

Causal Effect Identification by Adjustment under Confounding and Selection Biases

Juan D. Correa, Elias Bareinboim
Purdue University
{correagr,eb}@purdue.edu

January 24, 2017

Abstract

Controlling for selection and confounding biases are two of the most challenging problems in the empirical sciences as well as in artificial intelligence tasks. Covariate adjustment (or, Backdoor Adjustment) is the most pervasive technique used for controlling confounding bias, but the same is oblivious to issues of sampling selection. In this paper, we introduce a generalized version of covariate adjustment that simultaneously controls for both confounding and selection biases. We first derive a sufficient and necessary condition for recovering causal effects using covariate adjustment from an observational distribution collected under preferential selection. We then relax this setting to consider cases when additional, unbiased measurements over a set of covariates are available for use (e.g., the age and gender distribution obtained from census data). Finally, we present a complete algorithm with polynomial delay to find all sets of admissible covariates for adjustment when confounding and selection biases are simultaneously present and unbiased data is available.

1 Introduction

One of the central challenges in data-driven fields is to compute the effect of interventions – for instance, how increasing the educational budget will affect violence rates in a city, whether treating patients with a certain drug will help their recovery, or how increasing the product price will change monthly sales? These questions are commonly referred as the problem of identification of causal effects. There are two types of *systematic bias* that pose obstacles to this kind of inference, namely *confounding bias* and *selection bias*. The former refers to the presence of a set of factors that affect both the action (also known as treatment) and the outcome [Pearl, 1993], while the latter arises when the action, outcome, or other factors differentially affect the inclusion of subjects in the data sample [Bareinboim and Pearl, 2016].

The goal of our analysis is to produce an unbiased estimand of the *causal effect*, specifically, the probability distribution of the outcome when an action is performed by an autonomous agent (e.g., FDA, robot), regardless of how the decision would naturally occur [Pearl, 2000, Ch. 1]. For example, consider the graph in Fig. 1(a) in which X represents a treatment (e.g., taking or not a drug), Y represents an outcome (health status), and Z is a factor (e.g., gender, age) that affects both the propensity of being treated and the outcome. The edges (Z, X) and (Z, Y) may encode the facts “gender affects how the drug is being prescribed” and “gender affects recovery” respectively – for example, females may be more health conscious, so they seek for treatment more frequently than their male counterparts and at the same time are less likely to develop large complications for the particular disease. Intuitively, the causal effect represents the variations of X that bring about change in Y *regardless* of the influence of Z on X , which is graphically represented in Fig. 1(b). Mutilation is the graphical operation of removing arrows representing a decision made by an

autonomous agent of setting a variable to a certain value. The mathematical counterpart of mutilation is the $do()$ operator and the average causal effect of X on Y is usually written in terms of the do -distribution $P(y | do(x))$ [Pearl, 2000, Ch. 1].

The gold standard for obtaining the do -distribution is through the use of randomization, where the treatment assignment is selected by a randomized device (e.g., a coin flip) regardless of any other set of covariates (Z). In fact, this operation physically transforms the reality of the underlying population (Fig. 1(a)) into the corresponding mutilated world (Fig. 1(b)). The effect of Z on X is neutralized once randomization is applied. Despite its effectiveness, randomized studies can be prohibitively expensive, and even unattainable in certain cases, either for technical, ethical, or technical reasons – e.g., one cannot randomize the cholesterol level of a patient and record if it causes the heart to stop, when trying to assess the effect of cholesterol level on cardiac failure.

An alternative way of computing causal effects is trying to relate non-experimentally collected samples (drawn from $P(z, x, y)$) with the experimental distribution ($P(y | do(x))$). Non-experimental (often called observational) data relates to the model in Fig. 1(a) where subjects decide by themselves to take or not the drug (X) while influenced by other factors (Z). There are a number of techniques developed for this task, where the most general one is known as *do-calculus* [Pearl, 1995]. In practice, one particular strategy from do-calculus called *adjustment* is used the most. It consists of averaging the effect of X on Y over the different levels of Z , isolating the effect of interest from the effect induced by other factors. Controlling for confounding bias by adjustment is currently the standard method for inferring causal effects in data-driven fields, and different properties and enhancements have been studied in statistics [Rubin, 1974, Robinson and Jewell, 1991, Pirinen et al., 2012, Mefford and Witte, 2012] and AI [Pearl, 1993, Pearl, 1995, Pearl and Paz, 2010, Shpitser et al., 2010, Maathuis and Colombo, 2015, van der Zander et al., 2014].

Orthogonal to confounding, *sampling selection bias* is induced by preferential selection of units for the dataset, which is usually governed by unknown factors including treatment, outcome, and their consequences. It cannot be removed by a randomized trial and may stay undetected during the data gathering process, the whole study, or simply never be detected¹. Consider Fig. 1(e) where X and Y represent again treatment and outcome, but S represents a binary variable that indicates if a subject is included in the pool ($S=1$ means that the unit is in the sample, $S=0$ otherwise). The effect of X on Y in the entire population ($P(y | do(x))$) is usually not the same as in the sample ($P(y | do(x), S=1)$). For instance, patients that went to the hospital and were sampled are perhaps more affluent and have better nutrition than the average person in the population, which can lead to a faster recovery. This preferential selection of samples challenges the validity of inferences in several tasks in AI [Cooper, 1995, Cortes et al., 2008, Zadrozny, 2004] and Statistics [Little and Rubin, 1986, Kuroki and Cai, 2006] as well as in the empirical sciences [Heckman, 1979, Angrist, 1997, Robins, 2001].

The problem of selection bias can be addressed by removing the influence of the biased sampling mechanism on the outcome as if a random sample of the population was taken. For the graph in Fig. 1(d), for example, the distribution $P(y | do(x))$ is equal to $P(y | x, S=1)$ because there are not external factors that affect X and the selection mechanism S is independent of the outcome Y when the effect is estimated for the treatment X . There exists a complete non-parametric² solution for the problem of estimating statistical quantities from selection biased datasets [Bareinboim and Pearl, 2012], and also sufficient and algorithmic conditions for recovering from selection in the context of causal inference [Bareinboim et al., 2014, Bareinboim and Tian, 2015]. Whenever non-parametric recoverability is not feasible, a number of additional constraints over the model can be considered, including assumptions relative to the query (e.g., odds ratio), the type and dimensionality of the variables (e.g., discrete), and their topological relationships (e.g., IVs) [Didelez et al., 2010, Bareinboim and Pearl, 2012, Evans and Didelez, 2015].

Both confounding and selection biases carry extraneous “flow” of information between treatment and

¹[Zhang, 2008] noticed some interesting cases where detection is feasible in a class of non-chordal graphs.

²No assumptions about the about the functions that relates variables are made (i.e. linearity, monotonicity).

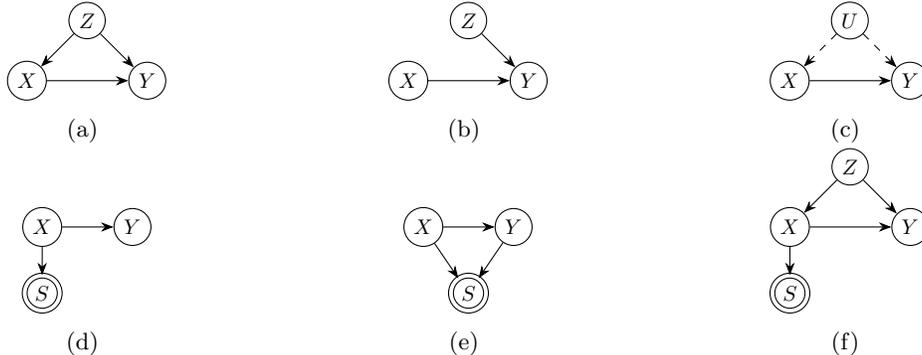


Figure 1: (a) and (d) give simple examples for confounding and selection bias respectively. (b) represents the model in (a) after an intervention is performed on X . (c) and (e) present examples where confounding and selection bias can not be removed respectively. In (f) we can control for either confounding or selection bias, but not for both unless we have external data on $P(z)$.

outcome, which is usually deemed “spurious correlation” since it does not correspond to the effect we want to compute on. Despite all the progress made in controlling these biases separately, we show that to estimate causal effects considering both problems requires a more refined analysis. First, note that the effect of X on Y can be estimated by blocking confounding and controlling for selection, respectively, in Figs. 1(a) and (d). On the other hand, confounding cannot be removed in Fig. 1(c) nor it can be recovered from selection bias in Fig. 1(e). Perhaps surprisingly, Fig. 1(f) presents a scenario where either confounding or selection can be addressed separately ($P(y|do(x)) = \sum_z P(y|x, z)P(z)$ and $P(z, y|do(x)) = P(z, y|do(x), S=1)$), but not simultaneously (without external data). As this example suggests, there is an intricate connection between these two biases that disallow the methods developed for these problems of being applied independently and then combined.

In this paper, we study the problem of estimating causal effects from models with an arbitrary structure that involve both biases. We establish necessary and sufficient conditions that a set of variables should fulfill so as to guarantee that the target effect can be unbiasedly estimated by adjustment. We consider two settings – first when only biased data is available, and then a more relaxed setting where additional unbiased samples of covariates are available for use (e.g., census data). Specifically, we solved the following problems:

1. **Identification and recoverability without external data:** The data is collected under selection bias, $P(\mathbf{v} | S=1)$, when does a set of covariates \mathbf{Z} allow $P(\mathbf{y} | do(\mathbf{x}))$ to be estimated by adjusting for \mathbf{Z} ?
2. **Identification and recoverability with external data:** The data is collected under selection bias $P(\mathbf{v} | S=1)$ and unbiased samples of $P(\mathbf{t}), \mathbf{T} \subseteq \mathbf{V}$, are available. When does a set of covariates $\mathbf{Z} \subseteq \mathbf{T}$ license the estimation of $P(\mathbf{y} | do(\mathbf{x}))$ by adjusting for \mathbf{Z} ?
3. **Finding admissible adjustment sets with external data:** How can we list all admissible sets \mathbf{Z} capable of identifying and recovering $P(\mathbf{y} | do(\mathbf{x}))$, for $\mathbf{Z} \subseteq \mathbf{T} \subseteq \mathbf{V}$?

2 Preliminaries

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 (SCM) [Pearl, 2000, pp. 204-207]. Formally, a SCM M is a 4-tuple $\langle U, V, F, P(u) \rangle$, where U is a set of exogenous (latent) variables and V is a set of endogenous (measured) variables. F represents a collection of functions $F = \{f_i\}$ such that each endogenous variable $V_i \in V$ 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 , $U_i \subseteq U$, $Pa_i \subseteq V \setminus V_i$ (where Pa_i is the set of endogenous variables that are arguments of f_i), and the entire set F forms a mapping from U to V . The uncertainty is encoded through a probability distribution over the exogenous variables, $P(u)$. Within the structural semantics, performing an action $X=x$ is represented through the do-operator, $do(X=x)$, which encodes the operation of replacing the original equation of X by the constant x and induces a submodel M_x . For a detailed discussion on the properties of structural models, we refer readers to [Pearl, 2000, Ch. 7].

Structural Causal Models are, by convention, represented economically using directed acyclic graphs with nodes for the measured variables and edges to represent the functional dependencies among them. Bidirected, dashed arrows are used to indicate the presence of an unobserved confounder between the variables connected by it. In this paper, the usual family notation over the graphs is used, so that Pa_X , An_X and De_X stand respectively for the set of parents, ancestors and descendants for a particular variable X . Moreover, sets of variables are represented in bold. The causal effect of a set \mathbf{X} when it is assigned a set of values \mathbf{x} , on a set \mathbf{Y} when it is instantiated as \mathbf{y} will be written as $P(\mathbf{y} \mid do(\mathbf{x}))$, which is a short hand notation for $P(\mathbf{Y}=\mathbf{y} \mid do(\mathbf{X}=\mathbf{x}))$. Mainly, the problems presented operate on $P(\mathbf{v})$, $P(\mathbf{v} \mid do(\mathbf{x}))$, $P(\mathbf{v} \mid S=1)$, respectively, the observational, experimental, and selection-biased distributions.

Formally, the task of estimating a probabilistic quantity from a selection-biased distribution is known as *recovering* from selection bias [Bareinboim and Pearl, 2012]. It is not uncommon for observations of a subset of the variables over the entire population (unbiased data) to be available for use. Therefore, our treatment consider two subsets of \mathbf{V} , $\mathbf{M}, \mathbf{T} \subseteq \mathbf{V}$, where \mathbf{M} contains the variables for which data was collected under selection bias, and \mathbf{T} encompasses the variables observed in the overall population, without bias. The absence of unbiased data is equivalent to have $\mathbf{T} = \emptyset$.

3 Selection Bias with Adjustment

The main justification for the validity of adjustment for confounding comes under a graphical condition called the “Backdoor criterion” [Pearl, 1993, Pearl, 2000], as shown below:

Definition 1 (Backdoor Criterion [Pearl, 2000]). *A set of variables \mathbf{Z} satisfies the Backdoor Criterion relative to a pair of variables (X, Y) in a directed acyclic graph G if:*

- (i) *No node in \mathbf{Z} is a descendant of X .*
- (ii) *\mathbf{Z} blocks every path between X and Y that contains an arrow into X .*

The heart of the criterion lies in cond. (ii), where the set \mathbf{Z} is required to block all the backdoor paths between X and Y that generate confounding bias. Furthermore, cond. (i) forbids the inclusion of descendants of X in \mathbf{Z} , which intends to avoid opening new non-causal paths. For example, the empty set is admissible for adjustment in Fig. 1(e), but adding S would not be allowed since it is a descendant of X and opens the non-causal path $X \rightarrow S \leftarrow Y$. On the other hand, even though S does not open any non-causal path in Fig. 1(f), the criterion does not allow it to be used for adjustment.

[Bareinboim et al., 2014] noticed that adjustment could be used for controlling for selection bias, in addition to confounding, which lead to a sufficient graphical condition called *Selection-Backdoor* criterion.

Definition 2 (Selection-Backdoor Criterion [Bareinboim and Tian, 2015]). *A set $\mathbf{Z} = \mathbf{Z}^+ \cup \mathbf{Z}^-$, with $\mathbf{Z}^- \subseteq De_X$ and $\mathbf{Z}^+ \subseteq \mathbf{V} \setminus De_X$ (where De_X is the set of variables that are descendants of X in G) satisfies the selection backdoor criterion (*s-backdoor*, for short) relative to X, Y and \mathbf{M}, \mathbf{T} in a directed acyclic graph G if:*

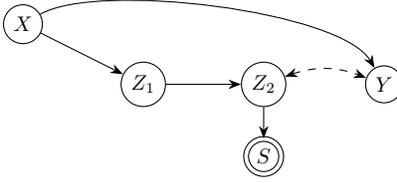


Figure 2: A graph that does not satisfy the s-backdoor criterion (respect to \mathbf{Z}), but the adjustment formula is recoverable and corresponds to desired causal effect.

- (i) \mathbf{Z}^+ blocks all back door paths from X to Y
- (ii) X and \mathbf{Z}^+ block all paths between \mathbf{Z}^- and Y , namely, $(\mathbf{Z}^- \perp\!\!\!\perp Y \mid X, \mathbf{Z}^+)$
- (iii) X and \mathbf{Z} block all paths between S and Y , namely, $(Y \perp\!\!\!\perp S \mid X, \mathbf{Z})$
- (iv) $\mathbf{Z} \cup \{X, Y\} \subseteq \mathbf{M}$ and $\mathbf{Z} \subseteq \mathbf{T}$

The first two conditions echo the extended-backdoor [Pearl and Paz, 2010]³, while cond. (iii) and (iv) guarantee that the resultant expression is estimable from the available datasets. If the S-Backdoor criterion holds for \mathbf{Z} relative to X, Y and \mathbf{M}, \mathbf{T} in G , then the effect $P(y \mid do(x))$ is identifiable, recoverable, and given by

$$P(y \mid do(x)) = \sum_{\mathbf{z}} P(y \mid x, \mathbf{z}, S=1)P(\mathbf{z}) \quad (1)$$

Note that the S-Backdoor is sufficient but not necessary for adjustment. To witness, consider the model in Fig. 2 where $\mathbf{Z} = \{Z_1, Z_2\}$, $\mathbf{M} = \{X, Y, Z_1, Z_2\}$, and $\mathbf{T} = \{Z_1, Z_2\}$. Here, $\mathbf{Z}^+ = \emptyset$, $\mathbf{Z}^- = \{Z_1, Z_2\}$. Condition (ii) in Def. 2 is violated, namely $(Z_1, Z_2 \not\perp\!\!\!\perp Y \mid X)$. Perhaps surprisingly, the effect $P(y \mid do(x))$ is identifiable and recoverable, as follows:

$$P(y \mid do(x)) = P(y \mid x) \quad (2)$$

$$= P(y \mid x) \sum_{z_1} P(z_1) \quad (3)$$

$$= \sum_{z_1} P(y \mid x, z_1)P(z_1) \quad (4)$$

$$= \sum_{z_1, z_2} P(y \mid x, z_1, z_2)P(z_2 \mid x, z_1)P(z_1) \quad (5)$$

$$= \sum_{z_1, z_2} P(y \mid x, z_1, z_2)P(z_2 \mid z_1)P(z_1) \quad (6)$$

$$= \sum_{z_1, z_2} P(y \mid x, z_1, z_2, S=1)P(z_1, z_2) \quad (7)$$

(2) follows from the application of the second rule of do calculus and the independence $(X \perp\!\!\!\perp Y)_{G_{\overline{X}}}$. Equations (5),(6),(7) use the independences $(Y \perp\!\!\!\perp Z_1 \mid X)$, $(Z_2 \perp\!\!\!\perp X \mid Z_1)$ and $(S \perp\!\!\!\perp Y \mid X, Z_1, Z_2)$ respectively. The final expression (7) is estimable from the available data.

Considering that $\mathbf{Z} = \emptyset$ controls for confounding, adjusting for $\mathbf{Z} = \{Z_1, Z_2\}$ seems unnecessary. As it turns out, covariates irrelevant for confounding control, could play a role when we compound this task with recovering from selection bias (where Y will need to be separated from S).

³The extended-backdoor augments the backdoor criterion to allow for descendants of X that could be harmless in terms of bias.

4 Generalized Adjustment without External Data

Let us consider the case when only biased data $P(\mathbf{v} \mid S=1)$ over \mathbf{V} is measured. Our interest in this section is on conditions that allow $P(\mathbf{y} \mid do(\mathbf{x}))$ to be computed by adjustment without external measurements.

Consider the model G in Fig. 4(a). Note that Y and S are marginally independent in $G_{\overline{X}}$ (the graph after an intervention on X where all edges into X are not present). As for confounding, Z needs to be conditioned on, but doing so opens a path between Y and S , letting spurious correlation from the bias to be included in our calculation. It turns out that with a careful manipulation of the expression, both biases can be controlled as follows:

$$P(y \mid do(x)) = P(y \mid do(x), S=1) \quad (8)$$

$$= \sum_Z P(y \mid do(x), z, S=1)P(z \mid do(x), S=1) \quad (9)$$

$$= \sum_Z P(y \mid x, z, S=1)P(z \mid S=1) \quad (10)$$

Eq. (8) follows from the independence ($Y \perp\!\!\!\perp S \mid X$) in the mutilated graph $G_{\overline{X}}$. Next we condition on Z and (10) is valid by the application of the second rule of do-calculus to the first term and the third rule to the second from (9). Note that every term in (10) is estimable from the biased distribution.

One important difference between the criterion presented in this section and the Backdoor criterion (both the standard and extended versions) is the explicit consideration of sets of outcome and, specially, treatment variables. This formulation requires a distinction between causal paths that contain a treatment variable only at the beginning, and those that are intercepted by another variable of the same kind. The former type of causal paths are called *proper* based in the following:

Definition 3 (Proper Causal Path [Shpitser et al., 2010]). *Let \mathbf{X} and \mathbf{Y} be sets of nodes. A causal path from a node in \mathbf{X} to a node in \mathbf{Y} is called proper if it does not intersect \mathbf{X} except at the end point.*

To illustrate the concept consider the graph in Fig. 3. In one hand, the path $X_1 \rightarrow W_1 \rightarrow X_2 \rightarrow W_2 \rightarrow Y$ is not proper because of the presence of X_2 in between the endpoints. On the other hand, $X_1 \rightarrow W_3 \rightarrow Y$ and $X_2 \rightarrow W_2 \rightarrow Y$ are proper. Note that, if \mathbf{X} is composed of a single variable, then all causal paths are proper. To build intuition around the necessity of the notion of proper causal path, consider the strategy

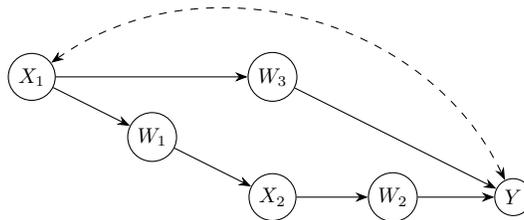


Figure 3: Graph with proper and non-proper causal paths

entailed in the backdoor adjustment: to block all non-causal paths and only the non-causal ones. In the case of the path $X_1 \rightarrow W_1 \rightarrow X_2 \rightarrow W_2 \rightarrow Y$ from before, the intervention on X_2 cuts the flow of causal information from X_1 to Y through that path. As a consequence adjusting for W_1 is not interfering with the causal influence of X_1 on Y , which is already nullified by the mentioned intervention on X_2 . Hence, a carefree characterization may disallow the use of W_1 for adjustment in similar scenarios, even when it may help with the selection bias problem.



Figure 4: Models where \mathbf{Z} satisfies Def. 4

In the sequel, we introduce a graphical criterion to determine whether a set of covariates is admissible for adjustment so as to simultaneously *identify and recover* a causal effect.

Definition 4 (Generalized Adjustment Criterion Type 1). *A set \mathbf{Z} satisfies the generalized criterion relative to the pair \mathbf{X} and \mathbf{Y} in a causal model with graph G , augmented with the selection mechanism S if:*

- (a) *No element of \mathbf{Z} is a descendant in $G_{\overline{\mathbf{X}}}$ of any $W \notin \mathbf{X}$ which lies on a proper causal path from \mathbf{X} to \mathbf{Y} .*
- (b) *All non-causal paths between \mathbf{X} and \mathbf{Y} in G are blocked by \mathbf{Z} .*
- (c) *\mathbf{Y} is d -separated from S given \mathbf{X} under the intervention $do(\mathbf{x})$, i.e., $(\mathbf{Y} \perp\!\!\!\perp S \mid \mathbf{X})_{G_{\overline{\mathbf{X}}}}$.*
- (d) *\mathbf{Z} can be partitioned into sets $\mathbf{Z}^+, \mathbf{Z}^-$ such that $\mathbf{Z}^+ = \{Z' \in \mathbf{Z} \mid (Z' \perp\!\!\!\perp \mathbf{X} \mid S)_{G_{\overline{\mathbf{X}}(S)}}\}$ and $(\mathbf{Y} \perp\!\!\!\perp \mathbf{Z}^- \mid \mathbf{X}, \mathbf{Z}^+, S)_{G_{\overline{\mathbf{X}}}}$.*

$G_{\overline{\mathbf{X}}(S)}$ is the graph where all edges into $X \in \mathbf{X} \setminus An_S$ are removed, where An_S is the set of ancestors of the variable S in G .

Conditions (a) and (b) echo the Extended Backdoor/Adjustment Criterion [Pearl and Paz, 2010, Shpitser et al., 2010] and guarantee that \mathbf{Z} is admissible for adjustment in the unbiased distribution. Condition (c) requires the outcome \mathbf{Y} to be independent of the selection mechanism S under intervention, without observing any covariate \mathbf{Z} , this is effectively saying that the causal effect is invariant to the selection mechanism. Condition (d) ensures that either conditioning on S is not introducing spurious information from \mathbf{X} to the weights given by the covariates \mathbf{Z} in the adjustment (for the \mathbf{Z}^+ subset), or that \mathbf{Y} is insensitive to \mathbf{Z} (as for \mathbf{Z}^-). The following theorem claims the completeness of this criterion:

Theorem 1 (Generalized Adjustment Formula Type 1). *Given disjoint sets of variables \mathbf{X}, \mathbf{Y} and \mathbf{Z} in a causal model with graph G . The effect $P(\mathbf{y} \mid do(\mathbf{x}))$ is given by*

$$P(\mathbf{y} \mid do(\mathbf{x})) = \sum_{\mathbf{z}} P(\mathbf{y} \mid \mathbf{x}, \mathbf{z}, S=1) P(\mathbf{z} \mid S=1) \quad (11)$$

in every model inducing G if and only if \mathbf{Z} satisfies the generalized adjustment criterion type 1 relative to the pair \mathbf{X}, \mathbf{Y} .

Proof. Suppose \mathbf{Z} satisfy the criterion relative to \mathbf{X}, \mathbf{Y} . Then it can be decomposed into \mathbf{Z}^- and \mathbf{Z}^+ as defined in condition (d). The causal effect can be derived as follows:

$$P(\mathbf{y} \mid do(\mathbf{x})) = P(\mathbf{y} \mid do(\mathbf{x}), S=1) \quad (12)$$

$$= \sum_{\mathbf{z}^+} P(\mathbf{y} \mid do(\mathbf{x}), \mathbf{z}^+, S=1) P(\mathbf{z}^+ \mid do(\mathbf{x}), S=1) \quad (13)$$

$$= \sum_{\mathbf{z}^+} P(\mathbf{y} \mid do(\mathbf{x}), \mathbf{z}^+, S=1) P(\mathbf{z}^+ \mid S=1) \quad (14)$$

$$= \sum_{\mathbf{z}^+} P(\mathbf{y} \mid do(\mathbf{x}), \mathbf{z}^+, S=1) \sum_{\mathbf{z}^-} P(\mathbf{z} \mid S=1) \quad (15)$$

$$\begin{aligned}
&= \sum_{\mathbf{z}} P(\mathbf{y} \mid do(\mathbf{x}), \mathbf{z}, S=1)P(\mathbf{z} \mid S=1) & (16) \\
&= \sum_{\mathbf{z}} P(\mathbf{y} \mid \mathbf{x}, \mathbf{z}, S=1)P(\mathbf{z} \mid S=1) & (17)
\end{aligned}$$

Eq. (12) follows from cond. (c). Conditioning on \mathbf{Z}^+ and applying the third rule of do-calculus using the definition of \mathbf{Z}^+ in cond. (d) yield eq. (14). Summing over \mathbf{Z}^- in the second factor and adding \mathbf{Z}^- to the first term using cond. (d) results in (16). Conditions (a) and (b) imply $(\mathbf{Y} \perp\!\!\!\perp \mathbf{X} \mid \mathbf{Z})_{G_{\overline{\mathbf{X}}}}$, furthermore, cond. (c) ensures that observing S will not open any path between \mathbf{X} and \mathbf{Y} , because such path will either violate (c) or have some $X \in \mathbf{X}$ as a collider which contradicts (b). Hence, $(\mathbf{Y} \perp\!\!\!\perp \mathbf{X} \mid \mathbf{Z}, S)_{G_{\overline{\mathbf{X}}}}$ holds and can be used together with rule 2 of do-calculus to remove of the *do* operator, which results in the adjustment formula in Eq. (11). The necessity part of the proof is presented in the appendix. \square

The set $\mathbf{Z} = \{Z\}$ for the model in Fig. 4(b) also satisfies Def. 4. Similarly to Fig. 4(a), if we control for confounding and then try to remove the do-operator, it appears that the second term of the adjustment expression cannot be estimated, because the independence $(\mathbf{Z} \perp\!\!\!\perp S)$ does not hold in G . Still, there exists a derivation strategy encapsulated in Def. 4 / Thm. 1 that allow one to recover from both selection and confounding biases.

5 Generalized Adjustment With External Data

A natural question that arises is whether additional measurements in the population level over the covariates can help identifying and recovering the desired causal effect. The following criterion relaxes the previous result by leveraging the unbiased data available.

Definition 5 (Generalized Adjustment Criterion Type 2). *A set \mathbf{Z} satisfies the generalized criterion relative to \mathbf{X}, \mathbf{Y} , a set of variables measured under selection bias \mathbf{M} and a set of variables observed in the overall population \mathbf{T} in a causal model with graph G augmented with the selection mechanism S if:*

- (a) *No element of \mathbf{Z} is a descendant in $G_{\overline{\mathbf{X}}}$ of any $W \notin \mathbf{X}$ which lies on a proper causal path from \mathbf{X} to \mathbf{Y} .*
- (b) *All non-causal paths between \mathbf{X} and \mathbf{Y} in G are blocked by \mathbf{Z} .*
- (c) *\mathbf{Y} is d-separated from the selection mechanism S given \mathbf{Z} and \mathbf{X} , i.e., $(\mathbf{Y} \perp\!\!\!\perp S \mid \mathbf{X}, \mathbf{Z})$.*
- (d) *The variables are measured with bias ($\mathbf{Z}, \mathbf{X}, \mathbf{Y} \subseteq \mathbf{M}$) and the covariates are available without bias ($\mathbf{Z} \subseteq \mathbf{T}$)*

As in Def. 4, conditions (a) and (b) ensure \mathbf{Z} is valid for adjustment without selection bias. Condition (c) requires that the influence of the selection mechanism in the outcome is nullified by conditioning on \mathbf{X} and \mathbf{Z} . Note that (c) is stated over the graph G and not $G_{\overline{\mathbf{X}}}$ as in the previous case. Actually, given cond. (b) they are interchangeable because they may differ only when \mathbf{X} is a collider in a path between S and \mathbf{Y} . But, this is also a non-causal path between \mathbf{X} and \mathbf{Y} that, by condition (b), should be blocked by \mathbf{Z} anyways. Given that, cond. (c) as in the definition for considering it simpler. Condition (d) guarantees that the adjustment expression can be estimated from the available data. The following theorem claims the completeness of the criterion for this case:

Theorem 2 (Generalized Adjustment Formula Type 2). *Given disjoint sets of variables \mathbf{X}, \mathbf{Y} and \mathbf{Z} , and sets \mathbf{M}, \mathbf{T} in a causal model with graph G . In every model inducing G , the effect $P(\mathbf{y} \mid do(\mathbf{x}))$ is given by*

$$P(\mathbf{y} \mid do(\mathbf{x})) = \sum_{\mathbf{z}} P(\mathbf{y} \mid \mathbf{x}, \mathbf{z}, S=1)P(\mathbf{z}) \quad (18)$$

if and only if the set \mathbf{Z} satisfies the generalized adjustment criterion type 2 relative to the pair \mathbf{X}, \mathbf{Y} .

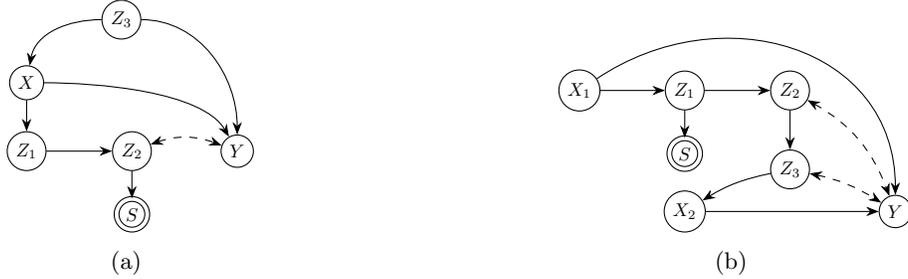


Figure 5: Models where the set \mathbf{Z} satisfies Def. 5.

Proof. Suppose \mathbf{Z} satisfy the conditions of the theorem relative to the pair \mathbf{X}, \mathbf{Y} and the sets \mathbf{M}, \mathbf{T} . By conditions (a) and (b), the effect can be written as:

$$P(\mathbf{y} \mid do(\mathbf{x})) = \sum_{\mathbf{z}} P(\mathbf{y} \mid \mathbf{x}, \mathbf{z})P(\mathbf{z})$$

Note that S can be introduced to the first term by cond. (c), which entail Eq. (18). Cond. (d) ensures that both terms in the expression are estimable from the available distributions. The necessity part of the proof is presented in the appendix. \square

Fig. 5 presents two causal models that satisfies the previous criterion if measurements over $\mathbf{Z} = \{Z_1, Z_2, Z_3\}$ are available. To witness how the expression can be reached using do-calculus and probability axioms, consider Fig. 5(a):

$$P(y \mid do(x)) = \sum_{Z_3} P(y \mid do(x), z_3)P(z_3 \mid do(x)) \quad (19)$$

$$= \sum_{Z_3} P(y \mid x, z_3)P(z_3) \quad (20)$$

$$= \sum_{Z_1, Z_3} P(y \mid x, z_1, z_3)P(z_1, z_3) \quad (21)$$

$$= \sum_{\mathbf{z}} P(y \mid x, \mathbf{z})P(z_2 \mid x, z_1, z_3)P(z_1, z_3) \quad (22)$$

$$= \sum_{\mathbf{z}} P(y \mid x, \mathbf{z}, S=1)P(\mathbf{z}) \quad (23)$$

First conditioning on Z_3 and removing $do(x)$ using rule 3 of the do-calculus from the second term. Then, conditioning the second term on Z_1 , moving the summation to the left, and introducing Z_1 into the first term results in (21). Eq. (22) follows from conditioning the first term on Z_2 , and finally removing X from the second term using the independence ($Z_2 \perp\!\!\!\perp X \mid Z_1, Z_3$). Combining the last two distributions over the Z 's and introducing the selection bias term using the independence ($Y \perp\!\!\!\perp S \mid X, \mathbf{Z}$) results in (23), which corresponds to the stated formula (18).

Model in Fig. 5(b) also satisfies the type 2 criterion and illustrates how this can be applied in cases where \mathbf{X} and \mathbf{Y} are sets of variables.

6 Finding Admissible Sets for Generalized Adjustment

An obvious extension to the problem is how to systematically list admissible sets for adjustment, using the criteria discussed in the previous sections. This is specially important in practice where factors such as feasibility, cost, and statistical power relate to the choosing of a covariate set.

In order to perform this kind of task efficiently, [van der Zander et al., 2014] introduced a transformation of the model called the *Proper Backdoor Graph* and formulate a criterion equivalent to the Adjustment Criterion:

Definition 6 (Proper Backdoor graph). *Let $G = (\mathbf{V}, \mathbf{E})$ be a DAG, and $\mathbf{X}, \mathbf{Y} \subseteq \mathbf{V}$ be pairwise disjoint subsets of variables. The proper backdoor graph, denoted as $G_{\mathbf{X}\mathbf{Y}}^{pbd}$, is obtained from G by removing the first edge of every proper causal path from \mathbf{X} to \mathbf{Y} .*

Definition 7 (Constructive Backdoor Criterion (CBD)). *Let $G = (\mathbf{V}, \mathbf{E})$ be a DAG, and $\mathbf{X}, \mathbf{Y} \subseteq \mathbf{V}$ be pairwise disjoint subsets of variables. The set \mathbf{Z} satisfies the Constructive Backdoor Criterion relative to (\mathbf{X}, \mathbf{Y}) in G if:*

- i) $\mathbf{Z} \subseteq \mathbf{V} \setminus Dpcp(\mathbf{X}, \mathbf{Y})$ and
- ii) \mathbf{Z} d -separates \mathbf{X} and \mathbf{Y} in the proper backdoor graph $G_{\mathbf{X}\mathbf{Y}}^{pbd}$.

Where $Dpcp(\mathbf{X}, \mathbf{Y}) = De((De_{\overline{\mathbf{X}}}(\mathbf{X}) \setminus \mathbf{X}) \cap An_{\underline{\mathbf{X}}}(Y))$

The set $Dpcp(\mathbf{X}, \mathbf{Y})$ is exactly the set of nodes forbidden by the first condition in both of the criteria given in this paper. Moreover, $G_{\mathbf{X}\mathbf{Y}}^{pbd}$ only contain \mathbf{X}, \mathbf{Y} paths that need to be blocked. The following lemmas and theorem allow the use of the algorithmic framework from [van der Zander et al., 2014] to solve instances of the question posed by the type 2 criterion.

Lemma 3 (Constructive Backdoor \implies Generalized Adjustment Type 2). *Any set \mathbf{Z} satisfying the CBD applied to $G_{(\mathbf{X} \cup S)\mathbf{Y}}^{pbd}$ and $Dpcp(\mathbf{X} \cup S, \mathbf{Y}) \cup (\mathbf{V} \setminus \mathbf{T})$ relative to \mathbf{X}, \mathbf{Y} in G also satisfies the Generalized Adjustment Criterion type 2.*

Proof. By the equivalence between the CBD criterion and the adjustment criterion, we have that $Dpcp(\mathbf{X}, \mathbf{Y})$ is exactly the set of nodes forbidden by cond. (a) of the type 2 criterion, so

$$Dpcp(\mathbf{X} \cup S, \mathbf{Y}) = De((De_{\overline{\mathbf{X}, S}}(\mathbf{X} \cup \{S\}) \setminus (\mathbf{X} \cup S)) \cap An_{\underline{\mathbf{X}, S}}(Y))$$

Since S has no descendants, $De_{\overline{\mathbf{X}, S}}(\mathbf{X} \cup \{S\}) = De_{\overline{\mathbf{X}}}(\mathbf{X}) \cup S$ and $An_{\underline{\mathbf{X}, S}}(Y) = An_{\underline{\mathbf{X}}}(Y)$. As a consequence $Dpcp(\mathbf{X} \cup S, \mathbf{Y}) = Dpcp(\mathbf{X}, \mathbf{Y})$ implying cond. (a) of Def. 5.

$G_{(\mathbf{X} \cup S)\mathbf{Y}}^{pbd}$ has all non-causal paths from \mathbf{X} to \mathbf{Y} present in $G_{\mathbf{X}\mathbf{Y}}^{pbd}$, therefore, if \mathbf{Z} block all non-causal paths in the former, it will do in the latter satisfying condition (b).

Every $S - \mathbf{Y}$ path may or may not contain \mathbf{X} . If not, \mathbf{Z} should block it in $G_{(\mathbf{X} \cup S)\mathbf{Y}}^{pbd}$. In the latter case, the subpath from \mathbf{X} to \mathbf{Y} is either causal or non-causal. If it is causal \mathbf{Z} will not block it, but the $S - \mathbf{Y}$ path will be blocked by \mathbf{X} . If the subpath is non-causal \mathbf{Z} should block it, therefore, the larger path is blocked too. This argument implies condition (c). Since CBD holds for $Dpcp(\mathbf{X} \cup S, \mathbf{Y}) \cup (\mathbf{V} \setminus \mathbf{T})$ every element in \mathbf{Z} must belong to \mathbf{T} satisfying condition (d). \square

Lemma 4 (Generalized Adjustment Type 2 \implies Constructive Backdoor). *Any set \mathbf{Z} satisfying the Generalized Adjustment Criterion type 2 relative to \mathbf{X}, \mathbf{Y} in G also satisfies the constructive backdoor criterion applied to $G_{(\mathbf{X} \cup S)\mathbf{Y}}^{pbd}$ and $Dpcp(\mathbf{X} \cup S, \mathbf{Y}) \cup (\mathbf{V} \setminus \mathbf{T})$.*

Proof. By lemma 3, $Dpcp(\mathbf{X} \cup S, \mathbf{Y}) = Dpcp(\mathbf{X}, \mathbf{Y})$, which combined with condition (d) implies condition (i) of the CBP.

By cond. (b) every non-causal path from \mathbf{X} to \mathbf{Y} is blocked by \mathbf{Z} and all paths from S to \mathbf{Y} (which are always non-causal when S is treated as an \mathbf{X}) are blocked by \mathbf{Z}, \mathbf{X} by cond. (c). Those two facts together imply cond. (ii) of the CBD. \square



Figure 6: (a) shows a causal model and (b) the proper backdoor graph associated with it relative to $\mathbf{X} \cup S$ and Y . The gray nodes in (b) represents variables in $Dpcp$.

Theorem 5 (Generalized Adjustment Type 2 \Leftrightarrow Constructive Backdoor). *A set \mathbf{Z} satisfies the Generalized Adjustment Criterion type 2 relative to \mathbf{X}, \mathbf{Y} in G if and only if it satisfies the CBC applied to $G_{(\mathbf{X} \cup S)\mathbf{Y}}^{pbd}$ and $Dpcp(\mathbf{X} \cup S, \mathbf{Y}) \cup (\mathbf{V} \setminus \mathbf{T})$.*

Proof. It follows immediately from lemmas 3,4. □

Thm. 5 allows us to use the LISTSEP procedure [van der Zander et al., 2014] to list all the valid sets for the generalized adjustment type 2. The algorithm guarantees $O(n(n+m))$ polynomial delay, where n is the number of nodes and m is the number of edges in G (see [Takata, 2010]). That means that the time needed to output the first solution or indicate failure, and the time between the output of consecutive solutions, is $O(n(n+m))$.

To provide the reader an intuition of how the algorithm works, consider the graph in Fig. 6(a) and its associated constructive backdoor graph in (b). W is a “forbidden node” in the sense that it cannot be used for adjustment and for this example is the only element in $Dpcp(\mathbf{X}, Y)$ assuming that unbiased measurement on the covariates Z_1, Z_2 and Z_3 are available (i.e. $\{Z_1, Z_2, Z_3\} \subseteq \mathbf{T}$). The algorithm LISTSEP will output every set of variables that d-separates $X \cup S$ from Y in the proper backdoor graph that does not contain any node in $Dpcp(\mathbf{X}, Y)$.

7 Conclusions

We provide necessary and sufficient conditions for identification and recoverability from selection bias of causal effects by adjustment, applicable for data-generating models with latent variables and arbitrary structure in non-parametric settings. Def. 4 and Thm. 1 provide a complete characterization of identification and recoverability by adjustment when no external information is available. Def. 5 and Thm. 2 provide a complete graphical condition for when external information on a set of covariates is available. Thm. 5 allowed us to list all sets that satisfies the last criterion in polynomial-delay time, effectively helping in the decision of what covariates need to be measured for recoverability. This is especially important when measuring a variable is associated with a particular cost or effort. Despite the fact that adjustment is neither complete nor the only method to identify causal effects, it is in fact the most used tool in the empirical sciences. The methods developed in this paper should help to formalize and alleviate the problem of sampling selection and confounding biases in a broad range of data-intensive applications.

A Appendix: Proof of the theorems

In order to prove the necessity of the criteria presented in the paper, it is imperative to construct Structural Causal Models (SCM) that serve as counter-examples to the identifiability or recoverability of the causal effect, whenever the set of covariates \mathbf{Z} fails to satisfy the conditions relative to the pair \mathbf{X}, \mathbf{Y} . The following lemmata will be useful to construct such models. The first one, lemma 6 licenses the the direct specification of the conditional distributions of any variable given its parents, in accordance to the causal diagram G .

Lemma 6 (Family Parametrization). *Let G be a causal diagram over a set \mathbf{V} of n variables. Consider also, a set of conditional distributions $P(v_i | pa_{V_i}), 1 \leq i \leq n$ such that Pa_{V_i} is the set of nodes in G from which there are outgoing edges pointing into V_i . Then, there exists a model M compatible with G that induces $P(\mathbf{v}) = \prod_{i=1}^n P(v_i | pa_{V_i})$.*

Proof. (By construction) For every V_i define any ordering on the values of its domain, and let $v_i^{(j)}$ refer to the j^{th} value in that order. Also, define a continuous unobservable variable $U_i \sim U[0, 1]$ (uniformly distributed in the interval $[0, 1]$) for every variable $V_i \in \mathbf{V}$. Then, construct a SCM $M = \langle \mathbf{U}, \mathbf{V}, \mathcal{F}, P(\mathbf{u}) \rangle$ where:

- \mathbf{V} is the same set of observables in G
- $\mathbf{U} = \bigcup_{i=1}^n U_i'$
- $\mathcal{F} = \left\{ f_i(pa_{V_i}, u_i) = \inf_j \left\{ \sum_{k=1}^j P(v_i^{(k)} | pa_{V_i}) \geq u_i \right\}, 1 \leq i \leq n \right\}$
- $U_i \sim U[0, 1], 1 \leq i \leq n$

At every variable V_i , given a particular configuration of Pa_{V_i} , M simulates its value using the distribution $P(v_i | pa_{V_i})$. By the Markov property, the joint distribution will be equal to the product of those distributions. \square

The following lemma, permits the construction of a SCM M compatible with a causal diagram G , using another model compatible with a related, but different, causal diagram G' where some arrows in a chain of variables have the reverse direction.

Lemma 7 (Chain Reversal). *Let G be a causal diagram containing a chain $R_\ell \rightarrow R_{\ell-1} \rightarrow \dots \rightarrow R_1 \rightarrow T \rightarrow W_1 \rightarrow \dots \rightarrow W_{k-1} \rightarrow W_k$, for two constants k, ℓ , where the only edge incoming to $R_{\ell-1}, \dots, R_1, T, W_1, \dots, W_k$ is the one in the chain, and R_ℓ has no parents. Then, for any SCM M compatible with G , there exists a model M' compatible with a causal diagram G' where the chain of variables mentioned before, is replaced by a chain of the form $R_\ell \leftarrow R_{k-1} \leftarrow \dots \leftarrow R_\ell \leftarrow T \rightarrow W_1 \rightarrow \dots \rightarrow W_{k-1} \rightarrow W_k$. Moreover, M' compatible with any observational distribution $P(\mathbf{v})$ induced by M .*

Proof. (By construction) Given M and any probability distribution $P(\mathbf{v})$ induced by it, compute the joint distribution $P(r_1, \dots, r_\ell, t)$. Construct a new model M' with the same set of observable variables and identical functions for all variables but for R_1, \dots, R_ℓ, T . For those, assign the functions $f_{R_i}(r_{i-1}, U_{R_i}), 1 \leq i \leq \ell - 1$ as in lemma 6. Also, let $f_{R_\ell}(U_{R_\ell}) = U_{R_\ell}, P(U_{R_\ell}) = P(r_\ell)$. By lemma 6 the sub-models composed of R_1, \dots, R_ℓ, T in M' and M produce the exact same distribution and since the set of parents and function for every other part of the model are exactly the same the overall distribution is identical. \square

Finally, the following lemma allows to simplify the parametrization of an arbitrarily long chain of binary variables.

Lemma 8 (Collapsible Path Parametrization). *Consider a causal diagram G and a probability distribution $P(\mathbf{v})$ induced by any SCM compatible with G . If G contains a chain $W_0 \rightarrow W_1 \rightarrow \dots \rightarrow W_k$, where each W_i represents a binary random variable, for every $1 \leq i \leq k$ the only incoming edge to W_i is from W_{i-1} , and every conditional distribution $P(w_i | w_{i-1}) = p$, $P(w_i | \bar{w}_{i-1}) = q$, for some $0 < p, q < 1$. Then, the conditional distribution $P(w_k | w_0) = \frac{q-(p-1)(p-q)^k}{q-p+1}$, $P(w_k | \bar{w}_0) = \frac{q-q(p-q)^k}{q-p+1}$.*

Proof. Since W_0, \dots, W_k is a chain, the value of W_k is a function of W_0 when all other W_1, \dots, W_{k-1} are marginalized. All W_i , $1 \leq i \leq k$ are independent of any other variable given W_0 . Therefore, the distribution $P(w_k | w_0)$ is equal to $\sum_{i=1}^{k-1} \prod_{i=1}^k P(w_i | w_{i-1})$, because any other variable can be removed from any product in this expression and summed out. This distribution can be calculated as the product of 2x2 matrices corresponding to the conditional distributions $P(w_i | w_{i-1})$ when encoded as $W_M = \begin{bmatrix} p & q \\ 1-p & 1-q \end{bmatrix}$. The product of k of such matrices is readily available if W_M is decomposed using its eigenvalues $\{1, p-q\}$ and eigenvectors $\left\{ \begin{bmatrix} q \\ (1-p) \end{bmatrix}, [1, 1] \right\}$:

$$P(w_k | w_0) = \sum_{i=1}^{k-1} \prod_{i=1}^k P(w_i | w_{i-1}) = (W_M)^k = \begin{bmatrix} \frac{q-(p-1)(p-q)^k}{q-p+1} & \frac{q-q(p-q)^k}{q-p+1} \\ 1 - \frac{q-(p-1)(p-q)^k}{q-p+1} & 1 - \frac{q-q(p-q)^k}{q-p+1} \end{bmatrix} \quad (24)$$

□

A.1 Proof for the First Criterion

Below, the first criterion is restated, and the proof of the associated theorem is given in full:

Definition 4 (Generalized Adjustment Criterion Type 1). *A set \mathbf{Z} satisfies the generalized criterion relative to the pair \mathbf{X} and \mathbf{Y} in a causal model with graph G , augmented with the selection mechanism S if:*

- (a) *No element of \mathbf{Z} is a descendant in $G_{\bar{\mathbf{X}}}$ of any $W \notin \mathbf{X}$ which lies on a proper causal path from \mathbf{X} to \mathbf{Y} .*
- (b) *All non-causal paths between \mathbf{X} and \mathbf{Y} in G are blocked by \mathbf{Z} .*
- (c) *\mathbf{Y} is d -separated from S given \mathbf{X} under the intervention $do(\mathbf{x})$, i.e., $(\mathbf{Y} \perp\!\!\!\perp S | \mathbf{X})_{G_{\bar{\mathbf{X}}}}$.*
- (d) *\mathbf{Z} can be partitioned into sets $\mathbf{Z}^+, \mathbf{Z}^-$ such that $\mathbf{Z}^+ = \{Z' \in \mathbf{Z} \mid (Z' \perp\!\!\!\perp \mathbf{X} | S)_{G_{\bar{\mathbf{X}}(S)}}\}$ and $(\mathbf{Y} \perp\!\!\!\perp \mathbf{Z}^- | \mathbf{X}, \mathbf{Z}^+, S)_{G_{\bar{\mathbf{X}}}}$.*

$G_{\bar{\mathbf{X}}(S)}$ is the graph where all edges into $X \in \mathbf{X} \setminus An_S$ are removed, where An_S is the set of ancestors of the variable S in G .

Theorem 1 (Generalized Adjustment Formula Type 1). *Given disjoint sets of variables \mathbf{X}, \mathbf{Y} and \mathbf{Z} in a causal model with graph G . The effect $P(\mathbf{y} | do(\mathbf{x}))$ is given by*

$$P(\mathbf{y} | do(\mathbf{x})) = \sum_{\mathbf{z}} P(\mathbf{y} | \mathbf{x}, \mathbf{z}, S=1) P(\mathbf{z} | S=1) \quad (11)$$

in every model inducing G if and only if \mathbf{Z} satisfies the generalized adjustment criterion type 1 relative to the pair \mathbf{X}, \mathbf{Y} .

Proof. (if) Suppose \mathbf{Z} satisfy the the conditions of the theorem relative to the pair \mathbf{X}, \mathbf{Y} . Then \mathbf{Z} can be partitioned into \mathbf{Z}^- and \mathbf{Z}^+ as in Def. 4. The causal effect is derived as follows:

First, by condition (c), $(\mathbf{Y} \perp\!\!\!\perp S \mid X)_{G_{\overline{\mathbf{X}}}}$ and S can be introduced into the expression:

$$P(\mathbf{y} \mid do(\mathbf{x})) = P(\mathbf{y} \mid do(\mathbf{x}), S=1) \quad (25)$$

Conditioning on \mathbf{Z}^+ , it becomes:

$$P(\mathbf{y} \mid do(\mathbf{x})) = \sum_{\mathbf{z}^+} P(\mathbf{y} \mid do(\mathbf{x}), \mathbf{z}^+, S=1) P(\mathbf{z}^+ \mid do(\mathbf{x}), S=1) \quad (26)$$

Cond. (d), $(\mathbf{Z}^+ \perp\!\!\!\perp \mathbf{X} \mid S)_{G_{\overline{\mathbf{X}(S)}}$, and rule 3 of the do-calculus allow the removal of $do(\mathbf{x})$ from the second term of the previous expression

$$P(\mathbf{y} \mid do(\mathbf{x})) = \sum_{\mathbf{z}^+} P(\mathbf{y} \mid do(\mathbf{x}), \mathbf{z}^+, S=1) P(\mathbf{z}^+ \mid S=1) \quad (27)$$

Summing over \mathbf{Z}^- in the second factor yields

$$P(\mathbf{y} \mid do(\mathbf{x})) = \sum_{\mathbf{z}^+} P(\mathbf{y} \mid do(\mathbf{x}), \mathbf{z}^+, S=1) \sum_{\mathbf{z}^-} P(\mathbf{z} \mid S=1) \quad (28)$$

The summations can be put together, and by cond. (d), $(\mathbf{Y} \perp\!\!\!\perp \mathbf{Z}^- \mid \mathbf{X}, \mathbf{Z}^+, S)_{G_{\overline{\mathbf{X}}}}$ holds true. Then, the variables in \mathbf{Z}^- can be introduced into the first factor,

$$P(\mathbf{y} \mid do(\mathbf{x})) = \sum_{\mathbf{z}} P(\mathbf{y} \mid do(\mathbf{x}), \mathbf{z}, S=1) P(\mathbf{z} \mid S=1) \quad (29)$$

Conditions (a) and (b) imply $(\mathbf{Y} \perp\!\!\!\perp \mathbf{X} \mid \mathbf{Z})_{G_{\overline{\mathbf{X}}}}$, furthermore condition (c) ensures that observing S will not open any path between \mathbf{X} and \mathbf{Y} , because such path will either violate (c) or have some $X \in \mathbf{X}$ as a collider which contradicts (b). Hence, $(\mathbf{Y} \perp\!\!\!\perp \mathbf{X} \mid \mathbf{Z}, S)_{G_{\overline{\mathbf{X}}}}$ holds and can be used together with rule 2 of do-calculus to remove of the $do()$ operator from the first factor of eq. (29), which results in the adjustment formula

$$P(\mathbf{y} \mid do(\mathbf{x})) = \sum_{\mathbf{z}} P(\mathbf{y} \mid \mathbf{x}, \mathbf{z}, S=1) P(\mathbf{z} \mid S=1) \quad (30)$$

(Only if) This direction is proved using the contrapositive. Suppose conditions (a) and (b) do not hold, then for any model compatible with $G_{\overline{\mathbf{S}}}$, which is also compatible with G , the adjustment formula is equal to $\sum_{\mathbf{z}} P(\mathbf{y} \mid \mathbf{x}, \mathbf{z}) P(\mathbf{z})$. But by the adjustment criterion [Shpitser et al., 2010] this expression will not be equal to $P(\mathbf{y} \mid do(\mathbf{x}))$.

For conditions (c) and (d), counter examples for the identifiability and recoverability of the causal effect are shown. In every case, let \mathbf{V} represent all variables in the graph except for the selection mechanism S , and Q refer to the adjustment formula as in (30). We construct two Structural Causal Models (SCM) M_1 and M_2 , that induce probability distributions P_1 and P_2 respectively. Both, models will be compatible with G , and they will agree in the probability distribution under selection bias

$$P_1(\mathbf{v} \mid S=1) = P_2(\mathbf{v} \mid S=1) \quad (31)$$

but Q_1 in the first model provides a different distribution than Q_2 in the second model. Let M_1 be compatible with G and M_2 with $G_{\overline{\mathbf{S}}}$ (in M_2 , S is independent from all other variables) such that $(\mathbf{V} \perp\!\!\!\perp S)_{P_2}$. Recoverability should hold for any parametrization, hence without loss of generality, all variables are assumed to be binary. The construction parametrizes P_1 through its factors (as in lemma 6) and then parametrizes P_2 to enforce (31). As a consequence, (31) is also equals to $P_2(\mathbf{v})$.

Suppose condition (c) does not hold, then, there is an open path between \mathbf{Y} and S in $G_{\overline{\mathbf{X}}}$. Without loss of generality, our attention can be directed into the particular $Y' \in \mathbf{Y}$ not satisfying the condition, and on the causal effect for Y' . To do this, the constructed model will have every variable in $\mathbf{Y} \setminus \{Y'\}$ disconnected from the graph, more precisely $(\mathbf{Y} \setminus \{Y'\}) \perp\!\!\!\perp \mathbf{V}$ holds, so that:

$$\begin{aligned}
P(\mathbf{y} \mid do(\mathbf{x})) &= \sum_{\mathbf{z}} P(\mathbf{y} \mid \mathbf{x}, \mathbf{z})P(\mathbf{z}) \\
&= \prod_{\mathbf{Y}} \sum_{\mathbf{z}} P(y \mid \mathbf{x}, \mathbf{z})P(\mathbf{z}) \\
&= \left(\prod_{\mathbf{Y} \setminus Y'} P(y) \right) \sum_{\mathbf{z}} P(y' \mid \mathbf{x}, \mathbf{z})P(\mathbf{z}) \\
&= \gamma \sum_{\mathbf{z}} P(y' \mid \mathbf{x}, \mathbf{z})P(\mathbf{z})
\end{aligned}$$

where γ represents the product of the marginal distribution of the remaining $\mathbf{Y} \setminus \{Y'\}$.

The following are the cases for which Y' may violate cond. (c). Figure 7 illustrate every case for easier reference.

case 1: $Y' \in Pa_S$

Let \mathbf{W} be the set of nodes connecting \mathbf{X} and Y' with directed paths. Consider the induced subgraph G' where all nodes in $\mathbf{V} \setminus \{\mathbf{X}, \mathbf{W}, Y', S\}$ are disconnected from $\{\mathbf{X}, \mathbf{W}, Y', S\}$. It must be the case that \mathbf{Z} and \mathbf{W} are disjoint, else condition (a) is violated. Consequently, every \mathbf{Z} is disconnected, and $(\mathbf{Z} \perp\!\!\!\perp Y')_{P_1}$ holds. M_1 and M_2 are constructed from G' , the adjustment formula in the second model can be expressed as:

$$\begin{aligned}
Q_2 &= \gamma \sum_{\mathbf{z}} P_2(y' \mid \mathbf{x}, \mathbf{z})P_2(\mathbf{z}) \\
&= \gamma \sum_{\mathbf{z}} P_1(y' \mid \mathbf{x}, \mathbf{z}, S=1)P_1(\mathbf{z} \mid S=1) \\
&= \gamma \sum_{\mathbf{z}} P_1(y' \mid \mathbf{x}, S=1)P_1(\mathbf{z} \mid S=1) \\
&= \gamma P_1(y' \mid \mathbf{x}, S=1) \\
&= \gamma \frac{P_1(y', \mathbf{x}, S=1)}{\sum_{Y'} P_1(y', \mathbf{x}, S=1)} \\
&= \gamma \frac{P_1(S=1 \mid y')P_1(y' \mid \mathbf{x})}{P_1(S=1 \mid y')P_1(y' \mid \mathbf{x}) + P_1(S=1 \mid \overline{y'})P_1(\overline{y'} \mid \mathbf{x})}
\end{aligned}$$

Using lemma 6, let $P_1(S=1 \mid y') = \alpha$ and $P_1(S=1 \mid \overline{y'}) = \beta$ with $0 < \alpha, \beta < 1$ and $\alpha \neq \beta$. Proceed with lemma 8 ($p = q = 1/2$) to define $P(y' \mid \mathbf{x}) = 1/2$. The previous expression becomes:

$$Q_2 = \gamma \frac{\alpha}{\alpha + \beta}$$

Following a similar derivation it can be established that $Q_1 = \gamma/2$ which is never equal to Q_2 in this parametrization.

case 2: There is a directed path p from Y' to S without any \mathbf{Z} .

Let R be the parent of S in such path, let \mathbf{W}_1 be the set of nodes between \mathbf{X} and Y' as in the previous

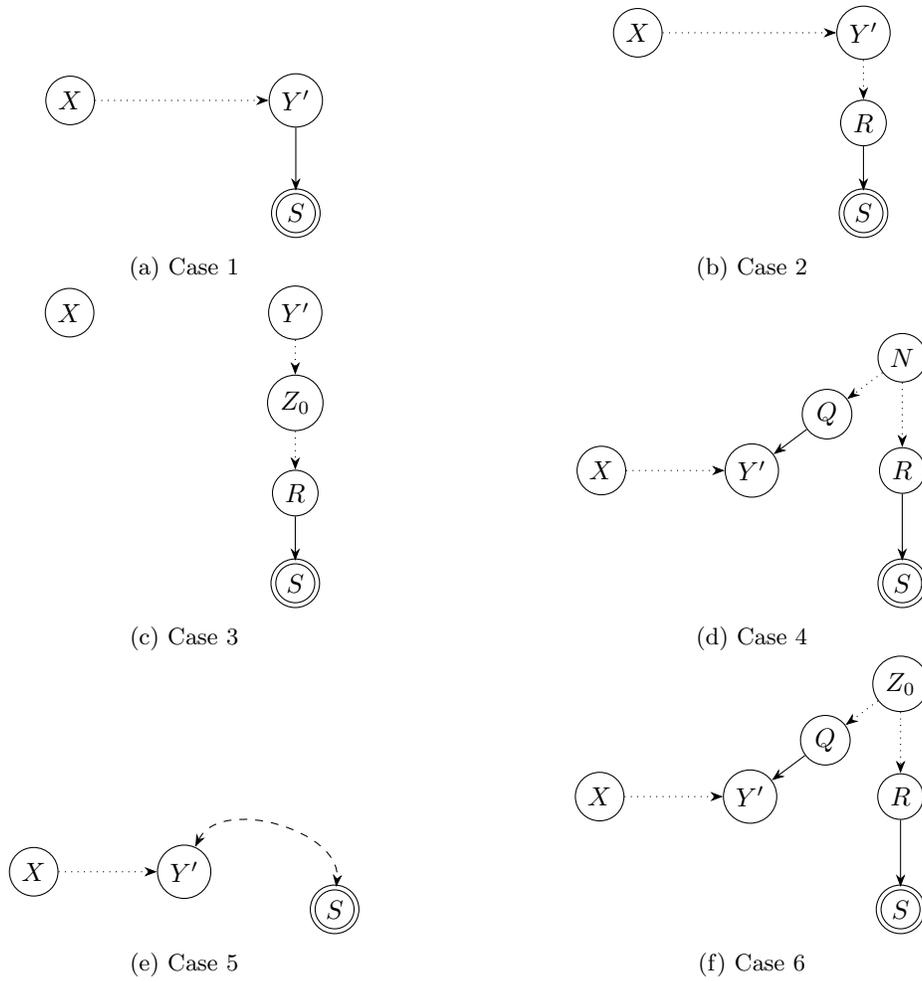


Figure 7: Cases considered for the necessity of condition (c) in the proof for Thm. 1. Dotted directed arrows indicate chains of arbitrary length in the graph.

case. Similarly let \mathbf{W}_2 be the variables in the path from Y' to R . Now, consider the graph G' where all nodes except for $\{\mathbf{X}, \mathbf{W}_1, Y', \mathbf{W}_2, R, S\}$ are disconnected from $\{\mathbf{X}, \mathbf{W}_1, Y', \mathbf{W}_2, R, S\}$. Proceeding as in the previous case, with the consideration that \mathbf{Z} is disconnected from the rest of the graph, yields:

$$Q_2 = \gamma P_1(y' | \mathbf{x}, S=1) = \gamma \frac{P_1(y', \mathbf{x}, S=1)}{P_1(\mathbf{x}, S=1)}$$

The numerator can be rewritten as:

$$\begin{aligned} P_1(y', \mathbf{x}, S=1) &= \sum_R P_1(y', \mathbf{x}, r, S=1) \\ &= \sum_R P_1(\mathbf{x}) P_1(y' | \mathbf{x}) P_1(r | y') P_1(S=1 | r) \end{aligned}$$

Factorizing the denominator analogously, the term $P_1(\mathbf{x})$ is the same and can be cancelled out, then Q_2 becomes:

$$Q_2 = \gamma \frac{P_1(y' | \mathbf{x}) \sum_R P_1(r | y') P_1(S=1 | r)}{\sum_{Y'} P_1(y' | \mathbf{x}) \sum_R P_1(r | y') P_1(S=1 | r)}$$

Using lemma 8 to set $P_1(r | y') = 1/2 + \epsilon/2$, $P_1(r | \bar{y}') = 1/2 - \epsilon/2$ where $\epsilon = (1/5)^k$ (using $p = 3/5$, $q = 2/5$), and defining $P(S=1 | r) = 2/3$ and $P(S=1 | \bar{r}) = 1/2$ leads to $Q_2 = \gamma(1/2 + \epsilon/14)$ and $Q_1 = \gamma/2$ which are never equal.

case 3: There is a directed path p from Y' to S that contains some $Z_0 \in \mathbf{Z}$

Let R be the parent of S in such path. It can be assured that, \mathbf{X} and Y' are not connected by any causal path, otherwise Z_0 violates condition (a). Let \mathbf{W}_1 be the nodes in the subpath between Y' and Z_0 , and \mathbf{W}_2 those in between Z_0 and R . Consider the graph G' where all nodes except for $\{\mathbf{X}, Y', \mathbf{W}_1, Z_0, \mathbf{W}_2, R, S\}$ are disconnected from $\{\mathbf{X}, Y', \mathbf{W}_1, Z_0, \mathbf{W}_2, R, S\}$.

Every $\mathbf{Z}' = \mathbf{Z} \setminus \{Z_0\}$ is disconnected from the rest of the graph, then:

$$\begin{aligned} Q_2 &= \gamma \sum_{\mathbf{z}} P_1(y' | \mathbf{x}, \mathbf{z}, S=1) P_1(\mathbf{z} | S=1) \\ &= \gamma \sum_{z_0} \sum_{\mathbf{z}'} P_1(y' | \mathbf{x}, z_0, \mathbf{z}', S=1) P_1(z_0, \mathbf{z}' | S=1) \\ &= \gamma \sum_{z_0} P_1(y' | \mathbf{x}, z_0, S=1) P_1(z_0 | S=1) \\ &= \gamma \sum_{z_0} P_1(y' | \mathbf{x}, z_0) P_1(z_0 | S=1) && (Y' \perp\!\!\!\perp S | \mathbf{X}, Z_0)_{P_1} \\ &= \gamma \sum_{z_0} \frac{P_1(y', \mathbf{x}, z_0)}{\sum_{Y'} P_1(y', \mathbf{x}, z_0)} P_1(z_0 | S=1) \end{aligned}$$

The numerator of the fraction in the last expression is equal to:

$$P_1(y', \mathbf{x}, z_0) = P_1(\mathbf{x}) P_1(y') P_1(z_0 | y')$$

A similar factorization can be employed for the denominator, as well as for Q_1 . The factor $P_1(\mathbf{x})$

appears in both parts of the fractions and can be canceled:

$$Q_1 = \gamma \sum_{z_0} \frac{P_1(y')P_1(z_0 | y')}{\sum_{Y'} P_1(y')P_1(z_0 | y')} P_1(z_0)$$

$$Q_2 = \gamma \sum_{z_0} \frac{P_1(y')P_1(z_0 | y')}{\sum_{Y'} P_1(y')P_1(z_0 | y')} P_1(z_0 | S=1)$$

Now, $P_1(z_0)$ and $P_1(z_0 | S=1)$ are derived in similar terms:

$$P_1(z_0) = \sum_{Y'} P_1(z_0, y') = \sum_{Y'} P_1(y')P_1(z_0 | y')$$

$$P_1(z_0 | S=1) = \frac{P_1(z_0, S=1)}{P_1(S=1)} = \frac{P_1(S=1 | z_0) \sum_{Y'} P_1(y')P_1(z_0 | y')}{P_1(S=1)}$$

Replacing $P_1(z_0)$ and $P_1(z_0 | S=1)$ in Q_1 and Q_2 , then simplifying:

$$Q_1 = \gamma \sum_{z_0} P_1(y')P(z_0 | y') = \gamma P_1(y')$$

$$Q_2 = \gamma \frac{\sum_{z_0} P_1(y')P_1(z_0 | y')P_1(S=1 | z_0)}{P_1(S=1)} = \gamma \frac{P_1(y', S=1)}{P_1(S=1)} = \gamma \frac{P_1(y')P_1(S=1 | y')}{\sum_{Y'} P_1(S=1 | y')P_1(y')}$$

The term $P_1(S=1 | y') = \sum_R P_1(S=1 | r)P_1(r | y')$. Lemma 8 can be employed exactly as in the previous case, and $P_1(y')$ can be defined directly since it has no parents, for instance $P_1(y') = 1/2$, then the queries end up as:

$$Q_1 = \gamma \frac{1}{2} \qquad Q_2 = \gamma \left(\frac{1}{2} + \frac{\epsilon}{14} \right)$$

Which are never equal.

case 4: There is a path p connecting Y' and S that goes through and ancestor of both, and does not contain any node in \mathbf{Z} .

Let N be the closest common ancestor of Y' and S . Let R be the parent of S and Q the parent of Y' in the mentioned path. Let \mathbf{W}_1 be the set of nodes between \mathbf{X} and Y' . Let \mathbf{W}_2 and \mathbf{W}_3 be the nodes in the paths from N to Q and from N to R respectively. Consider the graph G' where the arrows in the subpath from N to Q are reversed and all nodes except for $\{\mathbf{X}, \mathbf{W}_1, Y', Q, \mathbf{W}_2, N, \mathbf{W}_3, R, S\}$ are disconnected from $\{\mathbf{X}, \mathbf{W}_1, Y', Q, \mathbf{W}_2, N, \mathbf{W}_3, R, S\}$. Any model constructed for G' can be translated to a model compatible with G using lemma 7. Following the same derivation as in case 2 (taking into account that \mathbf{Z} is disconnected from the rest of the graph) yields:

$$Q_2 = \gamma \frac{P_1(y', \mathbf{x}, S=1)}{\sum_{Y'} P_1(y', \mathbf{x}, S=1)}$$

The numerator of the last expression can be rewritten as:

$$P_1(y', \mathbf{x}, S=1) = \sum_Q P_1(y', \mathbf{x}, q, S=1)$$

$$= \sum_Q P_1(\mathbf{x})P_1(y' | \mathbf{x}, q)P_1(q)P_1(S=1 | q)$$

By rewriting the denominator similarly, the term $P_1(\mathbf{x})$ appearing in both vanishes, then Q_2 becomes:

$$Q_1 = \gamma \frac{\sum_Q P_1(y' | \mathbf{x}, q) P_1(q)}{\sum_{Y', Q} P_1(y' | \mathbf{x}, q) P_1(q)}$$

$$Q_2 = \gamma \frac{\sum_Q P_1(y' | \mathbf{x}, q) P_1(q) P_1(S=1 | q)}{\sum_{Y', Q} P_1(y' | \mathbf{x}, q) P_1(q) P_1(S=1 | q)}$$

Lemma 8 can be employed to set $P_1(r | q) = 1/2 + \epsilon/2, P_1(r | \bar{q}) = 1/2 - \epsilon/2$ where $\epsilon = (1/5)^k$ (using $p = 3/5, q = 2/5$). Define $P(S=1 | r) = 2/3$ and $P(S=1 | \bar{r}) = 1/2$. Calculate $P(S=1 | q)$ as $\sum_R P_1(r | q) P_1(S=1 | r)$. Also by lemma 8 let $P_1(y' | q, x) = P_1(y' | q, \bar{x}) = 3/4, P_1(y' | \bar{q}, x) = P_1(y' | \bar{q}, \bar{x}) = 1/2$. It leads to:

$$Q_1 = \gamma \frac{3}{8} \qquad Q_2 = \gamma \left(\frac{3}{8} + \frac{\epsilon}{56} \right)$$

which are never equal.

case 5: There is confounding path between Y' and S consisting of unobservable variables.

The models for this case can be constructed as in case 4, then moving the variables in the in the path from Q to R (included) from the set of observables to the set of unobservables.

case 6: There is a path p connecting Y' and S that goes through an ancestor of both, and contains some $Z_0 \in \mathbf{Z}$.

Let N, Q, R be defined as in the previous case, also construct G' the same way. Following the same derivation strategy as in case 3, the query expressions become:

$$Q_1 = \gamma \sum_{z_0} \frac{\sum_Q P_1(y' | \mathbf{x}, q) P_1(q) P_1(z_0 | q)}{\sum_{Y', Q} P_1(y' | \mathbf{x}, q) P_1(q) P_1(z_0 | q)} P_1(z_0)$$

$$Q_2 = \gamma \sum_{z_0} \frac{\sum_Q P_1(y' | \mathbf{x}, q) P_1(q) P_1(z_0 | q)}{\sum_{Y', Q} P_1(y' | \mathbf{x}, q) P_1(q) P_1(z_0 | q)} P_1(z_0 | S=1)$$

Now, $P_1(z_0)$ and $P_1(z_0 | S=1)$ are derived in similar terms:

$$P_1(z_0) = \sum_Q P_1(z_0, q) = \sum_Q P_1(q) P_1(z_0 | q)$$

$$P_1(z_0 | S=1) = \frac{\sum_R P_1(S=1, z_0, r)}{\sum_{z_0, R} P_1(S=1, z_0, r)} = \frac{P_1(z_0) \sum_R P_1(r | z_0) P_1(S=1 | r)}{\sum_{z_0} P_1(z_0) \sum_R P_1(r | z_0) P_1(S=1 | r)}$$

Use lemma 8 to parametrize $P_1(r | z_0) = 1/2 + \epsilon_1/2, P_1(r | \bar{z}_0) = 1/2 - \epsilon_1/2, P(z_0 | q) = 1/2 + \epsilon_2/2, P(z_0 | \bar{q}) = 1/2 - \epsilon_2/2$ where $\epsilon_i = (1/5)^{k_i}, i = \{1, 2\}$ (using $p = 3/5, q = 2/5$ in both cases). Define $P(S=1 | r) = 2/3$ and $P(S=1 | \bar{r}) = 1/2$. Also by lemma 8 let $P_1(y' | q, x) = P_1(y' | q, \bar{x}) = 3/4, P_1(y' | \bar{q}, x) = P_1(y' | \bar{q}, \bar{x}) = 1/2$. The queries end up as:

$$Q_1 = \gamma \frac{3}{8} \qquad Q_2 = \gamma \left(\frac{3}{8} + \frac{\epsilon_1 \epsilon_2}{56} \right)$$

Which are never equal.

Now, suppose condition (d) does not hold. It should be the case that some $Z_0 \in \mathbf{Z}^-$ is connected to \mathbf{Y} given $\mathbf{X}, \mathbf{Z}^+, S$ in $G_{\bar{\mathbf{X}}}$. There two possible cases (depicted in Fig. 8) not contradicting previous conditions:

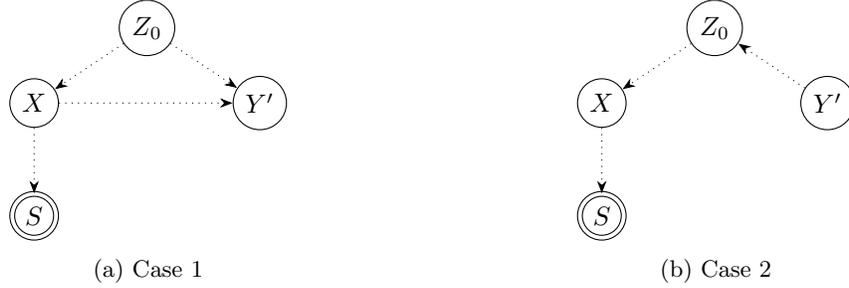


Figure 8: Cases considered for the necessity of condition (d) in the proof for Thm. 1. Dotted directed arrows indicate chains of arbitrary length in the graph.

case 1: Z_0 is an ancestor of \mathbf{X} and Y' , and S is a descendant of \mathbf{X} .

Let \mathbf{W}_1 be the nodes in the path between Z_0 and \mathbf{X} , \mathbf{W}_2 those between Z_0 and Y' , and \mathbf{W}_3 those between \mathbf{X} and S . As in previous cases, consider the graph G' where all nodes but $\{\mathbf{X}, Z_0, Y', \mathbf{W}_1, \mathbf{W}_2, \mathbf{W}_3\}$ are disconnected from this set. Also, suppose \mathbf{X} and Y' are not connected by any path not going through Z_0 . The queries in the corresponding models can be expressed as:

$$Q_1 = \gamma \sum_{Z_0} P_1(y' | \mathbf{x}, z_0) P_1(z_0) = \gamma \sum_{Z_0} P_1(y' | z_0) P_1(z_0) = \gamma P_1(y')$$

$$Q_2 = \gamma \sum_{Z_0} P_1(y' | \mathbf{x}, z_0, S=1) P_1(z_0 | S=1) = \gamma \sum_{Z_0} P_1(y' | z_0) P_1(z_0 | S=1)$$

The term $P_1(z_0 | S=1)$ is available as:

$$P_1(z_0 | S=1) = \frac{P_1(z_0) \sum_{\mathbf{x}} P_1(\mathbf{x} | z_0) P(S=1 | \mathbf{x})}{\sum_{\mathbf{x}, z_0} P_1(z_0) P_1(\mathbf{x} | z_0) P(S=1 | \mathbf{x})}$$

Let $P_1(z_0) = 1/2$, $P_1(y' | z_0) = 1/2 + \epsilon_1/2$, $P(y' | \bar{z}_0) = 1/2 - \epsilon_1/2$, $P_1(\mathbf{x} | z_0) = 1/2 + \epsilon_2/2$, $P(\mathbf{x} | \bar{z}_0) = 1/2 - \epsilon_2/2$. Also $P_1(S=1 | \mathbf{x}) = 1/2 + \epsilon_3$, $P_1(S=1 | \bar{\mathbf{x}}) = 1/2 - \epsilon_3$ where $\epsilon_i = (1/5)^{k_i}$, $i = 1, 2, 3$ (using lemma 8 with $p = 3/5$, $q = 2/5$ in all cases):

$$Q_1 = \frac{1}{2} \gamma \qquad Q_2 = \gamma \left(\frac{1}{2} + \frac{\epsilon_1 \epsilon_2 \epsilon_3}{2} \right)$$

Which are never equal.

case 2: Z_0 is an ancestor of \mathbf{X} and a descendant of Y' and S is a descendant of \mathbf{X} .

In this case there are no causal paths between \mathbf{X} and Y' otherwise Z_0 violates condition (a). Lemma 7 can be used to change the direction of the edges in the path from Y' to Z_0 while staying in the same equivalence class, then the same parametrization from the previous case applies.

□

A.2 Proof for the Second Criterion

Here the definition and theorem is stated again and then proved in full:

Definition 5 (Generalized Adjustment Criterion Type 2). *A set \mathbf{Z} satisfies the generalized criterion relative to \mathbf{X}, \mathbf{Y} , a set of variables measured under selection bias \mathbf{M} and a set of variables observed in the overall population \mathbf{T} in a causal model with graph G augmented with the selection mechanism S if:*

- (a) No element of \mathbf{Z} is a descendant in $G_{\overline{\mathbf{X}}}$ of any $W \notin \mathbf{X}$ which lies on a proper causal path from \mathbf{X} to \mathbf{Y} .
- (b) All non-causal paths between \mathbf{X} and \mathbf{Y} in G are blocked by \mathbf{Z} .
- (c) \mathbf{Y} is d -separated from the selection mechanism S given \mathbf{Z} and \mathbf{X} , i.e., $(\mathbf{Y} \perp\!\!\!\perp S \mid \mathbf{X}, \mathbf{Z})$.
- (d) The variables are measured with bias $(\mathbf{Z}, \mathbf{X}, \mathbf{Y} \subseteq \mathbf{M})$ and the covariates are available without bias $(\mathbf{Z} \subseteq \mathbf{T})$

Theorem 2 (Generalized Adjustment Formula Type 2). *Given disjoint sets of variables \mathbf{X}, \mathbf{Y} and \mathbf{Z} , and sets \mathbf{M}, \mathbf{T} in a causal model with graph G . In every model inducing G , the effect $P(\mathbf{y} \mid do(\mathbf{x}))$ is given by*

$$P(\mathbf{y} \mid do(\mathbf{x})) = \sum_{\mathbf{z}} P(\mathbf{y} \mid \mathbf{x}, \mathbf{z}, S=1)P(\mathbf{z}) \quad (18)$$

if and only if the set \mathbf{Z} satisfies the generalized adjustment criterion type 2 relative to the pair \mathbf{X}, \mathbf{Y} .

Proof. (if) Suppose \mathbf{Z} satisfy the conditions of the theorem relative to the pair \mathbf{X}, \mathbf{Y} and the sets \mathbf{M}, \mathbf{T} . By conditions (a) and (b), the effect can be written as:

$$P(\mathbf{y} \mid do(\mathbf{x})) = \sum_{\mathbf{z}} P(\mathbf{y} \mid \mathbf{x}, \mathbf{z})P(\mathbf{z})$$

Note that S can be introduced to the first term by cond. (c), which entail Eq. (18). Cond. (d) ensures that both terms in the expression are estimable from the available distributions.

(Only if) This direction is proved using the contrapositive: Assume conditions (a) or (b) do not hold, then, the same argument as in the proof of the previous type applies here. That is, even if recoverable the adjustment formula does not rely the correct causal effect.

As in the previous theorem, the prove of necessity for condition (c), required counter-examples for the identifiability and recoverability of the causal effect. Let \mathbf{V} represents all variables in the graph except for the selection mechanism S , and Q refer to the adjustment formula as in (18) we construct two SCMs, compatible with G such that they agree in the probability distribution under selection bias:

$$P_1(\mathbf{v} \mid S=1) = P_2(\mathbf{v} \mid S=1) \quad (32)$$

and on the non-biased distribution over \mathbf{Z} ,

$$P_1(\mathbf{z}) = P_2(\mathbf{z}) \quad (33)$$

but Q_1 in the first model provides a different distribution than Q_2 in the second model. Let P_1 be compatible with G and P_2 compatible with $G_{\overline{S}}$ (a model where S is independent from all other variables) such that $(\mathbf{V} \perp\!\!\!\perp S)_{P_2}$. Recoverability should hold for any parametrization, hence, without loss of generality, all variables are assumed to be binary. Every construction parametrizes P_1 through its factors (as in lemma 6) and then parametrize P_2 to enforce (32) and (33). Moreover, (32) also equals to $P_2(\mathbf{v})$.

Now, lets assume condition (c) does not hold, then there is an open path between \mathbf{Y} and S not blocked when \mathbf{Z} is observed. As for the previous proof, a variable $Y' \in \mathbf{Y}$ not satisfying the condition is fixed and the query is studied as:

$$Q = \gamma \sum_{\mathbf{z}} P(y' \mid \mathbf{x}, \mathbf{z})P(\mathbf{z})$$

where γ represents the product of the marginal distribution of the remaining \mathbf{Y} .

Fig. 9 illustrates the cases in which condition (c) may be unsatisfied, they are described in the sequel:

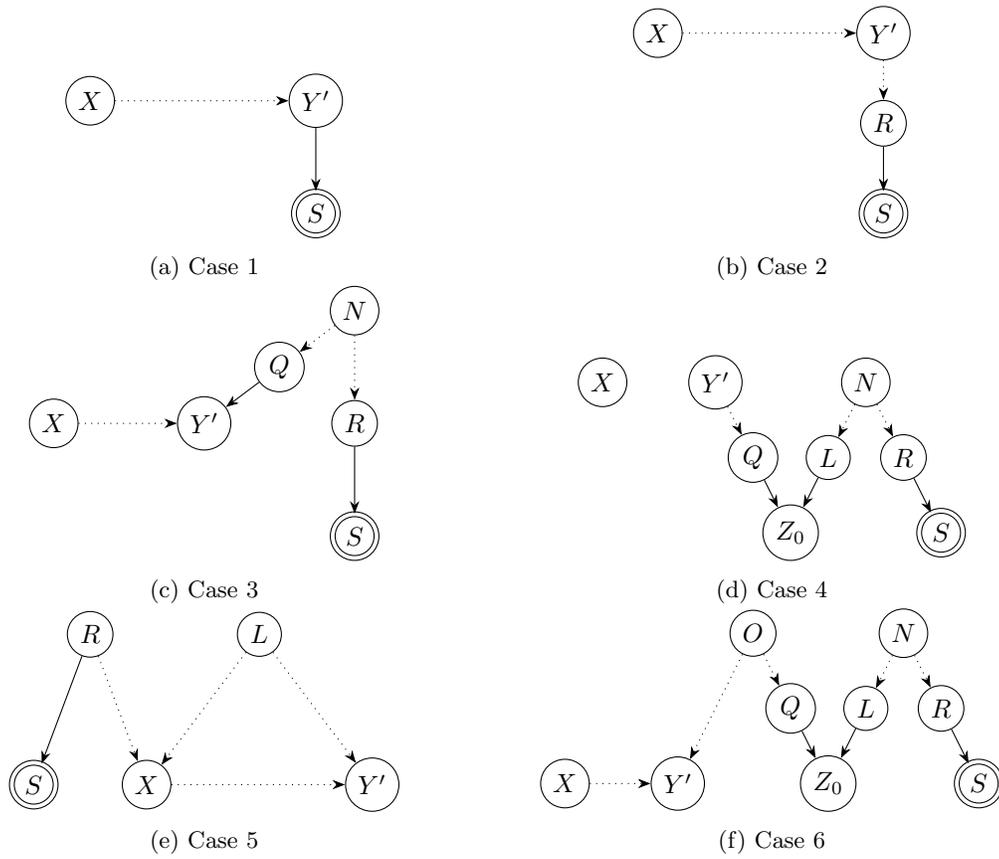


Figure 9: Cases considered for the necessity of condition (c) in the proof for Thm. 2. Dotted directed arrows indicate chains of arbitrary length in the graph.

case 1: $Y' \in Pa_S$

Proceed exactly as in case 1 of the proof for Thm. 1.

case 2: There is a directed path p from Y' to S not containing any node in \mathbf{Z}

Same as case 2 of Thm. 1.

case 3: There is a path from Y' to S passing through an ancestor of both, with no colliders or any node in \mathbf{Z} .

Same as case 4 in the proof for Thm. 1.

case 4: There is a non-directed path from Y' to S with some $Z_0 \in \mathbf{Z}$ as a collider.

\mathbf{X} must be disconnected from Y' otherwise condition (a) is violated since Z_0 is a descendant of Y' which belongs to every causal path from \mathbf{X} to Y' . Let N be the common ancestor of S and Z_0 , let L be the parent of Z_0 and R be the parent of S in that path. Consider the subgraph G' where all nodes but for those in the mentioned path between Y' and S are disconnected from all others, using lemma 7 in the portion between L and N to reverse the directionality of the arrows and still produce an equivalent model for the original graph. This case differs from the previous ones, here the causal effect $P(y' | do(\mathbf{x})) = P(y')$, however the adjustment formula Q will not be equal to this effect in every model compatible with the graph. To see it, it is enough to construct one model M as follows:

$$\begin{aligned} Q &= \gamma \sum_{\mathbf{z}} P(y' | \mathbf{x}, \mathbf{z}, S=1) P(\mathbf{z}) \\ &= \gamma \sum_{z_0} \sum_{\mathbf{z}'} P(y' | \mathbf{x}, z_0, \mathbf{z}', S=1) P(z_0, \mathbf{z}') \\ &= \gamma \sum_{z_0} P(y' | \mathbf{x}, z_0, S=1) P(z_0) \\ &= \gamma \sum_{z_0} \frac{P(y', \mathbf{x}, z_0, S=1)}{\sum_{Y'} P(y', \mathbf{x}, z_0, S=1)} P(z_0) \end{aligned}$$

The numerator of the fraction in the last expression can be decomposed as:

$$\begin{aligned} P(y', \mathbf{x}, z_0, S=1) &= \sum_{Q,L} P(y', \mathbf{x}, q, z_0, l, S=1) \\ &= P(\mathbf{x}) P(y') \sum_{Q,L} P(q | y') P(z_0 | q, l) P(l) P(S=1 | l) \end{aligned}$$

Performing a similar decomposition in the denominator, and considering that the factor $P(\mathbf{x})$ appears in both and can be cancelled out, the queries become:

$$Q = \gamma \sum_{z_0} \frac{P(y') \sum_{L,Q} P(q | y') P(z_0 | q, l) P(l) P(S=1 | l)}{\sum_{Y'} P(y') \sum_{L,Q} P(q | y') P(z_0 | q, l) P(l) P(S=1 | l)} P(z_0)$$

The term $P(z_0)$ can be computed as:

$$P(z_0) = \sum_{Y'} P(y') \sum_{L,Q} P(q | y') P(l) P(z_0 | q, l)$$

And $P(S=1 | l)$ as

$$P(S=1 | l) = \sum_R P(r | l) P(S=1 | r)$$

Using lemma 6 set $P(y') = P(l) = 1/2$, $P(z_0 | q, l) = 1/3$, $P(z_0 | q, \bar{l}) = 2/3$, $P(z_0 | \bar{q}, l) = 1/2$, $P(z_0 | \bar{q}, \bar{l}) = 1/2$. By lemma 8 let $P(q | y) = 1/2 - \epsilon$, $P(q | \bar{y}) = 1/2 + \epsilon$, $P(r | l) = 1/2 - \epsilon_2$, $P(r | \bar{l}) = 1/2 + \epsilon_1$, where $\epsilon_i = (1/5)^{k_i}$, $i = 1, 2$ (having $p = 3/5$, $q = 2/5$ in both cases). Let $P(S=1 | l) = \alpha$, $P(S=1 | \bar{l}) = \beta$, and pick any $0 < \alpha, \beta < 1$. This result in:

$$P(\mathbf{y} | do(\mathbf{x})) = \frac{\gamma}{2}$$

$$Q = \frac{\gamma}{2} \frac{1764 - \epsilon_2^2 - \epsilon_1 \epsilon_2^2}{1764 - \epsilon_2^2}$$

Q is always different that $P(\mathbf{y} | do(\mathbf{x}))$.

case 5: There is a path between Y' and S passing by an ancestor of Y' having some $X' \in \mathbf{X}$ as a collider. Such path would have arrows incoming to X' . But it is also an open non-causal path between Y' and X' , violating condition (b).

case 6: There is a path p between Y' and S passing by an ancestor of Y' having some $Z_0 \in \mathbf{Z}$ as a collider. Let O be the common ancestor of Y' and Z_0 and let N be the common ancestor of Z_0 and S in that path. Let R be the parent of S , Q the parent of Z_0 in the path to Y' and L the parent of Z_0 in the path to S . Consider the subgraph G' where \mathbf{X} is disconnected from Y' and lemma 7 is applied to all the edges in the path between O and Y' . Then, the model can be parametrized exactly as for case 4.

If condition (d) does not hold it is either because, one, \mathbf{Z}, \mathbf{X} or \mathbf{Y} are not available in the biased distribution, in which case the estimation is impossible. Second, $P(\mathbf{z})$ is not available nor recoverable, therefore the adjustment expression given is not estimable either. \square

References

- [Angrist, 1997] Angrist, J. D. (1997). Conditional independence in sample selection models. *Economics Letters*, 54(2):103–112.
- [Bareinboim and Pearl, 2012] Bareinboim, E. and Pearl, J. (2012). Controlling selection bias in causal inference. In Lawrence, N. and Girolami, M., editors, *Proceedings of the 15th International Conference on Artificial Intelligence and Statistics (AISTATS)*, pages 100–108. JMLR, La Palma, Canary Islands.
- [Bareinboim and Pearl, 2016] Bareinboim, E. and Pearl, J. (2016). Causal inference and the data-fusion problem. *Proceedings of the National Academy of Sciences*, 113:7345–7352.
- [Bareinboim and Tian, 2015] Bareinboim, E. and Tian, J. (2015). Recovering causal effects from selection bias. *Proceedings of the Twenty-Ninth AAAI Conference on Artificial Intelligence*, pages 3475–3481.
- [Bareinboim et al., 2014] Bareinboim, E., Tian, J., and Pearl, J. (2014). Recovering from selection bias in causal and statistical inference. In Brodley, C. E. and Stone, P., editors, *Proceedings of the Twenty-eighth AAAI Conference on Artificial Intelligence*, pages 2410–2416, Palo Alto, CA. AAAI Press.
- [Cooper, 1995] Cooper, G. (1995). Causal discovery from data in the presence of selection bias. In *Proceedings of the Fifth International Workshop on Artificial Intelligence and Statistics*, pages 140–150.
- [Cortes et al., 2008] Cortes, C., Mohri, M., Riley, M., and Rostamizadeh, A. (2008). Sample selection bias correction theory. In *International Conference on Algorithmic Learning Theory*, pages 38–53. Springer.
- [Didelez et al., 2010] Didelez, V., Kreiner, S., and Keiding, N. (2010). Graphical models for inference under outcome-dependent sampling. *Statistical Science*, 25(3):368–387.
- [Evans and Didelez, 2015] Evans, R. J. and Didelez, V. (2015). Recovering from selection bias using marginal structure in discrete models.
- [Heckman, 1979] Heckman, J. J. (1979). Sample selection bias as a specification error. *Econometrica*, 47(1):153–161.
- [Kuroki and Cai, 2006] Kuroki, M. and Cai, Z. (2006). On recovering a population covariance matrix in the presence of selection bias. *Biometrika*, 93(3):601–611.
- [Little and Rubin, 1986] Little, R. J. A. and Rubin, D. B. (1986). *Statistical Analysis with Missing Data*. John Wiley & Sons, Inc., New York, NY, USA.
- [Maathuis and Colombo, 2015] Maathuis, M. H. and Colombo, D. (2015). A generalized back-door criterion. *Ann. Statist.*, 43(3):1060–1088.
- [Mefford and Witte, 2012] Mefford, J. and Witte, J. S. (2012). The covariate’s dilemma. *PLoS Genet*, 8(11):e1003096.
- [Pearl, 1993] Pearl, J. (1993). Aspects of graphical models connected with causality. In *Proceedings of the 49th Session of the International Statistical Institute*, pages 391–401, Tome LV, Book 1, Florence, Italy.
- [Pearl, 1995] Pearl, J. (1995). Causal diagrams for empirical research. *Biometrika*, 82(4):669–688.
- [Pearl, 2000] Pearl, J. (2000). *Causality: Models, Reasoning, and Inference*. Cambridge University Press, New York. 2nd edition, 2009.

- [Pearl and Paz, 2010] Pearl, J. and Paz, A. (2010). Confounding equivalence in causal equivalence. In *Proceedings of the Twenty-Sixth Conference on Uncertainty in Artificial Intelligence*, pages 433–441. AUAI, Corvallis, OR.
- [Pirinen et al., 2012] Pirinen, M., Donnelly, P., and Spencer, C. C. (2012). Including known covariates can reduce power to detect genetic effects in case-control studies. *Nature genetics*, 44(8):848–851.
- [Robins, 2001] Robins, J. M. (2001). Data, design, and background knowledge in etiologic inference. *Epidemiology*, 12(3):313–320.
- [Robinson and Jewell, 1991] 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*, pages 227–240.
- [Rubin, 1974] Rubin, D. (1974). Estimating causal effects of treatments in randomized and nonrandomized studies. *Journal of Educational Psychology*, 66:688–701.
- [Shpitser et al., 2010] Shpitser, I., VanderWeele, T. J., and Robins, J. M. (2010). On the validity of covariate adjustment for estimating causal effects. In *Proceedings of UAI 2010*, pages 527–536.
- [Takata, 2010] Takata, K. (2010). Space-optimal, backtracking algorithms to list the minimal vertex separators of a graph. *Discrete Applied Mathematics*, 158(15):1660–1667.
- [van der Zander et al., 2014] van der Zander, B., Liskiewicz, M., and Textor, J. (2014). Constructing separators and adjustment sets in ancestral graphs. In *Proceedings of UAI 2014*, pages 907–916.
- [Zadrozny, 2004] Zadrozny, B. (2004). Learning and evaluating classifiers under sample selection bias. In *Proceedings of the twenty-first international conference on Machine learning*, page 114. ACM.
- [Zhang, 2008] Zhang, J. (2008). On the completeness of orientation rules for causal discovery in the presence of latent confounders and selection bias. *Artif. Intell.*, 172:1873–1896.