

Causal simulation experiments: Lessons from bias amplification

Tyrel Stokes¹ , Russell Steele¹ and Ian Shrier^{2,3}

Statistical Methods in Medical Research

2022, Vol. 31(1) 3–46

© The Author(s) 2021



Article reuse guidelines:

sagepub.com/journals-permissions

DOI: 10.1177/0962280221995963

journals.sagepub.com/home/smm



Abstract

Recent theoretical work in causal inference has explored an important class of variables which, when conditioned on, may further amplify existing unmeasured confounding bias (bias amplification). Despite this theoretical work, existing simulations of bias amplification in clinical settings have suggested bias amplification may not be as important in many practical cases as suggested in the theoretical literature. We resolve this tension by using tools from the semi-parametric regression literature leading to a general characterization in terms of the geometry of OLS estimators which allows us to extend current results to a larger class of DAGs, functional forms, and distributional assumptions. We further use these results to understand the limitations of current simulation approaches and to propose a new framework for performing causal simulation experiments to compare estimators. We then evaluate the challenges and benefits of extending this simulation approach to the context of a real clinical data set with a binary treatment, laying the groundwork for a principled approach to sensitivity analysis for bias amplification in the presence of unmeasured confounding.

Keywords

Causal simulation, bias amplification, sensitivity analysis, causal inference, simulation experiments

I Introduction

Causal identification strategies aim to condition on a sufficient set of observables such that the potential outcomes are conditionally independent of the treatment of interest.^{1–3} Causal variable selection procedures often assume that at least one subset of the observed variables forms such a sufficient set.⁴ The object in causal variable selection then becomes how to separate variables which are necessary for identification of the causal effect from those variables which are extraneous^{4,5} in the interest of reducing estimator variance or covariate dimensionality.^{4,6}

In non-experimental observational studies, we do not have full access to a sufficient set in many realistic settings, and important confounding pathways remain unblocked.^{7,8} This is referred to as unmeasured confounding or endogeneity in the statistics and econometrics literatures, respectively. However, applied researchers currently rely on variable selection techniques such as lasso, step-wise, change-in-estimator selection, and outcome and/or treatment oriented approaches⁹ despite violating their underlying assumptions.

Often, variable selection techniques are used to avoid conditioning on negligible confounding pathways without introducing meaningful bias to the estimator. In this paper, we explore how this intuition can break down under even mild violations of the underlying assumptions. In particular, we build on the work of Pearl¹⁰ and explore how treatment prediction-oriented approaches may inadvertently amplify bias if unmeasured

¹Department of Mathematics and Statistics, McGill University, Montreal, QC, Canada

²Department of Family Medicine, McGill University, Montreal, QC, Canada

³Centre for Clinical Epidemiology, Lady Davis Institute, Montreal, QC, Canada

Corresponding author:

Russell Steele, Department of Mathematics and Statistics, McGill University, Burnside Hall, Room 1005 805, Sherbrooke Street, West Montreal, Quebec H3A 0B9, Canada.

Email: russell.steele@mcgill.ca

confounding pathways remain. We use the lessons from this failure of intuition to better understand the consequences of variable selection in both linear models and partially linear orthogonalized models in the presence of important unmeasured confounding. We then build a principled approach to simulate from such systems of equations which may be used to compare the theoretical properties of different estimators, such as accurately quantifying the change in bias when we modify the causal relationship between a confounder and the treatment, or for the purpose of sensitivity analysis.

First consider data generated from the following directed acyclic graph (DAG) (Figure 1) and set of structural equations

$$Y = \alpha_y + A\beta_a + U\beta_u + \epsilon_1 \quad (1)$$

$$A = \alpha_a + U\gamma_u + \sum_{i=1}^{10} \mathbf{BAV}_i \gamma_{bav_i} + \epsilon_2 \quad (2)$$

$$U = \alpha_u + \sum_{i=1}^{10} \mathbf{BAV}_i \psi_{bav_i} + \epsilon_3 \quad (3)$$

where Y is the outcome, A is the treatment of interest, U is an unmeasured variable, \mathbf{BAV} refers to 10 different potential bias amplifying variables that are measured and affect both A and U but have no direct effect on Y , γ_x is the coefficient for the effect of the variable on A , β_x is the coefficient for the effect of the variable on Y , and ϵ_{number} is an error term for the corresponding equation. This model contains one confounding path that cannot be blocked ($A \leftarrow U \rightarrow Y$) and 10 confounding paths ($A \leftarrow \mathbf{BAV}_1 \rightarrow U \rightarrow Y$; ...; $A \leftarrow \mathbf{BAV}_{10} \rightarrow U \rightarrow Y$) that can be blocked by including the measured \mathbf{BAV}_i variables in the model. However, including any of these \mathbf{BAV}_i might also increase bias (potential bias amplifying variables). Our goal is to find the least biased estimator of the average causal effect (ACE) of treatment (β_a).

By including more of the \mathbf{BAV} variables in the model, intuition suggests the remaining unmeasured confounding bias should decrease because more potential confounders have been included in the conditioning set. However, as demonstrated in the bias amplification literature,^{10–12} conditioning on confounders may still increase bias. For example, suppose further the 10 observable variables account for 90% of the variance in the variable U responsible for unmeasured confounding. The blue violin plot in Figure 2 represents the sampling distribution of the estimator from the true outcome model that includes the treatment and both measured/unmeasured confounding variables included as regressors. As expected, the estimates are approximately normally distributed around the true value $\beta_a = 0.7$. The green violin plot represents the biased estimates from the naive model, the simple regression of the outcome Y on the treatment A , which does not include any of the confounders (measured or unmeasured). The red violin plot represents the linear model adjusted for all 10 measured confounders which account for 90% of the unmeasured confounding. The adjusted model performs much worse than the naive model both in terms of bias (0.73 compared to 0.43, interpretable as standard deviations) and variance (standard deviation of 0.1 compared to 0.02). In fact, in 4990 of 5000 simulations the adjusted estimate was farther from the truth than the naive estimate and nearly 65% of the adjusted estimates had the incorrect effect sign.

The purpose of this paper is to explain why model selection intuition fails us in this case and how we can use a combination of data and simulation approaches to improve model selection. We build upon an emerging theoretical literature exploring a class of variables which can amplify existing unmeasured confounding bias,^{10,13,14,30,31,32,35} called bias amplifiers. Throughout this text, we will refer to (potential) bias amplifiers as those variables which, upon their inclusion in a model, (may) increase the absolute bias in the estimation of particular target parameters relative to a smaller, nested model. In the above example, all 10 \mathbf{BAV} 's increase the bias due to the path $A \leftarrow U \rightarrow Y$ and in fact could also increase the bias of any other variable \mathbf{BAV} if it were not included in the model, since they are all confounders.

The class of bias-amplifying variables is potentially very large and common in practical applications. Bias amplification occurs when the absolute bias of an estimator for a target parameter is larger than the absolute bias of an estimator for a competing, nested model. In this text, we define the natural model comparison to be the naive regression of the outcome (Y) on the treatment (A) and bias amplification will be the additional relative bias that results from adding additional predictors. Throughout, the goal of model selection in this context will be to

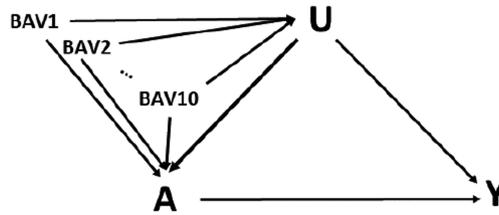


Figure 1. Directed acyclic graph (DAG): Meyers (2011) extended to 10 possible bias amplifying variables BAVs.

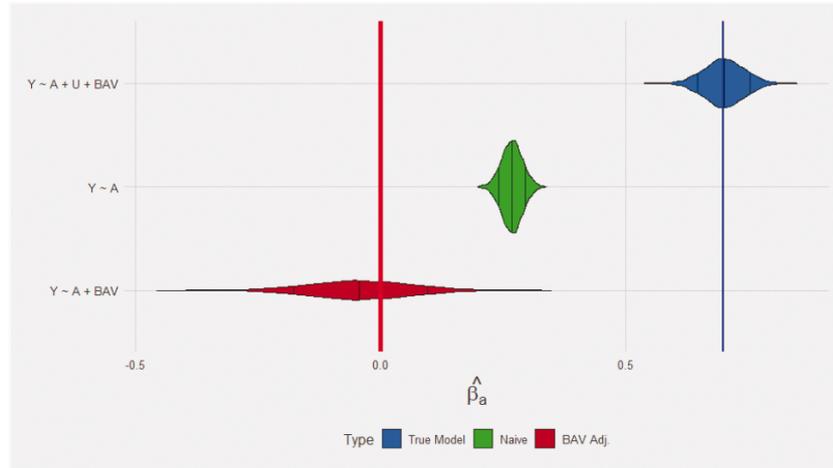


Figure 2. Violin plot from simulations from equations (1) to (3). There were 5000 replication with $n = 5000$. The true effect of interest was $\beta_a = 0.7$, represented by the blue line. The confounding effects were $\beta_u = -0.5$ and $\gamma_u = 0.59$. The vector of coefficients for **BAV** on **U** was $\psi_{bav} = \{-0.55, -0.45, -0.3, 0.30, .25, 0.20, -0.20, 0.20, -0.15, 0.10\}$, which in general were larger than the impact of **BAV** on **A**, $\gamma_{bav} = \{-.1, -.15, -.1, .21, -.2, .3, -.2, -.15, -.2, .075\}$. **U**, **BAV**, **A**, and **Y** are all standard normal variables. The error terms were normally distributed with standard deviations: $\sigma_1 = 0.93$, $\sigma_2 = 0.07$, and $\sigma_3 = 0.32$. All intercepts were set to 0. The code to reproduce the plot can be found in the supplementary materials. The red vertical line represents $\hat{\beta}_a = 0$.

choose the variables and thus the model which minimizes bias in estimating the true average treatment of effect in the presence of unmeasured confounding pathways.

We adopt a matrix notation framework to characterize this problem because (1) we can more easily generalize to a much larger class of directed acyclic graphs and structural equations than previously studied; (2) it offers a unifying geometric explanation for the amplified bias in the context of least squares estimation; and (3) it offers a solid foundation for how to build data-informed model selection procedures. Finally, we develop a procedure for simulating from a more complete parameter space in a way that respects the underlying amplification process. In addition to lending itself better to articulating and answering causal simulation questions, this procedure helps explain why some previous studies have incorrectly concluded that applied investigators need not worry about amplification in practice.¹⁵ We evaluate the challenges of implementing this approach with a real clinical example with binary treatment.

2 Problem formulation

Figure 3 shows a basic directed acyclic graph (DAG) which has both measured and unmeasured confounding. Let **Y** represent the outcome and let **A** be the treatment or variable of interest. Let **U** be an unmeasured confounding variable that we cannot include in a regression model, but which has a functional relationship with both **Y** and **A**. The bias amplifying variable (**BAV**) in this DAG is analogous to **U** in that it is a cause of **Y** and a cause of **A**; however, we are able to measure **BAV** and not **U**. Intuition from currently recommended causal variable selection techniques would tell us to include **BAV** in the regression to reduce bias because it is the root of a confounding

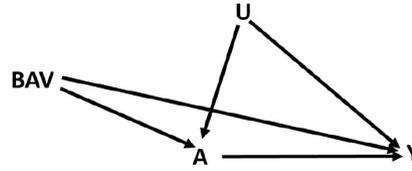


Figure 3. DAG: Two confounding paths, where A is the treatment of interest, Y is the outcome, U is an unmeasured variable and BAV is a measured variable.

path ($A \leftarrow BAV \rightarrow Y$). However, as has been demonstrated,^{10–12,14} blocking this confounding path can actually increase or amplify the bias relative to the naive estimator that depends only on A .

Assume that we now restrict the possible models for the DAG in Figure 3 to only linear associations amongst variables. The data generating model (or causal structural equations) representing Figure 3 under strictly linear association can be written as

$$Y = \alpha_y + A\beta_a + U\beta_u + BAV\beta_{bav} + \epsilon_1 \quad (4)$$

$$A = \alpha_a + U\gamma_u + BAV\gamma_{bav} + \epsilon_2 \quad (5)$$

where α_y and α_a are the intercept terms for Y and A , respectively. As noted above, we use the form β_x throughout this paper to denote true structural coefficients for some variable X on the outcome Y . For example, the true structural coefficient for U on Y is β_u . Analogously, the true regression parameter for some variable X on the treatment A is represented by γ_x . The estimates of these parameters by OLS are denoted by $\hat{\beta}_x$, $\hat{\gamma}_x$ with additional superscripts to clarify which set of estimating equations the estimator is derived from. By assumption ϵ_1 and ϵ_2 are error terms independent of each other we further assume that $E[\epsilon_1|A, U, BAV] = 0$ and $E[\epsilon_2|A, BAV] = 0$, and that the error terms have some finite variance $\sigma_{\epsilon_{1,2}}^2$. In simulation experiments, we typically simulate from normal distributions which are independent from all other variables, but only the above assumptions are necessary for the theoretical results to hold.

We also assume that the target estimand is the ACE. In the linear model case, this is simply the β_a above in equation (4) (See Appendix A.3.I for derivation of the ACE). U is unmeasured and thus we cannot identify the ACE from the observed data if $\beta_u \neq 0$, but we are interested in estimating the quantity with as little bias as possible.

2.1 Matrix notation and probability limits

To tackle the question of model selection, we must derive properties of different $\hat{\beta}_a$ estimators, with the restriction that the estimators be functions of only observed variables, under different assumed conditional regression models. To this aim, we propose expressing OLS estimates using matrix notation and ideas borrowed from the partial regression literature. Further, we propose considering also the probability limits of the estimators to extend our results to more general and realistic cases of bias amplification (see Appendix A.2). For the naive estimator, we use a simple linear regression to estimate a conditional expectation of the form

$$E[Y|A] = \alpha_y^{naive} + A\beta_a^{naive} + v_I \quad (6)$$

where v_I is the error of the regression term. Unbiased estimation of the naive model by OLS, and thus of the true conditional expectation $E[Y|A]$, requires the assumption that $E[v_I|A] = 0$. However, this of course is not true since according to the data generating equation (4), the error will be a function of the confounding terms, resulting in non-zero bias. The naive estimator bias is a special case of the classic omitted variables problem, where we have two omitted variables which are related to both the treatment and the exposure, U and BAV .

Let $\hat{\beta}_a^{naive}$ be the estimate of β_a from the naive model (6). Throughout this paper, we will consider the matrix Z to be a matrix of all the variables that we include in a regression that are not the variable of interest A , in other words control variables in a selection on observables approach. In the naive model, $Z = \mathbf{1}$ where throughout $\mathbf{1}$ will

denote an $n \times 1$ vector of 1s. In matrix notation, applying the Frisch-Waugh-Lovell (FWL) theorem (see Appendix A.1), we can write $\widehat{\beta}_a^{naive}$ as

$$\widehat{\beta}_a^{naive} = \frac{\mathbf{A}^T \mathbf{M}_I \mathbf{Y}}{\mathbf{A}^T \mathbf{M}_I \mathbf{A}} \quad (7)$$

where \mathbf{M}_I is a centering projection matrix, defined and described in detail in Appendix section A.1. In the case of linear relationships between all the variables, following Pearl,¹⁰ this estimator has the following expectation

$$E[\widehat{\beta}_a^{naive}] = \beta_a + (\beta_u \gamma_u \sigma_u^2 + \beta_{bav} \gamma_{bav} \sigma_{bav}^2) \frac{1}{\sigma_a^2} \quad (8)$$

The absolute bias for the ACE then clearly is $|(\beta_u \gamma_u \sigma_u^2 + \beta_{bav} \gamma_{bav} \sigma_{bav}^2) \frac{1}{\sigma_a^2}|$. Now consider the estimates resulting from further conditioning on the observable \mathbf{BAV} variable, i.e. estimating a conditional expectation of the form

$$E[\mathbf{Y}|\mathbf{A}, \mathbf{BAV}] = \alpha_y^{naive} + \mathbf{A} \beta_a^{bav} + \mathbf{BAV} \beta_{bav}^{bav} + v_2 \quad (9)$$

We will denote the resulting estimator $\widehat{\beta}_a^{bav}$ which can be written as follows by again applying the FWL theorem

$$\widehat{\beta}_a^{bav} = \frac{\mathbf{A}^T \mathbf{M}_z \mathbf{Y}}{\mathbf{A}^T \mathbf{M}_z \mathbf{A}} \quad (10)$$

where $\mathbf{Z} = [\mathbf{I}, \mathbf{BAV}]$ and \mathbf{M}_z is the annihilator projection matrix of the matrix \mathbf{Z} (see Appendix A.1 for details and properties). Again following Pearl,¹⁰ the expectation of $\widehat{\beta}_a^{bav}$ is

$$E[\widehat{\beta}_a^{bav}] = \beta_a + (\beta_u \gamma_u \sigma_u^2) \frac{1}{\sigma_a^2 - \gamma_{bav}^2 \sigma_{bav}^2} \quad (11)$$

In Appendix A.4, we explicitly show Pearl's derivation and how it relies on the conditional expectation $E[\mathbf{U}|\mathbf{A}, \mathbf{BAV}]$ being linear in both \mathbf{A} and \mathbf{BAV} . Pearl's derivation is limited in that it is cumbersome and does not generalize well to a broad class of DAGs and functional forms. A simple example where we are unable to use Pearl's method is the case of an interaction term in the exposure structural equation between \mathbf{U} and \mathbf{BAV} . This implies that $E[\mathbf{U}|\mathbf{A}, \mathbf{BAV}]$ is nonlinear in \mathbf{A} and \mathbf{BAV} . This cannot be represented by an unbiased least squares projection of the form $\mathbf{U} = \alpha_u + \mathbf{A} \zeta_a + \mathbf{BAV} \zeta_{bav} + \epsilon_3$ as required by Pearl's derivation method (see Appendix A.4 for details), where ζ_i represents the true regression coefficient for variable i . If we impose further strict distributional assumptions over all the variables, we may still be able to directly solve the conditional expectation and find an expression for bias in terms of the underlying parameters. However, in many applied cases, these distributional assumptions will not be justified, particularly assuming a distribution for the unmeasured confounding which will always be untestable.

In contrast, if we consider the probability limits, we do not need to assume that $E[\mathbf{U}|\mathbf{A}, \mathbf{BAV}]$ is linear, nor do we have to make any additional distributional assumptions to find meaningful limiting expressions for our estimators in a broad class of clinically relevant circumstances. In addition to giving rise to a meaningful interpretation, the closed form asymptotics we derive allow us to more easily harness domain knowledge about the underlying causal process for the purpose of model selection.

Since we are still interested in the finite sample expectation of the estimators and the bias directly, we report the expectations when appropriate and feasible. The probability limit facilitates insight under weaker assumptions than those necessary to derive exact forms of the expectations. Additionally, in some cases, like the linear model of Pearl,¹⁰ the probability limits for $\widehat{\beta}_a^{naive}$ and $\widehat{\beta}_a^{bav}$ are precisely equal to their expectations (see Appendix A.7).

3 Treatment variance explained and amplifying effects for strictly linear models

In this section we examine the probability limits of the estimators to better understand the mechanics and root causes of bias amplification. First, notice the denominator of equation (11) ($\sigma_a^2 - \gamma_{bav}^2 \sigma_{bav}^2$) gets smaller as the

causal edge $BAV \rightarrow A$ increases in strength. This is because when we specify the functional form of a system of random variables and conditional independence assumptions, we are also determining a formula for its variance. Under the structural equation we specified for the exposure (equation (5), with the corresponding DAG shown later in Figure 7(a)), $\sigma_a^2 - \gamma_{bav}^2 \sigma_{bav}^2$ is precisely equal to the remaining residual treatment variance in A after having adjusted for BAV in a regression model. When we assume that all the variables are standardized, this term becomes $1 - \gamma_{bav}^2$ as presented in Pearl¹⁰ (under the assumption of standard normal, here we drop the distributional requirement), because the variance of standard normal variables is equal to 1 ($\sigma_a^2, \sigma_{bav}^2 = 1$).

In order to visualize this phenomenon, we use ideas from partial regression plots.¹⁶ By the FWL theorem, we can always pre-multiply an estimating equation by the residual-making variables of a set of regressors and get the same estimates (see Appendix A.1 for further details). For example, the following two regression equations produce the same numerical ordinary least squares estimates of $\hat{\beta}_a^{naive}$

$$Y = \alpha_y + A\beta_a + v_I \quad (12)$$

$$M_I Y = M_I A \beta_a + v_I \quad (13)$$

Equation (13) is the model for a simple linear regression of a modified outcome, $M_I Y$, on a modified treatment, $M_I A$ (see Appendix A.1). There is no intercept term as the mean of the modified treatment must be equal to zero.

$$Y = \alpha_y + A\beta_a + BAV\beta_{bav} + v_2 \quad (14)$$

$$M_z Y = M_z A \beta_a + v_2 \quad (15)$$

Similarly equations (14) and (15) above produce equivalent ordinary least squares estimates of $\hat{\beta}_a$, where $Z = [I \ BAV]$ is a column of 1s and the BAV variable. Equation (15) is a single variable regression on a transformed set of variables. The modified Y is produced by taking the residuals from regressing Y on a column of 1s and BAV . In other words, the dependent variable is the residual vector resulting from regressing Y on an intercept column and BAV . The independent variable is the residual vector resulting from regressing A on a column of 1s and BAV . Another way to think of the residuals is in the context of orthogonalization techniques, where $A^T M_z$ for example is the part of A orthogonal to linear combinations of the control variables A . We then regress the orthogonalized outcome on the orthogonalized treatment, which is implicitly the mechanics of what happens whenever we use least squares estimation. See section 4.2 for more general orthogonalization techniques.

Since we have reduced the multi-variable regression problem to a simple linear regression for two linearly transformed variables, we can easily visualize the amplification process via a partial regression plot. In Figure 4 the left subplot visualizes the naive regression equation (13), whereas the blue subplot on the right visualizes the regression equation (15) that includes BAV . The data were simulated from a special case of equations (4) and (5), with $n = 1000$. Details can be found in the appendix (Appendix section A.13).

The unbiased ACE is the slope of the black line ($\beta_a = 0.2$) in these plots. The slope of the blue line (equal to the OLS estimator from the amplifying model) is clearly farther away from the true slope (in black) compared to the slope of the red line from the naive model, and thus the conditional estimator is more biased.

Note first that including BAV in the model reduces the variance in the adjusted treatment, which can be seen by comparing the relative sparsity of points along the x-axis in red compared to the relative density of points along the x-axis in blue. However, if we inspect the spread of points vertically along the y-axis, we can see that the red and blue samples are similarly dispersed in this dimension because conditional on the treatment, linear combinations of BAV explain very little of the variance in the outcome. Most importantly, including only BAV does not change the variance in Y due to U , the unmeasured confounder.

When we add the BAV to the regression model, the bias is 0.14 larger in absolute terms (or approximately 65% greater in relative terms) than the naive estimate, even though it blocks a confounding path between the treatment A and the outcome Y . More simply, trying to block a confounding path with weak response association can amplify bias in causal effect estimation because it increases the proportion of treatment association due to unmeasured confounding on unblocked paths.

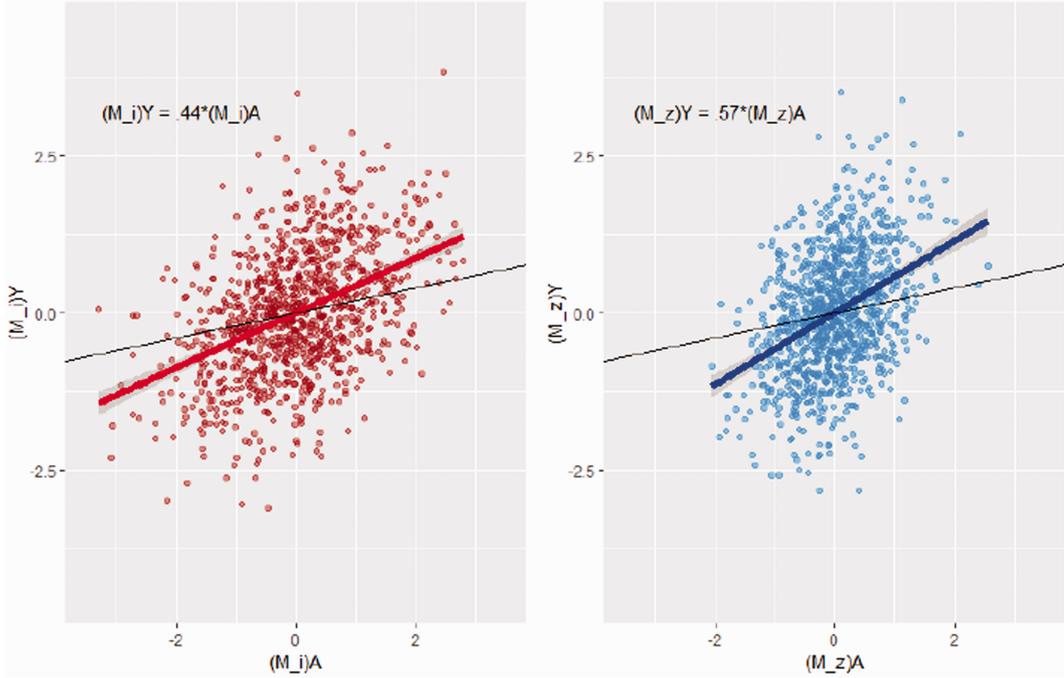


Figure 4. In both panels, the unbiased ACE ($\beta_a = 0.2$) is shown by the dotted black line. In the left panel, the red dots represent the centered treatment (M_iA) plotted against the centered outcome (M_iY) and the estimated slope $\hat{\beta}_a^{naive}$ is shown with the bolded red line. This represents the equivalent regressions in equations (12) and (13). In the right panel, the blue dots represent the modified treatment (M_zA) plotted against the modified outcome (M_zY). The solid blue line represents the treatment estimate from the equivalent regressions (14) and (15).

The magnitude of bias amplification can be potentially very large. From equations (8) and (9), the absolute bias of the **BAV** estimator will be larger than the absolute bias of the naive estimator whenever

$$\frac{|\text{Bias}(\hat{\beta}_a^{bav})|}{|\text{Bias}(\hat{\beta}_a^{naive})|} = \left(\frac{\sigma_a^2}{\sigma_a^2 - \gamma_{bav}^2 \sigma_{bav}^2} \right) \left(\frac{|\beta_u \gamma_u \sigma_u^2|}{|\beta_u \gamma_u \sigma_u^2 + \beta_{bav} \gamma_{bav} \sigma_{bav}^2|} \right) > 1 \quad (16)$$

We can rewrite the first term on the right side of the equal sign as $(1 - \mathcal{R}_{A|BAV}^2)^{-1}$, which is the inverse of 1 minus percentage variance explained in A adjusting for **BAV**

$$\frac{|\text{Bias}(\hat{\beta}_a^{bav})|}{|\text{Bias}(\hat{\beta}_a^{naive})|} = \left(\frac{1}{1 - \mathcal{R}_{A|BAV}^2} \right) \left(\frac{|\beta_u \gamma_u \sigma_u^2|}{|\beta_u \gamma_u \sigma_u^2 + \beta_{bav} \gamma_{bav} \sigma_{bav}^2|} \right) > 1 \quad (17)$$

Notice that the first term must always be greater than or equal to 1. The more strongly the control variables predict the treatment, the larger the magnitude of the term, increasing monotonically¹¹ and hyperbolically in the treatment variance explained by **BAV**.

We now examine the second term after the equal sign. If $\beta_u = 0$ or $\gamma_u = 0$, then there is no unmeasured confounding and $\hat{\beta}_{bav}$ will be unbiased. If $\beta_u \neq 0$ and $\gamma_u \neq 0$ then there is unmeasured confounding, so by manipulating equation (17), we can write the second term as

$$\frac{|\text{Bias}(\hat{\beta}_a^{bav})|}{|\text{Bias}(\hat{\beta}_a^{naive})|} = \left(\frac{1}{1 - \mathcal{R}_{A|BAV}^2} \right) \left(\frac{1}{|1 + \eta \text{sgn}(\beta_u) \text{sgn}(\gamma_u) \text{sgn}(\beta_{bav}) \text{sgn}(\gamma_{bav})|} \right) > 1 \quad (18)$$

where $\text{sgn}(\cdot)$ is equal to $+1$, -1 , and 0 if the argument is positive, negative, and equal to 0 , respectively, and η is defined as the ratio of the absolute strength of the confounding path through **BAV** ($|\beta_{bay}\gamma_{bay}\sigma_{bay}^2|$) to the absolute strength of the confounding path through **U** ($|\beta_u\gamma_u\sigma_u^2|$).

$$\eta = \frac{|\beta_{bay}\gamma_{bay}\sigma_{bay}^2|}{|\beta_u\gamma_u\sigma_u^2|}$$

The second term depends on the signs of the structural coefficients for the two confounding paths. This is a dimensionless quantity that does not depend on the scale of the data, which can make it useful for sensitivity analysis. If $\gamma_{bay} = 0$, then **BAV** does not affect **A**, so it has no effect on the bias because $\eta = 0$ and $\mathcal{R}_{A|BAV}^2 = 0$. In the special case that $\beta_{bay} = 0$, **BAV** is a true instrumental variable under the DAG in Figure 3. If there are no interactions as in the structural equations (4) and (5), then including an instrumental variable will always increase bias in the presence of unmeasured confounding.^{3,10,14} We can see this clearly in equation (16) since by definition an instrumental variable does not explain variance in the outcome except through the treatment and must correlate with the treatment. A strong instrument in this case must be a strong amplifier since it will strongly predict variance in the treatment making the first term in equation (16) large and the second term exactly 1.

In the case that none of the structural coefficients are equal to 0, note that $\text{sgn}(\beta_u)\text{sgn}(\gamma_u)\text{sgn}(\beta_{BAV})\text{sgn}(\gamma_{BAV})$ will be equal to $+1$ if there is an even number of positive signs for the structural coefficients and it will be equal to -1 if there is an odd number of signs. We have defined η in terms of absolute values, so it must be non-negative. If $0 \leq \eta < 1$, then the confounding path through **U** is stronger than the path through **BAV**. If $\eta > 1$, then the confounding path through **BAV** is stronger than the path through **U**. When $\eta = 1$, the strengths of the two confounding paths are equal.

We can use equation (18) to characterize all possible confounding structures that lead to bias amplification. First assume that there are an even number of positively signed structural coefficients. This implies that the second term of equation (18) is equal to $(1 + \eta)^{-1}$ and bias amplification will occur when $(1 - \mathcal{R}_{A|BAV}^2)^{-1} > (1 + \eta)$ or equivalently when $\mathcal{R}_{A|BAV}^2(1 - \mathcal{R}_{A|BAV}^2)^{-1} > \eta$. So if there is an even number of positively signed structural coefficients, the larger amount of treatment variance explained by **BAV**, the greater the range of possible η values that will lead to bias amplification.

Next assume that there is an odd number of positively signed structural coefficients, so that the second term will be equal to $|1 - \eta|^{-1}$. In this case, if $0 < \eta \leq 1$ (i.e. if the confounding path through **U** is stronger than the confounding path through **BAV**), then there will always be bias amplification. If $\eta > 1$, then there will be bias amplification if $(1 - \mathcal{R}_{A|BAV}^2)^{-1} > (\eta - 1)$, similar to the case when $\eta > 1$.

In Appendix A.5 we further re-express η in terms of only correlations (or partial correlations if desired) and the coefficients of determination and free of model parameters. This is a dimensionless quantity that does not depend on the scale of the data, which can make it useful for sensitivity analysis. Additionally, the different ways of expressing the quantity can reveal portions of the unknown quantity which are estimable, i.e. functions of only observed data. Depending on the application parameters, it may be easier to express domain knowledge through correlations or percentage variance explained. This allows one to move more easily from one to the other, which should be very useful to applied researchers. We expand on these ideas in sections 5 and develop principled ways to simulate from such systems of equation. In Appendix A.5, we discuss exploiting the fact that correlations and variances are sufficient for determining the whole system of equations.

In summary, variable selection approaches which aggressively target confounding paths with strong associations with treatment and weak associations with outcome are at greater risk of bias amplification as they are much more sensitive to the assumption that a full sufficient set is measurable. Adding controlling variables in proportion to their ability to predict the treatment in linear models is only guaranteed to be bias reducing if the resulting selected variables satisfy ignorability assumptions. When this is not the case, we have shown that unmeasured confounding bias may be severely amplified since the treatment variance term increases both monotonically and hyperbolically. It is often not possible or extremely unlikely to select a sufficient set in many non-experimental settings, i.e. most investigators are not willing to assume their observational study is equivalent to a randomized trial. Thus the potential for bias amplification must be considered carefully.

4 Generalizing to a larger class of causal models

In sections (1) to (3), we showed examples of bias amplification under two different DAGs, each with linear structural equations in both the coefficients and variables and probability limit results under the DAG in Figure 3 and linear structural equations. In section 4.1 we extend bias amplification to the class of structural functions that are additive in the outcome and whose target causal effect is represented by β_a . In other words, we consider structural equations of the form

$$Y = \alpha_y + A\beta_a + f_1(U) + f_2(\mathbf{BAV}) + \epsilon_I \quad (19)$$

$$A = \alpha_a + g(U, \mathbf{BAV}) + \epsilon_2 \quad (20)$$

We require minimal assumptions of treatment structural equation (20), in fact it is completely general except that it must be decomposable into a function of U and \mathbf{BAV} (i.e. the $E[A|U, \mathbf{BAV}] = \alpha_a + g(U, \mathbf{BAV})$) and an error term orthogonal to this function ($\epsilon_2 : E[\epsilon_2|U, \mathbf{BAV}] = 0$). For the outcome equation, we require more structure and assumptions on top of an orthogonal decomposition into the conditional expectation and error term ($\epsilon_I : E[\epsilon_I|U, A, \mathbf{BAV}] = 0$). Specifically we require that the outcome equation (19) is additive in functions of A , U , and \mathbf{BAV} (ruling out interaction effects for example) and that the causal effect of interest is a single parameter β_a ($\frac{\partial E[Y|A, U, \mathbf{BAV}]}{\partial A} = \beta_a$).

In section 4.2, we again consider the same structural setting, but extend the results and intuitions developed in sections 3 and 4.1 to more general orthogonalization techniques beyond least squares, namely a specific case of Neyman-Orthogonalization in partially linear models as used in double debiased approaches like in Chernozhukov et al.¹⁷ We show that in both settings, we can decompose bias amplification into a component due to covariance between the treatment and the outcome through the confounding paths and a second term which is a ratio of marginal treatment variance to residual treatment variance after adjusting for other covariates. We can always estimate the residual treatment variance and should report these estimates in observational studies, particularly if variable selection was in part determined by strength of treatment prediction. As the remaining treatment variance decreases, the potential for bias amplification increases and any inferences about the causal effect depend more strongly on the assumption that unmeasured confounding is negligible. Further, similar to the earlier setting, we can directly estimate correlations, coefficients of determination, and conditional expectations composed only of observable variables. While the entire confounding pathways cannot be identified, some of the components can be. These estimable quantities, in turn, place restrictions on what the total confounding pathways can be. In section 5 we show how to use these principles and perform simulations to test the performance of competing estimators. In section 6 we apply these ideas to a real clinical data set.

4.1 Generalized bias amplification in least squares

If $f_2(\mathbf{BAV})$ is known or can be well approximated (say by an appropriate basis expansion that grows in dimension as $n \rightarrow \infty$) up to a constant of proportionality and intercept, then in the limit there will not be any bias due to misspecification of $f_2(\mathbf{BAV})$ (for a look at the case when $f_2(\mathbf{BAV})$ is misspecified see Appendix section A.9.1). In such a case, we again can compare two feasible OLS estimators, the first being the Naive estimator ($\hat{\beta}_a^{naive}$) as before and the second the OLS estimator resulting from including $f_2(\mathbf{BAV})$ or its approximation ($\hat{\beta}_a^{f_2(bav)}$). Under any DAG, it can be shown that under structural equations (19) and (20), there will be bias amplification when

$$\frac{|\text{Bias}(\hat{\beta}_a^{bav})|}{|\text{Bias}(\hat{\beta}_a^{naive})|} = \left(\frac{1}{1 - \mathcal{R}_{A|f_2(\mathbf{BAV})}^2} \right) \times \left(\frac{|COR(A, f_1(U)) - COR(A, f_2(\mathbf{BAV}))COR(f_1(U), f_2(\mathbf{BAV}))|}{|COR(A, f_1(U)) + COR(A, f_2(\mathbf{BAV}))\kappa|} \right) > 1 \quad (21)$$

where $\kappa = \sigma_{f_2(\mathbf{BAV})}\sigma_{f_1(U)}^{-1}$. In the above expression, we assume the case that there is in fact bias in the naive model, i.e. $\text{Bias}(\hat{\beta}_a^{naive}) \neq 0$.

Bias amplification is decomposed into two terms, the ratio of remaining variance in the treatment once the control variables are projected out, and the remaining covariance between the outcome and treatment through the

confounding pathways. Just like in the linear case, the first term $(1 - \mathcal{R}_{A|f_2(\mathbf{BAV})}^2)^{-1}$ is the ratio of remaining treatment variance once we have projected out linear combinations of the control variables, which is just a constant in the naive case and $f_2(\mathbf{BAV})$ in the adjusted case. This first term must be greater than 1 and is a hyperbolic function increasing as $f_2(\mathbf{BAV})$ linearly explains more variance in the treatment.

Now consider the second bias term $\left(\frac{|COR(A, f_1(U)) - COR(A, f_2(\mathbf{BAV}))COR(f_1(U), f_2(\mathbf{BAV}))|}{|COR(A, f_1(U)) + COR(A, f_2(\mathbf{BAV}))\kappa|}\right)$ which is the ratio of the covariance between the outcome and treatment remaining through the uncontrolled confounding pathways. The numerator of this expression $(|COR(A, f_1(U)) - COR(A, f_2(\mathbf{BAV}))COR(f_1(U), f_2(\mathbf{BAV}))|)$ is the magnitude of the correlation between the treatment and $f_1(U)$ once we have projected out linear combinations of \mathbf{BAV} . When $f_1(U)$ and $f_2(\mathbf{BAV})$ are uncorrelated, or independent as they would be under the DAG in Figure 3 and equation (17) in section 3, then the numerator reduces to simply $|COR(A, f_1(U))|$ since \mathbf{BAV} does not linearly explain any of the shared variance in the treatment A and $f_1(U)$.

In general, when there is correlation between $f_1(U)$ and $f_2(\mathbf{BAV})$, adjusting for $f_2(\mathbf{BAV})$ reduces the confounding bias of $f_2(U)$ when the signs of the correlations $COR(A, f_2(\mathbf{BAV}))$ and $COR(f_1(U), f_2(\mathbf{BAV}))$ are the same. Otherwise, this adjustment may further increase the confounding bias. As a simple example, if all of the correlations are positive (i.e. $COR(A, f_1(U)), COR(A, f_2(\mathbf{BAV})), COR(f_1(U), f_2(\mathbf{BAV})) > 0$) then the second term in the amplification expression will be less than or equal to 1 $\left(\frac{|COR(A, f_1(U)) - COR(A, f_2(\mathbf{BAV}))COR(f_1(U), f_2(\mathbf{BAV}))|}{|COR(A, f_1(U)) + COR(A, f_2(\mathbf{BAV}))\kappa|}\right) \leq 1$). In this case, we reduce bias in two ways. First we remove the bias due to $COR(A, \mathbf{BAV})$ directly and then we reduce the remaining covariance between A and $f_1(U)$ through the part of this correlation explained by \mathbf{BAV} .

There will be amplification, on the other hand, if the first term (which is always greater than one) is bigger than the reciprocal of the second term (which in this case is less than one). If $COR(A, f_1(U))$ and $COR(A, f_2(\mathbf{BAV}))$ are opposite signed, then when we adjust for \mathbf{BAV} and remove bias due to the confounding path $A \leftarrow \mathbf{BAV} \rightarrow Y$ we may increase bias since $COR(A, f_1(U))$ and $COR(A, f_2(\mathbf{BAV}))$ would no longer be partially cancelling each other out as they do in the naive model. Similarly, projecting out linear combinations of \mathbf{BAV} may increase or decrease the remaining covariance between the treatment and $f_1(U)$ as seen in the numerator of the second term. Overall, this second term may be greater than or less than 1, meaning it may contribute to increasing or decreasing the bias relative to the naive estimator. By estimating the correlations and variances which are functions of observables and eliciting domain knowledge about the remaining unobserved quantities in the relative bias expression, we can make informed predictions about the plausibility of bias amplification.

Even in the most general case, where there are interactions and non-separability of U and \mathbf{BAV} in the data generating model equation, the second term of the bias term will always be proportional to $1 - \mathcal{R}_{A|Z}^2$, where $Z = [U, h(\mathbf{BAV})]$ and $h(\mathbf{BAV})$ are whatever function of \mathbf{BAV} that we include in a regression. As can be seen in the appendix (section A.9.1), this is even true when we have misspecified $f_2(\mathbf{BAV})$. However, there will be an additional bias term due to misspecification, but this will also be amplified with respect to the denominator term. The overall lesson is that when using OLS for causal effect estimation, unmeasured confounding bias (and remaining misspecification bias) will be amplified hyperbolically with respect to how well the controlling variables linearly explain variance in the treatment, i.e. the magnitude of $\mathcal{R}_{A|Z}$, where Z is all of the variables that are not A that we include in our OLS regression. Further, it is important to note that the amplifying factor $((1 - \mathcal{R}_{A|f_2(\mathbf{BAV})}^2)^{-1})$ is always a function of only observables and can always be estimated by running the regression A on $f_2(\mathbf{BAV})$. We can use these estimates for sensitivity analysis or in a model selection procedure. In a case like the simulation in Figure 2, our potential controls explain the large majority of the treatment variance. When this occurs, we should require more confidence that there is no unmeasured confounding after adjusting for our controls before we trust these estimates.

When there are causal pathways between U and \mathbf{BAV} , we showed that this may help make the first term smaller. However, this will open up an additional pathway for $f_2(\mathbf{BAV})$ to explain variance in the treatment (for example $\mathbf{BAV} \rightarrow U \rightarrow A$ as in the DAG in Figure 1, but in general the path could be indirect or the direction of causation may be reversed) and thus make the second hyperbolic term larger. This was the source of the dramatic bias amplification in Figure 2 simulated under linear structural equations from the DAG in Figure 1. In that specific case, \mathbf{BAV} was a very good proxy for the unmeasured confounding U and this has two effects. First, it means that including \mathbf{BAV} eliminates the large majority of the unmeasured confounding through its path to the outcome. Second, this means that \mathbf{BAV} linearly explains nearly as much variation in the treatment as including both \mathbf{BAV} and U , which is nearly all of the variation in the treatment in that simulation case. So the little unmeasured confounding remaining was amplified significantly, more than 25 times $\left(\left((1 - \mathcal{R}_{A|\mathbf{BAV}}^2)\sigma_a^2\right)^{-1} =$

$\left(\sigma_a^2 - \sum_{i=1}^{10} (\gamma_{bavi} + \gamma_{ui}\psi_{bavi})^2 \sigma_{bavi}^2\right)^{-1} = (1 - 0.96015)^{-1} = 25.09725$) under the DAG in Figure 1 and linear

structural equations (1) to (3)). This is because the potential amplifier linearly explained more than 96% of the variance in the treatment.

4.2 Bias amplification in more general semi-parametric models

In the beginning of this paper, we considered bias amplification in the context of linear expectations and linear models. In section 4.1 we extended this to when the structural models are non-linear when the estimator is ordinary least squares. Here we show proof of concept for bias amplification in more general semi-parametric orthogonalization methods. We plan to further develop this work in the future. We again consider the underlying structural equations to be represented by equations (19) and (20).

In the case that $f_2(\mathbf{BAV})$ is unknown, one might turn to more general semi-parametric estimators for the causal effect estimation. A reasonable semi-parametric estimator for these structural equations is Robinson's Double Residual Regression (DRR hereforth)¹⁸ when the outcome is hypothesized to be linear in the treatment (or well approximated by a partial linear function). In section 3 we showed via the FWL theorem that we can think of multiple regression as simple regression of the outcome and treatment once linear combinations of the control variables (which may be non-linear functions themselves as shown in section 4.1) have been projected out. In other words, we orthogonalized the outcome and treatment with respect to the subspace spanned by the linear combinations of all the control variables we included in our regression model (\mathbf{Z}) and then perform simple linear regression on the orthogonalized variables. The coefficient of interest in DRR, just like in multiple linear regression, is the result of a least squares estimate of an orthogonalized outcome and treatment, but we orthogonalize with respect to a more general subspace using the expectation operator.

Considering some arbitrary random variable Y and collection of controls X_1, \dots, X_n (each with finite variance), the conditional expectation $E[Y|X_1, \dots, X_n]$ is the projection of Y onto a closed subspace of \mathcal{L}^2 consisting of all Borel functions $\phi(X_1, \dots, X_n) : \mathcal{R}^n \rightarrow \mathcal{R}$, see for example Brockwell and Davis.¹⁹ This is the heart of the familiar result that the conditional expectation is the function of X_1, \dots, X_n which minimizes the mean-squared error ($\text{argmin}_{f: \mathcal{R}^n \rightarrow \mathcal{R}, \text{Borel}} E[(Y - f(X_1, X_2, \dots, X_n))^2] = E[Y|X_1, X_2, \dots, X_n]$). The set of linear functions is a strict subset of all Borel functions, making conditional expectation a more general projection than the projections implicit in least squares estimation. In other words, if it turns out that the conditional expectations for both Y and A are in fact linear in the same variables, the DRR estimator will be equivalent to OLS in the probability limit.

The first step of DRR is to estimate the conditional expectations, i.e. the projections, $E[Y|\mathbf{BAV}]$ and $E[A|\mathbf{BAV}]$. We can use our favorite non-parametric estimator(s) as long as it is consistent as $n \rightarrow \infty$. It is common to use a conditional density estimator, but one may choose to use modern regression tree models or machine learning techniques like neural networks as in the context of Chernozhukov et al..¹⁷ Now we construct our orthogonalized outcome and treatment, $\tilde{Y} = Y - \hat{E}[Y|\mathbf{BAV}]$ and $\tilde{A} = A - \hat{E}[A|\mathbf{BAV}]$. Finally, we perform least squares of the modified outcome (\tilde{Y}) on the modified treatment (\tilde{A}) to get our semi-parametric estimator denoted $\hat{\beta}^{\text{semi}}$.

Bias amplification will occur in the probability limit when

$$\frac{|\text{Bias}(\hat{\beta}^{\text{semi}})|}{|\text{Bias}(\hat{\beta}_a^{\text{naive}})|} = \left(\frac{1}{1 - \frac{\text{VAR}(E[A|\mathbf{BAV}])}{\sigma_a^2}} \right) \times \left(\frac{|E[\text{COV}(A, f_1(\mathbf{U})|\mathbf{BAV})]|}{|\text{COV}(A, f_1(\mathbf{U})) + \text{COV}(A, f_2(\mathbf{BAV}))|} \right) > 1$$

Once again, we can decompose this relative bias expression into two components. The first term is the ratio of remaining treatment variance once the controlling variables have been projected out of the treatment and the second term $\left(\frac{|E[\text{COV}(A, f_1(\mathbf{U})|\mathbf{BAV})]|}{|\text{COV}(A, f_1(\mathbf{U})) + \text{COV}(A, f_2(\mathbf{BAV}))|} \right)$ is the ratio of confounding covariances remaining through the outcome paths $\mathbf{U} \rightarrow \mathbf{A} \rightarrow \mathbf{Y}$ and $\mathbf{BAV} \rightarrow \mathbf{A} \rightarrow \mathbf{Y}$. The first term again always contributes to amplification, since $(1 - \text{VAR}(E[A|\mathbf{BAV}])/\sigma_a^2)^{-1} \geq 1$. This is the direct analogue of the results developed in the previous sections. In this case, $\mathcal{R}_{A|\mathbf{BAV}}^2$ is the variation of the treatment explained by linear combinations of \mathbf{BAV} (or $h(\mathbf{BAV})$ if we have included a different function of \mathbf{BAV} in our regression), and $\text{VAR}(E[A|\mathbf{BAV}])\sigma_a^{-2}$ is the total variation explained by fluctuations in \mathbf{BAV} .²⁰ As before, the first term increases hyperbolically as \mathbf{BAV} explains more treatment variance.

When we project out \mathbf{BAV} , the term $\text{COV}(A, f_2(\mathbf{BAV}))$ drops from the second component and $\text{COV}(A, f_1(\mathbf{U}))$ becomes $E[\text{COV}(A, f_1(\mathbf{U})|\mathbf{BAV})]$. Depending on the signs of the bias and the direction of the correlation $\text{COV}(E[A|\mathbf{BAV}], E[f_1(\mathbf{U})|\mathbf{BAV}])$, the second term may be greater or less than 1. This is similar to the first term

in the expression in equation (21) since the biases from the two confounding paths may partially cancel each other out and thus by projecting them both out, one may increase bias.

In other words, the underlying mechanics of bias amplification extends to more general orthogonalization methods, both explicit as in the case of DRR and implicit like that in OLS. In particular, the partially linear regression set-up is the basis for more complicated causal machine learning techniques such as Double Machine Learning residuals on residuals regression,¹⁷ which has become a popular approach for dealing with regularization bias in high dimensional settings. Even in these more complicated settings where we orthogonalize in more sophisticated ways, variables which predict large amounts of variance in the treatment may significantly increase bias if confounding pathways remain. Once again, the hyperbolic remaining treatment variance term is entirely a function of observed variables and thus can be estimated (non-parametrically if desired) to see if one is at risk of bias amplification with respect to priors or sensitivity analysis over expected levels of unmeasured confounding. Similarly, we can use the mean-squared error of our resulting model to evaluate the extent to which unmeasured confounding is possible through the outcome path since no interaction implies that remaining variation must either be independent error or unmeasured confounding. A full treatment of this framework is beyond the scope of this article and will be the subject of future work.

5 Causal simulation experiments: the case of bias amplification

In our experience, simulating bias amplification is challenging in a number of subtle, but important ways. Our context of interest is assessing the potential for bias amplification in an analysis of an observational study in which we have measured several independent variables and the outcome but there might be an unmeasured confounder. We are interested in evaluating the feasible estimators we have developed in the previous sections, $\hat{\beta}_a^{naive}$ and $\hat{\beta}_a^{bavi}$ for example, with respect to possible data sets generated by a class of DAGs and structural equations. In this section, we show that if we constrain certain aspects of the simulated data (in particular, the marginal variances of observed quantities), we are better able to articulate and answer causal questions about the effect of bias amplification on proposed estimators. While we discuss the example of bias amplification simulations specifically, this section has implications for simulating data to test causal estimators more broadly. Now, consider the challenge of determining the effect of increasing unmeasured confounding on bias amplification in Figure 5(a). We might, for example, be interested in how large an unmeasured confounder must be, with fixed amplifying variables, to cross some threshold of bias in the adjusted model as part of a sensitivity analysis. To answer such a question, we must define clearly what is meant by the strength of an unmeasured confounder. In Figure 5(a), there are two edges which determine the overall bias due to the unmeasured confounding path through U : the edge from U to A and the edge from U to Y . The bias due to the unmeasured confounding path through U in the naive model is simply the product of the weight of these two edges, scaled by the variance of the treatment as shown in equation (7). The extent to which bias can become amplified, however, is not symmetric with respect to the weight of the edges $U \rightarrow A$ and $U \rightarrow Y$, since amplification is the result of variance explained in the treatment as discussed in section 3. There is more potential for amplification of a strong unmeasured confounder (in the sense the product of the confounding edges is large) when the strength is due to U being a strong cause of Y compared to a strong cause of A . This is because when U is a strong cause of A , the BAV can only explain a small amount of the variance of A , limiting the possible amount of bias amplification. Thus, to answer a causal question about the effect of increased unmeasured confounding on bias amplification, we should only vary one of the confounding edges, holding all other edges fixed.

As an example, suppose we are interested in the change in bias amplification when we increase the strength of the edge from U to A , holding all else constant. This notion of intervening on a single edge of our DAG while holding the others fixed should be familiar to causal inference practitioners since it is the principle behind counterfactual analysis more broadly. Here we want to ensure that our results from varying a single edge are not confounded by variations in other edges as the result of unintended consequences or induced associations.

Because the goal is to increase the strength of a single edge, holding all else constant, we must specify a metric by which we measure the strength of the edge. In a fully linear system, we might consider the strength of the edge as the regression coefficient itself, γ_u , or the proportion of variance explained by U , $\frac{\gamma_u \sigma_u}{\sigma_a^2}$ and the sign of γ_u . It is tempting to see the two measures as equivalent with different scalings, but this is only true in the context of simulating a single equation. In the context of a system of linear equations, especially with the potential for bias amplification, we argue the relevant quantity is the proportion of variance explained by each child node of the parent variable. This can be seen most easily by examining the bias formula in equation (11), where the denominator is the remaining variation in A unexplained by the potential bias amplifying variables.

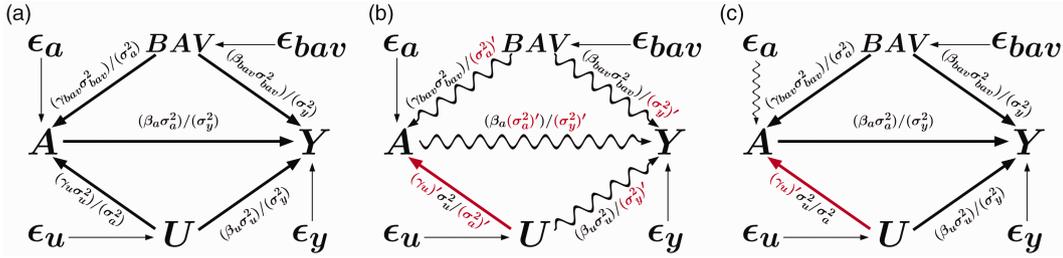


Figure 5. Causal diagrams where (a) represents the underlying causal structure. In addition to the usual causal pathways, we explicitly show the pathway of the independent error terms, ϵ_a , ϵ_u , ϵ_{bav} , ϵ_y . (b) The simulation experiment where we change strength of the edge $U \rightarrow A$ (shown in bold red) and do not renormalize the variance of the treatment (**A**). All edges which have been inadvertently modified are symbolized as squiggly arrows. All parameters which have been modified (inadvertently or intentionally) are shown in red along the edges. (c) Intervening on $U \rightarrow A$, where the variance of the treatment is fixed. To do so, we modify the variance of the noise term ϵ_a , visualized by the squiggly arrow. Notice no other edges are inadvertently modified.

Consider the implications of treating the coefficients themselves as the relevant measure of edge strength in a simulation trying to determine the effect of increasing the causal association along the path from U to A . If we want to increase γ_u to $\gamma_u' > \gamma_u$ without changing any other parameters, we must also increase the total variance in the treatment, A , since $\sigma_a^2 = \gamma_u^2 \sigma_u^2 + \gamma_{bav}^2 \sigma_{bav}^2 + \sigma_{\epsilon_a}^2$. A treatment with a larger variance is in some sense a different intervention, and thus this simulation is not compatible with the class of experiments which generated the original data with parameter γ_u . Further, from the previous sections we know this implies the total amount of variance explained from the bias amplifier BAV is reduced, since $\frac{\gamma_{bav}^2 \sigma_{bav}^2}{(\sigma_a^2)}$ has been reduced. Although we have not changed the parameter γ_{bav} we have decreased the extent to which BAV amplifies the bias as seen by examining equation (11). The increased variance in A in turn modifies the total variance of Y . Therefore, the relative proportion of variance of Y that is explained by BAV is modified by changing the causal effect of $U \rightarrow A$, as are the measured proportion of variance of $BAV \rightarrow A$, $A \rightarrow Y$, $U \rightarrow Y$, and $BAV \rightarrow Y$ and their associated covariance terms.

We can see in Figure 5(b) that by modifying a single coefficient and leaving all other coefficients unchanged we have inadvertently modified the relative proportion of variance explained by the four other edges ($BAV \rightarrow A$, $BAV \rightarrow Y$, $U \rightarrow Y$, and $A \rightarrow Y$) represented by the wavy arrows. Data generated by the second set of structural equations are not compatible with the constraints of the experiment which generated the first data and by intervening on a single edge we have modified all of the competing effects of interest. Comparing the distribution of estimates produced under γ_u and γ_u' gives us a confounded and thus biased estimate of the impact of increasing the unmeasured confounding through its causal pathway to the treatment on the estimators or functions thereof. We will show that this bias can result in under-estimating the impact of bias amplifying variables.

In general, when we vary one of the regression coefficients along a causal pathway, this has upstream and downstream effects on the proportion of variance explained by all variables going into or out of the varied node. In order to keep the proportional effects of the other edges constant, we need to use the error terms of the structural equations (ϵ_1 and ϵ_2) to absorb the shocks to the marginal variances.

In Figure 5(c), if we change γ_u and simply adjust the structural error term ϵ_a such that the total variance in A remains constant, we can isolate the effect of modifying $U \rightarrow A$. In Figure 6(a) we visualize the consequences of failing to hold the variance of the treatment when we modify γ_u . In red, for Figure 6(a), we simulate bias amplification where $\gamma_u = 0.3$. In green, we simulate bias amplification where γ_u is increased to 0.55 holding all other parameters constant, thus allowing the total variance of the treatment to grow from 1 to 1.21. This has the downstream effect of also increasing the variance of the outcome from 1 to 1.02. This also then impacts the relative proportions of variance explained of the treatment and the outcome that are explained by U and BAV , respectively. Notice that the bias increases from 20% ($\frac{0.36-0.3}{0.3}$) to 43% ($\frac{0.43-0.3}{0.3}$). In blue, we increase γ_u from 0.3 to 0.55, but re-normalize the variance in the treatment to remain constant at 1. The bias now increases further to 67% ($\frac{0.5-0.3}{0.3}$) with respect to the original simulation in red. We do this by decreasing the variance of the independent noise term, ϵ_a to ϵ'_a , allowing it to absorb the increase in variation from U . When we do not fix the variance, we underestimate the impact of the amplifier on both the bias and the variance *because the unfixed variance case simulates a different kind of intervention due to the change in variance of the treatment variable*. In the simulation

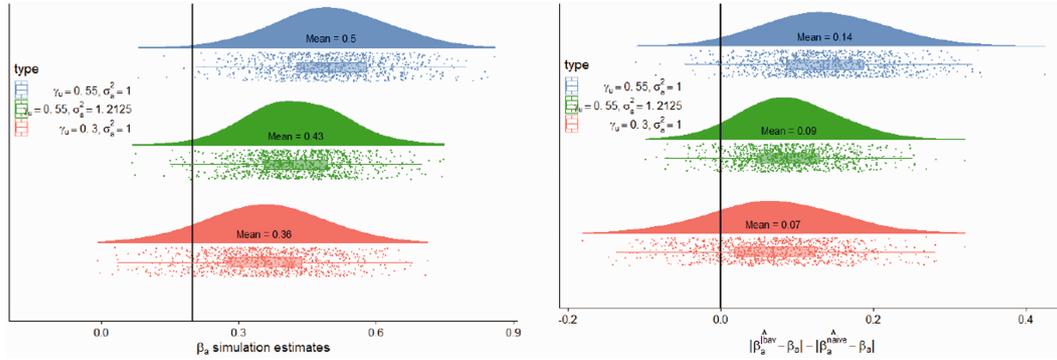


Figure 6. (a) Simulation results from the experiment of intervening on the edge $U \rightarrow A$. The ground truth, $\beta_a = 0.2$, is visualized by the black vertical line. In red we visualize the baseline bias amplification. Green shows the results of the conditional estimator in the case where we increase the weight of the edge, but fail to fix the variance of the treatment, allowing it to grow. In blue, we show the bias when we increase the weight of the edge but now hold the variance constant so the variances remain compatible with the original data. (b) Simulation results from the same experiment as (a) but the outcome is the difference in absolute bias between the conditional estimator $\hat{\beta}_a^{bav}$ and the naive estimator $\hat{\beta}_a^{naive}$. A value greater than zero (black vertical line) indicates that the conditional estimator was more biased than the naive estimator. Parameter values: $\beta_a = 0.2$, $\beta_u = 0.3$, $\beta_{bav} = -0.05$, $\gamma_{bav} = 0.6$.

above, by not keeping the variance fixed in A we implicitly reduced the amount of variance that **BAV** accounts for in the treatment from 36% to 30%. In effect, we were comparing the distribution of

$$P(\hat{\beta}_a | \gamma_u', \gamma_{bav}, \beta_u, \beta_a, \beta_{bav})$$

to

$$P(\hat{\beta}_a | \gamma_u, \gamma_{bav}, \beta_u, \beta_a, \beta_{bav})$$

when a more fair causal counterfactual would be to compare the distribution of

$$P(\hat{\beta}_a | (U \rightarrow A)', U \rightarrow Y, BAV \rightarrow A, BAV \rightarrow Y, A \rightarrow Y)$$

to

$$P(\hat{\beta}_a | U \rightarrow A, U \rightarrow Y, BAV \rightarrow A, BAV \rightarrow Y, A \rightarrow Y)$$

Therefore, our simulation *experiment* results in green are distorted because when we increased the unmeasured confounding through $U \rightarrow A$, we also decreased the strength of the bias amplifying variable through the pathway $BAV \rightarrow A$. Notice that this bias will impact decisions and conclusions we might make about the merits of different estimators in this context. For example, below we compare the conditional estimator, $\hat{\beta}_a^{bav}$, to the naive estimator, $\hat{\beta}_a^{naive}$ with respect to their bias in the same three simulation set ups.

In Figure 6(b) we show the direct comparison of the bias for the conditional and the naive estimators. When we increase the unmeasured confounding through γ_u but fail to renormalize the treatment variance, we do not capture the full extent to which the conditional estimator amplifies the bias. If we compare the green and red plot, it would seem that nearly doubling the unmeasured confounding coefficient only has a small impact on the relative bias of the naive and conditional estimator, since the relative bias only increased from 0.07 to 0.09 (29%). By comparing the green density plot to the blue, we see the relative bias doubles (from 0.07 to 0.14). Therefore, the decision to use the naive or conditional estimator is in fact much more sensitive to the amount of unmeasured confounding than it would appear under the improper simulation with floating variance. It is extremely important to do these kinds of simulations properly particularly in the context of sensitivity analysis where we are testing the performance of estimators with respect to untestable assumptions such as unmeasured confounding.

To properly simulate bias amplification and answer questions of clinical concern with respect to the merits of potential estimators, we must think of the structural equations as an interconnected system. While we typically specify such equations from the perspective of determining their conditional means, the structural equations along with our independence assumptions determine the variances of the variables in the system. Above, this necessitates increasing the strength of the edge $U \rightarrow A$ while holding all other edges constant, which requires us to re-normalize the variances to maintain the strength of the edge $BAV \rightarrow A$.

In Appendix section A.10, we consider the properties of a simulation experiment aiming to vary the strength of the edge $BAV \rightarrow A$. We show that in the case that we fail to fix the variance of the treatment that the bias of the conditional estimator $\hat{\beta}_a^{bav}$ is invariant to $\gamma_{bav}, \forall \gamma_{bav} \in (-\infty, \infty)$, but that the naive estimator is strictly increasing in γ_{bav} . It is clear from the theory we developed in section 3 that if we increase the edge from $BAV \rightarrow A$ that amplification should strictly increase, but if we allow the variance in A to increase as the parameter increases, the amplification effect is precisely cancelled out.

In general terms, simulating linear systems of location-scale family random variables requires first fixing the variances of the variables in the DAG. The relevant quantity determining the strength of the various edges are ratios of variances and covariances of the upstream parent nodes to the variance of the child node in determining the edge's strength. Since the effects are relative, in a simulation context we can normalize the variances to 1 or set them to the expected/observed variances of the data in a particular context. For simplicity we will demonstrate the normalized approach. In Figure 3, this means that $\sigma_u^2 = \sigma_a^2 = \sigma_{bav}^2 = \sigma_y^2 = 1$.

The second step is to be explicit about independence and conditional independence assumptions. Given the independence assumptions, we can specify the covariance matrix of each child variable Y_{child} in terms of the matrix of k arbitrary parent variables which form the edges going into the child variable, Y_{child} .

$$Y_{child} = X_{parent}\beta_{parent} + \epsilon_{child}$$

$$Var(Y_{child}) = \beta_{parent}^T Var(X_{parent}) \beta_{parent} + Var(\epsilon_{child})$$

$$\sigma_{y_{child}}^2 = 1 = [\beta_1 \quad \beta_2 \quad \dots \quad \beta_k] \begin{bmatrix} 1 & \sigma_{1,2} & \dots & \sigma_{1,k} \\ \sigma_{1,2} & 1 & \dots & \sigma_{2,k} \\ \vdots & \sigma_{j,2} & \ddots & \vdots \\ \sigma_{k,1} & \sigma_{k,2} & \dots & 1 \end{bmatrix} \begin{bmatrix} \beta_1 \\ \beta_2 \\ \vdots \\ \beta_k \end{bmatrix} + \sigma_{\epsilon}^2$$

The diagonal of all the parent covariance matrices is 1 since we have normalized all variables pictured in the DAG. The covariances themselves will be determined by the independence assumptions, the edges connecting the child nodes, and their structural equations. Essentially we are choosing the proportion of the child variation that the variances and the covariances of the parent variances explain. The error terms, ϵ 's are the only non-normalized variances, and they absorb the shocks when we increase and decrease the strength of the edges of the non-error variables. This maintains the strength of all other relations visualized on the DAG.

Since all variance terms must be non-zero (or equivalently that $\beta_{parent}^T Var(X_{parent}) \beta_{parent} \leq 1 = \sigma_{child}^2$), the variance equations define bounds on the simulation parameter space. In the above example, conditional on holding the strength of the edges $U \rightarrow Y, BAV \rightarrow A, BAV \rightarrow Y, A \rightarrow Y$, $\gamma_u \in (-0.893, 0.893)$ defines the feasible range. That is, the edge $U \rightarrow A$ can explain up to 79.75% of the variation in A ($\frac{\gamma_u^2 \sigma_u^2}{\sigma_a^2} = \frac{\gamma_u^2}{1}$) since the edge $BAV \rightarrow A$ explains 20.25% of the variation already. In general, the extent to which an edge can explain variation in the child node is constrained by the other child nodes and the covariance structure between those variables. A parameter, however, such as γ_u may be constrained by more than one set of inequalities. In this particular case, γ_u has to satisfy the following inequalities

$$|\gamma_u| \leq (1 - \gamma_{bav}^2)^{\frac{1}{2}}$$

$$\gamma_u \leq \frac{1 - \beta_a^2 - \beta_u^2 - \beta_{bav}^2 - 2\beta_a\beta_{bav}}{2\beta_a\beta_u}$$

where conditional on the strength of the particular edges ($U \rightarrow Y$, $BAV \rightarrow A$, $BAV \rightarrow Y$, $A \rightarrow Y$) in the above simulation, only the first inequality was binding.

The nuance here is that the extent to which we can simulate unmeasured confounding depends upon not only how much amplifying we have simulated, but also on the true effect of the treatment on the outcome $A \rightarrow Y$. Since this is an interdependent system of equations, all of the parameters are competing for shares of fixed variances. If the treatment, independent of U and BAV , explains the large majority of the outcome variance (i.e. the edge $A \rightarrow Y$), it means the weight of the edge $U \rightarrow Y$ must be relatively small, opposite signed, or the structural equations contain an effect modifier. This in turn constrains γ_u .

Consider again the above simulation experiment where we are interested in varying the strength of $U \rightarrow A$ conditional on all other pathways. Suppose that the pathway $A \rightarrow Y$ explains 64% of the variance in Y , i.e. that $\beta_a = 0.8$. Now both constraints on γ_u are binding and the simulation parameter space is $\gamma_u \in (-0.893, 0.3916)$.

In summary, when simulating linear location-scale family systems of equations, we start by identifying the DAG and the independence assumptions between variables. Second, our simulation experiment should attempt to answer a causal question about how a proposed estimator behaves in response to an intervention on the weights of causal DAG. Just like experimental design, properly estimating the relevant counterfactual requires that the difference in distributions between our intervention(s) and the control is the effect of the intervention(s) themselves. As demonstrated in this section, simulating linear systems of equations requires varying one of the edges of the DAG holding all else constant, and matching the means and variances of the simulated variables with that of the target observational study we are trying to mimic. This allows us to generate simulations whose distributions are proper counterfactuals. Third, conditional on the other edges, the covariance matrices impose bounds for the parameter space that we can simulate and thus the extent to which we can vary the edge of interest. For a specific realization of the experiment and accompanying valid parameters, the variables are constructed in the downstream direction, that is from parent nodes to child.

In the example of simulating the proper intervention in Figure 5(c), we first simulate U and BAV independently with variance 1, respectively. Given γ_u and γ_{bav} , the variance of the error term ϵ_2 from equation (5) is implied and can be simulated. Having U , BAV and ϵ_2 allows us to simulate the treatment A . Conditional on the already simulated variables, their associated parameters, and β_a , β_{bav} , and β_u , the variance of the error term ϵ_1 is implied and can be simulated. Finally, since all of the child variables for the outcome have been simulated, we can simulate the outcome. To be clear, we can fix proportions of variance explained by each edge in any order we'd like as long as we respect the underlying constraints. However, given an admissible set of weights of the edges we must proceed from parent to child nodes to conduct the simulation.

While this method requires us to calculate inequalities and make explicit the implied variance formulas for our variables, the benefits are that we can view our simulation as a well-defined causal experiment matching the constraints of our target study and we get sets of parameter bounds. When we do not keep the variance fixed, there are no defined bounds beyond heuristics, and more importantly, we are no longer matching the data to our target observational study. In many small systems, such as the one in Figure 5(a), it is often computationally inexpensive to simulate a discretized approximation to all possible parameter configurations. In extremely large systems, we can use domain knowledge to make refinements on these bounds and simulate a reasonable subset of the parameter space. This method allows us to make refinements over edges with strong priors while simulating the entirety of edges with greater uncertainty.

An alternative, but equivalent way of developing a causal simulation for linear systems of equations would be to work with correlation (or covariance) matrices and variances directly as opposed to parameters. In the Appendix section A.5, we discuss representations of regression coefficients and linear structural parameters in terms of partial correlations, correlations, and coefficients of determination. In some cases, it may be easier to elicit domain knowledge using partial correlations or correlations, compared to the parameters themselves, depending on the clinician and the application at hand. There are still constraints on the system, which are implied by the constraints of a non-singular correlation matrix. One of the potential benefits of working directly with the correlation matrix is that, as shown in section A.5 much of the pairwise correlation matrix can be directly estimated from the observed data. Further, there exist methods to augment the estimable portion of the matrix to a non-singular full correlation matrix as well as methods to add noise to a valid correlation matrix to represent the uncertainty inherent in the estimation process. Additionally, sometimes causal DAGs may imply further restrictions on correlation matrices (or partial correlation matrices) and thus the underlying regression parameters and causal effects.

6 Simulating bias amplification from a real data set

Here we conduct a data simulation for an observational study. We want to consider a medical example with realistic amounts of variance in the treatment and the outcome. Further, we specifically consider the case of a binary treatment which is common in medical applications, biostatistics, and epidemiology. The difficulty, in general, when working with real data in the spirit of plasmode simulation^{21,22} is that you do not know the true underlying parameter values. In this section, we start with a randomized controlled trial (RCT) and modify it appropriately, so that we can take the intention to treat (ITT) estimate as the true underlying effect for the foundation of our simulations.

In our simulation experiment, we keep the treatment data unchanged (thus fixing their variance), and then simulate unmeasured confounding (U) and bias amplifiers (BAV) in order to modify selected covariates (X) and the outcome (Y) to produce a synthetic observational experiment. In order to precisely control the relationships between the simulated variables and the real variables, we treat the binary treatment, A , as though it comes from a latent probit model.

$$\begin{aligned} A &= 1(A^* > 0) \\ &= 1(\alpha_a + U\gamma_u + \tilde{X}\gamma_{\tilde{X}} + \epsilon_2 > 0) \end{aligned}$$

where $\tilde{X} = \frac{X}{\sigma'} + BAV$, and σ' is a scaling variable such that X and \tilde{X} have the same population variance. All of the latent variables (U , $BAV_{n \times k}$, ϵ_2 , and hence $A^* = \alpha_a + U\gamma_u + BAV\gamma_{bav} + \epsilon_2 > 0$) are set to come from normal distributions. The details of the how the simulation is performed are in Appendix section A.11.

For this paper, we use data from Helle et al.,^{23,38,39} a published RCT with 294 participants and relatively balanced distribution of covariates. While the researchers examined many outcomes, we will focus on the effects of an e-Health intervention in infants on child eating behaviors. The researchers gave the parents in the treatment group access to a “monthly age-appropriate video addressing infant feeding topics together with corresponding cooking films/recipes,” and the outcome was eating habits of the child at a later point in time. In the observational study that we want to create (target observational study), we want to estimate the effect of the treatment on emotional overeating as measured by the Child Eating Behavior Questionnaire (CEBQ).

6.1 Unbiased ITT model

Our foundation is the unbiased ITT effect from the RCT data regressing the treatment on the outcome ($Y \sim A$) shown in the first column of Table 1.

In column 1 of Table 1, we see that the ITT estimate is 0.14. As this is an RCT, we do not expect baseline covariates [Child Food Neophobia Score ($CFNS$), Child Feeding Questionnaire (CFQ) subscale pressure, and Age of mother (Age_{mother})] to be associated with exposure. We thus assume that the experimental data are generated from the causal DAG in Figure 7(a), where X represents the matrix of all three covariates ($CFNS$, CFQ , and Age_{mother}) after they have been individually standardized to have mean 0 and variance 1. To verify that these variables are not bias amplifiers, that is explain only a negligible proportion of the treatment variance, we also present the results of the regression of the treatment on the three covariates in column 3 in Table 1. We can see that jointly and individually the three covariates explain very little of the variance in the treatment, $\mathcal{R}^2 = 0.009$. This should be expected in a truly randomized experiment set-up since proper randomization breaks the causal association from the covariates to the treatment. Since these covariates do not cause A and we have assumed that the ITT estimator is unbiased, when we estimate $Y \sim A + CFNS_{score} + CFQ_{pressure} + Age_{mother}$ the expectation and probability limit of $\hat{\beta}_a$ remains unchanged regardless of the strength of association between the covariates and the outcome. However, actual results may vary due to final sample variation. In our RCT data, the unadjusted model estimates a treatment effect of 0.144 and the adjusted model estimates 0.137. Since simulation experiments performed in section 6.2 all condition on covariates, we consider the covariate adjusted results from the RCT as the gold standard for determining bias due to unmeasured confounding in our simulated data.

6.2 Biased model simulations

Our objective is to simulate data according to the DAG in Figure 7(b). To produce the simulations, we took 10,000 bootstrap replications of the original outcome, treatment and covariates. From each bootstrap sample of

Table 1. Clinical data regression table.

Model	ITT	ITT Cond.	A
A	0.144 (0.053)	0.137 (0.053)	–
CFNS	–	0.007 (0.005)	–005 (0.007)
CFQ	–	0.058 (0.036)	0.036 (0.040)
Age_{mother}	–	0.008 (0.006)	0.009 (0.007)
R^2	0.018	0.045	0.009

Note: Each column represents a different regression. Column 1 is the unbiased Intention to Treat (ITT) model, regressing the outcome on the treatment. Column 2 represents the unbiased conditional ITT model with three covariates conditioned on. The third column represents the regression of the treatment on the three regressors. The row names are the independent variables, where A is the treatment, CFNS is the Child Food Neophobia Score, CFQ is the Child Feeding Questionnaire pressure subscale, and Age_{mother} is the age of the infant's mother. The relevant standard errors are displayed in brackets.

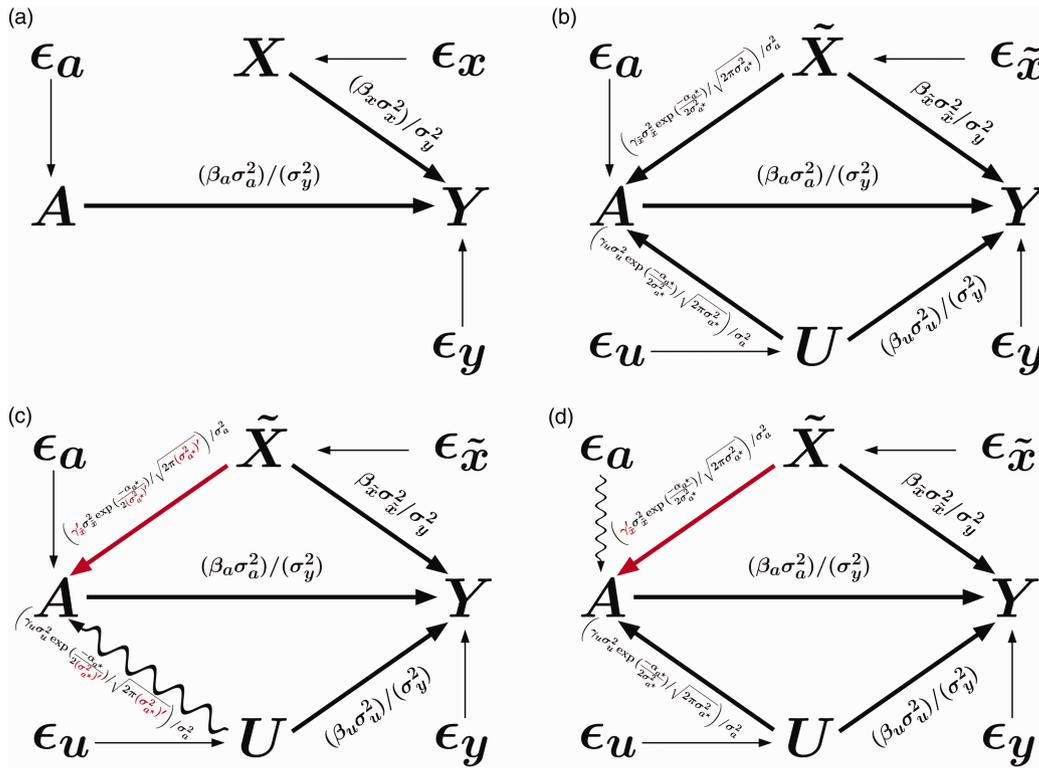


Figure 7. Causal diagrams. (a) DAG representing the original experiment data. (b) DAG with the modified data (see Appendix A.11 for details). (c) represents intervening on the causal DAG in (b) by changing the strength of the edge $\tilde{X} \rightarrow A$ without holding the latent treatment variance A^* constant. The edges which have been modified inadvertently are shown as squiggly arrows. Parameters which have been changed are shown in red. (d) The intervention on the edge $\tilde{X} \rightarrow A$ while holding the latent variance constant.

the treatment, $A_{bootstrap}$, of size $n = 294$ we simulated the latent variable A^* using the procedure outlined in the Appendix (section A.11). Next, conditional on the drawn latent samples of A^* and the bootstrapped covariates, we drew samples for the unmeasured confounding, U , and bias amplifying variable, BAV . The modified random control variables, $\tilde{X} = \frac{X}{\sigma} + BAV$, were produced by adding the bias amplifying variables to a scaled version of the original control variables. Linear combinations of the unmeasured confounding and modified covariates were then added with reasonable values to the outcome such that the following DAG and equations hold (see simulation results).

$$\tilde{Y} = \alpha_y + A\beta_a + \tilde{X}\beta_{\tilde{x}} + U\beta_u + \epsilon_t \quad (22)$$

$$\mathbf{A}^* = \alpha_a + \mathbf{U}\gamma_u + \tilde{\mathbf{X}}\gamma_{\tilde{\mathbf{X}}} + \epsilon_2 \quad (23)$$

$$\mathbf{A} = 1\{\mathbf{A}^* > 0\} \quad (24)$$

In section 6.1 we showed that the true treatment effect was 0.137 conditional on the covariates \mathbf{X} . In the bootstrap simulation pictured below, the unbiased model conditional on both the modified covariates, $\tilde{\mathbf{X}}$, and the unmeasured confounding \mathbf{U} is 0.136 as expected. The naive model estimator had an average estimate of 0.234 in the simulations and thus an absolute estimated bias of 0.097, or a relative bias of 1.8 standard deviations ($\frac{0.234-0.137}{0.053}$) with respect to the unbiased estimate in section (6.1).

When we further condition on the modified covariates, the absolute bias ($E[|\hat{\beta}_a^{\tilde{\mathbf{X}}} - 0.137|]$) more than doubles to 0.225, and the relative bias increases to 4.3 standard deviations ($\frac{0.36-0.137}{0.053}$) with respect to the unbiased estimate in section 6.1. The simulations confirm that bias amplification can be significant even when constrained to problems of realistic variance. Further, we see that bias amplification is potentially a problem for binary outcomes. This underscores the theoretical points made in sections (3) and (4) where we showed that the phenomenon behind bias amplification does not require specific distributional assumptions of the variables in the model.

More importantly, by combining the methodology outline in the appendix (see Appendix A.11) to simulate measured confounding using real data and the principles for simulating systems of equations in section 5, we can produce realistic and complete simulations of parameter spaces which match the underlying characteristics of the data. Investigators who choose covariates based on the assumption of no unmeasured confounding can now evaluate the amount of bias amplification that would occur if this assumption does not hold.

Finally, in the Appendix (section A.12) we consider an example of a causal simulation experiment with a binary treatment variable under the DAG in Figure 7(b) and structural equations (22) to (24). The experiment involves modifying the strength of the edge $\tilde{\mathbf{X}}_I \rightarrow \mathbf{A}$ and evaluating the impact on the naive and conditional estimators. With binary treatment (\mathbf{A}), we show that if we fail to hold the variance of the latent treatment (\mathbf{A}^*) constant and increase $\gamma_{\tilde{\mathbf{X}}_I}$, then it is possible to decrease the amount of observed treatment variance (σ_a^2) explained by $\tilde{\mathbf{X}}_I$. Further, the increased treatment variance also decreases the strength of the edge $\mathbf{U} \rightarrow \mathbf{A}$. As a result of performing the causal simulation experiment improperly, it appears as though that varying the strength of the potential amplifiers has a negligible or negative impact on the resulting bias amplification. The improper and proper approaches to intervention are shown in Figure 7(c) and (d), respectively and the results from these simulations are visualized in Figure 10 of the Appendix section A.12. This of course leads to improper inferences regarding the relative merits of the naive and conditional estimators as well. This highlights once again the importance of comparing simulations with comparable properties and ensuring that when we intervene on the edges of our causal diagram that we are not inadvertently varying the edges we mean to keep fixed. Just as in the experimental context, our simulation results become muddled or meaningless if we are not evaluating well-articulated counterfactuals.

7 Discussion

Causal model selection techniques have largely been developed under the assumption that a sufficient set of variables is available to create ignorability. When a sufficient set is not available or when a causal variable selection technique does not correctly identify the sufficient set, we are at risk of bias amplification. In the first simulation in section 1, we showed that even under mild perturbations of the usual assumptions, conditioning on a set of jointly strong proxy variables for \mathbf{A} in OLS led to a very biased estimator (0.73 standard deviations on average). Further, most current causal variable selection techniques are likely to include this set of variables since they are significant predictors of the outcome and the treatment as well as variables which cause large changes in estimates when included sequentially.

Under threat of bias amplification, treatment-oriented selection techniques for regression analyses using continuous exposure regimes should be used cautiously unless one has strong priors that a sufficient set is available and likely to be identified. We showed in section 3 that it is precisely the amount of variance in the treatment explained by the observables in our model which is responsible for bias amplification. Similarly, we can see that a significant change in estimate is not sufficient to suggest that overall bias is decreasing since this could be the result of further bias amplification.

These results call for new techniques to be developed for observational studies which can accommodate unmeasured confounding to help researchers choose reasonable and least-biased methods. We suggest to first

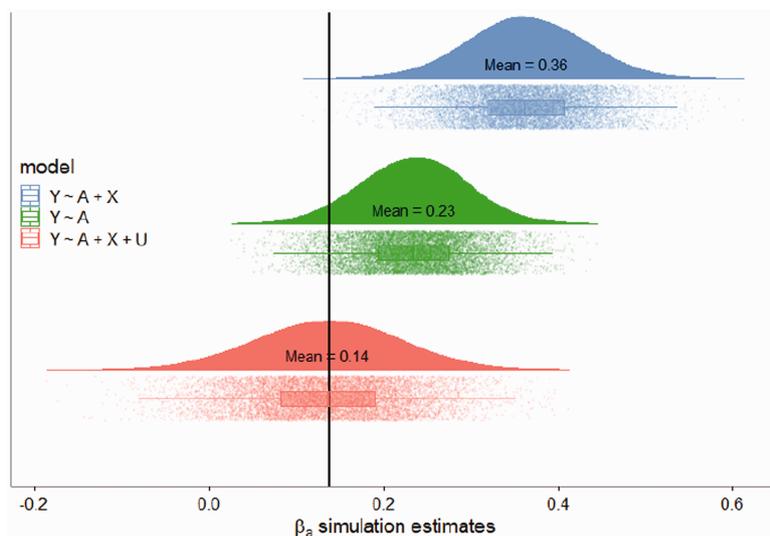


Figure 8. Here we compare three estimators for β_a from the structural equations in equations (22), (23), and (24). In red, the results of the unbiased and infeasible estimator are shown, centered at the true value $\beta_a = 0.137$ (shown by the vertical black line). In green, the replications for the naive estimator ($\hat{\beta}_a^{naive}$) is shown and in blue the replications for the conditional estimator ($\hat{\beta}_a^X$) is shown. Simulation details: Bootstrap replications = 10,000. $\beta_a = 0.137$, $\beta_u = 0.15$, $\beta_{\tilde{x}} = (0.10, -0.15, -0.10)$, $COV(A, U) = 0.25$, $COV(A, \tilde{X}) = (0.22, 0.15, 0.13)$.

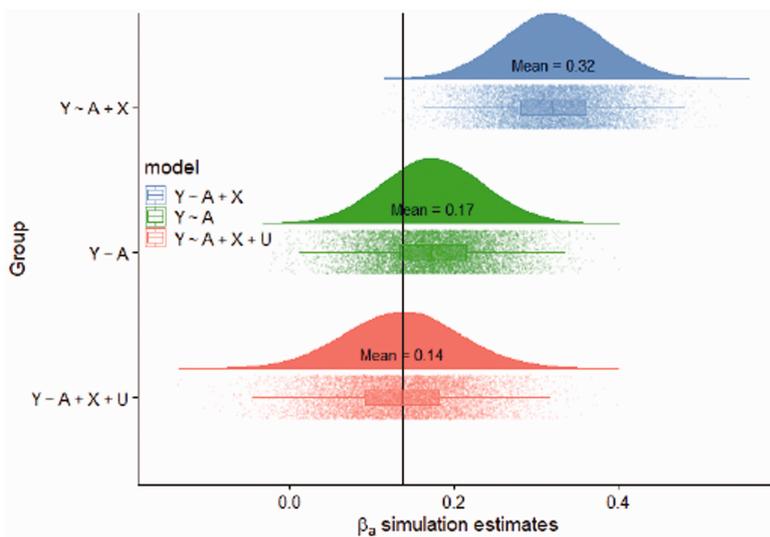


Figure 9. Control treatment estimators. Black line represents the true underlying parameter $\beta_a = .1377$. $N = 10,000$ simulation replications.

identify the most plausible causal DAG. From the DAG and basic structural equation assumptions, an expression for asymptotic bias can often be derived. Further, we suggest to estimate the always-identifiable amplification term in observational settings and to assess the risk of bias amplification. With a measure for amplification and a limiting bias expression, a sensitivity analyses can be performed. One reasonable sensitivity analysis approach would be to estimate the amount of unmeasured confounding required in the spirit of E-values⁷ to determine the strength of confounding associations required to “explain away the treatment effect”⁷ and to make principled inferences from the data. This would require, as we have shown, properly simulating the unmeasured confounder so that (1) the properties of the original data are respected and (2) other competing effects, i.e. edges of the DAG, are not inadvertently altered. In such a set-up, large effects when the controls jointly explain little of the treatment variance lend credibility to results as being robust to unmeasured confounding, particularly in cases when suitable

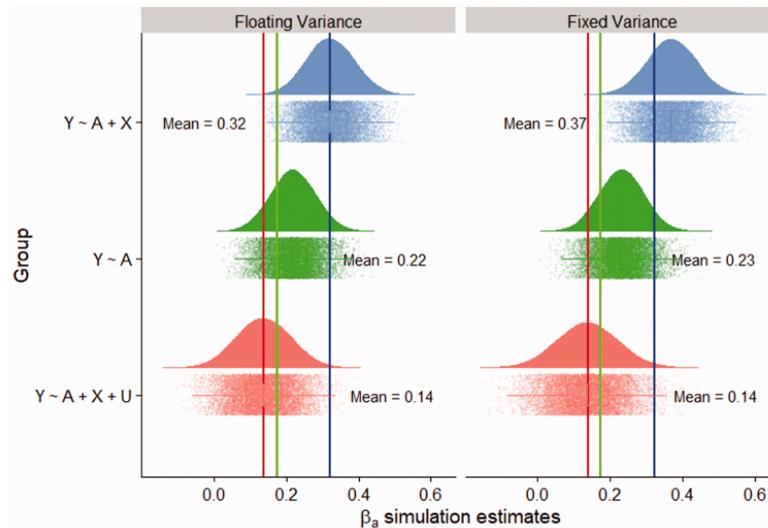


Figure 10. The vertical lines represent the means from the control experiment: red representing the mean of unbiased control estimator, green the mean of the naive estimator, and blue the mean of the amplified estimator.

priors can be placed on the variables along the unmeasured confounding pathway. Alternatively, one could follow the approach of Carnegie et al.⁸ and use the underlying structural equations and the data to generate candidate values of the unmeasured confounding. As we showed in section 5, it is important that any such simulation method take into account the asymmetry of bias amplification with respect to the weight of the edge $U \rightarrow A$ and $U \rightarrow Y$.

Ultimately, simulation experiments must aim to produce data from which we can draw causal conclusions to questions about estimators or functions. This means having well-defined interventions on the edges of the causal graphs and holding the other edges constant. In linear systems of equations, this requires keeping the moments of the variables, in particular variance, fixed when modifying the weight of the DAG's edges. If we allow the treatment variance to vary incidentally as we increase confounding effects, the intervention arm of our simulations will no longer match the target observational study in the control arm. As a further consequence, the additional variance in the exposure may absorb much of the amplifying effect. This leads to systematic underestimation of bias amplification and may be an explanation for why the threat of bias amplification has not been appreciated as a concern for applied researchers.¹⁵ Fixing the variance of the variables has the additional benefit of defining the feasible parameter space. By constraining the underlying parameters by the implied variance equations, it is computationally and conceptually easier to simulate the entire range of plausible treatment effects and biases. This leads to more representative simulations and more principled inferences.

Declaration of conflicting interests

The author(s) declared no potential conflicts of interest with respect to the research, authorship, and/or publication of this article.

Funding

The author(s) disclosed receipt of the following financial support for the research, authorship, and/or publication of this article: The authors received funding from the Canadian Institute of Health Research (CIHR) through the Collaborative Health Research Projects (NSERC partnered) for the research and publication of this article (grant number: CPG-140204).

ORCID iD

Tyrel Stokes  <https://orcid.org/0000-0001-9305-1859>

Supplemental Material

Supplementary material for this article is available online.

References

1. Rubin DB. Estimating causal effects of treatments in randomized and nonrandomized studies. *J Educ Psychol* 1974; **66**: 688.
2. Rosenbaum PR and Rubin DB. The central role of the propensity score in observational studies for causal effects. *Biometrika* 1983; **70**: 41–55.
3. Wooldridge JM. *Econometric analysis of cross section and panel data*. Cambridge, MA: MIT Press, 2010.
4. Witte J and Didelez V. Covariate selection strategies for causal inference: classification and comparison. *Biometric J* 2018; **61**: 1270–1289.
5. Hernn MA, Hernndez-Daz S, Werler MM, et al. Causal knowledge as a prerequisite for confounding evaluation: an application to birth defects epidemiology. *Am J Epidemiol* 2002; **155**: 176–184, <https://doi.org/10.1093/aje/155.2.176>
6. Greenland S and Pearce N. Statistical foundations for model-based adjustments. *Ann Rev Public Health* 2015; **36**: 89–108.
7. VanderWeele TJ and Ding P. Sensitivity analysis in observational research: introducing the E-value introducing the E-value. *Ann Intern Med* 2017; **167**: 268–274, <https://doi.org/10.7326/M16-2607>
8. Carnegie NB, Harada M and Hill JL. Assessing sensitivity to unmeasured confounding using a simulated potential confounder. *J Res Educ Effectiveness* 2016; **9**: 395–420, <https://doi.org/10.1080/19345747.2015.1078862>
9. Talbot D and Massamba VK. A descriptive review of variable selection methods in four epidemiologic journals: there is still room for improvement. *Eur J Epidemiol* 2019; **34**: 725–730. Retrieved on 4 March 2021 from <https://doi.org/10.1007/s10654-019-00529-y>
10. Pearl J. On a class of bias-amplifying variables that endanger effect estimates. *arXiv preprint arXiv:12033503*. 2012.
11. Pearl J. Invited commentary: understanding bias amplification. *Am J Epidemiol* 2011; **174**: 1223–1227.
12. Middleton JA, Scott MA, Diakow R, et al. Bias amplification and bias unmasking. *Political Analys* 2016; **24**: 307–323.
13. Wooldridge JM. Should instrumental variables be used as matching variables? *Res Econom* 2016; **70**: 232–237. Retrieved on 4 March 2021 from <http://www.sciencedirect.com/science/article/pii/S1090944315301678>
14. Ding P, Vanderweele T and Robins J. Instrumental variables as bias amplifiers with general outcome and confounding. *Biometrika* 2017; **104**: 291–302.
15. Myers JA, Rassen JA, Gagne JJ, et al. Effects of adjusting for instrumental variables on bias and precision of effect estimates. *Am J Epidemiol* 2011; **174**: 1213–1222. Retrieved on 4 March 2021 from <http://dx.doi.org/10.1093/aje/kwr364>
16. Velleman PF and Welsch RE. Efficient computing of regression diagnostics. *Am Stat* 1981; **35**: 234–242.
17. Chernozhukov V, Chetverikov D, Demirer M, et al. *Double/debiased machine learning for treatment and structural parameters*. Oxford, UK: Oxford University Press, 2018.
18. Robinson PM. Root-N-consistent semiparametric regression. *Econometrica: Econometric* 1988; 931–954.
19. Brockwell PJ, Davis RA and Fienberg SE. *Time series: theory and methods: theory and methods*. Berlin, Germany: Springer Science & Business Media, 1991.
20. Bowsher CG and Swain PS. Identifying sources of variation and the flow of information in biochemical networks. *Proc Natl Acad Sci* 2012; **109**: E1320–E1328. Retrieved on 4 March 2021 from <https://www.pnas.org/content/109/20/E1320>
21. Vaughan LK, Divers J, Padilla MA, et al. The use of plasmodes as a supplement to simulations: a simple example evaluating individual admixture estimation methodologies. *Computat Statistics Data Analys* 2009; **53**: 1755–1766.
22. Cattell RB and Jaspers J. A general plasmode (No. 30-10-5-2) for factor analytic exercises and research. *Multivariate Behavior Res Monograph* 1967; **67**: 1–212.
23. Helle C, Hillesund ER, Wills AK, et al. Replication data for examining the effects of an eHealth intervention from infant age 6 to 12 months on child eating behaviors and maternal feeding practices one year after cessation: the Norwegian randomized controlled trial Early Food for Future Health. *Data Set* 2019. Retrieved on 4 March 2021 from <https://doi.org/10.18710/R2KJHK>
24. Davidson R and MacKinnon JG. *Econometric theory and methods*. vol. 5. New York, NY: Oxford University Press, 2004.
25. Chung KL. *A course in probability theory*. Cambridge, MA: Academic Press, 2001.
26. Whittaker J. *Graphical models in applied multivariate statistics*. Hoboken, NJ: Wiley Publishing, 2009.
27. Artner R, Wellingerhof PP, Lafit G, et al. The shape of partial correlation matrices. *Commun Stat-Theory Meth* 2020. DOI: 10.1080/03610926.2020.1811338.
28. Budden M, Hadavas P, Hoffman L. On the generation of correlation matrices. *Appl Math E-Notes* 2008; **8**: 279–282.
29. Hardin J, Garcia SR and Golan D. A method for generating realistic correlation matrices. *Ann Appl Stat* 2013; **7**: 1733–1762.
30. Vander Weele TJ and Shpitser I. A new criterion for confounder selection. *Biometrics* 2011; **67**: 1406–1413.
31. Vansteelandt S, Bekaert M and Claeskens G. On model selection and model misspecification in causal inference. *Stat Meth Med Res* 2012; **21**: 7–30.
32. Vander Weele TJ. Principles of confounder selection. *Eur J Epidemiol* 2019; **34**: 211–219. Retrieved on 4 March 2021 from <https://doi.org/10.1007/s10654-019-00494-6>
33. Helle C, Hillesund E, Omholt M, et al. Early food for future health: A randomized controlled trial evaluating the effect of an eHealth intervention aiming to promote healthy food habits from early childhood. *BMC Public Health* 2017; **17**: 1–12.

34. Helle C, Hillesund ER, Wills AK, et al. Examining the effects of an eHealth intervention from infant age 6 to 12 months on child eating behaviors and maternal feeding practices one year after cessation: The Norwegian randomized controlled trial Early Food for Future Health. *PLoS One* 2019; **14**: e0220437.
35. Jiang L, Oualkacha K, Didelez V, et al. Constrained instruments and their application to Mendelian randomization with pleiotropy. *Genetic Epidemiol* 2019; **43**: 373–401.
36. Allen M, Poggiali D, Whitaker K, et al. Raincloud plots: a multi-platform tool for robust data visualization. *Wellcome Open Res* 2019; **4**.
37. Wickham H, Chang W and Wickham MH. GGplot2 package create elegant data visualisations using the grammar of graphics version. ■ 2016; **2**: 1–189.

Appendix I

A.1 Matrix notation and FWL theorem

Throughout this paper, we make use of matrix notation to concisely represent estimates and as a way of considering the geometry of the least squares. Here is a quick guide for understanding the notation in this paper.

Let A be the $n \times p$ matrix of treatment variables. For illustrative purposes, consider that A is a single binary $n \times 1$ vector.

$$A_{n \times 1} = \begin{bmatrix} 1 \\ 1 \\ 0 \\ 1 \\ \vdots \\ 0 \end{bmatrix}$$

There are n rows of data, each with a 1 or 0 representing the observation being treated or not.

Another piece of notation that is used is annihilator and orthogonal projection matrices. Let P_X be the orthogonal projection matrix of X , an $n \times k$ matrix, and M_X the annihilator or residual-making matrix of X

$$P_X = X(X^T X)^{-1} X^T;$$

$$M_X = I - P_X = I - X(X^T X)^{-1} X^T$$

A projection matrix maps each point to the nearest point in the subspace spanned by the columns in X , $S(X)$. The annihilator matrix maps each point to the orthogonal complement of $S(X)$, $S^\perp(X)$. The predicted outcome in ordinary least squares is $\hat{Y} = X\hat{\beta} = X(X^T X)^{-1} X^T Y = P_X Y$, which we can think of geometrically as “dropping a perpendicular”²⁴ from the outcome vector into the subspace spanned by the covariates in the regression. The orthogonal complement to the space spanned by the regressors $S^\perp(X)$ is where the fitted residual vector, $\hat{\epsilon}$, lives. We can see that the residual vector $\hat{\epsilon} = Y - \hat{Y} = Y - X\hat{\beta} = Y - P_X Y = (I - P_X)Y = M_X Y$ is just the projection of Y into the subspace orthogonal to $S(X)$.

By definition, we can always then decompose Y uniquely into its projection onto $S(X)$ and $S^\perp(X)$

$$Y = P_X Y + M_X Y$$

Orthogonal projection matrices have two important properties, they are symmetric and idempotent. This means that $P_X = P_X^T$ and $P_X P_X = P_X$ and that these same two properties are equally enjoyed by M_X . Further, any matrix in the subspace spanned by X is *annihilated* when operated on by M_X , since it is by definition orthogonal to $S(X)$.

We also appeal to the Frisch-Waugh-Lovell (FWL) theorem to construct the matrix notation regression estimates as well as for visualizing the 2 dimensional plot of a single regression in the context of a multivariable regression. Suppose we construct an arbitrary partition of $X_{n \times k} = [X_{1 \times k_1}, X_{2 \times k_2}]$, where $k_1 + k_2 = k$. The FWL

theorem states that the following two regressions, equations (25) and (27), produce numerically equivalent estimates of the vector $\hat{\beta}_2$ as well as numerically equivalent residuals.

$$Y = X_1\beta_1 + X_2\beta_2 + \epsilon \quad (25)$$

$$M_{X_1}Y = M_{X_1}X_1\beta_1 + M_{X_1}X_2\beta_2 + M_{X_1}\epsilon \quad (26)$$

$$= M_{X_1}X_2\beta_2 + M_{X_1}\epsilon \quad (27)$$

What this says in words is that it is numerically equivalent to regress Y on the columns of X_1 and X_2 simultaneously as it is to first regress both Y and X_2 on the columns of X_1 separately, then save the respective residuals, $M_{X_1}Y$ and $M_{X_1}X_2$, and regress the former on the later. By simply pre-multiplying both sides of equation (27) by X_2^T and rearranging, we get the general matrix notation formulation for the vector of coefficient estimates $\hat{\beta}_2 = (X_2^T M_{X_1} X_2)^{-1} X_2^T M_{X_1} Y$. Using the idempotency and symmetry properties, we can rewrite the coefficient estimate

$$\hat{\beta}_2 = (X_2^T M_{X_1} X_2)^{-1} X_2^T M_{X_1} Y \quad (28)$$

$$= \frac{(M_{X_1} X_2)^T (M_{X_1} Y)}{(M_{X_1} X_2)^T (M_{X_1} X_2)} \quad (29)$$

$$= \frac{\hat{v}_{x_2} \hat{v}_y}{\hat{v}_{x_2}^2} \quad (30)$$

where \hat{v}_{x_2} and \hat{v}_y are the residuals from the regression of X_2 and Y on X_1 , respectively.

A special case of the above result is when we have a $n \times k_1$ matrix of treatment variables, A and a $n \times k_2$ matrix of controlling variables. For example, we are trying to estimate the causal effect of the matrix A on the outcome Y using a selection on observables strategy by conditioning on Z . The FWL theorem tells us that the estimates of the causal effect, $\hat{\beta}_a$ can be obtained by the two following regression equations

$$Y = A\beta_a + Z\beta_z + v_1$$

$$M_Z Y = M_Z A \beta_a + v_1$$

The second regression is a simple linear regression, with only one dependent variable $M_Z Y$ and a single regressor, $M_Z A$. The error term remains unchanged by the projection into $S^\perp(Z)$ since it can be represented as $M_{A,Z} Y$ which is already contained in the subspace $S^\perp(Z)$. Another way we can write the estimate $\hat{\beta}_a$ is to apply the well-known $\hat{\beta} = (X^T X)^{-1} X^T Y$ to the second regression equation above.

$$\hat{\beta}_a = \frac{(M_Z A)^T (M_Z Y)}{(M_Z A)^T M_Z A}$$

$$= \frac{\frac{1}{n} (M_Z A)^T (M_Z Y)}{\frac{1}{n} (M_Z A)^T M_Z A}$$

The numerator $\frac{1}{n} (M_Z A)^T (M_Z Y)$ can be seen as the dot product of the residuals from the regression of the treatment on the control variables, $A \sim Z$, and the residuals from the regression of the outcome on the control variables, $Y \sim Z$, scaled by $\frac{1}{n}$. If a column of ones is included in the matrix Z , both sets of residuals will be centered. We can then think of the dot product in the numerator as an estimator for the covariance of the two residuals,

$COV(\epsilon_a, \epsilon_y)$. In general terms, we will have bias due to unmeasured confounding if the covariance is a function of U . The denominator can be seen as numerically equal to the sum of squared residuals.

An important special case of the annihilator matrix is \mathbf{M}_I , where \mathbf{I} is a $n \times 1$ vector of ones. This is sometimes called the centering matrix because it de-means the matrix it operates on, since $\mathbf{M}_I \mathbf{X} = \mathbf{X} - \mathbf{I}(\mathbf{I}^T \mathbf{I})^{-1} \mathbf{I}^T \mathbf{X} = \mathbf{X} - \frac{1}{n} \sum_{i=1}^n x_i = \mathbf{X} - \bar{\mathbf{X}}$

Using the symmetry and idempotency properties combined with the convergence in probability properties discussed in the next section, this implies

$$\frac{1}{n} \mathbf{X}_I \mathbf{M}_I \mathbf{X}_2 = \frac{1}{n} (\mathbf{X}_I - \bar{\mathbf{X}}_I)(\mathbf{X}_2 - \bar{\mathbf{X}}_2) \quad (31)$$

$$= \frac{1}{n} \sum_{i=1}^n (x_{1i} - \bar{X}_I)(x_{2i} - \bar{X}_2) \quad (32)$$

$$\xrightarrow{p} E[(\mathbf{X}_I - E[\mathbf{X}_I])(\mathbf{X}_2 - E[\mathbf{X}_2])] \quad (33)$$

$$= COV(\mathbf{X}_I, \mathbf{X}_2) \quad (34)$$

When $\mathbf{X}_I = \mathbf{X}_2$, the last line becomes $\mathbf{Var}(\mathbf{X}_I)$.

For a more complete and technical treatment of projection and annihilation matrices as it pertains to OLS, see *Econometric theory and methods* by Davidson and MacKinnon.²⁴

A.2 Convergence in probability

Throughout the paper, we use the notation $\text{plim}_{n \rightarrow \infty}$ to mean the limit in probability. Specifically, if $\text{plim}_{n \rightarrow \infty} \mathbf{Y}_n = \mathbf{Y}$, then

$$\lim_{n \rightarrow \infty} P(|\mathbf{Y}_n - \mathbf{Y}| > \epsilon) = 0, \forall \epsilon > 0$$

Alternatively we can write $\text{plim}_{n \rightarrow \infty} \mathbf{Y}_n = \mathbf{Y}$ as $\mathbf{Y}_n \xrightarrow{p} \mathbf{Y}$ or simply plim . Throughout this paper, all probability limits are as $n \rightarrow \infty$.

Below are a few important properties of Probability limits used throughout the paper. Suppose $\mathbf{X}_n \xrightarrow{p} \mathbf{X}$ and $\mathbf{Y}_n \xrightarrow{p} \mathbf{Y}$, then

$$\mathbf{X}_n + \mathbf{Y}_n \xrightarrow{p} \mathbf{X} + \mathbf{Y} \quad (35)$$

$$\mathbf{X}_n \mathbf{Y}_n \xrightarrow{p} \mathbf{X} \mathbf{Y} \quad (36)$$

$$\frac{\mathbf{X}_n}{\mathbf{Y}_n} \xrightarrow{p} \frac{\mathbf{X}}{\mathbf{Y}} \quad (37)$$

where the third line is just a special case of the second and holds whenever the denominator is well defined. These properties are well known and follow from the Continuous Mapping Theorem.

Another useful theorem we use in this paper is the Weak Law of Large Numbers (WLLN). Here we consider a set of standard assumptions. Suppose we take the sample average of random variables $\{\mathbf{X}_1, \mathbf{X}_2, \dots, \mathbf{X}_n\}$, such that \mathbf{X}_i 's are independent and identically distributed (iid) and $E[\mathbf{X}_i] = m < \infty$, i.e. the expectation is finite then

$$\frac{1}{n} \sum_{i=1}^n \mathbf{X}_i \xrightarrow{p} E[\mathbf{X}_i] = m. \quad (38)$$

In the paper, whenever specified we assume that the error terms are coming from a normal distribution. In light of the WLLN, we can see that normal error terms are not required for the results to hold, and that in fact we just need the error terms to come from an identical and independent distribution. The above results can be weakened further such that we can replace the iid condition with pairwise independence (see Chung²⁵ for details).

Further, some probability limit results do not always have a closed-form expression, for example the probability limit of a sum of iid cauchy random variables since they do not have finite expectations. Sometimes we express the resulting limit as a function of random variables. These random variables tend to their respective probability limits, provided they exist.

A.3 Derivations continued

Below are derivations, extensions, proofs, and alternate forms of the equations presented in the main text.

A.3.1 Average causal and average partial effects

Throughout this paper, we consider linear models with continuous exposures and as such a natural causal estimand of interest is the average partial effect (APE). Under the linearity assumptions, the APE coincides with the ACE. Below we show the derivation of the APE under the various DAG and structural equation assumptions. Implicitly, we further assume standard regularity conditions such as existence and boundedness of the estimators in \mathcal{L}_1 , so that the derivative operator can freely move inside the expectation integral.

Average partial effects for equations (4) and (5)

$$\begin{aligned} \text{APEs} &= \frac{\partial E[Y|A, U, BAV]}{\partial A} \\ &= \frac{\partial E[\alpha_y + A\beta_a + U\beta_u + BAV\beta_{bav} + \epsilon_I|A, U, BAV]}{\partial A} = \beta_a \end{aligned}$$

When we allow for BAV to be a $n \times k$ vector and for β_{BAV} to be potentially a zero vector, we can see that the above derivation holds for all of the DAG's and structural equations which assume there is no interaction term. $\hat{\beta}_a^{naive}$ in equation (7)

$$\hat{\beta}_a^{naive} = \frac{\mathbf{A}^T \mathbf{M}_z \mathbf{Y}}{\mathbf{A}^T \mathbf{M}_z \mathbf{A}} \quad (39)$$

$$= \beta_a + \beta_u \frac{\mathbf{A}^T \mathbf{M}_I \mathbf{U}}{\mathbf{A}^T \mathbf{M}_I \mathbf{A}} + \beta_{bav} \frac{\mathbf{A}^T \mathbf{M}_I \mathbf{BAV}}{\mathbf{A}^T \mathbf{M}_I \mathbf{A}} + \frac{\mathbf{A}^T \mathbf{M}_I \epsilon_I}{\mathbf{A}^T \mathbf{M}_I \mathbf{A}} \quad (40)$$

$$= \beta_a + \beta_u \frac{\frac{1}{n} \sum_{i=1}^n (a_i - \bar{a})(u_i - \bar{u})}{\frac{1}{n} \sum_{i=1}^n (a_i - \bar{a})^2} + \beta_{bav} \frac{\frac{1}{n} \sum_{i=1}^n (a_i - \bar{a})(bav_i - \bar{bav})}{\frac{1}{n} \sum_{i=1}^n (a_i - \bar{a})^2} + \frac{\frac{1}{n} \sum_{i=1}^n (a_i - \bar{a})(\epsilon_{1i} - \bar{\epsilon}_1)}{\frac{1}{n} \sum_{i=1}^n (a_i - \bar{a})^2} \quad (41)$$

A.4 Derivation of Pearl (2011) result

Here we will explicitly follow Pearl's derivation,¹⁰ to show the advantages of considering the probability limit over strictly expectations. We will derive the expectation for $\hat{\beta}_a^{bav}$ from estimating (9), which is the APE conditional on the treatment, A , and the BAV variable.

$$\begin{aligned} E[\hat{\beta}_a^{bav}] &= \frac{\partial E[Y|A, BAV]}{\partial A} \\ &= \frac{\partial E[\alpha_y + A\beta_a + U\beta_u + BAV\beta_{bav} + \epsilon_I|A, BAV]}{\partial A} \\ &= \beta_a + \beta_u \frac{\partial E[U|A, BAV]}{\partial A} \end{aligned}$$

We must find the expectation of U conditional on A and BAV . Pearl solves this challenge by supposing the true underlying relationship between U , A , and BAV is linear and writing this functional form as a linear regression equation

$$U = \alpha_u + A\zeta_a + BAV\zeta_{bav} + \epsilon_3$$

Using this equation in addition to the two structural equations for Y and A , respectively, we can express the regression coefficients, ζ_a and ζ_{bav} in terms of the structural coefficients β_a , β_u , γ_a , γ_u by equating expressions for the covariances under the two sets of structural equations.

$$\begin{aligned} COV(U, A) &= E[AU] - E[A]E[U] \\ &= E[(\alpha_u + U\gamma_u + BAV\gamma_{bav} + \epsilon_2)U] - E[(\alpha_u + U\gamma_u + BAV\gamma_{bav} + \epsilon_2)]E[U] \\ &= \gamma_u(E[U^2] - E[U]^2) \\ &= \gamma_u\sigma_u^2 \end{aligned}$$

Equivalently

$$\begin{aligned} COV(U, A) &= E[AU] - E[A]E[U] \\ &= E[A(\alpha_u + A\zeta_a + BAV\zeta_{bav})] - E[A]E[(\alpha_u + A\zeta_a + BAV\zeta_{bav})] \\ &= \zeta_a(E[A^2] - E[A]^2) + \zeta_{bav}(E[ABAV] - E[A]E[BAV]) \\ &= \zeta_a\sigma_a^2 + \zeta_{bav}(COV(A, BAV)) \\ &= \zeta_a\sigma_a^2 + \zeta_{bav}\gamma_{bav}\sigma_{bav}^2 \end{aligned}$$

where the last line follows analogously from our derivation of $COV(A, U)$. Putting these together, we have

$$\zeta_a = \frac{\gamma_u\sigma_u^2 - \zeta_{bav}\gamma_{bav}\sigma_{bav}^2}{\sigma_a^2}. \quad (42)$$

Similarly, putting the two steps together for succinctness

$$COV(U, BAV) = 0 \text{ (independence)} \quad (43)$$

$$= E[UBAV] - E[U]E[BAV] \quad (44)$$

$$= E[(\alpha_u + A\zeta_a + BAV\zeta_{bav})BAV] - E[(\alpha_u + A\zeta_a + BAV\zeta_{bav})]E[BAV] \quad (45)$$

$$= \zeta_a COV(A, BAV) + \zeta_{bav} Var(BAV) \quad (46)$$

$$= \zeta_a\gamma_{bav}\sigma_{bav}^2 + \zeta_{bav}\sigma_{bav} \quad (47)$$

Now we have two equations for the two new regression coefficients in terms of the structural equations. Combining

$$\zeta_a = \frac{\gamma_u\sigma_u^2}{\sigma_a^2 - \gamma_{bav}^2\sigma_{bav}^2} \quad (48)$$

$$\zeta_{bav} = \frac{-\gamma_u\gamma_{bav}\sigma_u^2}{\sigma_a^2 - \gamma_{bav}^2\sigma_{bav}^2} \quad (49)$$

Returning to the task of finding the partial effect of A on $E[Y|A, \mathbf{BAV}]$ and thus $E[\widehat{\beta}_a^{bav}]$.

$$\begin{aligned}
 E[\widehat{\beta}_a^{bav}] &= \frac{\partial E[Y|A, \mathbf{BAV}]}{\partial A} & (50) \\
 &= \frac{\partial E[\alpha_y + A\beta_a + U\beta_u + \mathbf{BAV}\beta_{bav} + \epsilon_I|A, \mathbf{BAV}]}{\partial A} \\
 &= \beta_a + \frac{\partial \beta_u E[U|A, \mathbf{BAV}]}{\partial A} \\
 &= \beta_a + \beta_u \frac{\partial E[(\alpha_u + \zeta_a A + \zeta_{bav} + \epsilon_I)|A, \mathbf{BAV}]}{\partial A} \\
 &= \beta_a + \beta_u \zeta_a \\
 &= \beta_a + \beta_u \frac{\gamma_u \sigma_u^2}{\sigma_a^2 - \gamma_{bav}^2 \sigma_{bav}^2} & (51)
 \end{aligned}$$

The approach of Pearl is limited in that it only works when the true underlying form of the conditional expectation, $E[U|A, \mathbf{BAV}]$ is linear in both A and \mathbf{BAV} . As a result, the derivation is cumbersome and does not easily generalize to more complicated cases with more variables or different functional forms. Similarly, we can find the expectation for the naive estimator $\widehat{\beta}_a^{naive}$ from (6) using this method

$$E[\widehat{\beta}_a^{naive}] = \frac{\partial E[Y|A]}{\partial A} \quad (52)$$

$$= \frac{\partial E[\alpha_y + A\beta_a + U\beta_u + \mathbf{BAV}\beta_{bav} + \epsilon_I|A]}{\partial A} \quad (53)$$

$$= \beta_a + \frac{\partial \beta_u E[U|A]}{\partial A} + \frac{\partial \beta_{bav} E[\mathbf{BAV}|A]}{\partial A} \quad (54)$$

We assume that the true underlying relationship between U and A is linear, while also assuming a linear relationship between \mathbf{BAV} and A

$$U = \alpha_u + A\tau_a + \epsilon_4 \quad (55)$$

$$\mathbf{BAV} = \alpha_{bav} + A\eta_a + \epsilon_5 \quad (56)$$

Using these two equations (55) and (56), we arrive at the following expressions for the $COV(A, U)$ and $COV(A, \mathbf{BAV})$

$$COV(A, U) = \tau_a \sigma_a^2 \quad (57)$$

$$COV(A, \mathbf{BAV}) = \eta_a \sigma_a^2 \quad (58)$$

Following an analogous process, we can show that $COV(A, \mathbf{BAV}) = \gamma_{bav} \sigma_{bav}^2$ in terms of the original structural coefficients. Combining the four covariance expressions and solving for τ_a and η_a in terms of the structural equations yields

$$\tau_a = \frac{\gamma_u \sigma_u^2}{\sigma_a^2} \quad (59)$$

$$\eta_a = \frac{\gamma_{bav}\sigma_{bav}^2}{\sigma_a^2} \quad (60)$$

Substituting equations (55) and (56) in equation (54) and using equations (59) and (60) yields the following expectation for $\widehat{\beta}_a^{naive}$:

$$E[\widehat{\beta}_a^{naive}] = \beta_a + \beta_u \frac{\gamma_u \sigma_u^2}{\sigma_a^2} + \beta_{bav} \frac{\gamma_{bav} \sigma_{bav}^2}{\sigma_a^2} \quad (61)$$

Note in this derivation we needed to assume two linear relationships in order to derive the expectation, $E[U|A]$ and $E[BAV|A]$. These assumptions are not necessary when using probability limits to define limiting expressions in terms of the structural parameters.

A.5 Correlations, partial correlations, \mathcal{R}^2 and dimensionless quantities for sensitivity analysis and simulation

In section 3 we discussed the quantity $\eta = \frac{|\beta_{bav}\gamma_{bav}\sigma_{bav}^2|}{|\beta_u\gamma_u\sigma_u^2|}$ and mentioned that it is possible to express as a dimensionless quantity which is useful for sensitivity analysis and simulation studies.

$$\eta = \frac{|\beta_{bav}\gamma_{bav}\sigma_{bav}^2|}{|\beta_u\gamma_u\sigma_u^2|} \quad (62)$$

$$= \frac{\rho_{Y,BAV|A,U}\rho_{A,BAV|U} (1 - \mathcal{R}_{Y|A,BAV}^2)^{\frac{1}{2}} (1 - \mathcal{R}_{U|A,BAV}^2)^{\frac{1}{2}} (1 - \mathcal{R}_{A|U}^2)^{\frac{1}{2}}}{\rho_{U,Y|A,BAV}\rho_{A,U|BAV} (1 - \mathcal{R}_{Y|A,U}^2)^{\frac{1}{2}} (1 - \mathcal{R}_{BAV|A,U}^2)^{\frac{1}{2}} (1 - \mathcal{R}_{A|BAV}^2)^{\frac{1}{2}}} \quad (63)$$

where a partial correlation is defined as the following for general variables Z_1 , Z_2 and p dimensional X

$$\rho_{Z_1,Z_2|X} = \frac{E[\widehat{\epsilon}_1 \widehat{\epsilon}_2]}{\sqrt{E[\widehat{\epsilon}_1^2]E[\widehat{\epsilon}_2^2]}} \quad (64)$$

and $\widehat{\epsilon}_i = \operatorname{argmin}_{(\beta_0, \beta) \in \mathcal{R}^{(p+1)}} (Z_i - \beta_0 - X\beta)^2$, for $i \in \{1, 2\}$. Partial correlations can be expressed as lower correlations by the following recursion formula.²⁶

Consider the partition of $X_{n \times p} = [X_{1n \times (p-1)}, X_{2n \times 1}]$.

$$\rho_{Z_1,Z_2|X_1,X_2} = \frac{\rho_{Z_1,Z_2|X_1} - \rho_{Z_1,X_2|X_1}\rho_{X_2,Z_2|X_1}}{\sqrt{(1 - \rho_{Z_1,X_2|X_1}^2)(1 - \rho_{X_2,Z_2|X_1}^2)}} \quad (65)$$

Further, it can be shown that a non-singular partial correlation matrix uniquely determines the correlation matrix and vice versa.²⁷ To directly express the partial correlations in terms of the correlation matrix, the following formula is convenient, where \tilde{R} is a partial correlation matrix, where element $\tilde{r}_{i,j} = \rho_{X_i,X_j|X \setminus \{X_i,X_j\}}$ is the partial correlation of two variables X_i , X_j partialling out all other variables. Analogously, let R be the correlation matrix, where elements correspond to the usual correlations.²⁷

$$\tilde{R} = -D_{R^{-1}}^{-1} R^{-1} D_{R^{-1}}^{-1} \quad (66)$$

where $D_{R^{-1}}^{-1}$ is a diagonal matrix, where the elements of the diagonal equal the vector $\sqrt{\operatorname{diag}(R^{-1})}$. Equivalently we can construct the correlation matrix from a specified partial correlation matrix.

$$R = D_{-\tilde{R}^{-1}}^{-1} (-\tilde{R})^{-1} D_{-\tilde{R}^{-1}}^{-1} \quad (67)$$

For example, using this relation, or iterative application of the recursion relation we can show that

$$\rho_{Y,BAV|A,U} = \frac{\left\{ \begin{array}{l} \rho_{A,U}\rho_{U,BAV}\rho_{Y,BAV} + \rho_{A,BAV}\rho_{U,BAV}\rho_{Y,U} - \rho_{A,U}\rho_{Y,U} \\ - \rho_{A,BAV}\rho_{Y,BAV} + \rho_{U,BAV}^2\rho_{Y,U} + \rho_{Y,A} - \rho_{U,BAV}^2 \end{array} \right\}}{\sqrt{\left\{ \begin{array}{l} (1 - \rho_{A,U}^2 - \rho_{A,BAV}^2 - \rho_{U,BAV}^2 + 2\rho_{A,U}\rho_{A,BAV}\rho_{U,BAV}) \times \\ (1 - \rho_{Y,BAV}^2 - \rho_{Y,U}^2 - \rho_{U,BAV}^2 + 2\rho_{U,BAV}\rho_{Y,U}\rho_{Y,BAV}) \end{array} \right\}}} \quad (68)$$

in the case the BAV is a single variable. When U and BAV are independent as is the case in the DAG considered in section 3, this simplifies

$$\rho_{Y,BAV|A,U} = \frac{\rho_{Y,A} - \rho_{A,U}\rho_{Y,U} - \rho_{A,BAV}\rho_{Y,BAV}}{\sqrt{(1 - \rho_{A,U}^2 - \rho_{A,BAV}^2)(1 - \rho_{Y,BAV}^2 - \rho_{Y,U}^2)}} \quad (70)$$

Similarly, when U and BAV are independent we have the following

$$\rho_{A,U|BAV} = \frac{\rho_{A,U}}{1 - \rho_{A,BAV}^2} \quad (72)$$

$$\rho_{A,BAV|U} = \frac{\rho_{A,BAV}}{1 - \rho_{A,U}^2} \quad (73)$$

On the other hand, $\rho_{U,Y|A,BAV}$ does not simplify nicely and is easier to handle numerically with the aid of the matrix relation rather than symbolically.

Two things of note. First, correlations and the variances of the variables determine the entire expression η . We showed this explicitly with the partial correlations, but it is also true that it is sufficient to fix a correlation matrix and the variances of the variables to determine the \mathcal{R}^2 values. It is sufficient to consider the space of valid correlation matrices, or equivalently a valid partial correlation matrix, and fix the variances of U and BAV (The variances for Y and A cancel) in order to simulate the entire parameter space for η . Moreover, if we have a real data set, many of the correlations are estimable and this greatly restricts the parameter space.

$$R = \begin{vmatrix} 1 & \rho_{Y,A} & \rho_{Y,BAV} & \rho_{Y,U} \\ \rho_{Y,A} & 1 & \rho_{A,BAV} & \rho_{A,U} \\ \rho_{Y,BAV} & \rho_{A,BAV} & 1 & \rho_{BAV,U} \\ \rho_{Y,U} & \rho_{A,U} & \rho_{BAV,U} & 1 \end{vmatrix}$$

The bolded correlations are strictly functions of observable and can be estimated from the data. By fixing the known values, or only considering variables within an estimated range, we can heavily restrict the space of plausible correlation matrices and thus the range of possible values of η . Similarly, we can estimate σ_{BAV}^2 from the data.

Notice that the estimable correlations form a valid correlation matrix themselves, call this $R_{Y,A,BAV}$. Assuming the treatment and outcome are univariate, $R_{Y,A,BAV}$ is of dimension $(2+p) \times (2+p)$ where p is the dimension of the control variables p . It is well known how to estimate $\hat{R}_{Y,A,BAV}$, and in order to do a proper simulation or sensitivity analysis it suffices to extend this correlation matrix to a valid $(2+p+1) \times (2+p+1)$. There are existing algorithms that can be implemented to extend $(n-1)$ dimension correlation matrices to n dimension correlation matrices²⁸ which can be used instead of solving parameter inequalities. Further, there is work looking at how to simulate realistic correlation matrices around a given structure²⁹ and finding ways to add noise to correlation matrices while ensuring the result is still positive definite. Such methods could allow for researchers to better represent the underlying uncertainty associated with the estimated matrix $\hat{R}_{Y,A,BAV}$.

A.6 Partial correlations and regression coefficients

Here we show the relationship between partial correlations and regression coefficients. Like in the text we distinguish between true structural coefficients β_x and estimates from various regressions. First let us consider an estimate.

$$\widehat{\beta}_x = \underset{\beta \in \mathcal{R}^p}{\operatorname{argmin}} (Y - X\beta)^2 \quad (74)$$

$$\widehat{\beta}_{x_i} = \frac{\widehat{\epsilon}_{y \setminus \{x_i\}} \widehat{\epsilon}_{x \setminus \{x_i\}}}{\widehat{\epsilon}_{x \setminus \{x_i\}}^2} \quad (75)$$

where $\widehat{\epsilon}_{y \setminus \{x_i\}} = \underset{\beta \in \mathcal{R}^{p-1}}{\operatorname{argmin}} (Y - (X \setminus \{X_i\})\beta)^2$ and $\widehat{\epsilon}_{x \setminus \{x_i\}} = \underset{\beta \in \mathcal{R}^{p-1}}{\operatorname{argmin}} (X_i - (X \setminus \{X_i\})\beta)^2$ are the vector of residuals from the regression of Y on X excluding the variable X_i and the vector of residuals from the regression of X_i on X excluding X_i , respectively. This follows directly from the FWL theorem as discussed in section A.1.

$$\widehat{\beta}_{x_i} \xrightarrow{p} \frac{E[\widehat{\epsilon}_{y \setminus \{x_i\}} \widehat{\epsilon}_{x \setminus \{x_i\}}]}{E[\widehat{\epsilon}_{x \setminus \{x_i\}}^2]} \quad (76)$$

$$= \frac{\rho_{Y, X_i | X \setminus \{X_i\}} \sqrt{E[\widehat{\epsilon}_{y \setminus \{x_i\}}^2]}}{\sqrt{E[\widehat{\epsilon}_{x \setminus \{x_i\}}^2]}} \quad (77)$$

$$= \frac{\rho_{Y, X_i | X \setminus \{X_i\}} (1 - \mathcal{R}_{Y | X \setminus \{X_i\}}^2)^{\frac{1}{2}} \sigma_y}{(1 - \mathcal{R}_{X_i | X \setminus \{X_i\}}^2)^{\frac{1}{2}} \sigma_{x_i}} \quad (78)$$

So here we see that in the probability limit we can express any regression coefficient estimator in terms of population partial correlations and limits of coefficients of determinations. The sample regression coefficient can be expressed in terms of the appropriate sample statistics.

Similarly, if we have a data generating process that can be expressed as a linear function in parameters over the parameter space $\beta \in \mathcal{R}^p$, then we can similarly express the true regression coefficients in terms of population partial correlations and coefficients of determination. Let

$$Y = X_1 \beta_1 + X_2 \beta_2 + \epsilon \quad (79)$$

be the true data generating process. As in the text we assume that $E[\epsilon | X] = 0$. Here we suppose that $[X_1, X_2]$ is an arbitrary partition of $X_{n \times p}$ such that X_1 is a vector and X_2 is $(p-1) \times n$ dimensional.

$$\beta_1 = \frac{\rho_{Y, X_1 | X_2} (1 - \mathcal{R}_{Y | X_2}^2)^{\frac{1}{2}} \sigma_y}{(1 - \mathcal{R}_{X_1 | X_2}^2)^{\frac{1}{2}} \sigma_{x_1}} \quad (80)$$

Notice that this does not depend in any way on the true data process for X_1 .

If we think about X_1 as our treatment of interest, then we can see, using the tools from above, that we can express the treatment effect estimator $\widehat{\beta}_{x_1}$ directly in terms of partial correlation and coefficients of determinations. This could be further decomposed into a function of only pairwise correlations and variances, since the coefficients of determination can themselves be decomposed into further functions of correlations and variances. This tells us that a known covariance matrix is sufficient to completely determine any OLS estimator. As we say in section 3, we were able to express the relative bias of the naive and adjusted estimators entirely in terms of correlations and variances. This can be extremely useful for sensitivity analysis or for simulating estimators. As mentioned in the section above with a real data set, much of the matrix can be estimated and the structure

of admissible covariance and correlation matrices is well known and there exist many algorithms to simulate them and even priors such as the LKJ prior which can be placed over them. Correlation matrices especially may be easy to elicit domain knowledge which may further inform or restrict potential simulation studies or sensitivity analysis.

If the true outcome generating process is conditionally linear, we can similarly decompose and understand the true regression coefficient and this is independent of any assumptions on the treatment assignment mechanism and does not require any assumptions on independence between the variables. Of course, if the DAG or partial DAG is known and it implies independencies or conditional independencies, this can help fill unknown partial correlations, correlations, or coefficients of determination. The ability to be able to easily convert correlations to partial correlations is particularly useful with DAGs which imply conditional independence, where there may be obvious restrictions on partial correlations which are harder to discern in terms of pairwise correlations. Here, the way that we model variables other than the outcome may be more important however, as partial correlations are defined in terms of linear projections which is only equivalent to conditioning on variables when the conditional expectations are linear, a special case of which is when the variables are multivariable normal. Here, careful asymptotic arguments may still be used at times to simplify the partial correlation structures in some non-linear or non-multivariate normal cases.

A.7 Probability limit calculations

Now we will show the generality of probability limits for generating meaningful expressions of estimator behavior and again we will use $\widehat{\beta}_a^{bav}$ from equation (9), where $\mathbf{Z} = [\mathbf{I}, \mathbf{BAV}]$.

From the FWL theorem, $\widehat{\beta}_a^{bav} = \frac{\mathbf{A}^T \mathbf{M}_z \mathbf{Y}}{\mathbf{A}^T \mathbf{M}_z \mathbf{A}}$.

$$\widehat{\beta}_a^{bav} = \frac{\mathbf{A}^T \mathbf{M}_z \mathbf{Y}}{\mathbf{A}^T \mathbf{M}_z \mathbf{A}} \quad (81)$$

$$= \frac{\mathbf{A}^T \mathbf{M}_z (\alpha_y + \mathbf{A} \beta_a + \mathbf{U} \beta_u + \mathbf{BAV} \beta_{bav} + \epsilon_1)}{\mathbf{A}^T \mathbf{M}_z \mathbf{A}} \quad (82)$$

$$= \beta_a + \frac{\mathbf{A}^T \mathbf{M}_z (\mathbf{U} \beta_u + \epsilon_1)}{\mathbf{A}^T \mathbf{M}_z \mathbf{A}} \quad (83)$$

$$= \beta_a + \beta_u \frac{\mathbf{A}^T \mathbf{M}_z \mathbf{U}}{\mathbf{A}^T \mathbf{M}_z \mathbf{A}} + \frac{\mathbf{A}^T \mathbf{M}_z \epsilon_1}{\mathbf{A}^T \mathbf{M}_z \mathbf{A}} \quad (84)$$

$$= \beta_a + \beta_u \frac{\frac{1}{n} \mathbf{A}^T \mathbf{M}_z \mathbf{U}}{\frac{1}{n} \mathbf{A}^T \mathbf{M}_z \mathbf{A}} + \frac{\frac{1}{n} \mathbf{A}^T \mathbf{M}_z \epsilon_1}{\frac{1}{n} \mathbf{A}^T \mathbf{M}_z \mathbf{A}} \quad (85)$$

which follows by simply substituting in the true structural equations for \mathbf{Y} and \mathbf{A} , 4 and 5, respectively, and then applying the annihilating properties of \mathbf{M}_z to set linear combinations of constants and \mathbf{BAV} to 0. Notice that we have not used any information about structural equation for the treatment. Thus the numerical form in equation (85) holds for any treatment structural equation, $\mathbf{A} = f(\mathbf{U}, \mathbf{BAV}, \epsilon_2)$. Further, since this is written in general matrix notation, \mathbf{U} and \mathbf{BAV} can be extended to be arbitrary p_1 and p_2 dimensional variables.

We can now solve for the probability limits of the three remaining expressions separately ($\frac{1}{n} \mathbf{A}^T \mathbf{M}_z \mathbf{A}$, $\frac{1}{n} \mathbf{A}^T \mathbf{M}_z \mathbf{U}$, and $\frac{1}{n} \mathbf{A}^T \mathbf{M}_z \epsilon_1$) and combine them due to the properties of probabilities limits, namely equations (35), (36), and (37). We begin with deriving $\frac{1}{n} \mathbf{A}^T \mathbf{M}_z \epsilon_1$,

$$\frac{1}{n} \mathbf{A}^T \mathbf{M}_z \epsilon_1 = \xrightarrow{p} E[\mathbf{A}^T \mathbf{M}_z \epsilon_1] \quad (86)$$

$$= E[\mathbf{A}^T \mathbf{M}_z E[\epsilon_1 | \mathbf{A}, \mathbf{Z}]] \quad (87)$$

$$= E[\mathbf{A}^T \mathbf{M}_z E[\epsilon_2]] \quad (88)$$

$$= 0 \quad (89)$$

where the first line follows from the WLLN, line two from the Law of Iterated Expectations, and the third from the independence of ϵ_2 from \mathbf{A} and \mathbf{BAV} . Next we consider $\frac{1}{n} \mathbf{A}^T \mathbf{M}_z \mathbf{A}$

$$\frac{1}{n} \mathbf{A}^T \mathbf{M}_z \mathbf{A} = \frac{1}{n} (\alpha_a + \mathbf{U}\gamma_u + \mathbf{BAV}\gamma_{bav} + \epsilon_2)^T \mathbf{M}_z (\alpha_a + \mathbf{U}\gamma_u + \mathbf{BAV}\gamma_{bav} + \epsilon_2) \quad (90)$$

$$= \frac{1}{n} (\mathbf{U}\gamma_u + \epsilon_2)^T \mathbf{M}_z (\mathbf{U}\gamma_u + \epsilon_2) \quad (91)$$

$$\xrightarrow{p} \text{plim} \frac{1}{n} (\mathbf{U}\gamma_u + \epsilon_2)^T \mathbf{M}_I (\mathbf{U}\gamma_u + \epsilon_2) \quad (92)$$

$$= \text{plim} \frac{1}{n} (\mathbf{U}\gamma_u)^T \mathbf{M}_I (\mathbf{U}\gamma_u) + \text{plim} \frac{1}{n} 2(\mathbf{U}\gamma_u)^T \mathbf{M}_I \epsilon_2 + \text{plim} \frac{1}{n} \epsilon_2^T \mathbf{M}_I \epsilon_2 \quad (93)$$

$$= \gamma_u^2 E[(\mathbf{U} - \bar{\mathbf{U}})^T (\mathbf{U} - \bar{\mathbf{U}})] + 2\gamma_u E[\mathbf{U} E[\epsilon_2 | \mathbf{U}]] + E[\epsilon_2^T \epsilon_2] \quad (94)$$

$$= \gamma_u^2 \sigma_u^2 + 0 + \sigma_{\epsilon_2}^2 \quad (95)$$

$$= \sigma_a^2 - \gamma_{bav}^2 \sigma_{bav}^2 \quad (96)$$

where line 92 follows from the fact that \mathbf{BAV} is independent of both \mathbf{U} and ϵ_2 . Since \mathbf{M}_z is a residual making vector, we can compare the residuals in the probability limit from the following two regressions

$$(\mathbf{U}\gamma_u + \epsilon_2) = \alpha_{\mathbf{U}\gamma_u + \epsilon_2} + \mathbf{BAV}\eta_{bav} + \mathbf{v}_I \quad (97)$$

$$(\mathbf{U}\gamma_u + \epsilon_2) = \alpha_{\mathbf{U}\gamma_u + \epsilon_2} + \mathbf{v}_2 \quad (98)$$

Due to independence, $\widehat{\eta}_{bav} \xrightarrow{p} 0$ and thus the residuals from the two regressions will be equivalent asymptotically. Thus we can replace \mathbf{M}_z with \mathbf{M}_I in equation (92), which as the centering projection matrix enjoys favorable properties as discussed in section A.1.

Finally we need to find the probability limit of $\frac{1}{n} \mathbf{A}^T \mathbf{M}_z \mathbf{U}$.

$$\frac{1}{n} \mathbf{A}^T \mathbf{M}_z \mathbf{U} = \frac{1}{n} (\alpha_a + \mathbf{U}\gamma_u + \mathbf{BAV}\gamma_{bav} + \epsilon_2)^T \mathbf{M}_z \mathbf{U} \quad (99)$$

$$= \frac{1}{n} (\mathbf{U}\gamma_u + \epsilon_2)^T \mathbf{M}_z \mathbf{U} \quad (100)$$

$$\xrightarrow{p} \frac{1}{n} (\mathbf{U}\gamma_u + \epsilon_2)^T \mathbf{M}_I \mathbf{U} \quad (101)$$

$$= \frac{1}{n} (\gamma_u (\mathbf{U}^T \mathbf{M}_I \mathbf{U}) + \gamma_u \mathbf{U} \mathbf{M}_I \epsilon_2) \quad (102)$$

$$\xrightarrow{p} \gamma_u E[(\mathbf{U} - \bar{\mathbf{U}})^2] + E[(\mathbf{U} - \bar{\mathbf{U}})(\epsilon_2 - \bar{\epsilon}_2)] \quad (103)$$

$$= \gamma_u \sigma_u^2 \quad (104)$$

Putting this altogether, this implies

$$\widehat{\beta}_a^{bav} \xrightarrow{p} \beta_u \frac{\gamma_u \sigma_u^2}{\sigma_a^2 - \gamma_{bav}^2 \sigma_{bav}^2} \quad (105)$$

This is equivalent to the expectation in this case. The benefit is that it is more robust to functional form assumptions and by using properties (35) to (37) and the FWL theorem, we can find asymptotic bias expressions by partitioning the estimator into a series of functions of residuals from simpler regressions. Further, we can always find the limiting expression for the numerator and the denominator separately. Expectations cannot be split up in such a manner and ratios of variables can be very difficult to find closed-form expressions for the expectation without imposing restrictive assumptions.

A.8 Probability limit derivations for section 4

Derivation of bias amplification limits of estimators in equation (21) under the structural equations in equations (19) and (20). Here we assume without loss of generality that $E[f_1(U)] = 0$ and $E[f_2(\mathbf{BAV})] = 0$, since the intercept can absorb the means of such function if they are non-zero.

$$\widehat{\beta}_a^{naive} = \frac{\mathbf{A}^T \mathbf{M}_1 \mathbf{Y}}{\mathbf{A}^T \mathbf{1}} \quad (106)$$

$$= \frac{\mathbf{A}^T \mathbf{M}_1 (\alpha_y + \mathbf{A} \beta_a + f_1(U) + f_2(\mathbf{BAV}) + \epsilon_1)}{\mathbf{A}^T \mathbf{1}} \quad (107)$$

$$= \beta_a + \frac{\mathbf{A}^T \mathbf{M}_1 f_1(U)}{\mathbf{A}^T \mathbf{1}} + \frac{\mathbf{A}^T \mathbf{M}_1 f_2(\mathbf{BAV})}{\mathbf{A}^T \mathbf{1}} + \frac{\mathbf{A}^T \mathbf{M}_1 \epsilon_1}{\mathbf{A}^T \mathbf{1}} \quad (108)$$

$$\xrightarrow{p} \beta_a + \frac{COV(\mathbf{A}, f_1(U))}{\sigma_a^2} + \frac{COV(\mathbf{A}, f_2(\mathbf{BAV}))}{\sigma_a^2} + \frac{E[\mathbf{A} \epsilon_1]}{\sigma_a^2} \quad (109)$$

$$= \beta_a + \frac{COV(\mathbf{A}, f_1(U))}{\sigma_a^2} + \frac{COV(\mathbf{A}, f_2(\mathbf{BAV}))}{\sigma_a^2} + \frac{E[\mathbf{A} E[\epsilon_1 | \mathbf{A}, U, \mathbf{BAV}]]}{\sigma_a^2} \quad (110)$$

$$= \beta_a + \frac{COV(\mathbf{A}, f_1(U))}{\sigma_a^2} + \frac{COV(\mathbf{A}, f_2(\mathbf{BAV}))}{\sigma_a^2} \quad (111)$$

Since for example

$$\frac{\mathbf{A}^T \mathbf{M}_1 f_1(U)}{\mathbf{A}^T \mathbf{1}} = \frac{\frac{1}{n} \sum_{i=1}^n (A_i - \bar{A})(f_1(U)_i - f_1(\bar{U}))}{\frac{1}{n} \sum_{i=1}^n (A_i - \bar{A})^2} \quad (112)$$

$$\xrightarrow{p} \frac{E[(\mathbf{A} - E[\mathbf{A}])(f_1(U) - E[f_1(U)])]}{E[(\mathbf{A} - E[\mathbf{A}])^2]} \quad (113)$$

$$= \frac{COV(\mathbf{A}, f_1(U))}{VAR(\mathbf{A})} \quad (114)$$

Where the probability limit follows from the weak law of large numbers and the continuous mapping theorem. For the adjusted estimator $\widehat{\beta}^{[f_2(\mathbf{BAV})]}$, let $\mathbf{Z} = [\mathbf{1}, f_2(\mathbf{BAV})]$.

$$\widehat{\beta}^{f_2(BAV)} = \frac{\mathbf{A}^T \mathbf{M}_Z \mathbf{Y}}{\mathbf{A}^T \mathbf{M}_Z \mathbf{A}} \quad (115)$$

$$= \frac{\mathbf{A}^T \mathbf{M}_Z (\alpha_y + \mathbf{A} \beta_a + f_1(\mathbf{U}) + f_2(\mathbf{BAV}) + \epsilon_1)}{\mathbf{A}^T \mathbf{M}_Z \mathbf{A}} \quad (116)$$

$$= \beta_a + \frac{\mathbf{A}^T \mathbf{M}_Z f_1(\mathbf{U})}{\mathbf{A}^T \mathbf{M}_Z \mathbf{A}} + \frac{\mathbf{A}^T \mathbf{M}_Z \epsilon_1}{\mathbf{A}^T \mathbf{M}_Z \mathbf{A}} \quad (117)$$

$$= \beta_a + \frac{\mathbf{A}^T \mathbf{M}_Z (f_1(\mathbf{U})^T \mathbf{M}_Z)^T}{\mathbf{A}^T \mathbf{M}_Z \mathbf{A}} + \frac{\mathbf{A}^T \mathbf{M}_Z \epsilon_1}{\mathbf{A}^T \mathbf{M}_Z \mathbf{A}} \quad (118)$$

By idempotence of \mathbf{M}_Z , $\mathbf{A}^T \mathbf{M}_Z$ is the residual from the regression of the treatment on \mathbf{Z} and likewise $(f_1(\mathbf{U})^T \mathbf{M}_Z)^T$ is the transpose of the residuals from the regression of $f_1(\mathbf{U})$ on \mathbf{Z} . To find the residuals, it suffices to find the coefficients from the above regressions.

Let $\mathbf{A} = \mathbf{1}\alpha_1 + f_2(\mathbf{BAV})\zeta_{f_2(bav)_1} + v_1$ and $f_1(\mathbf{U}) = \mathbf{1}\alpha_2 + f_2(\mathbf{BAV})\zeta_{f_2(bav)_2}$ represent the regression equations. Following the FWL theorem once more, we have the following estimators

$$\widehat{\alpha}_1 = \frac{\mathbf{1}^T \mathbf{M}_{f_2(\mathbf{BAV})} \mathbf{A}}{\mathbf{1}^T \mathbf{M}_{f_2(\mathbf{BAV})} \mathbf{1}} \xrightarrow{p} E[\mathbf{A}] - \frac{COV(\mathbf{A}, f_2(\mathbf{BAV}))}{Var(f_2(\mathbf{BAV}))} E[f_2(\mathbf{BAV})] = E[\mathbf{A}] \quad (119)$$

$$\widehat{\zeta_{f_2(bav)_1}} = \frac{f_2(\mathbf{BAV})^T \mathbf{M}_1 \mathbf{A}}{f_2(\mathbf{BAV})^T \mathbf{M}_1 f_2(\mathbf{BAV})} \xrightarrow{p} \frac{COV(\mathbf{A}, f_2(\mathbf{BAV}))}{Var(f_2(\mathbf{BAV}))} \quad (120)$$

$$\widehat{\alpha}_2 = \frac{\mathbf{1}^T \mathbf{M}_{f_2(\mathbf{BAV})} f_1(\mathbf{U})}{\mathbf{1}^T \mathbf{M}_{f_2(\mathbf{BAV})} \mathbf{1}} \xrightarrow{p} E[f_1(\mathbf{U})] - \frac{COV(f_1(\mathbf{U}), f_2(\mathbf{BAV}))}{Var(f_2(\mathbf{BAV}))} E[f_2(\mathbf{BAV})] = 0 \quad (121)$$

$$\widehat{\zeta_{f_2(bav)_2}} = \frac{f_2(\mathbf{BAV}) \mathbf{M}_1 f_1(\mathbf{U})}{f_2(\mathbf{BAV}) \mathbf{M}_1 f_2(\mathbf{BAV})} \xrightarrow{p} \frac{COV(f_2(\mathbf{BAV}), f_1(\mathbf{U}))}{Var(f_2(\mathbf{BAV}))} \quad (122)$$

Using the probability limits of the estimates, the continuous mapping theorem and the fact that $\mathbf{A}^T \mathbf{M}_Z \mathbf{A} \xrightarrow{p} (1 - \mathcal{R}_{A|Z}^2) \sigma_a^2$ where $\mathcal{R}_{A|Z}^2$ is understood to be the probability limit of the \mathcal{R}^2 from the regression of the treatment on the columns of \mathbf{Z} , we get

$$\widehat{\beta}^{f_2(BAV)} \xrightarrow{p} \frac{COV(\mathbf{A}^T \mathbf{M}_Z, (f_1(\mathbf{U})^T \mathbf{M}_Z)}{\sigma_a^2 (1 - \mathcal{R}_{A|Z}^2)} \quad (123)$$

$$= \beta_a + \frac{COV(\mathbf{A}, f_1(\mathbf{U})) - \frac{COV(\mathbf{A}, f_2(\mathbf{BAV})) COV(f_1(\mathbf{U}), f_2(\mathbf{BAV}))}{Var(f_2(\mathbf{BAV}))}}{\sigma_a^2 (1 - \mathcal{R}_{A|Z}^2)} \quad (124)$$

and the error term drops out since $E[\epsilon_1 | \mathbf{A}, \mathbf{1}, f_2(\mathbf{BAV})] = 0$.

A.8.1 Orthogonal semi-parametric models derivation

From the standard OLS formula we get

$$\widehat{\beta}^{semi} = \frac{\frac{1}{n} (\mathbf{A} - \widehat{E}[\mathbf{A} | \mathbf{BAV}])^T (\mathbf{Y} - \widehat{E}[\mathbf{Y} | \mathbf{BAV}])}{\frac{1}{n} (\mathbf{A} - \widehat{E}[\mathbf{A} | \mathbf{BAV}])^T (\mathbf{A} - \widehat{E}[\mathbf{A} | \mathbf{BAV}])} \quad (125)$$

$$\xrightarrow{p} \frac{E[(A - E[A|BAV])^T(Y - E[Y|BAV])]}{E[(A - E[A|BAV])^T(A - E[A|BAV])]} \quad (126)$$

$$= \frac{E[E[(A - E[A|BAV])^T(Y - E[Y|BAV])|BAV]]}{E[E[(A - E[A|BAV])^T(A - E[A|BAV])|BAV]]} \quad (127)$$

$$= \frac{E[COV(A, Y|BAV)]}{E(Var(A|BAV))} \quad (128)$$

$$= \frac{E[COV(A, Y|BAV)]}{\sigma_a^2 - var(E[A|BAV])} \quad (129)$$

This holds under any DAG and any set of structural equations. However, as discussed in section 4, without making restrictions on the underlying structural equations, the estimator may not be meaningful or easily comparable to a causal effect of interest. Under the structural equations in equations (19) and (20), the estimator probability limit is

$$\widehat{\beta}^{semi} \xrightarrow{p} \frac{E[E[(A - E[A|BAV])^T(Y - E[Y|BAV])|BAV]]}{E[E[(A - E[A|BAV])^T(A - E[A|BAV])|BAV]]} \quad (130)$$

$$= \frac{\left\{ \begin{array}{l} E[E[(A - E[A|BAV])^T((\alpha_y + A\beta_a + f_1(U) + f_2(BAV) + \epsilon_T))] \\ - E[E[(\alpha_y + A\beta_a + f_1(U) + f_2(BAV) + \epsilon_T)|BAV]|BAV] \end{array} \right\}}{E[E[(A - E[A|BAV])^T(A - E[A|BAV])|BAV]]} \quad (131)$$

$$= \beta_a + \frac{E[E[(A - E[A|BAV])^T(f_1(U) - E[f_1(U)|BAV])]}{\sigma_a^2 - var(E[A|BAV])} \quad (132)$$

$$= \beta_a + \frac{COV(A, f_1(U)|BAV)}{\sigma_a^2 - var(E[A|BAV])} \quad (133)$$

In the special case that U and BAV are independent, as they would be under the DAG in Figure 3, this reduces to

$$\widehat{\beta}^{semi} \xrightarrow{p} \beta_a + \frac{COV(A, f_1(U))}{\sigma_a^2 - var(E[A|BAV])} \quad (134)$$

A.9 Bias amplification under misspecification and unrestricted structural equations

A.9.1 Bias Amplification and misspecification of the control function

In section 4 of the main text, we considered the problem of bias amplification when $f_2(BAV)$ was specified correctly at least in the limit $n \rightarrow \infty$. Here we consider the additional possibility of misspecification. Suppose the structural equations were again those in equations (19) and (20) and instead of including the function $f_2(BAV)$ we include some other control function $h(BAV)$. For clarity, this control function could be just the linear term BAV or a polynomial function for example, but in general it might be anything the analyst believes is reasonable, but that is ultimately not quite correct. If $h(BAV)$ is unable to be written as a linear combination of $f_2(BAV)$ (i.e. $\nexists a, b \in \mathcal{R} : h(BAV) = a + bf_2(BAV)$), then there may be misspecification bias in addition to unmeasured confounding.

The naive estimator ($\widehat{\beta}_a^{naive}$) and the estimator conditional on $h(\mathbf{BAV})$ will have the following probability limits.

$$\widehat{\beta}_a^{naive} \xrightarrow{p} \beta_a + (Cov(\mathbf{A}, f_1(\mathbf{U})) + Cov(\mathbf{A}, f_2(\mathbf{BAV}))) \frac{1}{\sigma_a^2} \quad (135)$$

$$\widehat{\beta}^{h(\mathbf{BAV})} \xrightarrow{p} \beta_a + (Cov(\mathbf{A}^T \mathbf{M}_z, (f_1(\mathbf{U})^T \mathbf{M}_z)^T) + Cov(\mathbf{A}^T \mathbf{M}_z, (f_2(\mathbf{BAV})^T \mathbf{M}_z)^T)) \frac{1}{(1 - \mathcal{R}_{\mathbf{A}|h(\mathbf{BAV})}^2) \sigma_a^2} \quad (136)$$

$$= \beta_a + COV(\mathbf{A}, f_1(\mathbf{U})) - \frac{COV(\mathbf{A}, h(\mathbf{BAV})) COV(f_1(\mathbf{U}), h(\mathbf{BAV}))}{Var(h(\mathbf{BAV}))} \quad (137)$$

$$+ COV(\mathbf{A}, f_2(\mathbf{BAV})) - \frac{COV(\mathbf{A}, h(\mathbf{BAV})) COV(h(\mathbf{BAV}), f_2(\mathbf{BAV}))}{Var(h(\mathbf{BAV}))} \frac{1}{1 - \mathcal{R}_{\mathbf{A}|h(\mathbf{BAV})}^2} \quad (138)$$

When we compare the probability limit of the misspecified estimator ($\widehat{\beta}^{h(\mathbf{BAV})}$) in equation (136) and the adjusted estimator in equation (123), there is an additional term due to misspecification in the numerator ($Cov(\mathbf{A}^T \mathbf{M}_z, (f_2(\mathbf{BAV})^T \mathbf{M}_z)^T)$). This does not necessarily mean that it is more biased and will depend on the signs of the covariances between all of the relevant variables. The numerator of the bias term ($Cov(\mathbf{A}^T \mathbf{M}_z, (f_1(\mathbf{U})^T \mathbf{M}_z)^T) + Cov(\mathbf{A}^T \mathbf{M}_z, (f_2(\mathbf{BAV})^T \mathbf{M}_z)^T)$) is still the covariance of the treatment and outcome through the confounding pathways $f_1(\mathbf{U}) \rightarrow \mathbf{A} \rightarrow \mathbf{Y}$ and the part of $f_2(\mathbf{BAV}) \rightarrow \mathbf{A} \rightarrow \mathbf{Y}$ which is not controlled by linear combinations of $h(\mathbf{BAV})$. The denominator is again always greater than 1 and contributing to amplification, but increases hyperbolically with respect to how well $h(\mathbf{BAV})$ linearly explains variance in the treatment rather than $f_2(\mathbf{BAV})$. Notice that even though $h(\mathbf{BAV})$ is misspecified with respect to the true outcome equation (19), this does not mean that the amplification term will be less severe. It is possible that $h(\mathbf{BAV})$ explains more of the treatment covariance since $\mathbf{A} = \alpha_a + g(\mathbf{U}, \mathbf{BAV}) + \epsilon_2$ and $h(\mathbf{BAV})$ may better approximate $g(\mathbf{U}, \mathbf{BAV})$. Overall, misspecification does not qualitatively change the problem of bias amplification, but in the case that the analyst does not know $f_2(\mathbf{BAV})$ it may be more difficult to harness domain knowledge to accurately predict the sign and direction of the bias. The amplifying denominator can still be estimated using observed data and used in the process of model selection or sensitivity analysis.

A.9.2 Bias amplification under unrestricted structural equations, heterogeneous causal effects

In the main text, we restricted our focus to systems of structural equations which were partially linear in the treatment of interest such that the desired causal effect was a single parameter (β_a) in the true model. In practice, it may be that the causal effect of interest is heterogeneous and a more complicated functional and the treatment of this setting is beyond the scope of this manuscript. In the case that we assume that the effect of interest is β_a but it is not, we can say the following about the OLS and partially least squares estimators.

If we make no assumptions about the functional form, and allow for non-linearities, and interactions between all variables, we can show that the OLS adjusted estimator is always the expression below (see Appendix A.1)

$$\widehat{\beta}_a^{|\mathbf{z}} = \frac{(\mathbf{M}_z \mathbf{A})^T \mathbf{M}_z \mathbf{Y}}{(\mathbf{M}_z \mathbf{A})^T (\mathbf{M}_z \mathbf{A})} \quad (139)$$

When \mathbf{Z} includes an intercept column, both $(\mathbf{M}_z \mathbf{A})^T$ and $\mathbf{M}_z \mathbf{Y}$ will have mean zero and thus we can think of the numerator as an empirical estimate of the covariance between the residuals from the regression of \mathbf{A} on \mathbf{Z} and the residuals from the regression of \mathbf{Y} on \mathbf{Z} . Unmeasured confounding bias in OLS occurs when after projecting out linear combinations of the controlling variables, \mathbf{Z} , there remain linear associations between the outcome and the treatment due to unobserved variables. The part of the bias due to unmeasured confounding is amplified whenever the control variables explain variance in the treatment. Holding all else constant, as the residuals from the regression of the treatment on \mathbf{Z} decrease in magnitude, the absolute value of the estimator $\widehat{\beta}_a^{|\mathbf{z}}$ will increase in magnitude. This is a general form of the result we showed in the previous section which is extremely powerful in that it encaptures a very large class of structural equations and DAGs. However, the cost of this generality is that

without making more specific assumptions about the particular form of the model, and in particular the outcome model, it becomes more difficult to incorporate the knowledge of the amplification factor into our model selection and thus a priori know which of the two estimators, $\widehat{\beta}_a^{naive}$ or $\widehat{\beta}_a^{bav}$, will be less biased, especially when the nonlinearities imply heterogenous causal effects and thus that β_a itself is not the causal effect of interest. In particular, we cannot easily decompose the bias or relative bias into covariances through the confounding pathways and treatment variance. The amplifying denominator remains the same, but the numerator cannot be easily decomposed. Interaction terms, for example, are an additional difficulty. Pearl,¹⁰ for example, showed that under a simple interaction effect between the unmeasured confounding and some function of a pure instrument, the adjusted estimator can be less biased than the naive case. To properly evaluate estimators in the context of bias amplification requires appropriate simulations.

Similarly, we can show that the partially linear estimator conditional on BAV under any structural equations is

$$\widehat{\beta}_{partial}^{BAV} = \frac{(A - \widehat{E}[A|BAV])^T (Y - \widehat{E}[Y|BAV])}{(A - \widehat{E}[A|BAV])^T (A - \widehat{E}[A|BAV])} \quad (140)$$

$$\xrightarrow{p} \frac{E[COV(A, Y|BAV)]}{(1 - \frac{VAR(E[A|BAV])}{\sigma_a^2})\sigma_a^2} \quad (141)$$

Much like the OLS estimator, the numerator of the partially linear estimator is an expectation of the covariance of the treatment and outcome once we have controlled for BAV . The denominator, just like in the main text is proportional to the remaining treatment variance once we have controlled for fluctuations in BAV . As BAV explains more of the treatment variation, the magnitude of the estimator increases hyperbolically. Again, however, without a priori knowledge of the structural form for the outcome and thus the form of the causal effect of interest, it is not clear if this will be closer or farther away from desired estimand. As an estimator, however, we can see that both of these orthogonalization estimators, OLS and partially least squares, have the property that explaining variance in the treatment increases the magnitude of the estimate hyperbolically.

A.10 Additional details simulation

Suppose we want to simulate a system of linear equations from equations (4) and (5) based on the DAG in Figure 3. Now suppose we are interested in assessing the effect of modifying the edge $BAV \rightarrow A$ on the conditional estimator $\widehat{\beta}_a^{bav}$. If we incorrectly run this simulation simply by changing the parameter γ_{bav} to some (or some set of) γ_{bav}' and fail to fix the variance of the treatment A as discussed in section 5, we can show that the bias of estimator $\widehat{\beta}_a^{bav}$ will remain unchanged. In section A.14 we show that the variance of A in the above simulation design is equal to

$$\sigma_a^2 = \gamma_u^2 \sigma_u^2 + \gamma_{bav}^2 \sigma_{bav}^2 + \sigma_{\epsilon_2}$$

Further we showed that the expectation and probability limit of the estimator is

$$E[\widehat{\beta}_a^{bav}] = \beta_a + \beta_u \frac{\gamma_u \sigma_u^2}{\sigma_a^2 - \gamma_{bav}^2 \sigma_{bav}^2}$$

Thus, if we change γ_{bav} to γ_{bav}' holding all other parameters constant, it can be shown that the resulting expectation is unchanged. This is because the increased amplification is precisely cancelled out by increasing the variance of the treatment.

$$\begin{aligned} \sigma_a'^2 &= \gamma_u^2 \sigma_u^2 + \gamma_{bav}'^2 \sigma_{bav}^2 + \sigma_{\epsilon_2} \\ &= \sigma_a^2 - \gamma_{bav}^2 \sigma_{bav}^2 + \gamma_{bav}'^2 \sigma_{bav}^2 \end{aligned}$$

$$= \sigma_a^2 + (\gamma_{bav}'^2 - \gamma_{bav}^2)\sigma_{bav}^2$$

This implies the expectation of the estimator $\widehat{\beta}_a^{bav'}$ has the following expression

$$\begin{aligned} E[\widehat{\beta}_a^{bav'}] &= \beta_a + \beta_u \frac{\gamma_u \sigma_u^2}{\sigma_a^2 - \gamma_{bav}'^2 \sigma_{bav}^2} \\ &= \beta_a + \beta_u \frac{\gamma_u \sigma_u^2}{(\sigma_a^2 + (\gamma_{bav}'^2 - \gamma_{bav}^2)\sigma_{bav}^2) - \gamma_{bav}'^2 \sigma_{bav}^2} \\ &= \beta_a + \beta_u \frac{\gamma_u \sigma_u^2}{\sigma_a^2 - \gamma_{bav}^2 \sigma_{bav}^2} \\ &= E[\widehat{\beta}_a^{bav}] \end{aligned}$$

Thus the expectation of the estimator remains unchanged for any change of parameter, γ_{bav} . As a consequence, the comparison of this estimator with the naive estimator will seem favorable as the absolute magnitude of the parameter γ_{bav}' increases since the difference in absolute bias is

$$E[\widehat{\beta}_a^{naive}] - E[\widehat{\beta}_a^{bav'}] = \left| \frac{\beta_u \gamma_u \sigma_u^2}{\sigma_a^2} + \frac{\beta_{bav} \gamma_{bav} \sigma_{bav}^2}{\sigma_a^2} \right| - \left| \beta_u \frac{\gamma_u \sigma_u^2}{\sigma_a^2 - \gamma_{bav}^2 \sigma_{bav}^2} \right|$$

The bias of $\widehat{\beta}_a^{naive}$ is increasing in γ_{bav} for sufficiently large γ_{bav} and we showed above that the bias for $\widehat{\beta}_a^{bav'}$ is invariant to changes in γ_{bav} if we do not fix the variance of the treatment \mathbf{A} . Thus eventually the bias of the naive estimate is strictly increasing in γ_{bav} and will continue to appear worse relative to the conditional estimator. However, as discussed in section 5, this is a consequence of failing to conduct a proper causal simulation experiment comparing data sets plausibly generated from similar experiments and holding all other potentially confounding edges constant.

A.11 Real data simulation details

The goal of this section is to utilize the real randomized control trial data described in section 6, which comes from the DAG in Figure 7(a), and simulate modified covariates ($\tilde{\mathbf{X}}$) and a modified outcome (\tilde{Y}) such that they come from the DAG 7b and the equations (22), (23), and (24). This proceeds in two steps. First we need to simulate the latent variable \mathbf{A}^* and then conditionally simulate \mathbf{BAV} and \mathbf{U} .

For simplicity, we set \mathbf{U} and \mathbf{BAV} to be standard normal variables. The latent variable, \mathbf{A}^* , has variance of 1 but its mean, α_a is determined to ensure that $P(\mathbf{A}^* > 0) = P(\mathbf{U}\gamma_u + \mathbf{BAV}\gamma_{bav} + \epsilon_2 > \alpha_a) = p_a$, which is determined in our data by matching p_a to the observed quantity $\widehat{p}_a = \frac{1}{n} \sum_{i=1}^n A_i \approx 0.51$.

Under the assumptions above, this implies that $\alpha_a = -\Phi^{-1}(1 - p_a)$, where $\Phi(x)$ is the cumulative distribution function (CDF) of the standard normal distribution. As previously mentioned, the variance of the error term ϵ_2 is set precisely to ensure that variance of \mathbf{A}^* is 1, $\sigma_{\epsilon_2}^2 = 1 - \gamma_u^2 - \gamma_{bav}^2$.

The first step is to use the observed \mathbf{A} data to simulate the latent \mathbf{A}^* . Consider the CDF of the latent \mathbf{A}^* conditional on $\mathbf{A} = 1$. Let $P(\mathbf{A} = 1) = p_a$, so that

$$\begin{aligned} P(\mathbf{A}^* \leq a^* | \mathbf{A} = 1) &= P(\mathbf{A}^* \leq a^* | \mathbf{A}^* \geq 0) \\ &= \frac{P(\mathbf{A}^* \leq a^*, \mathbf{A}^* \geq 0)}{P(\mathbf{A}^* \geq 0)} \\ &= \frac{P(0 \leq \mathbf{A}^* \leq a^*)}{P(\mathbf{A}^* \geq 0)} \end{aligned}$$

$$\begin{aligned}
& P\left(\frac{-E[A^*]}{\sigma_{a^*}^2} \leq \frac{A^* - E[A^*]}{\sigma_{a^*}^2} \leq \frac{a^* - E[A^*]}{\sigma_{a^*}^2}\right) \\
&= \frac{P(A^* \geq 0)}{P(A^* \geq 0)} \\
&= \frac{\Phi\left(\frac{a^* - E[A^*]}{\sigma_{a^*}^2}\right) - \Phi\left(\frac{-E[A^*]}{\sigma_{a^*}^2}\right)}{p_a} \\
&= \frac{\Phi(a^* - \alpha_a) - \Phi(-\alpha_a)}{p_a} \\
&= \frac{\Phi(a^* - \alpha_a) - \Phi(\Phi^{-1}(1 - p_a))}{p_a} \\
&= \frac{\Phi(a^* - \alpha_a) - (1 - p_a)}{p_a}
\end{aligned}$$

Since $P(A^* \leq a^* | A = 1)$ is the CDF of a continuous random variable, it is distributed uniformly between 0 and 1. Let $X \sim U(0, 1)$ be a uniform random variable with support $[0, 1]$. Appealing to the probability inverse transform, X and $P(A^* \leq a^* | A = 1)$ are equivalent in distribution.

$$\begin{aligned}
\Rightarrow X &\equiv \frac{\Phi(a^* - \alpha_a) - (1 - p_a)}{p_a} \\
\Rightarrow p_a X + (1 - p_a) &\equiv \Phi(a^* - \alpha_a) \\
\Rightarrow \Phi^{-1}(p_a X + (1 - p_a)) + \alpha_a &\equiv a^*
\end{aligned}$$

Similarly, it can be shown that conditional on $A = 0$

$$a^* \equiv \Phi^{-1}((1 - p_a)X) + \alpha_a$$

Thus, in general

$$a^* \equiv \Phi^{-1}(p_a^A (1 - p_a)^{A-1} X + (1 - p_a)^A) + \alpha_a$$

Therefore, by conditioning on A and simulating a uniform random variable, we can take draws from the unobserved latent variable A^* . Once we have recovered the latent variable, we can jointly simulate U and BAV conditional on A^* and the observed covariates X . The observed covariates are centered and asymptotically multivariate normal. Since U , BAV , and A^* are univariate or multivariate normal variables and X are asymptotically normal, the conditional distribution will be asymptotically multivariate normal and proportional to the joint density. From standard multivariate normal theory

$$U = u, BAV = bav | A^* = a^*, X = x \sim N(\boldsymbol{\mu}^*, \boldsymbol{\Sigma})$$

As stated above, the conditional distribution is proportional to the joint model. Thus, we will define $\boldsymbol{\Sigma}$ and $\boldsymbol{\mu}$ for the joint density.

$$U = u, BAV = bav, A^* = a^*, X = x \sim N(\boldsymbol{\mu}, \boldsymbol{\Sigma})$$

$$\boldsymbol{\Sigma} = \begin{bmatrix} \Sigma_{(U, BAV), (U, BAV)} & \Sigma_{(U, BAV), (A^*, X)} \\ \Sigma_{(A^*, X), (U, BAV)} & \Sigma_{(A^*, X), (A^*, X)} \end{bmatrix}$$

$$\Sigma_{(U,BAV),(U,BAV)} = \begin{bmatrix} \sigma_u^2 & 0 & 0 & 0 \\ 0 & \sigma_{bav_1}^2 & 0 & 0 \\ 0 & 0 & \sigma_{bav_2}^2 & 0 \\ 0 & 0 & 0 & \sigma_{bav_3}^2 \end{bmatrix}$$

$$\Sigma_{(U,BAV),(A^*,X)} = \begin{bmatrix} \gamma_u \sigma_u^2 & 0 & 0 & 0 \\ \gamma_{\tilde{x}_1} \sigma_{bav_1}^2 & 0 & 0 & 0 \\ \gamma_{\tilde{x}_2} \sigma_{bav_2}^2 & 0 & 0 & 0 \\ \gamma_{\tilde{x}_3} \sigma_{bav_3}^2 & 0 & 0 & 0 \end{bmatrix}$$

$$\Sigma_{(A^*,X),(U,BAV)} = \Sigma_{(U,BAV),(A^*,X)}^T$$

$$\Sigma_{(A^*,X),(A^*,X)} = \begin{bmatrix} \sigma_{a^*}^2 & 0 & 0 & 0 \\ 0 & \sigma_{x_1}^2 & 0 & 0 \\ 0 & 0 & \sigma_{x_2}^2 & 0 \\ 0 & 0 & 0 & \sigma_{x_3}^2 \end{bmatrix}$$

$$\boldsymbol{\mu} = \begin{bmatrix} \mu_u \\ \mu_{bav_1} \\ \mu_{bav_2} \\ \mu_{bav_3} \\ \mu_{a^*} \\ \mu_{x_1} \\ \mu_{x_2} \\ \mu_{x_3} \end{bmatrix} = \begin{bmatrix} 0 \\ 0 \\ 0 \\ 0 \\ \alpha_a \\ 0 \\ 0 \\ 0 \end{bmatrix}.$$

Using standard multivariate normal theory and the matrices defined above, we can define the mean, $\boldsymbol{\mu}^*$ and variance, Σ^* of the conditional model

$$\Sigma^* = \Sigma_{(U,BAV),(U,BAV)} - \Sigma_{(U,BAV),(A^*,X)} \Sigma_{(A^*,X),(A^*,X)}^{-1} \Sigma_{(A^*,X),(U,BAV)}; \text{ and}$$

$$\boldsymbol{\mu}^* = \Sigma_{(U,BAV),(A^*,X)} \Sigma_{(A^*,X),(A^*,X)}^{-1} ([A^*, X] - [\mu_{a^*}])^T$$

Using the conditional distribution, we can thus take draws of U, BAV conditional on the particular values of A^* and X . Using the simulated BAV we add it to the covariates X to form the modified covariates, $\tilde{X} = \frac{X}{\sigma'} + BAV$, where σ' is a scaling factor chosen simultaneously with σ_{bav} such that the variance of \tilde{X} is precisely equal to $\sigma_x = 1$. This step is important if we would like to compare simulations with the modified and the unmodified covariates. In the particular simulations conducted in section 6, the scaling was chosen such that $Var(\frac{X_i}{\sigma'}) = 0.01, i = 1, 2, 3$ and thus $\sigma_{bav_i}^2 = 0.99, i = 1, 2, 3$.

Now that the modified covariates have been constructed, the modified outcome can be constructed. The original RCT data coming from Figure 7(a) is assumed to come from the linear model

$$Y = \alpha_y + A\beta_a + X\beta_x + \epsilon_1$$

where β_a and β_x are estimated unbiasedly in section 6. Next, we add the unmeasured confounding, $U\beta_u$ directly (where β_u is chosen) and then add $BAV\beta_x + \tilde{X}\beta_{adj}$, where $\beta_{adj} = \beta_{\tilde{x}} - \beta_x$ where $\beta_{\tilde{x}}$ is chosen to set the desired covariance between the modified covariates and the modified outcome

$$\tilde{Y} = Y + U\beta_u + BAV\beta_x + \tilde{X}(\beta_{\tilde{x}} - \beta_x)$$

$$\begin{aligned}
&= (\alpha_y + \mathbf{A}\beta_a + \mathbf{X}\beta_x + \epsilon_I) + \mathbf{U}\beta_u + \mathbf{BAV}\beta_x + \tilde{\mathbf{X}}(\beta_{\tilde{x}} - \beta_x) \\
&= \alpha_y + \mathbf{A}\beta_a + \mathbf{U}\beta_u + (\mathbf{X} + \mathbf{BAV})\beta_x + \tilde{\mathbf{X}}(\beta_{\tilde{x}} - \beta_x) + \epsilon_I \\
&= \alpha_y + \mathbf{A}\beta_a + \mathbf{U}\beta_u + (\tilde{\mathbf{X}})\beta_x + \tilde{\mathbf{X}}(\beta_{\tilde{x}} - \beta_x) + \epsilon_I \\
&= \alpha_y + \mathbf{A}\beta_a + \mathbf{U}\beta_u + \tilde{\mathbf{X}}\beta_{\tilde{x}} + \epsilon_I
\end{aligned}$$

This is precisely the outcome equation in equation (22) in Section 6.2. Thus following this method we can use the real data to create a data simulation using the original treatment data and matching many of the characteristics of the real data, but precisely control the causal structure and correlations between the variables. As with the other simulations, there will still be restrictions on the parameters and correlations that we set such as positive definiteness of all the variance matrices in the above simulation.

A.12 Real data simulation comparison of estimators

Consider a causal simulation experiment coming from a DAG and system of equations identical to the one considered in section 6.2 as described by Figure 7(b) and the system of equations (22), (23), and (24). The experiment uses the real data described in section 6 and the procedure detailed in section A.11. The simulation experiment involves intervening on the edge $\tilde{\mathbf{X}}_I \rightarrow \mathbf{A}$, that is increasing the covariance between $\tilde{\mathbf{X}}_I$ and \mathbf{A} . As in section 5 we will explore the consequences of failing to properly hold all non-intervention edges of the DAG.

Simulation Parameters	$\gamma_{\tilde{\mathbf{X}}}$	γ_u	β_u	$\beta_{\tilde{x}}$	β_a
Control	0.20, 0.38, 0.33	0.63	0.15	0.10, -0.15, -0.10	0.1377
Intervention	0.55, 0.38, 0.33	0.63	0.15	0.10, -0.15, -0.10	0.1377

Above, the parameters for the two simulation treatments are described. The only difference between the two is that in the control, $\gamma_{\tilde{x}_1} = 0.2$ and in the intervention $\gamma_{\tilde{x}_1} = 0.55$. Below we visualize the naive, adjusted, and unbiased estimators for the control treatment.

In the control treatment, we see that the unbiased estimator behaves as expected, centered on the true underlying parameter. The naive estimator $\hat{\beta}_a^{naive}$ is only slightly biased, since some unmeasured biases due to the vector $\tilde{\mathbf{X}}$ and \mathbf{U} happen to be of opposing signs and partially cancel each other out. If this is not the case, of course the naive estimator may be significantly more biased. The adjusted estimator behaves poorly with an average absolute bias of 0.18. Although the parameters in the latent space are relatively large, $\gamma_{\tilde{x}} = [0.2, 0.38, 0.33]$, the covariances in the observed space with respect to the treatment are relatively small, $COV(\mathbf{A}, \tilde{\mathbf{X}}) = [0.08, 0.15, 0.13]$, and yet the amplifying effect is quite large. In fact, the amplifying variables jointly explain only 18% of the variance of the treatment, but since the variance of the treatment was already quite small, $\sigma_a^2 \approx 0.25$, the amplifying variables had a more than proportional effect.

The bias attributed to the path $\mathbf{A} \leftarrow \mathbf{U} \rightarrow \mathbf{Y}$ for the naive estimator is $\frac{\beta_u COV(\mathbf{A}, \mathbf{U})}{\sigma_a^2} = \frac{0.15 \times 0.25}{0.25} = \frac{0.0375}{0.25} = 0.0375 \times 4 = 0.15$, whereas for the amplified estimator it is $\frac{0.15 \times 0.25}{0.25 - (0.08^2 + 0.15^2 + 0.13^2)} = \frac{0.0375}{.183} = 0.0375 \times 4.88 = 0.183$. Since $\mathbf{A}^T \mathbf{M} \mathbf{z} \mathbf{A} \leq \mathbf{A}^T \mathbf{M} \mathbf{A}$, $\forall \mathbf{z} : \mathbf{z} \subseteq \mathbf{Z}$, we can rewrite the bias due to bias amplification as $\frac{|\beta_u \times COV(\mathbf{A}, \mathbf{U})|}{(1-c) \times \sigma_a^2}$, $c \in [0, 1]$, where c is the proportion of treatment variance explained by the bias amplifiers jointly

$$\frac{\partial^2 \left(\frac{|\beta_u \times COV(\mathbf{A}, \mathbf{U})|}{(1-c) \times \sigma_a^2} \right)}{\partial \sigma_a^2 \partial c} = \begin{cases} \frac{-\beta_u \times COV(\mathbf{A}, \mathbf{U})}{(1-c)^2 (\sigma_a^2)^2}, & \beta_u \times COV(\mathbf{A}, \mathbf{U}) > 0 \\ \frac{\beta_u \times COV(\mathbf{A}, \mathbf{U})}{(1-c)^2 (\sigma_a^2)^2}, & \beta_u \times COV(\mathbf{A}, \mathbf{U}) < 0 \end{cases}$$

The derivative above shows us that as the variance, σ_a^2 , gets smaller, the marginal impact on absolute bias from an increase in the proportion of the variance explained by the amplifiers increases.

Now consider the intervention of increasing the proportion of variance explained by one of the potential amplifiers, \tilde{X}_I , by increasing $\gamma_{\tilde{x}_1}$ to 0.55. Again, we will consider the case of keeping all of the variances constant to the case where we simply change the parameter and allow the variances to float.

Notice in the left panel that although we have intentionally increased the amplification, the amplifying estimator has seemingly not changed. However, when we fix the variance, the bias amplification increases as we expected (the mean absolute bias increased from 0.18 to 0.23). The reason for this effect is that as the variance of the latent treatment A^* increases, the covariances of the variables of the unmeasured confounding and the treatment as well as the potential amplifiers, $COV(A, U)$ and $COV(A, \tilde{X})$, decrease. It can be shown that when A^* , \tilde{X} , and U are normal or multivariate normal

$$COV(A, U) = \frac{1}{p_a} E[U] + \frac{\gamma_u \sigma_u^2}{\sqrt{2\pi\sigma_{a^*}^2}} \exp\left(\frac{-\sigma_a^2}{2\sigma_{a^*}^2}\right) \quad (142)$$

$$COV(A, \tilde{X}) = \frac{1}{p_a} E[\tilde{X}] + \frac{\gamma_{\tilde{x}} \sigma_{\tilde{x}}^2}{\sqrt{2\pi\sigma_{a^*}^2}} \exp\left(\frac{-\sigma_a^2}{2\sigma_{a^*}^2}\right). \quad (143)$$

We can see in equation (142) that if we allow the variance of A^* to increase as $\gamma_{\tilde{x}_1}$ increases that $COV(A, U)$ decreases. In the case of this simulation, the covariance decreased from 0.25 in the control treatment to 0.22, since $\sigma_{a^*}^2$ increased from 1 to $1 + (0.55^2 - 0.2^2) = 1.26$. Thus we have decreased the strength of the edge $U \rightarrow A$ incidentally. Further by considering equation (143), we can see that if we increase $\gamma_{\tilde{x}_1}$ we do not necessarily increase the amount of variance explained by \tilde{X} since there are two opposing effects. First, consider the increase directly through $\gamma_{\tilde{x}_1}$ and the decrease through increasing $\sigma_{a^*}^2$. In the particular example, although our intended goal was to observe the effect of increasing the weight of the edge $\tilde{X}_I \rightarrow A$, we have in fact inadvertently decreased the covariance from 0.2 to 0.196.

Again we can see that when we fail to hold the variances constant, we are no longer comparing a controlled intervention on the weight of a particular set of nodes, but have modified the edges into and out of the intervened upon edge. This example shows that this is true in cases beyond fully linear systems of equations explored in section 5. Examining the simulation results we can see that this might lead to inappropriate conclusions about the effects of our interventions and the relative merits of particular estimators in contexts of interest to us.

A.13 Simulation for Figure 4

The figure was simulated from the general structural equations (4) and (5) with the particular values below.

$$Y = 2 + .2 \times A + .5 \times U + .05 \times BAV + v_I$$

$$A = 1 + .3 \times U + .75 \times BAV + v_2$$

$$BAV \sim N(0, 1), \quad U \sim N(0, 1)$$

$$v_I \sim N(0, \sigma_{v_1})$$

$$\sigma_{v_1} = (\sigma_y^2 - (\beta_a^2 \sigma_a^2 + \beta_u^2 \sigma_u^2 + \beta_{bav}^2 \sigma_{bav}^2$$

$$+ 2\beta_a \beta_u \gamma_u \sigma_u^2 + 2\beta_a \beta_{bav} \gamma_{bav} \sigma_{bav}^2 \sigma_{\epsilon_1}^2))^{1/2}$$

$$= 0.906$$

$$v_2 \sim N(0, \sigma_{v_2})$$

$$\begin{aligned}\sigma_{v_2} &= (\sigma_a^2 - (\gamma_u^2 \sigma_u^2 + \gamma_{bav}^2 \sigma_{bav}^2 + \sigma_{\epsilon_2}^2))^{\frac{1}{2}} \\ &= 0.809\end{aligned}$$

v_1 and v_2 had variances such that A and Y both have unit variance.

A.14 Variance derivations

A.14.1 Treatment variance for equation (5)

$$\begin{aligned}\mathbf{A} &= \alpha_a + \mathbf{U}\gamma_u + \mathbf{BAV}\gamma_{bav} + \epsilon_2 \\ \Rightarrow \text{Var}(\mathbf{A}) &= \gamma_u^2 \sigma_u^2 + \gamma_{bav}^2 \sigma_{bav}^2 + 2\gamma_u \gamma_{bav} \text{COV}(\mathbf{U}, \mathbf{BAV}) + \text{Var}(\epsilon_2) \\ \sigma_a^2 &= \gamma_u^2 \sigma_u^2 + \gamma_{bav}^2 \sigma_{bav}^2 + \sigma_{\epsilon_2}^2\end{aligned}\tag{144}$$

A.14.2 Outcome variance for equation (4)

$$\mathbf{Y} = \alpha_y + \mathbf{A}\beta_a + \mathbf{U}\beta_u + \mathbf{BAV}\beta_{bav} + \epsilon_1\tag{145}$$

$$\Rightarrow \text{Var}(\mathbf{Y}) = \beta_a^2 \text{Var}(\mathbf{A}) + \beta_u^2 \text{Var}(\mathbf{U}) + \beta_{bav}^2 \text{Var}(\mathbf{BAV}) + \sigma_{\epsilon_2}^2 +\tag{146}$$

$$2\beta_a \beta_u \text{Cov}(\mathbf{A}, \mathbf{U}) + 2\beta_a \beta_{bav} \text{Cov}(\mathbf{A}, \mathbf{BAV}) + 2\beta_u \beta_{bav} \text{Cov}(\mathbf{U}, \mathbf{BAV})\tag{147}$$

$$= \beta_a^2 \sigma_a^2 + \beta_u^2 \sigma_u^2 + \beta_{bav}^2 \sigma_{bav}^2 + \sigma_{\epsilon_2}^2 + 2\beta_a \beta_u \gamma_u \sigma_u^2 + 2\beta_a \beta_{bav} \gamma_{bav} \sigma_{bav}^2\tag{148}$$