

# Identification of Causal Effects in the Presence of Selection Bias

**Juan D. Correa**

Computer Science Department  
Purdue University  
correagr@purdue.edu

**Jin Tian**

Department of Computer Science  
Iowa State University  
jtian@iastate.edu

**Elias Bareinboim**

Computer Science Department  
Purdue University  
eb@purdue.edu

## Abstract

Cause-and-effect relations are one of the most valuable types of knowledge sought after throughout the data-driven sciences since they translate into stable and generalizable explanations as well as efficient and robust decision-making capabilities. Inferring these relations from data, however, is a challenging task. Two of the most common barriers to this goal are known as confounding and selection biases. The former stems from the systematic bias introduced during the treatment assignment, while the latter comes from the systematic bias during the collection of units into the sample. In this paper, we consider the problem of identifiability of causal effects when both confounding and selection biases are simultaneously present. We first investigate the problem of identifiability when all the available data is biased. We prove that the algorithm proposed by [Bareinboim and Tian, 2015] is, in fact, complete, namely, whenever the algorithm returns a failure condition, no identifiability claim about the causal relation can be made by any other method. We then generalize this setting to when, in addition to the biased data, another piece of external data is available, without bias. It may be the case that a subset of the covariates could be measured without bias (e.g., from census). We examine the problem of identifiability when a combination of biased and unbiased data is available. We propose a new algorithm that subsumes the current state-of-the-art method based on the back-door criterion.

## Introduction

One prominent challenge shared throughout the empirical disciplines is to infer cause and effect relationships – for instance, one may need to determine how increasing the state’s educational budget will bring about change in the average income of the population, whether exposing subjects to a new advertisement campaign would translate into additional sales revenue, or how patients will react to the decrease of the drug’s dosage, would they still recover in acceptable health conditions? Despite the disparate nature of these questions in terms of subject matter, they evoke the same set of principles and formal machinery, which comes under the rubric of *causal inference* (Pearl 2000; Spirtes, Glymour, and Scheines 2001).

Causal inference is concerned with the potential mismatch between the inferential power of the collected data

and the target inference. In practice, this is particularly relevant since data is almost invariably plagued with various biases, most prominently, confounding and selection. The former refers to the presence of a set of factors that affect both the action (also known as treatment) and the outcome, while the latter arises when the action, outcome, and other factors differentially affect the inclusion of subjects in the data sample (Bareinboim and Pearl 2016).

The problem of *identifiability* gives formal dressing to the issue of confounding (Pearl 2000, Ch. 3). Specifically, it is concerned with determining the effect of a treatment ( $X$ ) on an outcome ( $Y$ ), denoted  $P(y|do(x))$  (for short,  $P_x(y)$ ), based on the observational, non-experimental distribution  $P(\mathbf{v})$  (where  $\mathbf{V}$  represents observable variables) and causal assumptions commonly expressed as a directed acyclic graph. The difference between  $P(y|do(x))$  and its probabilistic counterpart,  $P(y|x)$ , is what is called *confounding bias* (Bareinboim and Pearl 2016). This problem has been extensively studied in the literature. A systematic treatment of this problem was given in (Pearl 1995), which introduced *do-calculus*. The do-calculus was shown to be complete for non-parametric identifiability from observational and experimental data (Tian and Pearl 2002a; Huang and Valtorta 2006; Shpitser and Pearl 2006; Bareinboim and Pearl 2012a).

The other source of disparities, *selection bias*, usually appears due to the preferential exclusion of units from the sample. For instance, in a typical study of the effect of grades on college admissions, subjects with higher achievement tend to report their scores more frequently than those who scored lower. In this case, the data-gathering process will reflect a distortion in the sample’s proportions and, since the data is no longer a faithful representation of the underlying population, biased estimates will be produced regardless of the number of samples collected (even when the treatment is controlled). The problem of selection bias can also be modeled graphically through the explicit articulation of the sampling mechanism,  $S$ . This mechanism can be seen as a binary indicator of entry into the data pool, such that  $S=1$  if a unit is included in the sample and  $S=0$  otherwise. Clearly, when the sampling process is entirely random,  $S$  is independent of all variables in the analysis. When samples are collected preferentially, the causal effects not only need to be identified but also *recovered* from the distribution

$P(\mathbf{v}|S=1)$ , instead of  $P(\mathbf{v})$  (Bareinboim and Pearl 2012b).

Selection bias has challenged inferences throughout a wide range of disciplines, including AI (Cooper 1995; Elkan 2001; Zadrozny 2004; Cortes et al. 2008), statistics (Whittemore 1978; Little and Rubin 1987; Robinson and Jewell 1991; Kuroki and Cai 2006; Evans and Didelez 2015), and the empirical sciences (e.g., genetics (Pirinen, Donnelly, and Spencer 2012; Mefford and Witte 2012), economics (Heckman 1979; Angrist 1997), and epidemiology (Robins 2001; Glymour and Greenland 2008)).

Even though selection and confounding biases appear together in most of the non-trivial, practical settings, they have been almost invariably treated independently in the literature. There are non-trivial interactions between them, however, which have not been investigated until recently. (Bareinboim, Tian, and Pearl 2014; Bareinboim and Tian 2015) provided sufficient conditions for the non-parametric recoverability of the causal effects from selection bias, and introduced a relaxation of this setting so that external (unbiased) data could be leveraged. (Evans and Didelez 2015) developed an approach for discrete models, where assumptions on the cardinality of the observable variables allow the estimation of the distribution over the sampling mechanism; in turn recovering the marginal distribution. (Correa and Bareinboim 2017) introduced a backdoor-like condition that controls for both biases, while (Correa, Tian, and Bareinboim 2018a) proved completeness for a more general backdoor criterion that allows for external data.

In this paper, we study the simultaneous effect of confounding and selection biases in general non-parametric settings. In particular, our contributions are as follow:

- We prove that the algorithm introduced in (Bareinboim and Tian 2015) is complete for the task of recoverability when all data available is biased. In other words, whenever the algorithm fails to recover a causal effect, the same is provable not recoverable by any other procedure.
- We relax the setting above and allow for the use of unbiased data in the form of a joint distribution over a subset of the observed variables. We develop a new algorithm for this task and prove that the approach is strictly more powerful than the current state-of-the-art method (Correa, Tian, and Bareinboim 2018a).

For the sake of space, the proofs not provided are available in the Appendix (Correa, Tian, and Bareinboim 2018b).

## Structural Models, Causal Effects, and Recoverability

The systematic analysis of confounding and selection biases requires a formal language where the characterization of the underlying data-generating model can be encoded explicitly. We use the language of Structural Causal Models (SCMs) (Pearl 2000, pp. 204-207). Formally, a SCM  $M$  is a 4-tuple  $\langle \mathbf{U}, \mathbf{V}, F, P(\mathbf{u}) \rangle$ , where  $\mathbf{U}$  is a set of exogenous (latent) variables and  $\mathbf{V}$  is a set of endogenous (measured) variables.  $F$  represents a collection of functions  $F = \{f_i\}$  such that each endogenous variable  $V_i \in \mathbf{V}$  is determined by a function  $f_i \in F$ , where  $f_i$  is a mapping from the respective domain of  $U_i \cup Pa_i$  to  $V_i$ ,  $U_i \subseteq \mathbf{U}$ ,  $Pa_i \subseteq \mathbf{V} \setminus V_i$ ,

and the entire set  $F$  forms a mapping from  $\mathbf{U}$  to  $\mathbf{V}$ . The uncertainty is encoded through a probability distribution over the exogenous variables,  $P(\mathbf{u})$ . Within the structural semantics, performing an action  $\mathbf{X}=\mathbf{x}$  is represented through the do-operator,  $do(\mathbf{X}=\mathbf{x})$ , which encodes the operation of replacing the original equation of  $\mathbf{X}$  by the constant  $\mathbf{x}$  and induces a submodel  $M_{\mathbf{x}}$ . For a detailed discussion of SCMs, causal inference and fusion, we refer readers to (Pearl 2000; Bareinboim and Pearl 2016).

Following the conventions in the field, we denote variables by capital letters and their realized values by small letters. Sets of variables are denoted in bold. We use typical graph-theoretic terminology with the abbreviations  $Pa(\mathbf{C})$ ,  $Ch(\mathbf{C})$ ,  $De(\mathbf{C})$ ,  $An(\mathbf{C})$ , which stand for the union of  $\mathbf{C}$  and respectively the parents, children, descendants, and ancestors of  $\mathbf{C}$ . The letter  $\mathcal{G}$  is used to refer to the causal graph, in which the unobserved common causes are encoded implicitly through the dashed bidirected arrows;  $\mathcal{G}_{\overline{\mathbf{XZ}}}$  denote the graph resulting from removing all incoming edges to  $\mathbf{X}$  and all outgoing edges from  $\mathbf{Z}$  in  $\mathcal{G}$ . For  $\mathbf{C} \subseteq \mathbf{V}$ , let  $\mathcal{G}_{\mathbf{C}}$  be the subgraph of  $\mathcal{G}$  composed only of variables in  $\mathbf{C}$ . Next, we formalize the notion of identifiability.

**Definition 1** (Effect Identifiability (Pearl 2000, pp.77)). *The causal effect of an action  $do(\mathbf{X}=\mathbf{x})$  on a set of variables  $\mathbf{Y}$  is said to be identifiable from  $P$  in  $\mathcal{G}$  if  $P(\mathbf{y}|do(\mathbf{x}))$  (for short,  $P_{\mathbf{x}}(\mathbf{y})$ ) is uniquely computable from  $P(\mathbf{v})$  in any model that induces  $\mathcal{G}$ . Formally, for every two models  $M_1$  and  $M_2$  compatible with  $\mathcal{G}$ ,  $P^{M_1}(\mathbf{v})=P^{M_2}(\mathbf{v})>0$  implies  $P^{M_1}(\mathbf{y}|do(\mathbf{x}))=P^{M_2}(\mathbf{y}|do(\mathbf{x}))$ .*

The systematic identification of causal effects calls for the ability to decompose them into easier-to-characterize quantities. For any set  $\mathbf{C} \subseteq \mathbf{V}$ , we then define  $Q[\mathbf{C}](\mathbf{v})$ , called *c-factor*, to denote the following function

$$Q[\mathbf{C}](\mathbf{v})=P_{\mathbf{v} \setminus \mathbf{c}}(\mathbf{c})=\sum_{\mathbf{U}} \prod_{\{i|V_i \in \mathbf{C}\}} P(v_i|pa_i, u_i)P(\mathbf{u}), \quad (1)$$

where  $pa_i$  is the set of observable parents of  $V_i$  and  $u_i$  is the set of unobserved parents. Of special interest are the c-factors associated with the elements of a partition on the observable variables induced by the presence of bidirected arrows, called C-Components (Tian and Pearl 2002a). The set  $\mathbf{V}$  is partitioned into c-components by assigning two variables to the same set if and only if they are connected by a path composed entirely of bidirected edges in  $\mathcal{G}$ .

While identification deals with the problem of controlling for confounding bias, an orthogonal problem arises when the observations are not a random sample from the population. This problem is what we referred to as *selection bias* (also called sampling selection bias).

**Definition 2** (Effect Recoverability (Bareinboim and Tian 2015)). *Given a causal graph  $\mathcal{G}$  augmented with the selection mechanism, represented by the  $S$  node, the causal effect  $P(\mathbf{y}|do(\mathbf{x}))$  is said to be recoverable from selection biased data if the assumptions embedded in  $\mathcal{G}$  render the effect expressible in terms of the distribution under selection,  $P(\mathbf{v}|S=1)$ . That is, for any models  $M_1$  and  $M_2$  compatible with  $\mathcal{G}$ ,  $P^{M_1}(\mathbf{v}|S=1)=P^{M_2}(\mathbf{v}|S=1)>0$  implies  $P^{M_1}(\mathbf{y}|do(\mathbf{x}))=P^{M_2}(\mathbf{y}|do(\mathbf{x}))$ .*

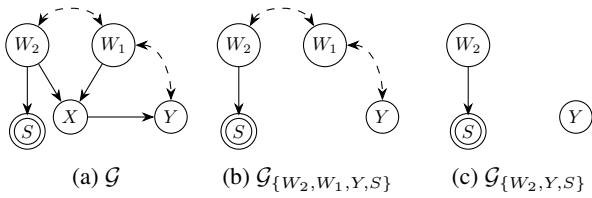


Figure 1: Subgraphs considered by RC while recovering  $P_x(y)$  for the model in (a).

Roughly speaking, the paths between an action  $\mathbf{X}$  and an outcome  $\mathbf{Y}$  in a causal graph can be partitioned into causal (directed paths) and non-causal (spurious). A path is called *proper* if it contains no variables in  $\mathbf{X}$  except at the starting point. The following construction graphically “disables” *proper* causal paths, by cutting the first arrow in such paths, leaving the spurious paths unperturbed.

**Definition 3** (Proper Backdoor Graph (van der Zander, Liskiewicz, and Textor 2014)). *Let  $\mathcal{G}$  be a causal diagram, and  $\mathbf{X}, \mathbf{Y}$  be disjoint subsets of variables. The proper backdoor graph, denoted as  $\mathcal{G}_{\mathbf{X}\mathbf{Y}}^{pbd}$ , is obtained from  $\mathcal{G}$  by removing the first edge of every proper causal path from  $\mathbf{X}$  to  $\mathbf{Y}$ .*

This transformation will allow us to characterize the failing condition for recoverability in the next section.

## Recoverability from Biased Data

In this section, we consider the problem of recovering the causal distribution when only biased data is available, namely, evaluating whether  $P_x(y)$  is computable from  $P(\mathbf{v}|S=1)$ . First, we consider the state-of-the-art sufficient procedure available in the literature, and then study the conditions under which it fails.

In order to recover a causal effect of the form  $P_x(y)$ , it is usually wise to express it as a product of c-factors associated with the c-components as follows:

$$\begin{aligned} P_x(y) &= \sum_{\mathbf{v} \setminus \mathbf{Y}} P_x(\mathbf{v} \setminus \mathbf{x}) = \sum_{\mathbf{v} \setminus \mathbf{Y}} Q[\mathbf{V} \setminus \mathbf{X}] \\ &= \sum_{\mathbf{D} \setminus \mathbf{Y}} Q[\mathbf{D}] = \sum_{\mathbf{D} \setminus \mathbf{Y}} \prod_{i=1}^l Q[D_i], \end{aligned} \quad (2)$$

where  $\mathbf{D} = An(\mathbf{Y})_{\mathcal{G}_{\mathbf{v} \setminus \mathbf{x}}}$ , and  $D_1, \dots, D_l$  are the c-components of  $\mathcal{G}_{\mathbf{D}}$ .

This factorization was employed as the basis for the algorithm RC (Bareinboim and Tian 2015), shown in Alg. 1. RC attempts to recover each  $Q[D_i]$ , by Lemma 3 in (Bareinboim and Tian 2015), every  $Q[C_i]$  in line 2 is recoverable, and the function IDENTIFY( $\mathbf{E}, C_i, Q[C_i]$ ) (Tian 2002) (line 4) can be used to determine the identifiability of  $Q[\mathbf{E}]$  from  $Q[C_i]$ , where  $\mathbf{E} \subseteq C_i$ . If all such factors are successfully recovered, then the effect  $P_x(y)$  is recoverable as (2). To understand the mechanics of the algorithm, we consider the model in Fig. 1(a) and assume our target distribution is  $P_x(y)$ . In this graph  $\mathbf{D} = \{Y\}$  hence  $P_x(y) = Q[Y]$ , consequently  $RC(\{Y\}, P(\mathbf{v}|S=1), \mathcal{G})$  will be invoked. Since all variables in  $\mathcal{G}$  are ancestors of  $Y$  or  $S$ , line 1’s condition does not apply, and RC iterates over each c-component of  $\mathcal{G}$ , adding

**Algorithm 1** Procedure in (Bareinboim and Tian 2015) for recovering  $Q[\mathbf{E}]$

**function** RC( $\mathbf{E}, P, \mathcal{G}$ )

*Input*  $\mathbf{E}$  a c-component,  $P$  a distribution and  $\mathcal{G}$  a causal diagram over variables  $\mathbf{V}$  and  $S$ .

*Output* Expression for  $Q[\mathbf{E}]$  in terms of  $P(\mathbf{v}|S=1)$  or FAIL

- 1: If  $\mathbf{V} \setminus (An(\mathbf{E}) \cup An(S)) \neq \emptyset$ ,  
return RC( $\mathbf{E}, \sum_{\mathbf{v} \setminus (An(\mathbf{E}) \cup An(S))} P, \mathcal{G}_{An(\mathbf{E}) \cup An(S)}$ )
- 2: Let  $C_1, \dots, C_k$  be the c-components of  $\mathcal{G}$  that contains no ancestors of  $S$ , and let  $\mathbf{C} = \bigcup_i C_i$
- 3: If  $\mathbf{C} = \emptyset$ , return FAIL
- 4: If  $\mathbf{E}$  is a subset of some  $C_i$ ,  
return IDENTIFY( $\mathbf{E}, C_i, Q[C_i]$ )
- 5: Return RC( $\mathbf{E}, \frac{P}{\prod_i Q[C_i]}, \mathcal{G}_{\mathbf{v} \setminus \mathbf{C}}$ )

those with no ancestor of  $S$  to the set  $\mathbf{C}$ . In this example,  $\mathcal{G}$  decomposes into three c-components:  $\{X\}$ ,  $\{W_1, W_2, Y\}$ , and  $\{S\}$ , where only  $\{X\}$  satisfies the condition to get into  $\mathbf{C}$  and  $Q[X]$  is recovered as  $P(x|w_1, w_2, S=1)$ . Since  $\{Y\}$  is not a subset of  $\{X\}$ , line 5 recursively calls  $RC(\{Y\}, P(\mathbf{v}|S=1)/P(x|w_1, w_2, S=1), \mathcal{G}_{\{W_2, W_1, Y, S\}})$ . This new graph is shown in Fig.1(b). Now that  $X$  is not in the graph, the variable  $W_1$  is no longer an ancestor of either  $Y$  or  $S$ , then line 1 performs a recursive call as  $RC(\{Y\}, \sum_{W_1} P(\mathbf{v}|S=1)/P(x|w_1, w_2, S=1), \mathcal{G}_{\{W_2, Y, S\}})$ . In the graph  $\mathcal{G}_{\{W_2, Y, S\}}$ , shown in Fig. 1(c), there are three c-components:  $\{W_2\}$ ,  $\{Y\}$ , and  $\{S\}$ . Since  $Y$  is not an ancestor of  $S$  in this graph, line 2 will recover  $Q[C_1]$  where  $C_1 = \{Y\}$  as  $\sum_{W_1} P(y|x, w_1, w_2|S=1)P(w_1, w_2|S=1)/P(w_2|S=1)$  and make  $\mathbf{C} = \{Y\}$ . Next, because our target  $\{Y\}$  is a subset of  $C_1 = \{Y\}$ , line 4 recovers  $Q[Y] = Q[C_1]$  and returns it, which, as noted before, corresponds to  $P_x(y)$ .

While RC was shown to be sound, it was not shown to be complete, that is, whether a FAIL triggered by line 3 implies that the target causal effect is not recoverable, or if the algorithm is not powerful enough to recover the expression.

In the following, we first present a necessary condition for the causal effect to be recoverable and then use it to show the completeness of the procedure RC.

**Theorem 1.** *Let  $\mathbf{X}, \mathbf{Y} \subset \mathbf{V}$  be two disjoint sets of variables and  $\mathcal{G}$  a causal diagram over  $\mathbf{V}$  and  $S$ . If  $(\mathbf{Y} \not\perp\!\!\!\perp S)_{\mathcal{G}_{\mathbf{X}\mathbf{Y}}^{pbd}}$ , then  $P_x(y)$  is not recoverable from  $P(\mathbf{v} | S=1)$  in  $\mathcal{G}$ .*

The necessary condition in Thm. 1 helps us to show that when RC fails,  $P_x(y)$  is not recoverable.

**Theorem 2.** *Let  $\mathbf{X}$  and  $\mathbf{Y}$  be two disjoint sets of variables and  $\mathcal{G}$  a causal diagram over  $\mathbf{V}$  and  $S$ . Let  $\mathbf{D} = An(\mathbf{Y})_{\mathcal{G}_{\mathbf{v} \setminus \mathbf{x}}}$  and  $D_1, \dots, D_\ell$  be the c-components of  $\mathcal{G}_{\mathbf{D}}$ . Then, the effect  $P_x(y)$  is recoverable from  $P(\mathbf{v}|S=1)$  if and only if each  $D_i, i = 1, \dots, \ell$  is recoverable by the function RC.*

Thm. 2 implies that the strategy employed by RC covers all recoverability scenarios, and all other algorithms concerned with this setting will be in some form or shape, at most, equivalent to it. In other words, the recoverability algorithm in (Bareinboim and Tian 2015) is complete.

## Recoverability with External Data

Whenever the conditions of Thm. 2 are not satisfied, the target effect is provably not inferable from  $P(\mathbf{v}|S=1)$ . One common strategy to circumvent this challenging situation is to try to find and leverage alternative sources of data. Popular baseline covariates such as age, sex, and ethnicity can be obtained without bias in many cases, for instance, using data from the census or smaller pilot studies.

We supplement Def. 2 to formally account for the availability of a new source of data, i.e.,

**Definition 4** (Recoverability from Selection Bias with External Data). *Given a causal graph  $\mathcal{G}$  augmented with the selection mechanism, represented by the  $S$  node, the causal effect  $P_{\mathbf{x}}(\mathbf{y})$  is said to be recoverable from selection bias with external data  $P(\mathbf{t})$  if for any two models  $M_1$  and  $M_2$  compatible with  $\mathcal{G}$ ,  $P^{M_1}(\mathbf{v}|S=1) = P^{M_2}(\mathbf{v}|S=1) > 0$  and  $P^{M_1}(\mathbf{t}) = P^{M_2}(\mathbf{t}) > 0$  implies  $P_{\mathbf{x}}^{M_1}(\mathbf{y}) = P_{\mathbf{x}}^{M_2}(\mathbf{y})$ .*

In other words, Def. 4 requires the causal effect to be uniquely computable from the available data (under selection bias and from the external source) and the assumptions embodied in the causal model.

We consider unbiased external data in the form of a distribution  $P(\mathbf{t}^0)$ , where  $\mathbf{T}^0 \subset \mathbf{V}$  is a set of variables measured (jointly) without bias. As shown in the next lemma, additional information can be inferred from the external data and model assumptions.

**Lemma 1.** *Given  $P(\mathbf{t}^0)$ , let  $\mathbf{T}'$  be a set of variables such that  $(S \perp\!\!\!\perp \mathbf{T}' \mid \mathbf{T}^0)$ , and let  $\mathbf{T} = \mathbf{T}^0 \cup \mathbf{T}'$ , then  $P(\mathbf{t})$  is recoverable.*

*Proof.*  $P(\mathbf{t}) = P(\mathbf{t}' \mid \mathbf{t}^0, S=1)P(\mathbf{t}^0)$ .  $\square$

From this point on, we will use  $\mathbf{T}$  to denote a set of variables such that  $P(\mathbf{t})$  is available (following from  $P(\mathbf{t}^0)$ ), and let  $\mathbf{R} = \mathbf{V} \setminus \mathbf{T}$  be the rest of the variables. Let  $PJ(\mathcal{G}, \mathbf{T})$  denote the graph derived from the original graph  $\mathcal{G}$  by representing the variables in  $\mathbf{R}$  as unobservables (with bidirected edges), known as the *projection* of  $\mathcal{G}$  on the set  $\mathbf{T}$  (Verma 1993) (see also Def. 1 in (Tian and Pearl 2002b)). Accordingly, we can define c-factors  $Q_R[\cdot]$  in this projection, denoting the following function

$$\begin{aligned} Q_R[\mathbf{C}] &= P_{\mathbf{t} \setminus \mathbf{c}}(\mathbf{c}) = P_{\mathbf{v} \setminus (\mathbf{c} \cup \mathbf{R})}(\mathbf{c}) \\ &= \sum_{\mathbf{U}, \mathbf{R}} \prod_{\{i \mid V_i \in \mathbf{C} \cup \mathbf{R}\}} P(v_i \mid pa_i, u_i) P(\mathbf{u}). \end{aligned} \quad (3)$$

In other words, the function  $Q_R[\cdot]$  represents a c-factor in  $\mathcal{G}$  when the variables in  $\mathbf{R}$  are treated as latent variables<sup>1</sup>.

The next result delineates the new c-factors that can be recovered from  $P(\mathbf{t})$ :

**Lemma 2.** *Let  $\mathbf{T} \subseteq \mathbf{V}$ ,  $\mathbf{R} = \mathbf{V} \setminus \mathbf{T}$ , and  $T_1, \dots, T_m$  be the c-components of  $PJ(\mathcal{G}, \mathbf{T})$ , then all  $Q_R[T_k]$  are recoverable from  $P(\mathbf{t})$ .*

<sup>1</sup>C-components with arbitrary variables as latent variables are defined in (Tian and Pearl 2002b).

*Proof.* We have that

$$P(\mathbf{t}) = \sum_{\mathbf{R}} P(\mathbf{v}) = \prod_{k=1}^m Q_R[T_k]. \quad (4)$$

By (Tian and Pearl 2002b, Lemma 2), all  $Q_R[T_k]$  are recoverable from  $P(\mathbf{t})$ .  $\square$

Building on Lemmas 1 and 2, we now state the main result of this section:

**Theorem 3.** *Let  $\mathbf{H} \subseteq \mathbf{V} \cup \{S\}$ , such that  $\mathbf{H}$  is partitioned into c-components  $H_1, \dots, H_l, H_s$  in the subgraph  $\mathcal{G}_{\mathbf{H}}$ , where  $S \in H_s$ . Assume*

$$f(P(\mathbf{v} \mid S=1)) = \frac{Q[H_s](\mathbf{v}, S=1)}{P(S=1)} \prod_i Q[H_i], \quad (5)$$

where  $f(P(\mathbf{v}|S=1))$  is some recoverable quantity, and  $P(\mathbf{t})$  is available. Let  $\mathbf{T}_{\mathbf{H}}^0 = \mathbf{T} \setminus De(\mathbf{V} \setminus \mathbf{H})_{\mathcal{G}}$  and  $\mathbf{T}'$  be the set of all variables in  $\mathbf{H}$  such that  $(\mathbf{T}' \perp\!\!\!\perp S \mid \mathbf{T}_{\mathbf{H}}^0)_{\mathcal{G}_{\mathbf{H}}}$ . Also, let  $\mathbf{T}_{\mathbf{H}} = \mathbf{T}_{\mathbf{H}}^0 \cup \mathbf{T}'$  and let  $\mathbf{R}_{\mathbf{H}} = \mathbf{H} \setminus \mathbf{T}_{\mathbf{H}}$ . Then, for  $j=1, \dots, l$ ,  $Q[H_j]$  is recoverable if  $H_j$  contains no variables that are both ancestors of  $H_s$  and belong to  $\mathbf{R}_{\mathbf{H}}$  or its children (i.e.  $H_j \cap An(H_s) \cap Ch(\mathbf{R}_{\mathbf{H}}) = \emptyset$ ) in  $\mathcal{G}_{\mathbf{H}}$ .

*Proof.* (sketch, see Appendix C for details) Let a topological order of the variables in  $\mathbf{H}$  be  $V_{h_1} < \dots < V_{h_k}$  in  $\mathcal{G}_{\mathbf{H}}$ . Let  $H^{\leq i} = \{V_{h_1}, \dots, V_{h_i}\}$  be the set of variables in  $\mathbf{H}$  ordered before  $V_{h_i}$  (including  $V_{h_i}$ ), and  $H^{> i} = \mathbf{H} \setminus H^{\leq i}$  for  $i = 1, \dots, k$ , and  $H^{\leq 0} = \emptyset$ . The assumptions of the theorem allow us to recover  $P_{\mathbf{v} \setminus \mathbf{h}}(\mathbf{t}_{\mathbf{H}})$  from  $f(P(\mathbf{v}|S=1))$  and  $P(\mathbf{t})$ . For any  $H_j$  that satisfies the condition of the theorem, the associated c-factor can be recovered as:

$$\begin{aligned} Q[H_j] &= \prod_{\{i \mid V_{h_i} \in H_j \cap An(H_s)\}} \frac{\sum_{h > i \cap \mathbf{T}_{\mathbf{H}}} P_{\mathbf{v} \setminus \mathbf{h}}(\mathbf{t}_{\mathbf{H}})}{\sum_{h > i - 1 \cap \mathbf{T}_{\mathbf{H}}} P_{\mathbf{v} \setminus \mathbf{h}}(\mathbf{t}_{\mathbf{H}})} \times \\ &\quad \prod_{\{i \mid V_{h_i} \in H_j \setminus An(H_s)\}} \frac{\sum_{h > i} f(P(\mathbf{v}|S=1))}{\sum_{h > i, V_{h_i}} f(P(\mathbf{v}|S=1))}. \end{aligned} \quad (6)$$

$\square$

Thm. 3 will be the main driving force for recovering causal effects from combined biased data  $P(\mathbf{v}|S=1)$  and unbiased data  $P(\mathbf{t})$ . To give an example of how this result can be used, consider the model in Fig 2(a) and assume we have external data over  $\mathbf{T}^0 = \{Z\}$ . Then,  $\mathbf{T} = \{Z, X, Y\}$  because  $(S \perp\!\!\!\perp X, Y \mid Z)$ ,  $\mathbf{R} = \mathbf{V} \setminus \mathbf{T} = \{R, W\}$ , and  $H_s = \{S, Z\}$  which is the c-component that contains  $S$ . Also, the biased distribution factorizes as follows:

$$P(\mathbf{v} \mid S=1) = \frac{Q[S, Z]}{P(S=1)} Q[W] Q[R] Q[X] Q[Y]. \quad (7)$$

Thm. 3 would allow us to recover  $Q[X]$  and  $Q[Y]$  since they do not contain any ancestor of  $H_s$ .

## Recovering Causal Effects Systematically

In order to recover the causal distribution  $P_{\mathbf{x}}(\mathbf{y})$  systematically, (Bareinboim and Tian 2015) proposed a strategy that recovers each  $Q[D_i]$  in Eq. (2) one by one. It turns out that

when external data  $P(\mathbf{t})$  is available, each  $Q[D_i]$  being recoverable is no longer necessary for the overall recoverability of  $P_{\mathbf{x}}(\mathbf{y})$ . To witness, let us follow up on the example from Fig 2(a), introduced at the end of the last section. Following the strategy dictated by Eq. (2), we note that

$$P_{\mathbf{x}}(\mathbf{y}) = \sum_{Z,R} Q[Y, Z, R] = \sum_{Z,R} Q[Y]Q[Z]Q[R]. \quad (8)$$

Thm. 3 licenses the recoverability of  $Q[Y]$ , but it is not difficult to show that neither  $Q[R]$  nor  $Q[Z]$  is recoverable. Perhaps surprisingly, however,  $P_{\mathbf{x}}(\mathbf{y})$  can be recovered as

$$\sum_Z Q[Y] \sum_R Q[R]Q[Z] = \sum_Z P(y|x, z, S=1)P(z). \quad (9)$$

The key observation here is that while  $Q[R]$  and  $Q[Z]$  are not recoverable individually,  $\sum_R Q[R]Q[Z]$  is, in fact, a function of  $Z$  and equal to  $Q_R[Z]$  (see Eq. 3), which can be recovered from  $P(\mathbf{t}) = P(z, x, y)$  as  $P(z)$  via Lemma 2.

To formally account for this situation, we re-write the causal effect in Eq. (2) by splitting  $\mathbf{D} \setminus \mathbf{Y}$  into two parts:  $\mathbf{A} = (\mathbf{D} \setminus \mathbf{Y}) \cap \mathbf{T}$  and  $\mathbf{B} = (\mathbf{D} \setminus \mathbf{Y}) \cap \mathbf{R}$  where  $\mathbf{R} = \mathbf{V} \setminus \mathbf{T}$ , and then we treat elements in  $\mathbf{B}$  as latent variables while defining c-factors  $Q_B[\cdot]$  in the resulting projected graph  $PJ(\mathcal{G}_{\mathbf{D}}, \mathbf{D} \setminus \mathbf{B})$ , as follows:

$$P_{\mathbf{x}}(\mathbf{y}) = \sum_{\mathbf{D} \setminus \mathbf{Y}} \prod_{i=1}^l Q[D_i] = \sum_{\mathbf{A}} \prod_{j=1}^{\ell} Q_B[C_j], \quad (10)$$

where  $\mathbf{D} = An(\mathbf{Y})_{\mathcal{G}_{\mathbf{V} \setminus \mathbf{X}}}$ ,  $D_1, \dots, D_l$  are the c-components of  $\mathcal{G}_{\mathbf{D}}$ ,  $C_1, \dots, C_{\ell}$  are the c-components of  $PJ(\mathcal{G}_{\mathbf{D}}, \mathbf{D} \setminus \mathbf{B})$ , and c-factors  $Q_B[C_j]$  are defined as

$$Q_B[C_j] = \sum_{\mathbf{U} \cup \mathbf{B}} \prod_{\{i|V_i \in C_j \cup \mathbf{B}\}} P(v_i|pa_i, u_i)P(\mathbf{u}). \quad (11)$$

$Q_B[C_j]$  could also be expressed in terms of  $Q[D_i]$  in the following form:

$$Q_B[C_j] = \sum_{B_j} \prod_{\{i|D_i \in F_j\}} Q[D_i], \quad (12)$$

where  $B_j$  are disjoint and possibly empty sets such that  $\cup_{j=1}^{\ell} B_j = \mathbf{B}$ , and  $F_1, \dots, F_{\ell}$  form a partition of  $\{D_1, \dots, D_l\}$ .

Under certain conditions, a c-factor  $Q_B[C_j]$  may be equal to the c-factor  $Q_R[C_j]$ , defined in  $PJ(\mathcal{G}, \mathbf{T})$ , which is potentially recoverable in terms of the unbiased distribution  $P(\mathbf{t})$ .

**Lemma 3.** *Let  $C_j$  be a c-component of  $PJ(\mathcal{G}_{\mathbf{D}}, \mathbf{D} \setminus \mathbf{B})$ . If  $\mathbf{B} \cap Pa(C_j) = \mathbf{R} \cap Pa(C_j)$ , then*

- (i)  $Q_B[C_j] = Q_R[C_j]$ , where  $Q_R[C_j]$  is a c-factor in  $PJ(\mathcal{G}, \mathbf{T})$  as defined in Eq. (3); and
- (ii) Let  $T_1, \dots, T_m$  be the c-components of  $PJ(\mathcal{G}, \mathbf{T})$ , then  $C_j$  must be a subset of some  $T_k$ .

*Proof.* (i) Let  $\hat{B}_j = \mathbf{B} \cap Pa(C_j)$ . Any variable in  $\mathbf{B}$  that is not in  $\hat{B}_j$  does not appear in the expression in (11), and can be summed out, leading to

$$Q_B[C_j] = \sum_{\mathbf{U} \cup \hat{B}_j} \prod_{\{i|V_i \in C_j \cup \hat{B}_j\}} P(v_i|pa_i, u_i)P(\mathbf{u}). \quad (13)$$

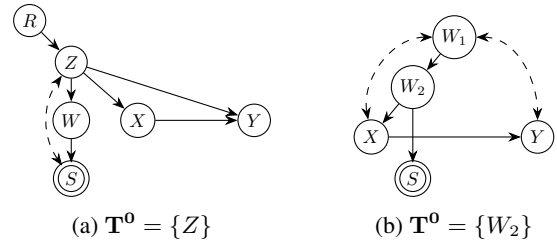


Figure 2: Examples of recoverability tasks for the effect  $P_{\mathbf{x}}(\mathbf{y})$ . Model in (a) can be recovered with external data on  $Z$ . Model in (b) is recoverable with external data on  $W_2$  or  $W_1$ .

Similarly, from (3) we have:

$$Q_R[C_j] = \sum_{\mathbf{U} \cup \mathbf{R}} \prod_{\{i|V_i \in C_j \cup \mathbf{R}\}} P(v_i|pa_i, u_i)P(\mathbf{u}).$$

Let  $\hat{R}_j = \mathbf{R} \cap Pa(C_j)$ . Then any variable in  $\mathbf{R}$  that is not in  $\hat{R}_j$  can be summed out, leading to

$$Q_R[C_j] = \sum_{\mathbf{U} \cup \hat{R}_j} \prod_{\{i|V_i \in C_j \cup \hat{R}_j\}} P(v_i|pa_i, u_i)P(\mathbf{u}). \quad (14)$$

It is clear that if  $\hat{R}_j = \hat{B}_j$ , then (14) is equal to (13).

(ii) Since  $\mathbf{D} \subseteq \mathbf{V}$  and  $\mathbf{B} \subseteq \mathbf{R}$ , a c-component of  $PJ(\mathcal{G}_{\mathbf{D}}, \mathbf{D} \setminus \mathbf{B})$  must be a subset of a c-component of  $PJ(\mathcal{G}, \mathbf{T})$ .  $\square$

The importance of Lemma 3 stems from the fact that  $Q_R[C_j]$  is potentially identifiable in  $PJ(\mathcal{G}, \mathbf{T})$  from the unbiased distribution  $P(\mathbf{t})$  based on Lemma 2. Specifically, we can use IDENTIFY( $C_j, T_k, Q_R[T_k]$ ) to try to recover  $Q_B[C_j] = Q_R[C_j]$ . If  $Q_R[C_j]$  is not identifiable from  $P(\mathbf{t})$ , then we further attempt to recover  $Q_B[C_j]$  by recovering each  $Q[D_i]$  in Eq. (12) factor by factor.

To recover an individual  $Q[D_i]$ , it turns out the RC algorithm (Alg. 1) is not complete anymore in our setting (even if line 2 of RC is enhanced with Thm. 3). Extending RC, we develop a new algorithm called RCE (Alg. 2) to recover any target c-component  $Q[\mathbf{E}]$ . RCE attempts to systematically recover  $Q[\mathbf{E}]$  by recovering, using Thm. 3 (line 2), the c-component  $Q[C_i]$  of  $\mathcal{G}$  that contains  $\mathbf{E}$ , and then call the function IDENTIFY to recover  $Q[\mathbf{E}]$  from  $Q[C_i]$  (line 3a). To facilitate this, RCE reduces the problem to simpler subgraphs, by removing irrelevant non-ancestors (line 1) or other recoverable c-components (line 3b and line 4) from the current graph. These other c-components are recovered either by Thm. 3 (line 2) or by recursively calling RCE (line 4). Due to the recursive nature of the process, RCE may try to compute a c-component more than once, which can be avoided by keeping track of the previous queries. For simplicity we omit these practical details.

Putting these results together, we develop a general, systematic procedure for recovering causal effects called IDSB. The function IDSB in Alg. 3 accepts as input two disjoint sets  $\mathbf{X}, \mathbf{Y}$ , distributions  $P(\mathbf{v}|S=1), P(\mathbf{t}^0)$ , and a causal diagram  $\mathcal{G}$ ; it outputs an expression for  $P_{\mathbf{x}}(\mathbf{y})$  in terms of the

---

**Algorithm 2** Recursive function used to recover an arbitrary c-component

---

**function** RCE( $\mathbf{E}, \mathcal{P}, \mathcal{G}$ )

*Input*  $\mathbf{E}$  a set of variables such that  $\mathbf{E}$  is a c-component,  $\mathcal{P}$  a distribution over  $\mathbf{V}$ ,  $\mathcal{G}$  a causal diagram over variables  $\mathbf{V}$  and  $S$ .  
 $P^*(\mathbf{t})$  a distribution over  $\mathbf{T}$  and  $\mathcal{G}^*$  the original graph over variables  $\mathbf{V}^*$  and  $S$  are defined globally.

*Output* Expression for  $Q[\mathbf{E}]$  or FAIL

- 1: Let  $\mathbf{W} = An(\mathbf{E}) \cup An(S)$ . If  $\mathbf{V} \setminus \mathbf{W} \neq \emptyset$ ,  
return RCE( $\mathbf{E}, \sum_{\mathbf{V} \setminus \mathbf{W}} \mathcal{P}, \mathcal{G}_{\mathbf{W}}$ )
  - 2: Let  $C_1, \dots, C_k$  be the c-components of  $\mathcal{G}$  that are recoverable by Thm. 3 (with  $f(P(\mathbf{v}|S=1)) = \mathcal{P}$  and  $P(\mathbf{t}) = P^*(\mathbf{t})$ ).  
Let  $\mathbf{C} = \bigcup_i C_i$
  - 3: If  $\mathbf{C} \neq \emptyset$ ,
    - (a) If  $\mathbf{E}$  is a subset of some  $C_i$ ,  
then return IDENTIFY( $\mathbf{E}, C_i, Q[C_i]$ )
    - (b) Return RCE( $\mathbf{E}, \frac{\mathcal{P}}{\prod_i Q[C_i]}, \mathcal{G}_{(\mathbf{V} \cup \{S\}) \setminus \mathbf{C}}$ )
  - 4: For each c-component  $B_i$  of  $\mathcal{G}$  that does not contain  $\mathbf{E}$  such that  $\mathbf{Z} = \mathbf{V} \setminus (An(S) \cup An(B_i)) \neq \emptyset$ :  
 $Q[B_i] = \text{RCE}(B_i, \sum_{\mathbf{Z}} \mathcal{P}, \mathcal{G}_{(\mathbf{V} \cup \{S\}) \setminus \mathbf{Z}})$   
If  $Q[B_i] \neq \text{FAIL}$ , return RCE( $\mathbf{E}, \frac{\mathcal{P}}{Q[B_i]}, \mathcal{G}_{(\mathbf{V} \cup \{S\}) \setminus B_i}$ )
  - 5: Return FAIL
- 

input distributions or FAIL. IDSB starts by simplifying the model via removing irrelevant non-ancestors (line 1) and recovering  $P(\mathbf{t})$  using Lemma 1 (lines 2, 3). IDSB then recovers  $P_{\mathbf{x}}(\mathbf{y})$  using Eq. (10) by recovering each  $Q_B[C_j]$  (line 5). For each  $Q_B[C_j]$ , IDSB first attempts to recover it from  $P(\mathbf{t})$  based on Lemma 2 by calling the function IDENTIFY if the condition in Lemma 3 is satisfied. If this fails, IDSB tries to recover  $Q_B[C_j]$  using (12) by calling RCE for each  $Q[D_i]$ . The next theorem states that IDSB is sound.

**Theorem 4.** *The procedure IDSB is sound.*

Due to page limits, we provide the proof of Thm. 4 in the Appendix D. Nevertheless, we will illustrate its mechanics using the example from Fig. 3(a) where we assume  $P(\mathbf{v}|S=1)$  and  $P(\mathbf{t}^0)$  are given, with  $\mathbf{T}^0 = \{V_2, V_3, V_6\}$ , and the goal is to recover  $P_{\mathbf{x}}(\mathbf{y})$ . Initially, in line 1  $\mathbf{W} = \mathbf{V}$ . Line 2 finds that  $(S \perp\!\!\!\perp X, V_1 | \mathbf{T}^0)$ , hence  $\mathbf{T} = \{X, V_1, V_2, V_3, V_6\}$  and line 3 recovers  $P(\mathbf{t}) = P(x, v_1 | v_2, v_3, v_6, S=1)P(v_2, v_3, v_6)$ . At line 4, we have  $\mathbf{D} = \{V_5, V_6, Y\}$ ,  $\mathbf{R} = \{V_4, V_5, Y\}$ ,  $\mathbf{A} = \{V_6\}$ , and  $\mathbf{B} = \{V_5\}$ . The graphs  $\mathcal{G}_{\mathbf{D}}$ ,  $PJ(\mathcal{G}, \mathbf{T})$ , and  $PJ(\mathcal{G}_{\mathbf{D}}, \mathbf{D} \setminus \mathbf{B})$  (Fig. 3(b), (c), and (d) respectively) are derived from  $\mathcal{G}$  (Fig. 3(a)). The table in Fig. 3(i) summarizes the decomposition of these graphs and recalls how each c-component and c-factor are denoted by IDSB in line 4. At this point, we know from Eq. (10) that  $P_{\mathbf{x}}(\mathbf{y}) = \sum_{V_6} Q_B[Y]Q_B[V_6]$ . Also,  $Q_B[Y] = Q[Y]$ ,  $Q_B[V_6] = \sum_{V_5} Q[V_5, V_6]$ , corresponding to  $B_1 = \emptyset$ ,  $B_2 = \{V_5\}$ ,  $F_1 = \{\{Y\}\}$  and  $F_2 = \{\{V_5, V_6\}\}$ . Clearly  $\mathbf{B} = B_1 \cup B_2$  and  $F_1, F_2$  constitute a partition over  $\{D_1, D_2\}$ .

Continuing with line 5, the algorithm considers the first c-component  $C_1 = \{Y\}$ , and since  $\mathbf{B} \cap Pa(Y) = \emptyset \neq \mathbf{R} \cap Pa(Y) = \{Y\}$ , it calls RCE to try to recover  $Q[Y]$

---

**Algorithm 3** Algorithm capable of recovering  $P_{\mathbf{x}}(\mathbf{y})$  from selection bias with external data

---

**function** IDSB( $\mathbf{X}, \mathbf{Y}, P, P(\mathbf{t}^0), \mathcal{G}$ )

*Input*  $\mathbf{X}, \mathbf{Y}$  disjoint sets of variables,  $P(\mathbf{v}|S=1)$  a distribution,  $P(\mathbf{t}^0)$  distribution over a set of variables  $\mathbf{T}^0$ , and  $\mathcal{G}$  a causal diagram over variables  $\mathbf{V}$  and  $S$

*Output* Expression for  $P_{\mathbf{x}}(\mathbf{y})$  in terms of  $P(\mathbf{v}|S=1)$  and  $P(\mathbf{t}^0)$  or FAIL

- 1: Let  $\mathbf{W} = An(\mathbf{Y}) \cup An(S)$ ,  $\mathcal{G} \leftarrow \mathcal{G}_{\mathbf{W}}$ ,  $P \leftarrow \sum_{\mathbf{V} \setminus \mathbf{W}} P$
  - 2: Let  $\mathbf{T}' \subset \mathbf{W}$  be the set of all the variables such that  $(S \perp\!\!\!\perp \mathbf{T}' | \mathbf{T}^0 \cap \mathbf{W})_{\mathcal{G}}$ , and  $\mathbf{T} = \mathbf{T}' \cup (\mathbf{T}^0 \cap \mathbf{W})$
  - 3: Recover  $P(\mathbf{t})$  by Lemma 1
  - 4: Let  $\mathbf{D} = An(\mathbf{Y})_{\mathcal{G}_{\mathbf{W} \setminus \mathbf{X}}}$ ,  
Let  $D_1, \dots, D_l$  be the c-components of  $\mathcal{G}_{\mathbf{D}}$ ,  
Let  $T_1, \dots, T_m$  be the c-components of  $PJ(\mathcal{G}, \mathbf{T})$ ,  
 $\mathbf{R} = \mathbf{W} \setminus \mathbf{T}$ ,  $\mathbf{A} = (\mathbf{D} \setminus \mathbf{Y}) \cap \mathbf{T}$ ,  $\mathbf{B} = (\mathbf{D} \setminus \mathbf{Y}) \cap \mathbf{R}$ ,  
Let  $C_1, \dots, C_\ell$  be the c-components of  $PJ(\mathcal{G}_{\mathbf{D}}, \mathbf{D} \setminus \mathbf{B})$ , such that  $Q_B[C_j]$  is given by Eq. (12).
  - 5: For each  $C_j$   
If  $\mathbf{B} \cap Pa(C_j) = \mathbf{R} \cap Pa(C_j)$  then  
Assume  $C_j$  is a subset of  $T_k$   
 $Q_B[C_j] = \text{IDENTIFY}(C_j, T_k, Q_R[T_k])$   
If  $\mathbf{B} \cap Pa(C_j) \neq \mathbf{R} \cap Pa(C_j)$  or  $Q_B[C_j] = \text{FAIL}$ , then  
 $Q_B[C_j] = \sum_{B_j} \prod_{i, D_i \in F_j} \text{RCE}(D_i, P, \mathcal{G})$   
If  $Q_B[C_j] = \text{FAIL}$ , then return FAIL
  - 6: Return  $\sum_{\mathbf{A}} \prod_{j=1}^{\ell} Q_B[C_j]$
- 

(which is equal to  $Q_B[Y]$ ) in the graph  $\mathcal{G}$ . The recursion induced by this call to RCE is depicted in Fig. 4, where each edge is annotated with the line number (in RCE) that initiates the call and Fig. 3(e)-(h) contain the relevant subgraphs. Each  $\mathcal{P}^{(i)}$ ,  $i=0, \dots, 4$  stands for the distributions associated with the corresponding subgraph, obtained as follows

$$\mathcal{P}^{(0)} = P(\mathbf{v}|S=1), \quad (15)$$

$$\mathcal{P}^{(1)} = \mathcal{P}^{(0)} / Q[V_2], \text{ where} \quad (16)$$

$$Q[V_2] = \sum_{X, V_3, V_6} P(\mathbf{t}) / \sum_{X, V_3, V_6, V_2} P(\mathbf{t}), \quad (17)$$

$$\mathcal{P}^{(2)} = \sum_{V_1} \mathcal{P}^{(1)}, \quad (18)$$

$$\mathcal{P}^{(3)} = \mathcal{P}^{(2)} / Q[X], \text{ where} \quad (19)$$

$$Q[X] = \sum_{V_4, V_5, Y} \mathcal{P}^{(2)} / \sum_{V_4, V_5, X, Y} \mathcal{P}^{(2)}, \text{ and} \quad (20)$$

$$\mathcal{P}^{(4)} = \sum_{V_3, V_4} \mathcal{P}^{(3)}. \quad (21)$$

Finally, the result returned by RCE is:

$$Q[Y] = \mathcal{P}^{(4)} / \sum_Y \mathcal{P}^{(4)}. \quad (22)$$

After  $Q[Y]$  is computed, IDSB moves on to  $C_2 = \{V_6\}$ . Since  $\mathbf{B} \cap Pa(V_6) = \{V_5\} = \mathbf{R} \cap Pa(V_6)$ , we have that  $Q_B[V_6]$  is equal to  $Q_R[V_6]$  which is potentially identifiable from  $Q_R[T_2]$  where  $T_2 = \{V_3, V_6\}$ . Next, IDSB calls IDENTIFY( $\{V_6\}, \{V_3, V_6\}, Q_R[T_2]$ ) to obtain  $Q_B[V_6] = P(v_6)$ .

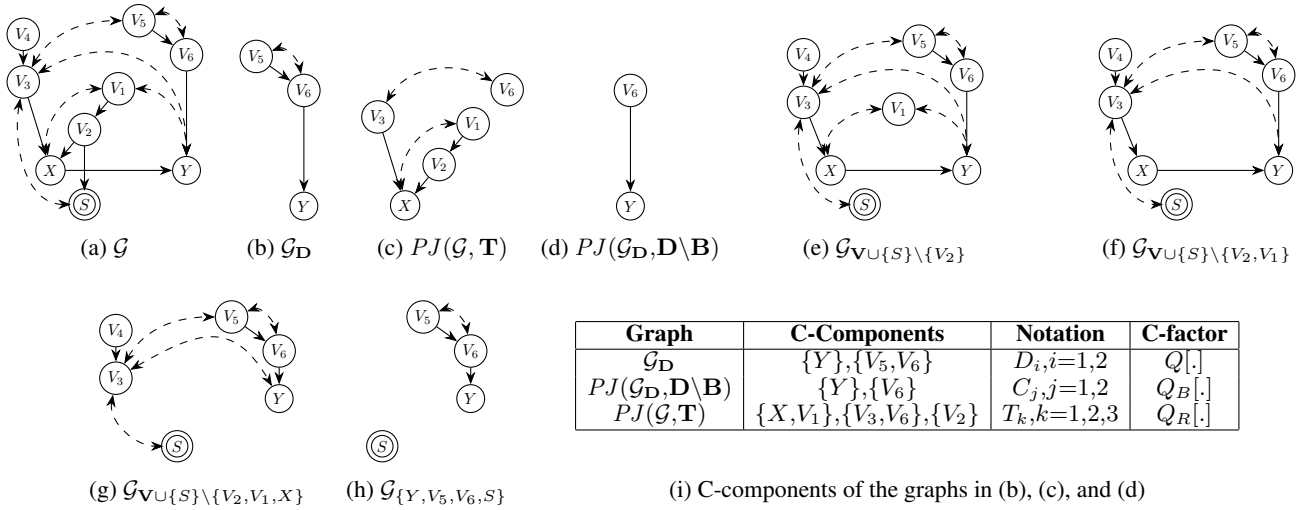


Figure 3: Example of a model and the transformations involved in recovering the target causal effect. We assume  $P(v|S=1)$  and  $P(v_2, v_3, v_6)$  are given.

$$\begin{aligned}
& \text{RCE}(\{Y\}, \mathcal{P}^{(0)}, \mathcal{G}) \\
& \quad | \text{3b} \\
& \text{RCE}(\{Y\}, \mathcal{P}^{(1)}, \mathcal{G}_{V \cup \{S\} \setminus \{V_2\}}) \\
& \quad | \text{1} \\
& \text{RCE}(\{Y\}, \mathcal{P}^{(2)}, \mathcal{G}_{V \cup \{S\} \setminus \{V_2, V_1\}}) \\
& \quad | \text{3b} \\
& \text{RCE}(\{Y\}, \mathcal{P}^{(3)}, \mathcal{G}_{V \cup \{S\} \setminus \{V_2, V_1, X\}}) \\
& \quad | \text{1} \\
& \text{RCE}(\{Y\}, \mathcal{P}^{(4)}, \mathcal{G}_{\{Y, V_5, V_6, S\}}) \\
& \quad | \text{3a} \\
& \text{IDENTIFY}(\{Y\}, \{Y\}, Q[Y])
\end{aligned}$$

Figure 4: Recursion of RCE when used to recover  $Q[Y]$  in the model in Fig. 3(a).

Despite IDSB’s generality, it is not clear at this point whether there are positive cases not covered by the algorithm – i.e., cases computable from  $P(\mathbf{t}^0)$  and  $P(\mathbf{v}|S=1)$ , but where IDSB returns “FAIL”. Still, the current state-of-the-art procedure that accepts external data, called Generalized Adjustment Criterion (GAC) (Correa, Tian, and Bareinboim 2018a), is constrained to backdoor-like expressions. The next proposition compares the power of the two approaches.

**Theorem 5.** *IDSB is strictly more powerful than the Generalized Adjustment Criterion for the task of recovering a causal effect  $P_{\mathbf{x}}(\mathbf{y})$  from a combination of biased distribution  $P(\mathbf{v}|S=1)$  and unbiased distribution  $P(\mathbf{t}^0)$  in  $\mathcal{G}$ .*

We outline how this statement can be proved (see Appendix D for the formal proof). We first show that whenever IDSB fails to recover  $P_{\mathbf{x}}(\mathbf{y})$ , then GAC is also unable to recover the effect. Then, to show that IDSB is strictly more general, we present an example where IDSB recovers  $P_{\mathbf{x}}(\mathbf{y})$  but GAC fails. Consider the problem of recovering  $P_{\mathbf{x}}(\mathbf{y})$  in the model in Fig. 2(b) with external data over  $\mathbf{T}^0 = \{W_2\}$ .

GAC asks for the following three conditions:

- Condition (iii) requires a set  $\mathbf{Z}^T$  to be available from external data such that the independence  $(S \perp\!\!\!\perp Y | \mathbf{Z}^T)_{\mathcal{G}_{\mathbf{X}\mathbf{Y}}^{pbd}}$  holds. For this model  $\mathbf{Z}^T = \{W_2\}$  suffices.
- Condition (i) requires that no covariate should be a descendant of a variable in a proper causal path from  $\mathbf{X}$  to  $\mathbf{Y}$ , which is also satisfied by  $\mathbf{Z} = \{W_2\}$ .
- However, condition (ii) requires the independence  $(X \perp\!\!\!\perp Y | \mathbf{Z}, S)_{\mathcal{G}_{\mathbf{X}\mathbf{Y}}^{pbd}}$  to hold, which cannot be satisfied in this model by  $\mathbf{Z} = \{W_2\}$ , or  $\mathbf{Z} = \{W_1, W_2\}$ , or any other  $\mathbf{Z}$ .

Since not all conditions are satisfiable, GAC fails. Nevertheless, IDSB is able to recover  $P_{\mathbf{x}}(\mathbf{y})$ . To witness, note that  $\mathbf{D} = An(\mathbf{Y})_{\mathcal{G}_{\mathbf{V} \setminus \mathbf{X}}} = \{Y\}$ , hence  $P_{\mathbf{x}}(\mathbf{y}) = Q[Y]$ . Also  $\mathbf{T} = \{W_1, W_2, X, Y\}$ ,  $\mathbf{R} = \emptyset$ . The set  $\{Y\}$  is a subset of c-component  $T_1 = \{W_1, X, Y\}$  in  $PJ(\mathcal{G}, \mathbf{T})$ . IDSB will call  $\text{IDENTIFY}(\{Y\}, T_1, Q_R[T_1])$ , where  $Q_R[T_1]$  is recoverable from  $P(\mathbf{t}) = P(y, x, w_1, w_2) = P(y, x, w_1 | w_2, S=1)P(w_2)$  by Lemma 2, and obtain

$$P_{\mathbf{x}}(\mathbf{y}) = \frac{\sum_{W_1} P(y, x | w_1, w_2) P(w_1)}{\sum_{W_1} P(x | w_1, w_2) P(w_1)}. \quad (23)$$

## Conclusions

We investigated the challenges arising due to confounding and selection biases, which come under the rubric of recoverability of causal effects. We first studied the algorithm RC (Alg. 1) (Bareinboim and Tian 2015), which takes as input a causal diagram and a biased distribution. We supplemented the algorithm with a necessary condition for recoverability (Thm. 1), and proved that RC is complete for this task, namely, it recovers all effects that are indeed recoverable (Thm. 2). We then relaxed the setting to allow the incorporation of unbiased data (Def. 4). We developed the algorithm

IDSB (Alg. 3), which takes as input a combination of biased and unbiased data. We proved that IDSB is strictly more powerful than the current state-of-the-art method available (Thm. 5). Since confounding and selection biases are problems pervasive across disciplines, we hope that the methods developed here should help to understand and alleviate this problem in a broad range of data-intensive applications.

**Acknowledgments** Juan D. Correa and Elias Bareinboim were in parts supported by grants from IBM Research, Adobe Research, NSF IIS-1704352, and IIS-1750807 (CA-REER). Jin Tian was partially supported by NSF grant IIS-1704352 and ONR grant N000141712140.

## References

- Angrist, J. D. 1997. Conditional independence in sample selection models. *Economics Letters* 54(2):103–112.
- Bareinboim, E., and Pearl, J. 2012a. Causal Inference by Surrogate Experiments: z-Identifiability. In *Proceedings of the Twenty-Eighth Conference on Uncertainty in Artificial Intelligence*, 113–120. AUAI Press.
- Bareinboim, E., and Pearl, J. 2012b. Controlling Selection Bias in Causal Inference. In *Proceedings of the 15th International Conference on Artificial Intelligence and Statistics (AISTATS)*. La Palma, Canary Islands: JMLR. 100–108.
- Bareinboim, E., and Pearl, J. 2016. Causal inference and the data-fusion problem. *Proceedings of the National Academy of Sciences* 113(27):7345–7352.
- Bareinboim, E., and Tian, J. 2015. Recovering Causal Effects from Selection Bias. *Proceedings of the Twenty-Ninth AAAI Conference on Artificial Intelligence* 3475–3481.
- Bareinboim, E.; Tian, J.; and Pearl, J. 2014. Recovering from Selection Bias in Causal and Statistical Inference. In *Proceedings of the Twenty-eighth AAAI Conference on Artificial Intelligence*, 2410–2416. Palo Alto, CA: AAAI Press.
- Cooper, G. 1995. Causal Discovery from Data in the Presence of Selection Bias. In *Proceedings of the Workshop on Artificial Intelligence and Statistics*, 140–150.
- Correa, J. D., and Bareinboim, E. 2017. Causal effect identification by adjustment under confounding and selection biases. In *AAAI*.
- Correa, J. D.; Tian, J.; and Bareinboim, E. 2018a. Generalized Adjustment Under Confounding and Selection Biases. In *Proceedings of the 32th Conference on Artificial Intelligence (AAAI 2018)*.
- Correa, J. D.; Tian, J.; and Bareinboim, E. 2018b. Identification of causal effects in the presence of selection bias. Technical report, R-38, Purdue AI Lab, Dep. of Computer Science, Purdue University.
- Cortes, C.; Mohri, M.; Riley, M.; and Rostamizadeh, A. 2008. Sample Selection Bias Correction Theory. In *Proceedings of the 19th International Conference on Algorithmic Learning Theory*, 38–53.
- Elkan, C. 2001. The foundations of cost-sensitive learning. In *International joint conference on artificial intelligence*, volume 17, 973–978.
- Evans, R. J., and Didelez, V. 2015. Recovering from Selection Bias Using Marginal Structure in Discrete Models. In *Proceedings of the UAI 2015 Conference on Advances in Causal Inference - Volume 1504, ACI'15*, 46–55.
- Glymour, M. M., and Greenland, S. 2008. Causal Diagrams. In *Modern Epidemiology*. Philadelphia, PA: Lippincott Williams & Wilkins, 3rd edition. 183–209.
- Heckman, J. J. 1979. Sample Selection Bias as a Specification Error. *Econometrica* 47(1):153.
- Huang, Y., and Valtorta, M. 2006. Pearl’s Calculus of Intervention Is Complete. In T.S. Richardson, R. D. a., ed., *Proceedings of the Twenty-Second Conference on Uncertainty in Artificial Intelligence*, 217–224. AUAI Press.
- Kuroki, M., and Cai, Z. 2006. On recovering a population covariance matrix in the presence of selection bias. *Biometrika* 93(3):601–611.
- Little, R., and Rubin, D. 1987. *Statistical analysis with missing data*. New York, NY, USA: John Wiley & Sons.
- Mefford, J., and Witte, J. S. 2012. The covariate’s dilemma. *PLoS Genet* 8(11):e1003096.
- Pearl, J. 1995. Causal diagrams for empirical research. *Biometrika* 82(4):669–688.
- Pearl, J. 2000. *Causality: models, reasoning, and inference*. Cambridge University Press, 2nd edition.
- Pirinen, M.; Donnelly, P.; and Spencer, C. 2012. Including known covariates can reduce power to detect genetic effects in case-control studies. *Nature Genetics* 44(8):848–851.
- Robins, J. M. 2001. Data, Design, and Background Knowledge in Etiologic Inference. *Epidemiology* 12(3):313–320.
- Robinson, L. D., and Jewell, N. P. 1991. Some Surprising Results about Covariate Adjustment in Logistic Regression Models. *International Statistical Review / Revue Internationale de Statistique* 59(2):227–240.
- Shpitser, I., and Pearl, J. 2006. Identification of Joint Interventional Distributions in Recursive semi-Markovian Causal Models. In *Proceedings of the Twenty-First AAAI Conference on Artificial Intelligence*, volume 21, 1219–1226.
- Spirites, P.; Glymour, C. N.; and Scheines, R. 2001. *Causation, Prediction, and Search*. MIT Press, 2nd edition.
- Tian, J., and Pearl, J. 2002a. A General Identification Condition for Causal Effects. In *Proceedings of the Eighteenth National Conference on Artificial Intelligence (AAAI 2002)*, 567–573. Menlo Park, CA: AAAI Press/The MIT Press.
- Tian, J., and Pearl, J. 2002b. On the Testable Implications of Causal Models with Hidden Variables. *Proceedings of the Eighteenth Conference on Uncertainty in Artificial Intelligence (UAI-02)* 519–527.
- Tian, J. 2002. *Studies in causal reasoning and learning*. Ph.D. Dissertation, Computer Science Department, University of California Los Angeles, CA.
- van der Zander, B.; Liskiewicz, M.; and Textor, J. 2014. Constructing separators and adjustment sets in ancestral graphs. In *Proceedings of UAI 2014*, 907–916.
- Verma, T. 1993. Graphical aspects of causal models. Technical report, Cognitive Systems Laboratory, Computer Science Department, University of California, Los Angeles.
- Whittemore, A. 1978. Collapsibility of multidimensional contingency tables. *Journal of the Royal Statistical Society. Series B* 40(3):328–340.
- Zadrozny, B. 2004. Learning and evaluating classifiers under sample selection bias. In *Twenty-first international conference on Machine learning - ICML '04*, 114–121.