlated due to the presence of unobserved confounder  $U$ . Compared  
155 with [3], the discretization method in Thm. 1 captures *all* constraints over counterfactual distributions  
156 while requiring only a factor of  $|U|$  increase in the cardinality of exogenous domains.

157 More recently, [15] proved a special case of Thm. 1 for interventional distributions in a specific  
158 class of causal diagrams that satisfy the running intersection property. When there is no direct arrow  
159 between endogenous variables, [38] showed that the observational distribution in a diagram could be  
160 represented using finite-state exogenous variables. Thm. 1 generalizes these results by showing that,  
161 for the first time, *all* counterfactual distributions in an *arbitrary* causal diagram could be generated  
162 using discrete exogenous variables taking values from a finite domain, without any loss of generality.

## 163 2.1 Partial identification of Counterfactual Distributions

164 To demonstrate the expressive power of discrete SCMs, we investigate the problem of partial iden-  
165 tification of counterfactual distributions. For an SCM  $M^* = \langle \mathbf{V}, \mathbf{U}, \mathbf{F}, P \rangle$ , we are interested in  
166 evaluating an arbitrary counterfactual probability  $P(\mathbf{y}_x, \dots, \mathbf{z}_w)$ . The detailed parametrization of  
167  $M^*$  is unknown. Instead, the learner only has access to the causal diagram  $\mathcal{G}$  and the observa-  
168 tional distribution  $P(\mathbf{v})$  induced by  $M^*$ . Our goal is to derive an informative bound  $[l, r]$  from the  
169 combination of  $\mathcal{G}$  and  $P(\mathbf{v})$  that contains the actual counterfactual probability  $P(\mathbf{y}_x, \dots, \mathbf{z}_w)$ .

<sup>2</sup>For every  $V \in \mathbf{V}$ ,  $\Omega_{Pa_V} \mapsto \Omega_V$  is the set of all functions mapping from domains  $\Omega_{Pa_V}$  to  $\Omega_V$ .

170 Let  $\mathcal{N}$  denote the family of discrete SCMs defined in Thm. [1](#) which are compatible with the causal  
 171 diagram  $\mathcal{G}$ . We derive a bound  $[l, r]$  over  $P(\mathbf{y}_x, \dots, \mathbf{z}_w)$  from the observational data  $P(\mathbf{v})$  by solving  
 172 the following optimization problem:

$$[l, r] = \min / \max \left\{ P_N(\mathbf{y}_x, \dots, \mathbf{z}_w) \mid \forall N \in \mathcal{N}, P_N(\mathbf{v}) = P(\mathbf{v}) \right\} \quad (6)$$

173 For instance, consider again the double-bow diagram  $\mathcal{G}$  in Fig. [1b](#). The observational distribution  
 174  $P(x, y, z)$  in any discrete SCM in  $\mathcal{N}$  could be written as:

$$P(x, y, z) = \sum_{u_1, u_2=1}^d \mathbb{1}_{\xi_Z^{(u_1)}=z} \wedge \mathbb{1}_{\xi_X^{(z, u_1, u_2)}=x} \wedge \mathbb{1}_{\xi_Y^{(x, u_2)}=y} \theta_{u_1} \theta_{u_2}. \quad (7)$$

175 One could derive a bound over the counterfactual distribution  $P(z, x_{z'}, y_{x'})$  from the observational  
 176 data  $P(x, y, z)$  by solving polynomial programs which optimize the objective Eq. [5](#) over parameters  
 177  $\theta_{u_1}, \theta_{u_2}, \xi_Z^{(u_1)}, \xi_X^{(z, u_1, u_2)}, \xi_Y^{(x, u_2)}$ , subject to the observational constraints Eq. [7](#).

178 As a corollary, it follows immediately from Thm. [1](#) that the solution  $[l, r]$  of the optimization problem  
 179 Eq. [6](#) is guaranteed to be a valid bound over the unknown counterfactual  $P(\mathbf{y}_x, \dots, \mathbf{z}_w)$ .

180 **Corollary 1** (Soundness). *Given a DAG  $\mathcal{G}$  and an observational distribution  $P(\mathbf{v})$ , let  $\mathcal{M}$  be the set  
 181 of all SCMs compatible with  $\mathcal{G}$  and let  $\mathcal{M}_o = \{\forall M \in \mathcal{M} \mid P_M(\mathbf{v}) = P(\mathbf{v})\}$ . For the solution  $[l, r]$   
 182 of Eq. [6](#),  $P_M(\mathbf{y}_x, \dots, \mathbf{z}_w) \in [l, r]$  for any SCM  $M \in \mathcal{M}_o$ .*

183 Since the underlying SCM  $M^* \in \mathcal{M}_o$ , Corol. [1](#) implies that the derived bound  $[l, r]$  must contain the  
 184 actual counterfactual probability  $P(\mathbf{y}_x, \dots, \mathbf{z}_w)$ . Our next result shows that such a bound  $[l, r]$  is  
 185 provably tight, i.e., it cannot be improved without additional assumptions.

186 **Corollary 2** (Tightness). *Given a DAG  $\mathcal{G}$  and an observational distribution  $P(\mathbf{v})$ , let  $\mathcal{M}$  be the set  
 187 of all SCMs compatible with  $\mathcal{G}$  and let  $\mathcal{M}_o = \{\forall M \in \mathcal{M} \mid P_M(\mathbf{v}) = P(\mathbf{v})\}$ . For the solution  $[l, r]$   
 188 of Eq. [6](#), there exist SCMs  $M_1, M_2 \in \mathcal{M}_o$  such that  $P_{M_1}(\mathbf{y}_x, \dots, \mathbf{z}_w) = l$ ,  $P_{M_2}(\mathbf{y}_x, \dots, \mathbf{z}_w) = r$ .*

189 Corol. [2](#) confirms the tightness of the bound  $[l, r]$  obtained from Eq. [6](#). Suppose there exists a valid  
 190 bound  $[l', r']$  strictly contained in  $[l, r]$ . One could construct from Corol. [2](#) an SCM  $M$  compatible  
 191 with the causal diagram  $\mathcal{G}$  and the observational distribution  $P(\mathbf{v})$ , but its counterfactual probability  
 192  $P(\mathbf{y}_x, \dots, \mathbf{z}_w)$  lies outside  $[l', r']$ , which is a contradiction.

193 The optimization problem of Eq. [6](#) is reducible to equivalent polynomial programs (see Appendix [E](#)).  
 194 Despite the soundness and tightness of derived bounds, solving such programs may take exponentially  
 195 long in the most general case [29](#). Our focus here is upon the causal inference aspect of the problem  
 196 and like earlier discussions we do not specify which solvers are used [3](#) [4](#). In some cases of  
 197 interest, effective approximate planning methods for polynomial programs do exist. Investigating  
 198 these methods is an ongoing subject of research [26](#) [31](#) [48](#) [28](#) [27](#).

### 199 3 Bayesian Approach for Partial Identification

200 This section describes an effective algorithm to approximate the optimal counterfactual bound in  
 201 Eq. [6](#), provided with finite samples  $\bar{\mathbf{v}} = \{\mathbf{v}^{(n)}\}_{n=1}^N$  drawn from the observational distribution  
 202  $P(\mathbf{v})$ , and prior distributions over parameters  $\theta_u$  and  $\xi_V^{(pa_V, u_V)}$  (possibly uninformative).

203 We first introduce Markov Chain Monte Carlo (MCMC) algorithms that sample the posterior distribu-  
 204 tion  $P(\theta_{\text{ctf}} \mid \bar{\mathbf{v}})$  over a counterfactual probability  $\theta_{\text{ctf}} = P(\mathbf{y}_x, \dots, \mathbf{z}_w)$ . More specifically, for every  
 205  $V \in \mathbf{V}$ ,  $\forall pa_V, u_V$ , parameters  $\xi_V^{(pa_V, u_V)}$  are drawn uniformly over the finite domain  $\Omega_V$ . For every  
 206  $U \in \mathbf{U}$ , exogenous probabilities  $\theta_u$  are drawn from a generalized Dirichlet distribution [12](#). We will  
 207 take the view of a stick-breaking construction [40](#) which successively breaks pieces off a unit-length  
 208 stick with size proportional to random draws from a Beta distribution. Parameters  $\theta_u$  are proportions  
 209 of each of the pieces relative to its original size. Formally,

$$\forall u = 1, 2, \dots, d_U, \quad \theta_u = \mu_u \prod_{i=1}^{u-1} (1 - \mu_i), \quad \mu_u \sim \text{Beta} \left( \alpha_U^{(u)}, \beta_U^{(u)} \right), \quad (8)$$

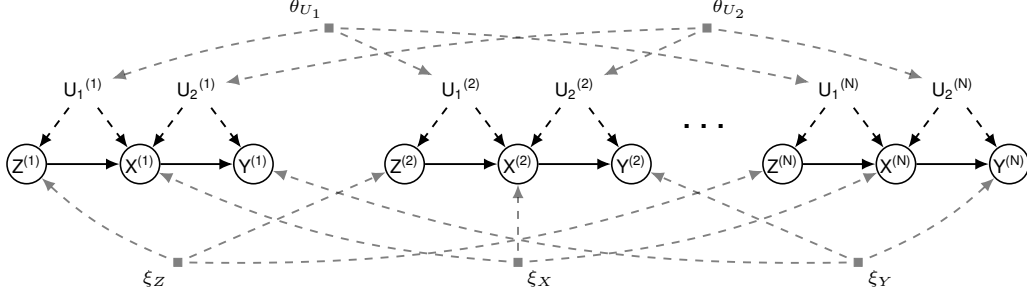


Figure 2: The data-generating process for the observational data  $\{X^{(n)}, Y^{(n)}, Z^{(n)}\}_{n=1}^N$  in an SCM associated with the causal diagram in Fig. 1b. For every exogenous variable  $U \in \mathcal{U}$ ,  $\theta_U = \{\theta_u \mid \forall u\}$ . For every endogenous variable  $V \in \mathcal{V}$ ,  $\xi_V = \{\xi_V^{(pa_V, u_V)} \mid \forall pa_V, u_V\}$ .

210 where  $d_U = \prod_{V \in \mathcal{C}_U} |\Omega_{Pa_V} \mapsto \Omega_V|$  and  $\alpha_U^{(u)}, \beta_U^{(u)} > 0$  are hyperparameters. Finally, we truncate  
 211 this construction by setting  $\mu_{d_U} = 1$ . Note from Eq. 8 that all parameters  $\theta_u$  for  $u > d_U$  are equal  
 212 to zero. As an example, Fig. 2 shows a graphical representation of the data-generating process over  
 213 parameters  $\theta_u$  and  $\xi_V^{(pa_V, u_V)}$  associated with SCMs in Fig. 1b, spanning over  $N$  observations.

214 Gibbs sampling is a well-known MCMC algorithm that allows one to sample posterior distributions.  
 215 For convenience, we introduce the following notations. Let parameters  $\theta = \{\theta_u \mid \forall U \in \mathcal{U}, \forall u\}$   
 216 and  $\xi = \{\xi_V^{(pa_V, u_V)} \mid \forall V \in \mathcal{V}, \forall pa_V, u_V\}$ . The set  $\bar{U} = \{U^{(n)}\}_{n=1}^N$  are exogenous variables  
 217 affecting  $N$  observations  $\bar{V} = \{V^{(n)}\}_{n=1}^N$ ; we use  $\bar{u}$  to represent their realizations. Our blocked  
 218 Gibbs sampler works by iteratively drawing values from the conditional distributions of variables as  
 219 follows [22]. Detailed derivations of complete conditional distributions are shown in Appendix F.

220 **Sampling  $P(\bar{u} \mid \bar{v}, \theta, \xi)$ .** Exogenous variables  $U^{(n)}$ ,  $n = 1, \dots, N$ , are mutually independent  
 221 given parameters  $\theta, \xi$ . We could draw each  $(U^{(n)} \mid \theta, \xi, \bar{V})$  corresponding to the  $n$ th observation  
 222 independently. The complete conditional for  $U^{(n)}$  is given by

$$P(u^{(n)} \mid v^{(n)}, \theta, \xi) \propto \prod_{V \in \mathcal{V}} \mathbb{1}_{\xi_V^{(pa_V^{(n)}, u_V^{(n)})} = v^{(n)}} \prod_{U \in \mathcal{U}} \theta_u. \quad (9)$$

223 **Sampling  $P(\xi, \theta \mid \bar{v}, \bar{u})$ .** Parameters  $\xi, \theta$  are independent given  $\bar{V}, \bar{U}$ . Therefore, we will derive  
 224 complete conditional  $\xi, \theta$  separately. Note that in discrete SCMs, the  $n$ th observation of variable  
 225  $V \in \mathcal{V}$  is decided by  $v^{(n)} \leftarrow \xi_V^{(pa_V, u_V)}$  given  $pa_V^{(n)} = pa_V, u_V^{(n)} = u_V$ . Thus, draw values of each  
 226  $\xi_V^{(pa_V, u_V)} \in \xi$  from the complete conditional defined as:

$$P(\xi_V^{(pa_V, u_V)} \mid \bar{v}, \bar{u}) = \begin{cases} \mathbb{1}_{\xi_V^{(pa_V, u_V)} = v^{(i)}} & \text{if } \exists i, \text{ s.t. } pa_V^{(i)} = pa_V, u_V^{(i)} = u_V, \\ 1/|\Omega_V| & \text{otherwise.} \end{cases} \quad (10)$$

227 Let  $n_u = \sum_{n=1}^N \mathbb{1}_{u^{(n)}=u}$  records the number of values in  $u^{(n)}$  that are equal to  $u$ . By the conjugacy  
 228 of the generalized Dirichlet distribution, the complete conditional of  $\theta_u$  is given by, for every  $U \in \mathcal{U}$ ,

$$\forall u = 1, 2, \dots, d_U, \quad \theta_u = \mu_u \prod_{i=1}^{u-1} (1 - \mu_i), \quad \mu_u \sim \text{Beta} \left( \alpha_U^{(u)} + n_u, \beta_U^{(u)} + \sum_{k=u+1}^{d_U} n_k \right). \quad (11)$$

229 Doing so eventually produces values drawn from the posterior distribution over  $(\theta, \xi, \bar{U} \mid \bar{V})$ . Given  
 230 parameters  $\theta, \xi$ , we compute the counterfactual probability  $\theta_{\text{ctf}} = P(\mathbf{y}_x, \dots, \mathbf{z}_w)$  following the  
 231 three-step algorithm in [33] which consists of abduction, action, and prediction. Thus computing  $\theta_{\text{ctf}}$   
 232 from each draw  $\theta, \xi, \bar{U}$  eventually gives us the draw from the posterior distribution  $P(\theta_{\text{ctf}} \mid \bar{v})$ .

### 233 3.1 Collapsed Gibbs Sampling

234 We also describe an alternative sampler that applies to stick-breaking priors with a known Pólya  
 235 urn characterization. Formally, consider stick-breaking priors in Eq. 8 with hyperparameters

236  $\alpha_U^{(u)} = \alpha_U/d_U$  and  $\beta_U^{(u)} = (d_U - u)\alpha_U/d_U$  for some real  $\alpha_U > 0$ . Let  $\bar{U}_{-n}$  denote the set  
 237 difference  $\bar{U} \setminus U^{(n)}$ ; so does  $\bar{V}_{-n} = \bar{V} \setminus V^{(n)}$ . Our collapsed Gibbs sampler first iteratively draws  
 238 values from the conditional distribution of  $(U^{(n)} | \bar{U}_{-n}, \bar{V})$ ,  $n = 1, \dots, N$ , as follows.

239 **Sampling**  $P(u^{(n)} | \bar{v}, \bar{u}_{-n})$ . At each iteration, draw  $U^{(n)}$  from the conditional given by

$$P(u^{(n)} | \bar{v}, \bar{u}_{-n}) \propto \prod_{V \in \mathcal{V}} P(v^{(n)} | pa_V^{(n)}, u_V^{(n)}, \bar{v}_{-n}, \bar{u}_{-n}) \prod_{U \in \mathcal{U}} P(u^{(n)} | \bar{v}_{-n}, \bar{u}_{-n}). \quad (12)$$

240 Among quantities in the above equation, for every  $V \in \mathcal{V}$ ,

$$P(v^{(n)} | pa_V^{(n)}, u_V^{(n)}, \bar{v}_{-n}, \bar{u}_{-n}) = \begin{cases} \mathbb{1}_{v^{(n)}=v^{(i)}} & \text{if } \exists i \neq n, pa_V^{(i)} = pa_V^{(n)}, u_V^{(i)} = u_V^{(n)}, \\ 1/|\Omega_V| & \text{otherwise.} \end{cases} \quad (13)$$

241 For every  $U \in \mathcal{U}$ , let  $\bar{u}_{-n}$  be a set of exogenous samples  $\{u^{(1)}, \dots, u^{(n-1)}, u^{(n+1)}, \dots, u^{(N)}\}$ . Let  
 242  $\{u_1^*, \dots, u_K^*\}$  denote  $K$  unique values that samples in  $\bar{u}_{-n}$  take on.

$$P(u^{(n)} | \bar{v}_{-n}, \bar{u}_{-n}) = \begin{cases} \frac{n_k^* + \alpha_U/d_U}{\alpha_U + N - 1} & \text{if } u^{(n)} = u_k^*, \text{ for } k = 1, \dots, K \\ \frac{\alpha_U(1 - K/d_U)}{\alpha_U + N - 1} & \text{if } u^{(n)} \notin \{u_1^*, \dots, u_K^*\} \end{cases} \quad (14)$$

243 where  $n_k^* = \sum_{i \neq n} \mathbb{1}_{u^{(i)}=u_k^*}$  records the number of values in  $u^{(i)} \in \bar{u}_{-n}$  that are equal to  $u_k^*$ .

244 Doing so eventually produces exogenous variables drawn from the posterior distribution of  $(\bar{U} | \bar{V})$ .  
 245 We then sample parameters from the posterior distribution of  $(\theta, \xi | \bar{U}, \bar{V})$ ; the complete conditional  
 246  $P(\xi, \theta | \bar{v}, \bar{u})$  are given in Eqs. (10) and (11). Finally, computing  $\theta_{\text{cf}}$  from each sample  $\theta, \xi$  gives  
 247 us a draw from the posterior distribution  $P(\theta_{\text{cf}} | \bar{v})$ .

248 When the cardinality  $d_U$  of exogenous domains is high, the collapsed Gibbs sampler described here is  
 249 more computational efficient than the blocked sampler, since it does not iteratively draw parameters  
 250  $\theta, \xi$  in the high-dimensional space. Instead, the collapsed sampler only draws  $\theta, \xi$  once after samples  
 251 drawn from the distribution of  $(\bar{U} | \bar{V})$  converge. On the other hand, when the cardinality  $d_U$  is  
 252 reasonably low, the blocked Gibbs sampler is preferable since it exhibits better convergence [22].

### 253 3.2 Credible Intervals over Counterfactual Probabilities

254 Given a MCMC sampler, one could bound the counterfactual probability  $\theta_{\text{cf}}$  by computing credible  
 255 intervals from the posterior distribution  $P(\theta_{\text{cf}} | \bar{v})$ .

256 **Definition 4.** Fix  $\alpha \in [0, 1)$ . A  $100(1 - \alpha)\%$  credible interval  $[l_\alpha, r_\alpha]$  for  $\theta_{\text{cf}}$  is given by

$$l_\alpha = \sup \{x | P(\theta_{\text{cf}} \leq x | \bar{v}) = \alpha/2\}, \quad r_\alpha = \inf \{x | P(\theta_{\text{cf}} \leq x | \bar{v}) = 1 - \alpha/2\}. \quad (15)$$

257 For a  $100(1 - \alpha)\%$  credible interval  $[l_\alpha, r_\alpha]$ , any counterfactual probability  $\theta_{\text{cf}}$  that is compatible  
 258 with observational data  $\bar{v}$  lies between the interval  $l_\alpha$  and  $r_\alpha$  with probability  $1 - \alpha$ . Credible  
 259 intervals have been widely applied for computing bounds over counterfactuals provided with finite  
 260 observations [20, 47, 37, 8, 46]. As the number of observational data  $N$  grows (to infinite), the 100%  
 261 credible interval  $[l_0, r_0]$  eventually converges to the optimal asymptotic bound  $[l, r]$  in Eq. (6) [11].

262 Let  $\{\theta^{(t)}\}_{t=1}^T$  be  $T$  samples drawn from  $P(\theta_{\text{cf}} | \bar{v})$ . One could compute the  $100(1 - \alpha)\%$  credible  
 263 interval for  $\theta_{\text{cf}}$  using the following consistent estimators [39]:

$$\hat{l}_\alpha(T) = \theta^{(\lceil (\alpha/2)T \rceil)}, \quad \hat{r}_\alpha(T) = \theta^{(\lceil (1 - \alpha/2)T \rceil)}, \quad (16)$$

264 where  $\theta^{(\lceil (\alpha/2)T \rceil)}, \theta^{(\lceil (1 - \alpha/2)T \rceil)}$  are the  $\lceil (\alpha/2)T \rceil$ th smallest and the  $\lceil (1 - \alpha/2)T \rceil$ th smallest of  
 265  $\{\theta^{(t)}\}_{t=1}^T$ <sup>3</sup>. Our next results establish non-asymptotic deviation bounds for the empirical estimates of  
 266 credible intervals defined in Eq. (16) for finite samples.

267 **Lemma 1.** Fix  $T > 0$  and  $\delta \in (0, 1)$ . Let function  $f(T, \delta) = \sqrt{2T^{-1} \ln(4/\delta)}$ . With probability at  
 268 least  $1 - \delta$ , estimators  $\hat{l}_\alpha(T), \hat{r}_\alpha(T)$  for any  $\alpha \in [0, 1)$  is bounded by

$$\hat{l}_\alpha(T) \in [l_{\alpha - f(T, \delta)}, l_{\alpha + f(T, \delta)}], \quad \hat{r}_\alpha(T) \in [r_{\alpha + f(T, \delta)}, r_{\alpha - f(T, \delta)}]. \quad (17)$$

<sup>3</sup>For any real  $\alpha \in \mathbb{R}$ ,  $\lceil \alpha \rceil$  denotes the smallest integer  $n \in \mathbb{Z}$  larger than  $\alpha$ , i.e.,  $\lceil \alpha \rceil = \min\{n \in \mathbb{Z} | n \geq \alpha\}$ .

269 We summarize our algorithm, CREDIBLEINTERVAL,  
 270 INTERVAL, in Alg. 1. It takes a credible level  
 271  $\alpha$  and tolerance levels  $\delta, \epsilon$  as inputs. In par-  
 272 ticular, CREDIBLEINTERVAL repeatedly draw  
 273  $T \geq \lceil 2\epsilon^{-2} \ln(4/\delta) \rceil$  samples from  $P(\theta_{\text{ctf}} | \bar{\mathbf{v}})$ .  
 274 It then computes estimates  $\hat{l}_\alpha(T), \hat{h}_\alpha(T)$  from  
 275 drawn samples following Eq. (16) and return  
 276 them as the output. It follows immediately from  
 277 Lem. 1 that such a procedure efficiently approx-  
 278 imates a  $100(1 - \alpha)\%$  credible interval.

---

**Algorithm 1: CREDIBLEINTERVAL**

---

- 1: **Input:** Credible level  $\alpha$ , tolerance level  $\delta, \epsilon$ .
  - 2: **Output:** An credible interval  $[l_\alpha, h_\alpha]$  for  $\theta_{\text{ctf}}$ .
  - 3: Let  $T = \lceil 2\epsilon^{-2} \ln(4/\delta) \rceil$ .
  - 4: Draw samples  $\{\theta^{(1)}, \dots, \theta^{(T)}\}$  from the posterior distribution  $P(\theta_{\text{ctf}} | \bar{\mathbf{v}})$ .
  - 5: Return interval  $[\hat{l}_\alpha(T), \hat{r}_\alpha(T)]$  (Eq. (16)).
- 

279 **Corollary 3.** Fix  $\delta \in (0, 1)$  and  $\epsilon > 0$ . With probability at least  $1 - \delta$ , the interval  $[\hat{l}, \hat{r}] =$   
 280 CREDIBLEINTERVAL( $\alpha, \delta, \epsilon$ ) for any  $\alpha \in [0, 1)$  is bounded by  $\hat{l} \in [l_{\alpha-\epsilon}, l_{\alpha+\epsilon}]$  and  $\hat{r} \in [r_{\alpha+\epsilon}, r_{\alpha-\epsilon}]$ .

281 Corol. 3 implies that any counterfactual parameter  $\theta_{\text{ctf}}$  compatible with observational data  $\bar{\mathbf{v}}$  falls  
 282 between  $[\hat{l}, \hat{r}] = \text{CREDIBLEINTERVAL}(\alpha, \delta, \epsilon)$  with probability  $P(\theta_{\text{ctf}} \in [\hat{l}, \hat{r}] | \bar{\mathbf{v}}) \approx 1 - \alpha \pm \epsilon$ . As  
 283 the tolerance rate  $\epsilon \rightarrow 0$ ,  $[\hat{l}, \hat{r}]$  converges to a  $100(1 - \alpha)\%$  credible interval with high probability.

## 284 4 Simulations and Experiments

285 We demonstrate our algorithms on various simulated SCM instances and a real world patient dataset  
 286 collected from the International Stroke Trial (IST) [10]. Overall, we found that simulation results sup-  
 287 port our findings and the proposed bounding strategy consistently dominates state-of-art algorithms.  
 288 When target distributions are identifiable (Experiment 1), our bounds collapse to the actual, unknown  
 289 counterfactual probabilities. For non-identifiable settings, our algorithm obtains sharp asymptotic  
 290 bounds when closed-form solutions already exist (Experiments 2 & 3); and improves over state-of-art  
 291 bounds in other more general cases where the optimal strategy is unknown (Experiment 4).

292 In all experiments, we evaluate our proposed bounding strategy based on credible intervals (*ci*). In  
 293 particular, we draw  $4 \times 10^3$  samples from the posterior distribution over the target counterfactual  
 294  $(\theta_{\text{ctf}} | \bar{\mathbf{V}})$ . This allows us to compute 100% credible interval over  $\theta_{\text{ctf}}$  within error  $\epsilon = 0.05$ , with  
 295 probability at least  $1 - \delta = 0.95$ . As the baseline, we also include the actual counterfactual probability  
 296  $\theta^*$ . For details on simulation setups and additional experiments, we refer readers to Appendix C

297 **Experiment 1: Frontdoor Graph** This experiment evaluates our sam-  
 298 pling algorithm on interventional probabilities that are identifiable from  
 299 the observational data. Consider the ‘‘Frontdoor’’ graph described in  
 300 Fig. 3 where  $X, Y, W$  are binary variables in  $\{0, 1\}$ ;  $U_1, U_2 \in \mathbb{R}$ . In this  
 301 case, the interventional distribution  $P(y_x)$  is identifiable from  $P(x, w, y)$   
 302 through the frontdoor adjustment [33, Thm. 3.3.4]. We collect  $N = 10^5$   
 303 observational samples  $\bar{\mathbf{V}} = \{X^{(n)}, Y^{(n)}, W^{(n)}\}_{n=1}^N$  from a randomly  
 304 generated SCM. Fig. 4a shows samples drawn from the posterior distribution of the target probability  
 305  $(P(Y_{x=0} = 1) | \bar{\mathbf{V}})$ . The analysis reveals that these samples collapse to the actual interventional  
 306 probability  $P(Y_{x=0} = 1) = 0.5085$ , which confirms the identifiability of  $P(y_x)$  in Fig. 3

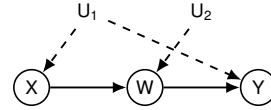


Figure 3: Frontdoor

307 **Experiment 2: Instrumental Variables (IV)** This experiment evaluates our bounding strategy in  
 308 non-identifiable settings, while closed-form solutions for the optimal bounds over target probabilities  
 309 already exist. Consider first the ‘‘IV’’ diagram in Fig. 1a where  $X, Y, Z \in \{0, 1\}$  and  $U_1, U_2 \in \mathbb{R}$ .  
 310 The non-identifiability of  $P(y_x)$  from the observational data  $P(x, y, z)$  with the instrument  $Z$  and the  
 311 unobserved confounding between  $X$  and  $Y$  has been acknowledged in [5]. For binary  $X, Y, Z$ , [2]  
 312 derived closed-form, sharp bounds over  $P(y_x)$  (labelled as *opt*). We collect  $N = 10^5$  observational  
 313 samples  $\bar{\mathbf{V}} = \{X^{(n)}, Y^{(n)}, Z^{(n)}\}_{n=1}^N$  from a randomly generated SCM instance. Fig. 4b shows  
 314 samples drawn from the posterior distribution of  $(P(Y_{x=0} = 1) | \bar{\mathbf{V}})$ . As a baseline, we also include  
 315 the optimal bound *opt*, and posterior samples obtained from the Gibbs sampler of [11], which utilizes  
 316 the canonical partitions of exogenous domains in [2] (*bp*). The analysis reveals that our algorithm  
 317 derives the valid bound over the actual probability  $P(Y_{x=0} = 1) = 0.3954$ ; the 100% credible  
 318 interval converges to the optimal IV bound  $l = 0.1468, r = 0.6617$ .



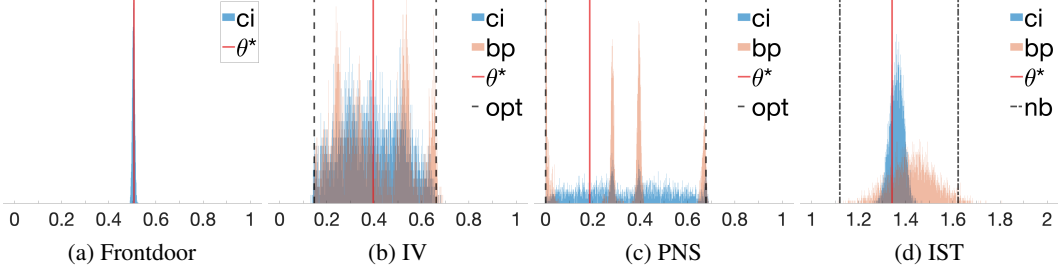


Figure 4: Histogram plots for samples drawn from the posterior distribution over target counterfactual probabilities. For all plots (a-d),  $ci$  represents our proposed algorithms;  $bp$  stands for Gibbs samplers using the representation of canonical partitions [2];  $\theta^*$  is the actual counterfactual probability. (b, c)  $opt$  represents the optimal asymptotic bound, if exists. (d)  $nb$  stands for the natural bounds [30].

319 **Experiment 3: Probability of Necessity and Sufficiency (PNS)** We now study the problem of  
 320 evaluating the *probability of necessity and sufficiency*  $P(Y_{x=1} = 1, Y_{x=0} = 0)$  from the observational  
 321 data  $P(x, y)$  in the “Bow” diagram of Fig. 1d where  $X, Y \in \{0, 1\}$  and  $U \in \mathbb{R}$ . The sharp bound for  
 322  $P(Y_{x=1} = 1, Y_{x=0} = 0)$  from  $P(x, y)$  was introduced in [44] (labelled as  $opt$ ). We collect  $N = 10^5$   
 323 observational samples  $\bar{V} = \{X^{(n)}, Y^{(n)}\}_{n=1}^N$  from an SCM instance. Fig. 4c shows samples drawn  
 324 from the posterior distribution of  $(P(Y_{x=1} = 1, Y_{x=0} = 0) | \bar{V})$ . As a baseline, we also include the  
 325 optimal bound  $opt$ , and posterior samples obtained from the Gibbs sampler which discretizes the  
 326 exogenous domains using canonical partitions [2] ( $bp$ ). The analysis reveals that our 100% credible  
 327 interval ( $ci$ ) matches the optimal PNS bound  $l = 0, r = 0.6775$ , i.e., the proposed strategy achieves  
 328 the sharp bound over the counterfactual probability  $P(Y_{x=1} = 1, Y_{x=0} = 0) = 0.1867$ .

329 **Experiment 4: International Stroke Trials (IST)** IST was a large, randomized, open trial of up  
 330 to 14 days of antithrombotic therapy after stroke onset [10]. In particular, the treatment  $X$  is a pair  
 331  $(i, j)$  where  $i = 0$  stands for no aspirin allocation, 1 otherwise;  $j = 0$  stands for no heparin allocation,  
 332 1 for median-dosage, and 2 for high-dosage. The primary outcome  $Y \in \{0, \dots, 3\}$  is the health  
 333 of the patient 6 months after the treatment, where 0 stands for death, 1 for being dependent on the  
 334 family, 2 for the partial recovery, and 3 for the full recovery.

335 To emulate the presence of unobserved confounding, we filter the experimental data with selection  
 336 rules  $f_X^{(Z)}$ ,  $Z \in \{0, \dots, 9\}$ , following a procedure in [49]. Doing so allows us to obtain  $N = 3 \times 10^3$   
 337 synthetic observational samples  $\bar{V} = \{X^{(n)}, Y^{(n)}, Z^{(n)}\}_{n=1}^N$  that are compatible with the “Double  
 338 bow” diagram of Fig. 1b. We are interested in evaluating the treatment effect  $E[Y_{x=(1,0)}]$  for  
 339 only assigning aspirin  $\bar{X} = (1, 0)$ . Fig. 4d shows samples drawn from the posterior distribution  
 340 of  $(E[Y_{x=(1,0)}] | \bar{V})$ . As a baseline, we also include a naïve generalization of the discretization  
 341 procedure ( $bp$ ) [2] (see Appendix D) and the natural bounds [36, 30] estimated at the 95% confidence  
 342 level ( $nb$ ) [49]. Posterior samples of  $ci$  and  $bp$  are drawn using our proposed collapsed sampler  
 343 due to the high-dimensional latent space. The analysis reveals that all algorithms achieve bounds  
 344 that contain the actual, target causal effect  $E[Y_{x=(1,0)}] = 1.3418$ . Our bounding strategy obtains a  
 345 100% credible interval  $l_{ci} = 1.2604, r_{ci} = 1.4687$ , which consistently improves over all the other  
 346 algorithms ( $l_{bp} = 1.1121, r_{bp} = 1.8073, l_{nb} = 1.1195, r_{nb} = 1.6221$ ).

## 347 5 Conclusion

348 This paper investigated the problem of partial identification of counterfactual distributions, which  
 349 concerns with bounding unknown counterfactual probabilities from the combination of the obser-  
 350 vational data and qualitative assumptions of the data-generating process, represented in the form of  
 351 a directed acyclic causal diagram. We studied a special family of SCMs with discrete exogenous  
 352 variables, taking values from a finite set of unobserved states, and showed that it could represent *all*  
 353 counterfactual distributions (over finite observed variables) in an arbitrary causal diagram. That is,  
 354 this new family of discrete SCMs is counterfactual equivalent to the original family of candidate  
 355 SCMs compatible with the causal diagram. Using this result, we developed a novel algorithm to  
 356 derive bounds over counterfactual probabilities from finite observations, which are provably tight.

357 **References**

- 358 [1] C. Avin, I. Shpitser, and J. Pearl. Identifiability of path-specific effects. In *Proceedings of the*  
359 *Nineteenth International Joint Conference on Artificial Intelligence IJCAI-05*, pages 357–363,  
360 Edinburgh, UK, 2005. Morgan-Kaufmann Publishers.
- 361 [2] A. Balke and J. Pearl. Counterfactual probabilities: Computational methods, bounds, and  
362 applications. In R. L. de Mantaras and D. Poole, editors, *Uncertainty in Artificial Intelligence*  
363 *10*, pages 46–54. Morgan Kaufmann, San Mateo, CA, 1994.
- 364 [3] A. Balke and J. Pearl. Counterfactuals and policy analysis in structural models. In P. Besnard  
365 and S. Hanks, editors, *Uncertainty in Artificial Intelligence 11*, pages 11–18. San Francisco,  
366 1995.
- 367 [4] A. Balke and J. Pearl. Bounds on treatment effects from studies with imperfect compliance.  
368 *Journal of the American Statistical Association*, 92(439):1172–1176, September 1997.
- 369 [5] E. Bareinboim and J. Pearl. Causal inference by surrogate experiments:  $z$ -identifiability.  
370 In N. de Freitas and K. Murphy, editors, *Proceedings of the Twenty-Eighth Conference on*  
371 *Uncertainty in Artificial Intelligence*, pages 113–120, Corvallis, OR, 2012. AUAI Press.
- 372 [6] H. Bauer. Probability theory and elements of measure theory. *Holt*, 1972.
- 373 [7] H. Bauer. *Measure and integration theory*, volume 26. Walter de Gruyter, 2011.
- 374 [8] F. A. Bugni. Bootstrap inference in partially identified models defined by moment inequalities:  
375 Coverage of the identified set. *Econometrica*, 78(2):735–753, 2010.
- 376 [9] C. Carathéodory. Über den variabilitätsbereich der fourier’schen konstanten von positiven har-  
377 monischen funktionen. *Rendiconti Del Circolo Matematico di Palermo (1884-1940)*, 32(1):193–  
378 217, 1911.
- 379 [10] A. Carolei et al. The international stroke trial (ist): a randomized trial of aspirin, subcutaneous  
380 heparin, both, or neither among 19435 patients with acute ischaemic stroke. *The Lancet*,  
381 349:1569–1581, 1997.
- 382 [11] D. Chickering and J. Pearl. A clinician’s tool for analyzing non-compliance. *Computing Science*  
383 *and Statistics*, 29(2):424–431, 1997.
- 384 [12] R. J. Connor and J. E. Mosimann. Concepts of independence for proportions with a generaliza-  
385 tion of the dirichlet distribution. *Journal of the American Statistical Association*, 64(325):194–  
386 206, 1969.
- 387 [13] J. Eckhoff. Helly, radon, and carathéodory type theorems. In *Handbook of convex geometry*,  
388 pages 389–448. Elsevier, 1993.
- 389 [14] R. J. Evans. Graphical methods for inequality constraints in marginalized dags. In *2012 IEEE*  
390 *International Workshop on Machine Learning for Signal Processing*, pages 1–6. IEEE, 2012.
- 391 [15] R. J. Evans et al. Margins of discrete bayesian networks. *The Annals of Statistics*, 46(6A):2623–  
392 2656, 2018.
- 393 [16] N. Finkelstein and I. Shpitser. Deriving bounds and inequality constraints using logical relations  
394 among counterfactuals. In *Conference on Uncertainty in Artificial Intelligence*, pages 1348–  
395 1357. PMLR, 2020.
- 396 [17] C. Frangakis and D. Rubin. Principal stratification in causal inference. *Biometrics*, 1(58):21–29,  
397 2002.
- 398 [18] D. Galles and J. Pearl. An axiomatic characterization of causal counterfactuals. *Foundation of*  
399 *Science*, 3(1):151–182, 1998.
- 400 [19] J. Halpern. Axiomatizing causal reasoning. In G. Cooper and S. Moral, editors, *Uncertainty*  
401 *in Artificial Intelligence*, pages 202–210. Morgan Kaufmann, San Francisco, CA, 1998. Also,  
402 *Journal of Artificial Intelligence Research* 12:3, 17–37, 2000.

- 403 [20] G. W. Imbens and C. F. Manski. Confidence intervals for partially identified parameters.  
404 *Econometrica*, 72(6):1845–1857, 2004.
- 405 [21] G. W. Imbens and D. B. Rubin. Bayesian inference for causal effects in randomized experiments  
406 with noncompliance. *The annals of statistics*, pages 305–327, 1997.
- 407 [22] H. Ishwaran and L. F. James. Gibbs sampling methods for stick-breaking priors. *Journal of the*  
408 *American Statistical Association*, 96(453):161–173, 2001.
- 409 [23] N. Kallus and A. Zhou. Confounding-robust policy improvement. In *Advances in neural*  
410 *information processing systems*, pages 9269–9279, 2018.
- 411 [24] N. Kallus and A. Zhou. Confounding-robust policy evaluation in infinite-horizon reinforcement  
412 learning. *Advances in Neural Information Processing Systems*, 2020.
- 413 [25] N. Kilbertus, M. J. Kusner, and R. Silva. A class of algorithms for general instrumental variable  
414 models. In *Advances in Neural Information Processing Systems*, 2020.
- 415 [26] J. B. Lasserre. Global optimization with polynomials and the problem of moments. *SIAM*  
416 *Journal on optimization*, 11(3):796–817, 2001.
- 417 [27] J. B. Lasserre. *Moments, positive polynomials and their applications*, volume 1. World Scientific,  
418 2009.
- 419 [28] M. Laurent. Sums of squares, moment matrices and optimization over polynomials. In *Emerging*  
420 *applications of algebraic geometry*, pages 157–270. Springer, 2009.
- 421 [29] H. R. Lewis. *Computers and intractability. a guide to the theory of np-completeness*, 1983.
- 422 [30] C. Manski. Nonparametric bounds on treatment effects. *American Economic Review, Papers*  
423 *and Proceedings*, 80:319–323, 1990.
- 424 [31] P. A. Parrilo. Semidefinite programming relaxations for semialgebraic problems. *Mathematical*  
425 *programming*, 96(2):293–320, 2003.
- 426 [32] J. Pearl. Causal diagrams for empirical research. *Biometrika*, 82(4):669–710, 1995.
- 427 [33] J. Pearl. *Causality: Models, Reasoning, and Inference*. Cambridge University Press, New York,  
428 2000. 2nd edition, 2009.
- 429 [34] J. Pearl. Principal stratification – a goal or a tool? *The International Journal of*  
430 *Biostatistics*, 7(1), 2011. Article 20, DOI: 10.2202/1557-4679.1322. Available at:  
431 <[http://ftp.cs.ucla.edu/pub/stat\\_ser/r382.pdf](http://ftp.cs.ucla.edu/pub/stat_ser/r382.pdf)>.
- 432 [35] A. Richardson, M. G. Hudgens, P. B. Gilbert, and J. P. Fine. Nonparametric bounds and  
433 sensitivity analysis of treatment effects. *Statistical science: a review journal of the Institute of*  
434 *Mathematical Statistics*, 29(4):596, 2014.
- 435 [36] J. Robins. The analysis of randomized and non-randomized aids treatment trials using a new  
436 approach to causal inference in longitudinal studies. In L. Sechrest, H. Freeman, and A. Mulley,  
437 editors, *Health Service Research Methodology: A Focus on AIDS*, pages 113–159. NCHSR,  
438 U.S. Public Health Service, Washington, D.C., 1989.
- 439 [37] J. P. Romano and A. M. Shaikh. Inference for identifiable parameters in partially identified  
440 econometric models. *Journal of Statistical Planning and Inference*, 138(9):2786–2807, 2008.
- 441 [38] D. Rosset, N. Gisin, and E. Wolfe. Universal bound on the cardinality of local hidden variables  
442 in networks. *Quantum Information & Computation*, 18(11-12):910–926, 2018.
- 443 [39] P. K. Sen and J. M. Singer. *Large sample methods in statistics: an introduction with applications*,  
444 volume 25. CRC press, 1994.
- 445 [40] J. Sethuraman. A constructive definition of dirichlet priors. *Statistica sinica*, pages 639–650,  
446 1994.

- 447 [41] I. Shpitser and J. Pearl. What counterfactuals can be tested. In *Proceedings of the Twenty-Third*  
 448 *Conference on Uncertainty in Artificial Intelligence*, pages 352–359. AUAI Press, Vancouver,  
 449 BC, Canada, 2007. Also, *Journal of Machine Learning Research*, 9:1941–1979, 2008.
- 450 [42] I. Shpitser and E. Sherman. Identification of personalized effects associated with causal  
 451 pathways. In *UAI*, 2018.
- 452 [43] J. Tian. *Studies in Causal Reasoning and Learning*. PhD thesis, Computer Science Department,  
 453 University of California, Los Angeles, CA, November 2002.
- 454 [44] J. Tian and J. Pearl. Probabilities of causation: Bounds and identification. *Annals of Mathematics*  
 455 *and Artificial Intelligence*, 28:287–313, 2000.
- 456 [45] J. Tian and J. Pearl. A general identification condition for causal effects. In *Proceedings of the*  
 457 *Eighteenth National Conference on Artificial Intelligence*, pages 567–573. AAAI Press/The  
 458 MIT Press, Menlo Park, CA, 2002.
- 459 [46] D. Todem, J. Fine, and L. Peng. A global sensitivity test for evaluating statistical hypotheses  
 460 with nonidentifiable models. *Biometrics*, 66(2):558–566, 2010.
- 461 [47] S. Vansteelandt, E. Goetghebeur, M. G. Kenward, and G. Molenberghs. Ignorance and uncer-  
 462 tainty regions as inferential tools in a sensitivity analysis. *Statistica Sinica*, pages 953–979,  
 463 2006.
- 464 [48] H. Waki, S. Kim, M. Kojima, and M. Muramatsu. Sums of squares and semidefinite program  
 465 relaxations for polynomial optimization problems with structured sparsity. *SIAM Journal on*  
 466 *Optimization*, 17(1):218–242, 2006.
- 467 [49] J. Zhang and E. Bareinboim. Bounding causal effects on continuous outcomes. In *Proceedings*  
 468 *of the 35nd AAAI Conference on Artificial Intelligence*, 2021.

## 469 Checklist

- 470 1. For all authors...
- 471 (a) Do the main claims made in the abstract and introduction accurately reflect the paper’s  
 472 contributions and scope? [Yes]
- 473 (b) Did you describe the limitations of your work? [Yes] “Throughout this paper, we  
 474 assume that endogenous variables  $V$  are discrete and finite; while exogenous variables  
 475  $U$  could take any (continuous) value.”
- 476 (c) Did you discuss any potential negative societal impacts of your work? [N/A] This work  
 477 does not present any foreseeable societal consequence.
- 478 (d) Have you read the ethics review guidelines and ensured that your paper conforms to  
 479 them? [Yes]
- 480 2. If you are including theoretical results...
- 481 (a) Did you state the full set of assumptions of all theoretical results? [Yes] See Sec. **L.1**
- 482 (b) Did you include complete proofs of all theoretical results? [Yes] See Appendices **A**  
 483 and **B**
- 484 3. If you ran experiments...
- 485 (a) Did you include the code, data, and instructions needed to reproduce the main experi-  
 486 mental results (either in the supplemental material or as a URL)? [Yes] See Appendix **C**
- 487 (b) Did you specify all the training details (e.g., data splits, hyperparameters, how they  
 488 were chosen)? [Yes] See Appendix **C**
- 489 (c) Did you report error bars (e.g., with respect to the random seed after running experi-  
 490 ments multiple times)? [N/A]
- 491 (d) Did you include the total amount of compute and the type of resources used (e.g., type  
 492 of GPUs, internal cluster, or cloud provider)? [Yes] See Appendix **C** “Experiments  
 493 were performed on a computer with 32GB memory.”

- 494 4. If you are using existing assets (e.g., code, data, models) or curating/releasing new assets...
- 495 (a) If your work uses existing assets, did you cite the creators? [Yes] “IST was a large,
- 496 randomized, open trial of up to 14 days of antithrombotic therapy after stroke onset
- 497 [10].” See also Appendix C
- 498 (b) Did you mention the license of the assets? [Yes] See Appendix C. The IST dataset is
- 499 shared under “Open Data Commons Attribution License (ODC-By) v1.0”.
- 500 (c) Did you include any new assets either in the supplemental material or as a URL? [N/A]
- 501
- 502 (d) Did you discuss whether and how consent was obtained from people whose data you’re
- 503 using/curating? [N/A]
- 504 (e) Did you discuss whether the data you are using/curating contains personally identifiable
- 505 information or offensive content? [N/A]
- 506 5. If you used crowdsourcing or conducted research with human subjects...
- 507 (a) Did you include the full text of instructions given to participants and screenshots, if
- 508 applicable? [N/A]
- 509 (b) Did you describe any potential participant risks, with links to Institutional Review
- 510 Board (IRB) approvals, if applicable? [N/A]
- 511 (c) Did you include the estimated hourly wage paid to participants and the total amount
- 512 spent on participant compensation? [N/A]