

A Crash Course in Good and Bad Controls

Carlos Cinelli* Andrew Forney† Judea Pearl ‡

March 21, 2022

Abstract

Many students of statistics and econometrics express frustration with the way a problem known as “bad control” is treated in the traditional literature. The issue arises when the addition of a variable to a regression equation produces an unintended discrepancy between the regression coefficient and the effect that the coefficient is intended to represent. Avoiding such discrepancies presents a challenge to all analysts in the data intensive sciences. This note describes graphical tools for understanding, visualizing, and resolving the problem through a series of illustrative examples. By making this “crash course” accessible to instructors and practitioners, we hope to avail these tools to a broader community of scientists concerned with the causal interpretation of regression models.

Introduction

Students, data analysts, and empirical social scientists have likely encountered the problem of “bad controls” (Angrist and Pischke, 2009, 2014). The problem arises when an analyst needs to decide whether or not the addition of a variable to a regression equation helps getting estimates closer to the parameter of interest. Analysts have long known that some variables, when added to the regression equation, can produce unintended discrepancies between the regression coefficient and the effect that the coefficient is expected to represent. Such variables have become known as “bad controls,” to be distinguished from “good controls” (also known as “confounders” or “deconfounders”) which are variables that must be added to the regression equation to eliminate what came to be known as “omitted variable bias” (Angrist and Pischke, 2009; Steiner and Kim, 2016; Cinelli and Hazlett, 2020a,b).

*Department of Statistics, University of Washington, Seattle. Email: cinelli@uw.edu

†Department of Computer Science, Loyola Marymount University, Los Angeles. Email: Andrew.Forney@lmu.edu

‡Department of Computer Science, University of California, Los Angeles. Email: judea@cs.ucla.edu. This research was supported in parts by grants from the National Science Foundation [#IIS-2106908], Office of Naval Research [#N00014-17-S-12091 and #N00014-21-1-2351], and Toyota Research Institute of North America [#PO000897].

The problem of “bad controls” however, has not received *systematic* attention in the standard statistics and econometrics literature. While most of the widely adopted textbooks discuss the problem of omitting “relevant” variables, they do not provide guidance on deciding which variables are relevant, nor which variables, if included in the regression, could induce, or worsen existing biases.¹ Researchers exposed only to this literature may get the impression that adding “more controls” to a regression model is always better. The few exceptions that do discuss the problem of “bad controls” unfortunately cover only a narrow aspect of the problem (e.g. Angrist and Pischke, 2009, 2014; Wooldridge, 2010; Imbens and Rubin, 2015; Gelman et al., 2020). Typical is the discussion found in Angrist and Pischke (2009, p.64)

Some variables are bad controls and should not be included in a regression model, even when their inclusion might be expected to change the short regression coefficients. Bad controls are variables that are themselves outcome variables in the notional experiment at hand. That is, bad controls might just as well be dependent variables too. Good controls are variables that we can think of having been fixed at the time the regressor of interest was determined.

Here, “good controls” are defined as variables that are thought to be unaffected by the treatment, whereas “bad controls” are variables that could be in principle affected by the treatment. Similar discussion can be found in Rosenbaum (2002) and Rubin (2009), for qualifying a variable for inclusion in propensity score analysis, as well as in Wooldridge (2005). Some authors (e.g, Wooldridge, 2010; Gelman et al., 2020) briefly warn about the potential of bias amplification of certain pre-treatment variables, but do not elaborate further. Although an improvement over an absence of discussion, these conditions are neither necessary nor sufficient for deciding whether a variable is a good control.

Recent advances in graphical models have produced simple criteria to distinguish “good” from “bad” controls; these range from necessary and sufficient conditions for deciding which set of variables should be adjusted for to identify the causal effect of interest (e.g, the back-door criterion and adjustment criterion in Pearl (1995) and Shpitser et al. (2012)), to deciding which, among a set of valid adjustment sets, would yield more precise estimates (Hahn, 2004; White and Lu, 2011; Henckel et al., 2019; Rotnitzky and Smucler, 2020; Witte et al., 2020). The purpose of this note is to provide practicing analysts a concise, simple, and *visual* summary of these criteria through illustrative examples.

¹See Chen and Pearl (2013) for a critical appraisal of econometrics textbooks, and Bollen and Pearl (2013) for eight misconceptions that still prevail in statistics and the social sciences. Warnings regarding the adjustment of post-treatment variables date back to at least Rosenbaum (1984), but a systematic solution to the problem of covariate selection was not available before the development of causal graphical models.

Preliminaries and basic terminology

Causal diagrams, and more specifically directed acyclic graphs (DAGs), have become popular in the social and health sciences for explaining and resolving difficult problems of causal inference in a rigorous, yet accessible manner. Many introductions to DAGs have now been published in a number of academic fields, such as sociology (Elwert, 2013; Morgan and Winship, 2015), economics (Hünermund and Bareinboim, 2019; Cunningham, 2021), psychology (Rohrer, 2018), epidemiology (Greenland et al., 1999; Hernán and Robins, 2020) and statistics (Pearl et al., 2009, 2016). Here we assume that readers are familiar with the basic notions of causal inference, DAGs, and in particular “path-blocking” as well as back-door paths. For those who need to refresh these notions, we provide a gentle introduction in the appendix. Still, given the simplicity of our illustrative examples, even the uninitiated reader will be able to understand and benefit from the main lessons of this crash course.

Briefly, causal DAGs provide a parsimonious representation of the qualitative aspects of the data generating process. Letters (e.g. X) represent random variables, and arrows, such as $X \rightarrow Y$, denote a (possible) direct causal effect of X on Y . No assumptions need to be made regarding the functional form of the causal relationships, nor about the distribution of variables. For this crash course, it is important to recall the three main sources of association that form the building blocks of a DAG, and when these are closed or opened:

1. *Mediators*, or chains, are patterns of the form $X \rightarrow Z \rightarrow Y$, meaning that X *causally* affects Y through the mediator Z . Conditioning on Z in a chain *blocks* (closes) this flow of association.
2. *Common causes*, or forks, are patterns of the form $X \leftarrow Z \rightarrow Y$, meaning X and Y share a common cause (a confounder) Z , thus inducing a *non-causal* association between both variables. Conditioning on Z in a fork *blocks* this flow of association.
3. *Common effects*, or colliders, are patterns of the form $X \rightarrow Z \leftarrow Y$, meaning that both X and Y share a common effect Z . Contrary to the other two variables, by default a common effect does not induce an association between X and Y . However, conditioning on Z induces a *non-causal* association between both variables.

Moreover, one important fact to keep in mind is that controlling for a descendant of a variable is equivalent to “partially” controlling for that variable. Any arbitrary path p from X to Y (consisting of a sequence of mediators, common causes, or colliders) will be blocked conditional on Z if, and only if, Z is a common cause or mediator along the path, or if p contains a collider and Z is *not* that collider, nor any of its descendants. We say that Z *d*-separates X from Y if Z blocks (closes) all paths from X to Y ; *d*-separation implies that Y and X are conditionally independent given Z .

Note that causal paths from X to Y are paths of the form $X \rightarrow \dots \rightarrow Y$, namely, those consisting of a sequence of (possibly empty) mediators. All other paths are non-causal, and may induce “spurious” associations between X and Y . In particular, for a

given variable X , we call “back-door” paths those confounding paths that begin with an arrow pointing into X . If we are interested, then, in estimating the causal effect of X on Y , our task is conceptually simple: we must block all spurious paths between X and Y , and we must not perturb any of the causal paths between them. This will be our guiding principle for deciding whether or not Z should be included in the regression equation, and it characterizes the essence of the graphical conditions known as the *back-door* criterion and the *adjustment* criterion (Pearl, 1995; Shpitser et al., 2012). Readers can find the formal statements of these graphical criteria in the appendix.

Illustrative examples

In the following set of models, the target of our analysis is the *average causal effect* (ACE) of a treatment X on an outcome Y , which stands for the expected increase of Y in response to a unit increase in X due to an *intervention*. Observed variables will be designated by black dots and unobserved variables by white empty circles. Variable Z , highlighted in red, will represent the variable whose inclusion in the regression equation is to be decided, with “good control” standing for bias reduction, “bad control” standing for bias increase, and “neutral control” when the addition of Z neither increases nor decreases the asymptotic bias. For this last case, we will also make brief remarks about how Z could affect the precision of the ACE estimate. Readers accustomed with the potential outcomes framework should know that deciding whether Z is a “good control” is equivalent to deciding whether ignorability of treatment assignment holds, conditional on Z . Readers who prefer to see algebraic derivations can find in the appendix analytical expressions for each graph, under the assumption of linearity² (the problem of “bad controls” is, however, non-parametric, i.e. it holds regardless of functional form assumptions).

Models 1, 2 and 3 – Good Controls (blocking back-door paths)

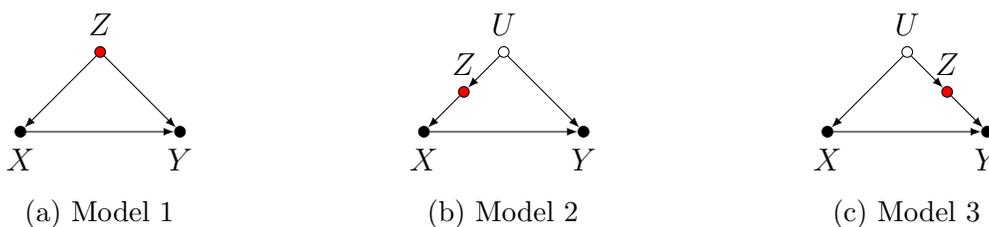


Figure 1: Models 1, 2, and 3

In Model 1, Z stands for a common cause of both X and Y . Once we control for Z , we block the back-door path from X to Y , producing an unbiased estimate of the ACE. In

²R code with numerical simulations for all examples can be found in: <https://www.kaggle.com/code/carloscinelli/crash-course-in-good-and-bad-controls-linear-r>.

Models 2 and 3, Z is not a common cause of both X and Y , and therefore, not a traditional “confounder” as in Model 1. Nevertheless, controlling for Z blocks the back-door path from X to Y due to the unobserved confounder U , and again, produces an unbiased estimate of the ACE.

Models 4, 5 and 6 – Good Controls (blocking back-door paths)

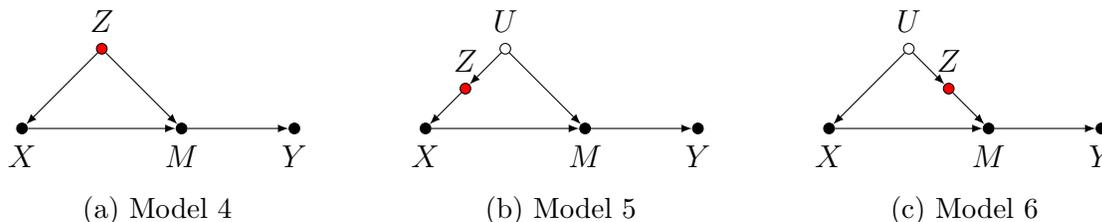


Figure 2: Models 4, 5, and 6

When thinking about possible threats of confounding, modelers need to keep in mind that common causes of X and any mediator (between X and Y) also confound the effect of X on Y . Therefore, Models 4, 5 and 6 are analogous to Models 1, 2 and 3—controlling for Z blocks the back-door path from X to Y and produces an unbiased estimate of the ACE.

Model 7 – Bad Control (M-bias)

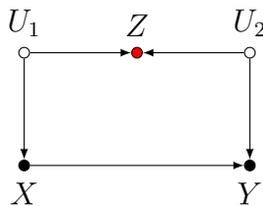


Figure 3: Model 7

We now encounter our first “bad control.” Here Z is correlated with the treatment and the outcome and it is also a “pre-treatment” variable. Traditional econometrics textbooks usually deem pre-treatment variables “good controls” (Angrist and Pischke, 2009, 2014; Imbens and Rubin, 2015). Careful analysis, however, reveals that Z is a “bad control.” Controlling for Z will induce bias by opening the back-door path $X \leftarrow U_1 \rightarrow Z \leftarrow U_2 \rightarrow Y$, thus spoiling a previously unbiased estimate of the ACE. This structure is known as the “M-bias,” and has spurred several controversies. Readers can find further discussion in Pearl (2009a, p. 186), Shrier (2009), Pearl (2009c,b), Sjölander (2009), Rubin (2009), Ding and Miratrix (2015), and Pearl (2015).

Undecidable—“Damned if you do, damned if you don’t.”

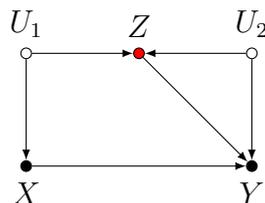


Figure 4: Variation of Model 7

Consider a variation of Model 7 such that Z has a direct effect on Y , as the one presented in Figure 4. Note that now we have an open back-door path, $X \leftarrow U_1 \rightarrow Z \rightarrow Y$, and the unadjusted estimate is no longer unbiased. While adjusting for Z closes this back-door path, it also opens back-door the path $X \leftarrow U_1 \rightarrow Z \leftarrow U_2 \rightarrow Y$, as we had in our previous example. In either case, the causal effect is not identified, and whether adjusting for Z reduces or increases the absolute value of the bias cannot be determined without further assumptions (see appendix). In this case, progress can be made with sensitivity analyses, by, for instance, positing plausible bounds on the the strength of the direct effect of Z on Y , or on the strength of the effects of the latent variables (Cinelli et al., 2019; Cinelli and Hazlett, 2020a).

Model 8 – Neutral Control (possibly good for precision)

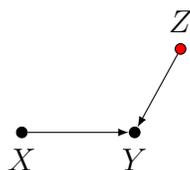


Figure 5: Model 8

In Model 8, Z is not a confounder nor does it block any back-door paths. Likewise, controlling for Z does not open any back-door paths from X to Y . Thus, in terms of asymptotic bias, Z is a “neutral control.” Analysis shows, however, that controlling for Z reduces the variation of the outcome variable Y , and helps to improve the precision of the ACE estimate in finite samples (Hahn, 2004; White and Lu, 2011; Henckel et al., 2019; Rotnitzky and Smucler, 2020).

Model 9 – Neutral Control (possibly bad for precision)

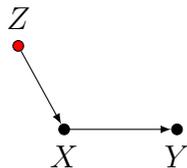


Figure 6: Model 9

Similar to the previous case, in Model 9 Z is “neutral” in terms of bias reduction. However, controlling for Z will reduce the variation of the treatment variable X and so may hurt the precision of the estimate of the ACE in finite samples (Henckel et al., 2019, Corollary 3.4). As a general rule of thumb, parents of X which are not necessary for identification are harmful for the asymptotic variance of the estimator; on the other hand, parents of Y which do not spoil identification are beneficial. See Henckel et al. (2019) for recent developments in graphical criteria for efficient estimation via adjustment in linear models. Remarkably, these conditions also have been shown to hold in non-parametric models for a broad class of non-parametric estimators (Rotnitzky and Smucler, 2020).

Model 10 – Bad Control (bias amplification)

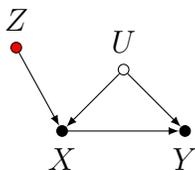


Figure 7: Model 10

We now encounter our second “pre-treatment” “bad control,” due to a phenomenon called “bias amplification” (Bhattacharya and Vogt, 2007; Wooldridge, 2009; Pearl, 2011, 2010, 2013; Middleton et al., 2016; Steiner and Kim, 2016). Naive control for Z in this model will not only fail to deconfound the effect of X on Y , but, in linear models, will amplify any existing bias.

Models 11 and 12 – Bad Controls (overcontrol bias)

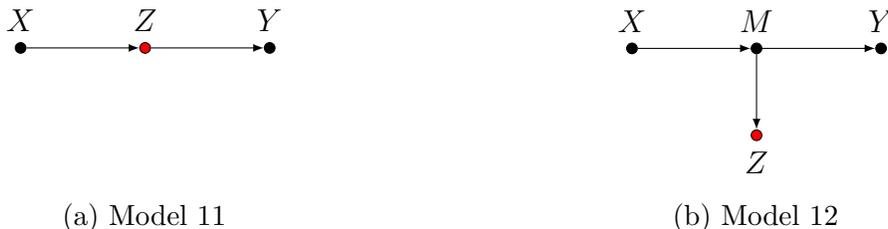


Figure 8: Models 11 and 12

If our target quantity is the ACE, we want to leave all channels through which the causal effect flows “untouched.” In Model 11, Z is a mediator of the causal effect of X on Y . Controlling for Z will block the very effect we want to estimate (the *total* effect of X on Y), thus biasing our estimates (this is usually known as “overcontrol bias”). In Model 12, although Z is not itself a mediator of the causal effect of X on Y , controlling for Z is equivalent to partially controlling for the mediator M , and will thus bias our estimates. Models 11 and 12 violate the back-door criterion (Pearl, 2009a), which excludes controls that are descendants of the treatment along paths to the outcome. Note that the same conclusions would hold if we had an extra direct causal path $X \rightarrow Y$.

Total *versus* direct effects

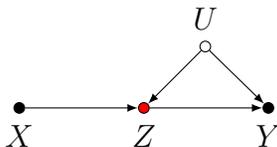


Figure 9: Variation of Model 11

The previous considerations assume the researcher is interested in the *total* effect of X on Y , as given by the ACE. If, instead, interest lies in the *controlled direct effect* (CDE) of X on Y (i.e, the effect of X while holding Z constant by intervention, see Pearl (2009a, 2011) as well as the appendix), then adjusting for Z in Model 11 (Figure 8a) would indeed be appropriate. However, consider a variation of Model 11 with an unobserved confounder of Z and Y , denoted by U , as shown in Figure 9. First notice that U *does not* confound the effect of X on Y , and thus our ACE estimate remains unbiased as it were in Model 11, so long as we do not adjust for Z . On the other hand, here adjusting for Z now opens the colliding path $X \rightarrow Z \leftarrow U \rightarrow Y$, thus biasing the CDE estimate.

Model 13 – Neutral Control (possibly good for precision)

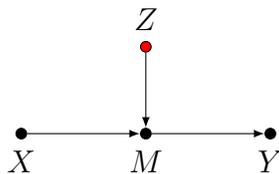


Figure 10: Model 13

At first look, Model 13 might seem similar to Model 12, and one may think that adjusting for Z would bias the effect estimate, by restricting variations of the mediator M . However, the key difference here is that Z is a cause, not an effect, of the mediator (and, consequently, also a cause of Y). Thus, Model 13 is analogous to Model 8, and so controlling for Z will be neutral in terms of bias and may increase the precision of the ACE estimate in finite samples. Readers can find further discussion of this case in Pearl (2013).

Models 14 and 15 – Neutral Controls (possibly helpful in the case of selection bias)



Figure 11: Models 14 and 15

Contrary to folklore, not all “post-treatment” variables are inherently bad controls. In Models 14 and 15 controlling for Z does not open any confounding paths between X and Y . Thus, Z is neutral in terms of bias. However, controlling for Z does reduce the variation of the treatment variable X and so may hurt the precision of the ACE estimate in finite samples. Additionally, in Model 15, suppose one has only samples with $W = w$ recorded (a case of selection bias³, which we explain next). In this case, controlling for Z can help to

³Some economists may denote confounding bias as “selection bias,” meaning preferential selection to treatment (Angrist and Pischke, 2009, 2014). Here selection bias means preferential selection into the available data.

obtain the W -specific effect of X on Y by blocking the colliding path due to W . In linear models, controlling for Z actually fully recovers the ACE (see appendix).

Models 16 and 17 – Bad Controls (selection bias)



Figure 12: Models 16 and 17

Contrary to Models 14 and 15, here controlling for Z is no longer harmless, and induces what is classically known as “selection bias” or “collider stratification bias.” Adjusting for Z in Model 16 opens the colliding path $X \rightarrow Z \leftarrow U \rightarrow Y$ and so biases the ACE. In Model 17, adjusting for Z not only opens the path $X \rightarrow Z \leftarrow Y$, but also the colliding path due to the latent parents of Y , thus biasing the ACE and motivating our final example.

Model 18 – Bad Control (case-control bias)

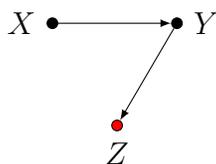


Figure 13: Model 18

In our last example, Z is not in the causal pathway from X to Y , Z is not a direct cause of X , and Z is connected to Y . Thus, one might surmise that, as in Model 8, controlling for Z is harmless for identification, and perhaps beneficial for finite sample efficiency. However, controlling for the effects of the outcome Y will induce bias in the estimate of the ACE, even without the direct arrow $X \rightarrow Z$, thus making Z a “bad control.” This happens because Z is in fact a descendant of a collider: the outcome Y itself. A visual explanation of this phenomenon using “virtual colliders” can be found in Pearl (2009a, Sec. 11.3). The same phenomenon can also be explained by explicitly drawing the potential outcomes on the DAG (see both explanations in the appendix). Model 18 is special case of selection bias

usually known as “case-control” bias. Finally, although controlling for Z will generally bias numerical estimates of the ACE, it does have an exception when X has no causal effect on Y . In this scenario, X is still d -separated from Y even after conditioning on Z . Thus, adjusting for Z is valid for testing whether the effect of X on Y is *zero*.

Bad controls in applied research

Despite their simplicity, these illustrative examples should provide practitioners with a principled framework to understand many problems found in real world applications. To demonstrate, we now briefly present three cases of bad controls discussed in applied research, coming from diverse areas such as epidemiology, sociology, and economics.

The birth-weight paradox (Hernández-Díaz et al., 2006). Infants born to smokers were found to have higher risks of mortality than infants born to non-smokers. However, among infants with low birth-weight (LBW), this relationship was reversed. This reversal of effects has created many controversies in epidemiology—does it mean that maternal smoking is beneficial for LBW infants? A plausible reason for such a finding could simply be collider stratification bias, as shown in Model 16. Here X is maternal smoking, Y infant mortality, Z birth-weight, and U stands for unobserved risk-factors (such as birth-defects and malnutrition), that could also affect birth-weight. Note that stratifying the analysis by birth-weight would induce a spurious association between smoking and mortality due to the competing risk-factors. LBW infants of non-smokers need to have alternative causes for their LBW (such as malnutrition), and such causes could also lead to higher mortality.

Homophily bias in social network analysis (Elwert and Winship, 2014). An important task in the causal inference of social networks is to estimate the causal effects of social contagion, also known as “interpersonal effects.” However, social ties in the analysis of social networks may be pre-treatment colliders as exemplified in the “M-bias” structure of Model 7. Suppose we are interested in assessing whether the civic engagement of individual 1 (X) leads to the civic engagement of individual 2 in the subsequent time period (Y). Let Z denote whether such individuals are friends, and U_1 and U_2 denote the personal characteristics (such as altruism) of individuals 1 and 2, respectively. Here, the social tie Z is a collider, and computing the association of Y and X between friends ($Z = 1$) would bias the interpersonal causal effects in civic engagement.

The Antebellum Puzzle (Schneider, 2020). An interesting puzzle of economic history is the fact that, during the nineteenth century in Britain and the United States, the average height of adult men fell even though the economic conditions of these countries improved alongside childhood nutrition. One possible explanation for such a paradoxical finding is selection bias in the forms of Models 17 and 18 wherein researchers using data

from individuals enlisted in the military or in prison are effectively conditioning on colliders. For military records, consider Model 18, and let X denote childhood nutrition, Y adult height, and Z an indicator of whether the individual was enlisted in the military. The causal path from Y to Z represents the fact that taller men may have better opportunities in the civilian market, and thus shorter men were more likely to enlist. Restricting the analysis to those enlisted in the military is therefore equivalent to controlling for Z , and leads to selection bias. Now for prison records, consider Model 17, and let Z be an indicator of whether the individual was arrested. Here one could argue that both childhood nutrition and adult height have pathways to committing a crime through socio-economic opportunities, thus again leading to selection bias.

These examples are by no means exhaustive. Readers can find other interesting cases across applied sciences, such as: the threats of collider bias in understanding risk factors of COVID-19 (Griffith et al., 2020); the “Obesity paradox,” in which obesity appears to benefit individuals who survive heart failure (Banack and Kaufman, 2013); and examples of “bad controls” due to adjustment of mediators and colliders in multigenerational mobility (Breen, 2018), anesthesiology research (Gaskell and Sleight, 2020) or animal science (Bello et al., 2018). Further discussion of bad controls in theoretical and applied works can be found in Pearl and Mackenzie (2018).

Multiple controls

When considering multiple controls, the status of a single control as “good” or “bad” may change depending on the context of the other variables under consideration. Nevertheless, the main lessons from our illustrative examples remain. A set of control variables \mathbf{Z} will be “good” if: (i) it blocks all non-causal paths from the treatment to the outcome; (ii) it leaves any mediating paths from the treatment to the outcome “untouched” (since we are interested in the total effect); and, (iii) it does not open new spurious paths between the treatment and the outcome (e.g., due to colliders). As to efficiency considerations, we should give preference to those variables “closer” to the outcome, in opposition to those closer to the treatment—so long as, of course, this does not spoil identification.

Finally, we remind readers that, when considering models with more complicated structures, one can always resort to specialized computer programs. Open-source software implementing algorithms for selecting adjustment sets can be found in the R packages `pcalg` (Kalisch et al., 2012), `dagitty` (Textor et al., 2016)⁴, and `causaleffect` (Tikka and Karvanen, 2017). Users familiar with the software `SAS` may find the procedure `CAUSALGRAPH` useful (Thompson, 2019). A web application implementing the methods discussed in Bareinboim and Pearl (2016) is also available.⁵ In other words, given a causal diagram, the problem of deciding which variables are good or bad controls has been automatized.

⁴Also available online in www.dagitty.net.

⁵Available online at www.causalfusion.net.

Beyond adjustment

Here we have focused on the identification of causal effects through simple covariate adjustment, classifying Z as a “good” or “bad” control according to this criterion. However, other identification opportunities may be available. For instance, going back again to Model 10, Z is what is usually known as an “instrumental variable.” In this case, while Z is indeed a “bad control,” it can still be used as an instrument to bound or point identify causal effects under certain parametric assumptions, albeit using a different formula (Wright, 1928; Bowden and Turkington, 1990; Balke and Pearl, 1994; Angrist et al., 1996; Balke and Pearl, 1997; Brito and Pearl, 2002). More generally, the *do*-calculus provides a complete solution for the task of non-parametric identification of treatment effects in causal DAGs, beyond the simple adjustment formula (Pearl, 1995; Shpitser and Pearl, 2008; Pearl, 2009a). In certain instances, such as the “front-door” criterion, this allows exploiting post-treatment variables for identification (Pearl, 2009a). Further details on the *do*-calculus should be the topic of a separate crash course.

Sensitivity analysis

In real world applications, it can be the case that the causal effect of X on Y cannot be identified from the DAG structure alone. When that happens, without further assumptions, it is usually not possible to determine whether including Z in the regression equation will reduce or increase the absolute value of the bias, as we have seen in the example of Figure 4. In such cases, claims about the causal effect of X on Y must rely on knowledge beyond the constraints of the DAG, such as plausibility judgments (i) on the direct effect of observed variables, (ii) on the strength of association of latent variables with X and Y , or (iii) on the relative importance of unobserved confounders as compared to observed confounders (Cinelli et al., 2018, 2019; Cinelli and Hazlett, 2020a,b; Zhang et al., 2021). A suite of sensitivity analysis tools to examine the robustness of linear regression estimates to omitted variable biases (OVB) can be found in the package `sensemkr` for R, Stata and Python (Cinelli et al., 2020; LaPierre et al., 2022). An interactive web application is also available.⁶ Generalization of OVB results to fully nonparametric models, using Debiased Machine Learning, is developed in Chernozhukov et al. (2021).

Concluding remarks

In this note, we demonstrated through illustrative examples how simple graphical criteria can be used to decide when a variable should (or should not) be included in a regression equation—and thus whether it can be deemed a “good” or “bad” control. Many of these examples act as cautionary notes against prevailing practices: for instance, Models 7 to 10

⁶Available online at https://carloscinelli.shinyapps.io/robustness_value/.

reveal that one should be cautious of the general recommendation, usually derived from propensity score logic, of conditioning on *all* pre-treatment predictors of the treatment assignment⁷; whereas Models 14 and 15 show that not all “post-treatment” variables are “bad-controls,” and some may even help with identification.

In all cases, structural knowledge is indispensable for deciding whether a variable is a good or bad control, and graphical models provide a natural language for articulating such knowledge, as well as efficient tools for examining its logical ramifications. We have found that an example-based approach to “bad controls,” such as the one presented here, can serve as a powerful instructional device to supplement more extended and formal discussions of the problem. By making this “crash course” accessible to instructors and practitioners, we hope to avail these tools to a broader community of scientists concerned with the causal interpretation of regression models.

References

- Angrist, J. and Pischke, J.-S. (2009). *Mostly harmless econometrics: an empiricists guide*. Princeton: Princeton University Press.
- Angrist, J. D., Imbens, G. W., and Rubin, D. B. (1996). Identification of causal effects using instrumental variables. *Journal of the American statistical Association*, 91(434):444–455.
- Angrist, J. D. and Pischke, J.-S. (2014). *Mastering ’metrics: The path from cause to effect*. Princeton University Press.
- Balke, A. and Pearl, J. (1994). Counterfactual probabilities: Computational methods, bounds and applications. In *Uncertainty Proceedings 1994*, pages 46–54. Elsevier.
- Balke, A. and Pearl, J. (1997). Bounds on treatment effects from studies with imperfect compliance. *Journal of the American Statistical Association*, 92(439):1171–1176.
- Banack, H. R. and Kaufman, J. S. (2013). The “obesity paradox” explained. *Epidemiology*, 24(3):461–462.
- Bareinboim, E. and Pearl, J. (2016). Causal inference and the data-fusion problem. *Proceedings of the National Academy of Sciences*, 113(27):7345–7352.
- Bello, N. M., Ferreira, V. C., Gianola, D., and Rosa, G. J. (2018). Conceptual framework for investigating causal effects from observational data in livestock. *Journal of animal science*, 96(10):4045–4062.
- Bhattacharya, J. and Vogt, W. B. (2007). Do instrumental variables belong in propensity scores? Technical report, National Bureau of Economic Research.

⁷For instance, examples of such recommendations can be found in Rosenbaum (2002, p.76), Rubin (2009), Imbens and Rubin (2015, p.265), Dorie et al. (2016, p.3453).

- Bollen, K. A. and Pearl, J. (2013). Eight myths about causality and structural equation models. In *Handbook of causal analysis for social research*, pages 301–328. Springer.
- Bowden, R. J. and Turkington, D. A. (1990). *Instrumental variables*. Cambridge university press.
- Breen, R. (2018). Some methodological problems in the study of multigenerational mobility. *European Sociological Review*, 34(6):603–611.
- Brito, C. and Pearl, J. (2002). Generalized instrumental variables. In *Proceedings of the Eighteenth conference on Uncertainty in artificial intelligence*, pages 85–93. Morgan Kaufmann Publishers Inc.
- Chen, B. and Pearl, J. (2013). Regression and causation: a critical examination of six econometrics textbooks. *Real-World Economics Review*.
- Chernozhukov, V., Cinelli, C., Newey, W., Sharma, A., and Syrgkanis, V. (2021). Omitted variable bias in machine learned causal models. *arXiv preprint arXiv:2112.13398*.
- Cinelli, C., Ferwerda, J., and Hazlett, C. (2020). sensemakr: Sensitivity analysis tools for OLS in R and Stata. *Available at SSRN 3588978*.
- Cinelli, C. and Hazlett, C. (2020a). Making sense of sensitivity: extending omitted variable bias. *Journal of the Royal Statistical Society: Series B (Statistical Methodology)*, 82(1):39–67.
- Cinelli, C. and Hazlett, C. (2020b). An omitted variable bias framework for sensitivity analysis of instrumental variables. *Work. Pap.*
- Cinelli, C., Kumor, D., Chen, B., Pearl, J., and Bareinboim, E. (2019). Sensitivity analysis of linear structural causal models. *International Conference on Machine Learning*.
- Cinelli, C., Pearl, J., and Chen, B. (2018). When confounders are confounded: Naive benchmarking in sensitivity analysis.
- Cramér, H. (1946). *Mathematical methods of statistics*. Princeton University Press.
- Cunningham, S. (2021). *Causal inference: The mixtape*. Yale University Press.
- Ding, P. and Miratrix, L. W. (2015). To adjust or not to adjust? sensitivity analysis of m-bias and butterfly-bias. *Journal of Causal Inference*, 3(1):41–57.
- Dorie, V., Harada, M., Carnegie, N. B., and Hill, J. (2016). A flexible, interpretable framework for assessing sensitivity to unmeasured confounding. *Statistics in medicine*, 35(20):3453–3470.

- Elwert, F. (2013). Graphical causal models. In *Handbook of causal analysis for social research*, pages 245–273. Springer.
- Elwert, F. and Winship, C. (2014). Endogenous selection bias: The problem of conditioning on a collider variable. *Annual review of sociology*, 40:31–53.
- Gaskell, A. L. and Sleigh, J. W. (2020). An introduction to causal diagrams for anesthesiology research. *Anesthesiology*, 132(5):951–967.
- Gelman, A., Hill, J., and Vehtari, A. (2020). *Regression and other stories*. Cambridge University Press.
- Greenland, S., Pearl, J., and Robins, J. M. (1999). Causal diagrams for epidemiologic research. *Epidemiology*, pages 37–48.
- Griffith, G. J., Morris, T. T., Tudball, M. J., Herbert, A., Mancano, G., Pike, L., Sharp, G. C., Sterne, J., Palmer, T. M., Smith, G. D., et al. (2020). Collider bias undermines our understanding of covid-19 disease risk and severity. *Nature communications*, 11(1):1–12.
- Hahn, J. (2004). Functional restriction and efficiency in causal inference. *Review of Economics and Statistics*, 86(1):73–76.
- Henckel, L., Perković, E., and Maathuis, M. H. (2019). Graphical criteria for efficient total effect estimation via adjustment in causal linear models. *arXiv preprint arXiv:1907.02435*.
- Hernán, M. and Robins, J. (2020). *Causal inference: What if*. Boca Raton: Chapman & Hill/CRC.
- Hernández-Díaz, S., Schisterman, E. F., and Hernán, M. A. (2006). The birth weight “paradox” uncovered? *American journal of epidemiology*, 164(11):1115–1120.
- Hünermund, P. and Bareinboim, E. (2019). Causal inference and data-fusion in econometrics. *arXiv preprint arXiv:1912.09104*.
- Imbens, G. W. and Rubin, D. B. (2015). *Causal inference in statistics, social, and biomedical sciences*. Cambridge University Press.
- Kalisch, M., Mächler, M., Colombo, D., Maathuis, M. H., and Bühlmann, P. (2012). Causal inference using graphical models with the R package pcalg. *Journal of Statistical Software*, 47(11):1–26.
- LaPierre, N., Hill, B. L., Zhang, Z., and Cinelli, C. (2022). sensemakr: Sensitivity analysis tools for OLS in python. <https://github.com/KennyZhang-17/PySensemakr>.
- Middleton, J. A., Scott, M. A., Diakow, R., and Hill, J. L. (2016). Bias amplification and bias unmasking. *Political Analysis*, 24(3):307–323.

- Morgan, S. L. and Winship, C. (2015). *Counterfactuals and causal inference*. Cambridge University Press.
- Pearl, J. (1995). Causal diagrams for empirical research. *Biometrika*, 82(4):669–688.
- Pearl, J. (2009a). *Causality*. Cambridge University Press.
- Pearl, J. (2009b). Letter to the editor: Remarks on the method of propensity score. *Statistics in Medicine*, 28:1420–1423. URL: <https://ucla.in/2NbS14j>.
- Pearl, J. (2009c). Myth, confusion, and science in causal analysis. *UCLA Cognitive Systems Laboratory*, Technical Report (R-348). URL: <https://ucla.in/2EihVyD>.
- Pearl, J. (2010). On a class of bias-amplifying variables that endanger effect estimates. In *Proceedings of the Twenty-Sixth Conference on Uncertainty in Artificial Intelligence*, pages 417–424. URL: <https://ucla.in/2N8mBMg>.
- Pearl, J. (2011). Invited commentary: understanding bias amplification. *American journal of epidemiology*, 174(11):1223–1227. URL: <https://ucla.in/2PORDX2>.
- Pearl, J. (2013). Linear models: A useful “microscope” for causal analysis. *Journal of Causal Inference*, 1(1):155–170. URL: <https://ucla.in/2LcpmHz>.
- Pearl, J. (2015). Comment on ding and miratrix: “to adjust or not to adjust?”. *Journal of Causal Inference*, 3(1):59–60. URL: <https://ucla.in/2PgOWNd>.
- Pearl, J. et al. (2009). Causal inference in statistics: An overview. *Statistics surveys*, 3:96–146.
- Pearl, J., Glymour, M., and Jewell, N. P. (2016). *Causal inference in statistics: a primer*. John Wiley & Sons.
- Pearl, J. and Mackenzie, D. (2018). *The Book of Why: The New Science of Cause and Effect*. Hachette UK.
- Rohrer, J. M. (2018). Thinking clearly about correlations and causation: Graphical causal models for observational data. *Advances in Methods and Practices in Psychological Science*, 1(1):27–42.
- Rosenbaum, P. R. (1984). The consequences of adjustment for a concomitant variable that has been affected by the treatment. *Journal of the Royal Statistical Society: Series A (General)*, 147(5):656–666.
- Rosenbaum, P. R. (2002). *Observational studies*. Springer.

- Rotnitzky, A. and Smucler, E. (2020). Efficient adjustment sets for population average causal treatment effect estimation in graphical models. *Journal of Machine Learning Research*, 21(188):1–86.
- Rubin, D. B. (2009). Should observational studies be designed to allow lack of balance in covariate distributions across treatment groups? *Statistics in Medicine*, 28(9):1420–1423.
- Schneider, E. B. (2020). Collider bias in economic history research. *Explorations in Economic History*, 78:101356.
- Shpitser, I. and Pearl, J. (2008). Complete identification methods for the causal hierarchy. *Journal of Machine Learning Research*, 9(Sep):1941–1979.
- Shpitser, I., VanderWeele, T., and Robins, J. M. (2012). On the validity of covariate adjustment for estimating causal effects. *arXiv preprint arXiv:1203.3515*.
- Shrier, I. (2009). Propensity scores. *Statistics in Medicine*, 28(8):1317–1318.
- Sjölander, A. (2009). Propensity scores and m-structures. *Statistics in medicine*, 28(9):1416–1420.
- Steiner, P. M. and Kim, Y. (2016). The mechanics of omitted variable bias: Bias amplification and cancellation of offsetting biases. *Journal of causal inference*, 4(2).
- Textor, J., van der Zander, B., Gilthorpe, M. S., Liškiewicz, M., and Ellison, G. T. (2016). Robust causal inference using directed acyclic graphs: the R package ‘dagitty’. *International journal of epidemiology*, 45(6):1887–1894.
- Thompson, C. (2019). Causal graph analysis with the causalgraph procedure. In *Proceedings of SAS Global Forum*.
- Tikka, S. and Karvanen, J. (2017). Identifying causal effects with the R package causaleffect. *Journal of Statistical Software*, 76.
- White, H. and Lu, X. (2011). Causal diagrams for treatment effect estimation with application to efficient covariate selection. *Review of Economics and Statistics*, 93(4):1453–1459.
- Witte, J., Henckel, L., Maathuis, M. H., and Didelez, V. (2020). On efficient adjustment in causal graphs. *Journal of Machine Learning Research*, 21:246.
- Wooldridge, J. (2009). Should instrumental variables be used as matching variables. Technical report, Citeseer.
- Wooldridge, J. M. (2005). Violating ignorability of treatment by controlling for too many factors. *Econometric Theory*, 21(5):1026–1028.
- Wooldridge, J. M. (2010). *Econometric analysis of cross section and panel data*. MIT press.

- Wright, P. G. (1928). *Tariff on animal and vegetable oils*. Macmillan Company, New York.
- Wright, S. (1921). Correlation and causation. *Journal of agricultural research*, 20(7):557–585.
- Zhang, C., Cinelli, C., Chen, B., and Pearl, J. (2021). Exploiting equality constraints in causal inference. In *International Conference on Artificial Intelligence and Statistics*, pages 1630–1638. PMLR.

A Appendix

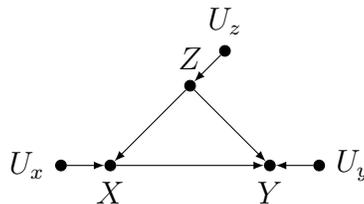
This appendix provides a short introduction to the notions of causal models, causal diagrams and “path-blocking” for the identification of causal effects via adjustment. Readers can find more extensive discussions in Pearl (2009a); Pearl et al. (2016) and Pearl and Mackenzie (2018).

Structural causal models and causal diagrams

In order to decide whether there is a discrepancy between a certain regression equation (an associational quantity), and a target “causal effect” (a causal quantity), we need to mathematically *define* what this causal effect is. And to do that, we first need the concept of a *causal model*. We briefly introduce *structural causal models* (SCM) (Pearl, 2009a) with an example.

$$M = \begin{cases} Z & \leftarrow f_z(U_z) \\ X & \leftarrow f_x(Z, U_x) \\ Y & \leftarrow f_y(X, Z, U_y) \\ \mathbf{U} & \sim P(\mathbf{U}) \end{cases}$$

(a) Structural causal model M



(b) Causal diagram G associated with M

Figure 14: Structural Causal Model M and its associated graph G

Consider the SCM M shown in Figure 14a. The variables $\mathbf{V} = \{Z, X, Y\}$ are called the *endogenous* variables, and stand for those variables that the investigator chose to model their cause-effect relationships; the variables $\mathbf{U} = \{U_z, U_x, U_y\}$ are called the *exogenous* variables and represent everything else that the investigator chose *not* to explicitly model (these are also usually called *disturbances*). The functions $\mathcal{F} = \{f_z, f_x, f_y\}$ are called *structural equations*, and each function represents a causal process that *assigns* to its respective endogenous variable a value based on the values of the other variables. We use the assignment symbol (\leftarrow) to emphasize the *asymmetry* in a causal relationship, flowing from cause to effect. Finally, the exogenous variables have an associated probability distribution $P(\mathbf{U})$ summarizing their uncertainty. In this particular example, we assume the exogenous variables are mutually independent (but in general, this need not be the case). The SCM M induces a joint distribution on the endogenous variables $P(\mathbf{V})$, which we denote by *observational distribution*. In observational studies, the investigator only has access to samples of $P(\mathbf{V})$.

Every SCM has an associated graph G , usually called its *causal diagram*. In the types of models we consider here, which do not exhibit cycles, the causal diagram will be a directed acyclic graph (DAG). The causal diagram of our example is shown in Figure 14b. The graph G contains one node for each variable in M , and a directed arrow $V_i \rightarrow V_j$ whenever

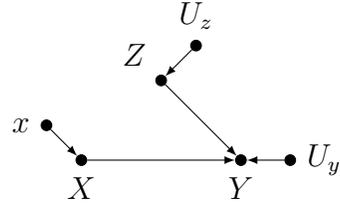
V_i appears in the structural equation of V_j , meaning that V_i is a *direct cause* of V_j . Here we explicitly show the exogenous variables, but, conventionally, these are omitted from the graph for brevity. When the exogenous variables are omitted from the diagram, a dashed bidirected arrow $V_i \leftrightarrow V_j$ should be added whenever the exogenous variables entering f_{v_i} and f_{v_j} are *not* independent.

Interventions and causal effects

Interventions are modeled by modifying mechanisms of the SCM. For example, the act $do(X = x)$ in the model of Figure 14a amounts to replacing the original mechanism $X \leftarrow f_x(Z, U_x)$ with a new mechanism in which X is externally forced to attain the value x , i.e., $X \leftarrow x$. This results in the modified SCM M_x of Figure 15a.

$$M_x = \begin{cases} Z & \leftarrow f_z(U_z) \\ X & \leftarrow x \\ Y & \leftarrow f_y(X, Z, U_y) \\ \mathbf{U} & \sim P(\mathbf{U}) \end{cases}$$

(a) Modified SCM M_x



(b) Modified causal diagram G_x

Figure 15: Effect of intervention $do(X = x)$

The model M_x induces an *interventional distribution* on the endogenous variables, denoted by $P(\mathbf{V} \mid do(X = x))$. With the concept of an intervention in mind, we can now define the average causal effect (i.e, the expected increase of Y in response to a unit increase in X due to an *intervention*) as the average contrast of Y under two distinct interventions:

$$ACE(x) = E[Y \mid do(x + 1)] - E[Y \mid do(x)]$$

In general the ACE varies depending on levels of x , but in linear models, as we show below, the ACE reduces to a single number. Other causal effects can be defined with the same model modification logic. For instance, the controlled direct effect (or CDE, i.e, the expected increase of Y per unit of a controlled increase in X , while holding Z constant) is defined as the difference:

$$CDE(x, z) = E[Y \mid do(x + 1), do(z)] - E[Y \mid do(x), do(z)]$$

Potential outcomes

Potential outcomes \mathbf{V}_x are defined as the solution of the endogenous variables \mathbf{V} in the modified model M_x . Thus, $P(\mathbf{V} \mid do(X = x))$ can be equivalently written as $P(\mathbf{V}_x)$ (likewise, we could have written all variables in M_x and G_x as Z_x , X_x and Y_x). As such, the ACE can be equivalently written as $ACE(x) = E[Y_{x+1}] - E[Y_x]$.

Causal and non-causal paths: chains, forks and colliders

Concretely, let us suppose that the structural equations of our example are linear, that is, $Z \leftarrow U_z$, $X \leftarrow \lambda_{zx}Z + U_x$, $Y \leftarrow \lambda_{xy}X + \lambda_{zy}Z + U_y$. Further assume that the disturbances U are normally distributed, and that the random variables X, Z, Y have mean zero and unit variance. Then the ACE evaluates to:

$$\text{ACE}(x) = E[Y \mid do(x+1)] - E[Y \mid do(x)] = \lambda_{xy}$$

To contrast, now let us compute the regression coefficient of Y on X , denoted by β_{yx}

$$\beta_{yx} = \frac{\text{Cov}(Y, X)}{\text{Var}(X)} = \lambda_{xy} + \lambda_{zx}\lambda_{zy}$$

Note how the regression coefficient $\beta_{yx} = \lambda_{xy} + \lambda_{zx}\lambda_{zy}$ differs from the ACE = λ_{xy} . This happens because the observed association of X and Y mixes both the *causal* association (the path $X \rightarrow Y$), and the *non-causal* association due to the confounder Z (the path $X \leftarrow Z \rightarrow Y$). We call such confounding paths, that start with an arrow pointing to X , “back-door paths.” Note, however, that the regression coefficient of Y on X adjusting for Z (denoted by $\beta_{yx.z}$) evaluates to (a derivation is provided later, in Equations 6 to 10)

$$\beta_{yx.z} = \lambda_{xy}$$

That is, controlling for Z in this model effectively blocks the back-door path, and recovers the ACE.

In general, how does path blocking work in a graphical model? To answer this question, we need to understand the three main patterns of a causal diagram, which help us characterize when paths (consisting of sequences of the following triplets) of the graph are blocked or open.

- **Chains (mediators).** Chains are patterns of the form $X \rightarrow Z \rightarrow Y$, meaning that X causally affects Y through the mediator Z . Conditioning on Z in a chain blocks this flow of association.
- **Forks (common causes).** Forks are patterns of the form $X \leftarrow Z \rightarrow Y$, meaning X and Y share a common cause (a confounder) Z , thus inducing a *non-causal* association between both variables. Conditioning on Z in a fork blocks this flow of association.
- **Colliders (common effects).** Colliders are patterns of the form $X \rightarrow Z \leftarrow Y$, meaning that both X and Y share a common effect Z . Contrary to the other two patterns, this path is closed by default—conditioning on Z *opens* the path and induces a *non-causal* association between X and Y .

A final rule to keep in mind is that controlling for a descendant of a variable is equivalent to “partially” controlling for that variable. Thus, controlling for a descendant of a mediator or

a confounder partially blocks the flow of association, whereas controlling for a descendant of a collider partially opens the flow of association.

We can now judge whether any path p in a graph, no matter how complicated, is blocked by a set \mathbf{Z} . This happens if, and only if: (i) p contains a chain or a fork, such that the middle node is in \mathbf{Z} ; *or*, (ii) p contains a collider, such that neither the middle node, nor any of its descendants, are in \mathbf{Z} .

The back-door and the adjustment criteria

Armed with these tools, the DAG reveals which set of variables \mathbf{Z} blocks the correct paths for valid estimation of the ACE. We would like to find a set \mathbf{Z} , such that,

- it blocks all spurious paths from X to Y ;
- it *does not* (partially) block any of the causal paths from X to Y ; and,
- it does not (partially) open other spurious paths.

The above conditions characterize the so-called *back-door* criterion, later generalized by the *adjustment* criterion (Pearl, 1995; Shpitser et al., 2012). If we can find such a set of controls $\mathbf{Z} = \{Z_1, \dots, Z_k\}$, then the interventional expectation of Y can be computed from the observational distribution as

$$E[Y \mid do(X = x)] = E[E[Y \mid X = x, \mathbf{Z}]] \quad (1)$$

Readers accustomed to potential outcomes should note that, if \mathbf{Z} satisfies the adjustment criterion, then conditional ignorability holds, ie., $Y_x \perp\!\!\!\perp X \mid \mathbf{Z}$.

Formal statements. For completeness, we now provide the formal statements of the back-door and adjustment criteria.

Definition 1 (Back-door criterion (Pearl, 2009a)) *A set of variables Z satisfies the back-door criterion relative to (X, Y) in a DAG G if:*

- *No node in Z is a descendant of X ; and,*
- *Z blocks every path between X and Y that contains an arrow into X .*

The adjustment criterion was later devised to explicitly handle cases in which Z may contain descendants of X .

Definition 2 (Adjustment criterion (Shpitser et al., 2012)) *A set of variables Z satisfies the adjustment criterion relative to (X, Y) in a DAG G if:*

- *In the dag $G_{\bar{x}}$ (the DAG G with the arrows incoming into X removed), no element in Z is a descendant of any $W \notin X$ which lies on a proper causal path from X to Y .*
- *All non-causal paths in G from X to Y are blocked by Z .*

Linear *versus* non-linear models

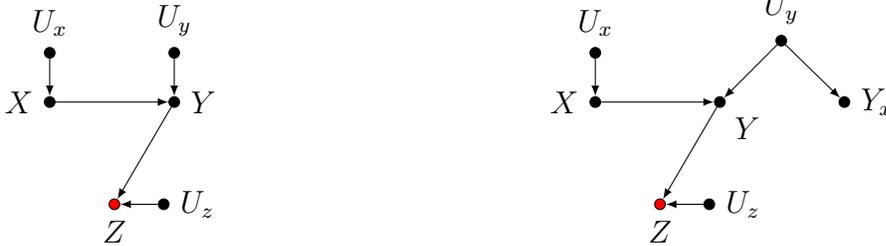
The previous identification result is non-parametric, and it involves *two* expectations. First we compute the conditional expectation $E[Y \mid X = x, \mathbf{Z} = \mathbf{z}]$, then we *average* this conditional expectation over $P(\mathbf{Z})$. If, however, the conditional expectation function $E[Y \mid X = x, \mathbf{Z} = \mathbf{z}]$ is linear, the expression simplifies to

$$E[E[Y \mid X = x, \mathbf{Z}]] = \alpha + \beta_{yx.z}x + \sum_{j=1}^k \beta_{yz_j.xz_{-j}}E[Z_j]$$

Where α is a constant and \mathbf{Z}_{-j} denotes the set \mathbf{Z} excluding Z_j . Therefore, under the parametric assumption of linearity, the ACE simply equals the regression coefficient $\beta_{yx.z}$, and no averaging over the distribution of \mathbf{Z} is necessary (similar result can be obtained if the conditional expectation is linearly separable on X). If, however, the conditional expectation is not linear, the regression coefficient $\beta_{yx.z}$ targets a different causal quantity, which may be an incomplete summary of the ACE (see, e.g., Angrist and Pischke, 2009). In such cases, users should resort back to the proper adjustment formula as given by Equation 1.

Virtual colliders and d-separation

Finally, we explain both d-separation and virtual colliders using the case of Model 18. Rewrite Model 18 showing the exogenous variables explicitly, as in Figure 16a.



(a) Model 18 showing exogenous variables

(b) Model 18 showing Y_x

Figure 16: Model 18 explained

We can now clearly see the colliding path $X \rightarrow Y \leftarrow U_y$. Conditioning on Z , a descendant of Y , thus partially opens this path, and creates a spurious association between X and U_y , the disturbance of Y , making Z a “bad control.” Another approach to see why Z is a bad control is to explicitly draw the potential outcome Y_x in the DAG, as shown in Figure 16b. As explained, recall that Y_x is the solution of Y in the modified model M_x . This results in $Y_x = f_Y(x, U_y)$, a function of the random variable U_y . Therefore, conditioning on Z , a descendant of Y , partially opens the path $X \rightarrow Y \leftarrow U_y \rightarrow Y_x$, and thus $Y_x \not\perp\!\!\!\perp X \mid \mathbf{Z}$.

Now let us consider the case in which the arrow $X \rightarrow Y$ is removed (zero causal effect of X on Y). First recall that two nodes X and Y are *d-separated* conditional on \mathbf{Z} if the

set \mathbf{Z} blocks every path from X to Y in the graph. If X and Y are d -separated conditional on \mathbf{Z} , this implies the conditional independence $Y \perp\!\!\!\perp X \mid \mathbf{Z}$. In Model 18, when there is no path from X to Y , conditioning on Z also does not open any other paths between these two variables. Hence, X is still d -separated from Y even after conditioning on Z , and the conditional independence $Y \perp\!\!\!\perp X \mid \mathbf{Z}$ holds.

Analytical expressions for linear models

Here we provide algebraic derivations for each illustrative model under the assumption of *linearity* of the structural equations. Before proceeding, we remind readers that adjusting for “bad controls” still lead to bias in non-parametric models. Furthermore, overt selection (rather than adjustment for) bad controls will also lead to bias (although the size and sign of the bias may differ).

Without loss of generality, we assume random variables have been standardized to have mean zero and unit variance. We use σ_{yx} to denote the covariance of Y and X and, like before, $\beta_{yx.z}$ to denote the partial regression coefficient of the regression of Y on X controlling for Z . The partial regression coefficient $\beta_{yx.z}$ can be written in terms of the covariances as (Cramér, 1946),

$$\beta_{yx.z} = \frac{\sigma_{yx} - \sigma_{xz}\sigma_{yz}}{1 - \sigma_{xz}^2} \quad (2)$$

In linear structural causal models, each edge $V_i \rightarrow V_j$ of a causal DAG can be mapped to a single structural coefficient $\lambda_{v_i v_j}$ representing the strength of the direct effect of V_i on V_j . We can use Wright’s path-tracing rules (Wright, 1921) to equate the covariance $\sigma_{v_i v_j}$ of any two variables V_i and V_j , to the sum of products of structural coefficients along *unblocked* paths between V_i and V_j .

For instance, in Model 1, path-tracing results in the following covariances,

$$\sigma_{yx} = \lambda_{xy} + \lambda_{zx}\lambda_{zy} \quad (3)$$

$$\sigma_{xz} = \lambda_{zx} \quad (4)$$

$$\sigma_{yz} = \lambda_{zy} + \lambda_{zx}\lambda_{xy} \quad (5)$$

With the help of Wright’s rules and Equation 2, one can easily proceed with the algebraic derivation for each model. For the sake of brevity, we provide the full derivation for Model 1, and for the remaining models (except Model 15) only the final result is presented (since the derivations would be very similar).

Model 1. In Model 1, the average causal effect of X on Y equals $\text{ACE} = \lambda_{xy}$. The unadjusted regression coefficient equals $\beta_{yx} = \sigma_{yx} = \lambda_{xy} + \lambda_{zx}\lambda_{zy}$, and the partial regression coefficient equals,

$$\beta_{yx.z} = \frac{\sigma_{yx} - \sigma_{xz}\sigma_{yz}}{1 - \sigma_{xz}^2} \quad (6)$$

$$= \frac{\lambda_{xy} + \lambda_{zx}\lambda_{zy} - \lambda_{zx}(\lambda_{zy} + \lambda_{zx}\lambda_{xy})}{1 - \lambda_{zx}^2} \quad (7)$$

$$= \frac{\lambda_{xy} + \lambda_{zx}\lambda_{zy} - \lambda_{zx}\lambda_{zy} - \lambda_{zx}^2\lambda_{xy}}{1 - \lambda_{zx}^2} \quad (8)$$

$$= \frac{(1 - \lambda_{zx}^2)\lambda_{xy}}{1 - \lambda_{zx}^2} \quad (9)$$

$$= \lambda_{xy} \quad (10)$$

Model 2. Path-tracing leads to the following covariances,

$$\sigma_{yx} = \lambda_{xy} + \lambda_{zx}\lambda_{uz}\lambda_{uy} \quad (11)$$

$$\sigma_{xz} = \lambda_{zx} \quad (12)$$

$$\sigma_{yz} = \lambda_{zx}\lambda_{xy} + \lambda_{uz}\lambda_{uy} \quad (13)$$

We have: $\text{ACE} = \lambda_{xy}$, $\beta_{yx} = \lambda_{xy} + \lambda_{zx}\lambda_{uz}\lambda_{uy}$, and $\beta_{yx.z} = \lambda_{xy}$.

Model 3. Path-tracing leads to the following covariances,

$$\sigma_{yx} = \lambda_{xy} + \lambda_{ux}\lambda_{uz}\lambda_{zy} \quad (14)$$

$$\sigma_{xz} = \lambda_{ux}\lambda_{uz} \quad (15)$$

$$\sigma_{yz} = \lambda_{zy} + \lambda_{uz}\lambda_{ux}\lambda_{uy} \quad (16)$$

We have: $\text{ACE} = \lambda_{xy}$, $\beta_{yx} = \lambda_{xy} + \lambda_{ux}\lambda_{uz}\lambda_{zy}$, and $\beta_{yx.z} = \lambda_{xy}$.

Model 4. Path-tracing leads to the following covariances,

$$\sigma_{yx} = \lambda_{xm}\lambda_{my} + \lambda_{zx}\lambda_{zm}\lambda_{my} \quad (17)$$

$$\sigma_{xz} = \lambda_{zx} \quad (18)$$

$$\sigma_{yz} = \lambda_{zm}\lambda_{my} + \lambda_{zx}\lambda_{xm}\lambda_{my} \quad (19)$$

We have: $\text{ACE} = \lambda_{xm}\lambda_{my}$, $\beta_{yx} = \lambda_{xm}\lambda_{my} + \lambda_{zx}\lambda_{zm}\lambda_{my}$, and $\beta_{yx.z} = \lambda_{xm}\lambda_{my}$.

Model 5. Path-tracing leads to the following covariances,

$$\sigma_{yx} = \lambda_{xm}\lambda_{my} + \lambda_{zx}\lambda_{uz}\lambda_{um}\lambda_{my} \quad (20)$$

$$\sigma_{xz} = \lambda_{zx} \quad (21)$$

$$\sigma_{yz} = \lambda_{uz}\lambda_{um}\lambda_{my} + \lambda_{zx}\lambda_{xm}\lambda_{my} \quad (22)$$

We have: $\text{ACE} = \lambda_{xm}\lambda_{my}$, $\beta_{yx} = \lambda_{xm}\lambda_{my} + \lambda_{zx}\lambda_{uz}\lambda_{um}\lambda_{my}$, and $\beta_{yx.z} = \lambda_{xm}\lambda_{my}$.

Model 6. Path-tracing leads to the following covariances,

$$\sigma_{yx} = \lambda_{xm}\lambda_{my} + \lambda_{ux}\lambda_{uz}\lambda_{zm}\lambda_{my} \quad (23)$$

$$\sigma_{xz} = \lambda_{ux}\lambda_{uz} \quad (24)$$

$$\sigma_{yz} = \lambda_{zm}\lambda_{my} + \lambda_{uz}\lambda_{ux}\lambda_{xm}\lambda_{my} \quad (25)$$

We have: $\text{ACE} = \lambda_{xm}\lambda_{my}$, $\beta_{yx} = \lambda_{xm}\lambda_{my} + \lambda_{ux}\lambda_{uz}\lambda_{zm}\lambda_{my}$, and $\beta_{yx.z} = \lambda_{xm}\lambda_{my}$.

Model 7. Path-tracing leads to the following covariances,

$$\sigma_{yx} = \lambda_{xy} \quad (26)$$

$$\sigma_{xz} = \lambda_{u_1x}\lambda_{u_1z} \quad (27)$$

$$\sigma_{yz} = \lambda_{u_2z}\lambda_{u_2y} + \lambda_{u_1z}\lambda_{u_1x}\lambda_{xy} \quad (28)$$

We have: $\text{ACE} = \lambda_{xy}$, $\beta_{yx} = \lambda_{xy}$, and $\beta_{yx.z} = \lambda_{xy} - \frac{\lambda_{u_1x}\lambda_{u_1z}\lambda_{u_2z}\lambda_{u_2y}}{1 - (\lambda_{u_1x}\lambda_{u_1z})^2}$.

Variation of Model 7. Path-tracing leads to the following covariances,

$$\sigma_{yx} = \lambda_{xy} + \lambda_{u_1x}\lambda_{u_1z}\lambda_{zy} \quad (29)$$

$$\sigma_{xz} = \lambda_{u_1x}\lambda_{u_1z} \quad (30)$$

$$\sigma_{yz} = \lambda_{zy} + \lambda_{u_2z}\lambda_{u_2y} + \lambda_{u_1z}\lambda_{u_1x}\lambda_{xy} \quad (31)$$

We have: $\text{ACE} = \lambda_{xy}$, $\beta_{yx} = \lambda_{xy} + \lambda_{u_1x}\lambda_{u_1z}\lambda_{zy}$, and $\beta_{yx.z} = \lambda_{xy} - \frac{\lambda_{u_1x}\lambda_{u_1z}\lambda_{u_2z}\lambda_{u_2y}}{1 - (\lambda_{u_1x}\lambda_{u_1z})^2}$. Thus, depending on the parameterization of the model, the absolute value of the bias of $\beta_{yx.z}$ can be greater than that of β_{yx} .

Model 8. Path-tracing leads to the following covariances,

$$\sigma_{yx} = \lambda_{xy} \quad (32)$$

$$\sigma_{xz} = 0 \quad (33)$$

$$\sigma_{yz} = \lambda_{zy} \quad (34)$$

We have: $\text{ACE} = \lambda_{xy}$, $\beta_{yx} = \lambda_{xy}$, and $\beta_{yx.z} = \lambda_{xy}$.

Model 9. Path-tracing leads to the following covariances,

$$\sigma_{yx} = \lambda_{xy} \quad (35)$$

$$\sigma_{xz} = \lambda_{zx} \quad (36)$$

$$\sigma_{yz} = \lambda_{zx}\lambda_{xy} \quad (37)$$

We have: $\text{ACE} = \lambda_{xy}$, $\beta_{yx} = \lambda_{xy}$, and $\beta_{yx.z} = \lambda_{xy}$.

Model 10. Path-tracing leads to the following covariances,

$$\sigma_{yx} = \lambda_{xy} + \lambda_{ux}\lambda_{uy} \quad (38)$$

$$\sigma_{xz} = \lambda_{zx} \quad (39)$$

$$\sigma_{yz} = \lambda_{zx}\lambda_{xy} \quad (40)$$

We have: $\text{ACE} = \lambda_{xy}$, $\beta_{yx} = \lambda_{xy} + \lambda_{ux}\lambda_{uy}$, and $\beta_{yx.z} = \lambda_{xy} + \frac{\lambda_{ux}\lambda_{uy}}{1-\lambda_{zx}^2}$. Since $0 < 1 - \lambda_{zx}^2 < 1$, the absolute value of the bias of $\beta_{yx.z}$ is always greater than that of β_{yx} .

Model 11. Path-tracing leads to the following covariances,

$$\sigma_{yx} = \lambda_{xz}\lambda_{zy} \quad (41)$$

$$\sigma_{xz} = \lambda_{xz} \quad (42)$$

$$\sigma_{yz} = \lambda_{zy} \quad (43)$$

We have: $\text{ACE} = \lambda_{xz}\lambda_{zy}$, $\text{CDE} = 0$, $\beta_{yx} = \lambda_{xz}\lambda_{zy}$, $\beta_{yx.z} = 0$.

Model 12. Path-tracing leads to the following covariances,

$$\sigma_{yx} = \lambda_{xm}\lambda_{my} \quad (44)$$

$$\sigma_{xz} = \lambda_{xm}\lambda_{mz} \quad (45)$$

$$\sigma_{yz} = \lambda_{mz}\lambda_{my} \quad (46)$$

We have: $\text{ACE} = \lambda_{xm}\lambda_{my}$, $\text{CDE} = 0$, $\beta_{yx} = \lambda_{xm}\lambda_{my}$, and $\beta_{yx.z} = \lambda_{xm}\lambda_{my} \times \left(\frac{1-\lambda_{mz}^2}{1-\lambda_{xm}^2\lambda_{mz}^2} \right)$.

Variation of Model 11. Path-tracing leads to the following covariances,

$$\sigma_{yx} = \lambda_{xz}\lambda_{zy} \quad (47)$$

$$\sigma_{xz} = \lambda_{xz} \quad (48)$$

$$\sigma_{yz} = \lambda_{zy} + \lambda_{uz}\lambda_{uy} \quad (49)$$

We have: $\text{ACE} = \lambda_{xz}\lambda_{zy}$, $\text{CDE} = 0$, $\beta_{yx} = \lambda_{xz}\lambda_{zy}$, and $\beta_{yx.z} = -\frac{\lambda_{xz}\lambda_{uz}\lambda_{uy}}{1-\lambda_{xz}^2}$.

Model 13. Path-tracing leads to the following covariances,

$$\sigma_{yx} = \lambda_{xm}\lambda_{my} \quad (50)$$

$$\sigma_{xz} = 0 \quad (51)$$

$$\sigma_{yz} = \lambda_{zm}\lambda_{my} \quad (52)$$

We have: $\text{ACE} = \lambda_{xm}\lambda_{my}$, $\beta_{yx} = \lambda_{xm}\lambda_{my}$, and $\beta_{yx.z} = \lambda_{xm}\lambda_{my}$.

Model 14. Path-tracing leads to the following covariances,

$$\sigma_{yx} = \lambda_{xy} \quad (53)$$

$$\sigma_{xz} = \lambda_{xz} \quad (54)$$

$$\sigma_{yz} = \lambda_{xz}\lambda_{xy} \quad (55)$$

We have: $\text{ACE} = \lambda_{xy}$, $\beta_{yx} = \lambda_{xy}$, and $\beta_{yx.z} = \lambda_{xy}$.

Model 15. Model 15 requires a more elaborate derivation since we need to consider four variables. Here we only have samples with $W = w$ recorded. In linear models, this is equivalent to having adjusted for W in a regression model. Thus, this means the researcher does not have access to the regression coefficients β_{yx} nor $\beta_{yx.z}$, but rather $\beta_{yx.w}$ and $\beta_{yx.wz}$ (that is, W is *always* conditioned on, due to sample selection). Path-tracing leads to the following covariances (there is no need to compute all of them to solve this problem, but we show them here for completeness),

$$\sigma_{yx} = \lambda_{xy} \quad (56)$$

$$\sigma_{xz} = \lambda_{xz} \quad (57)$$

$$\sigma_{xw} = \lambda_{xz}\lambda_{zw} \quad (58)$$

$$\sigma_{yz} = \lambda_{xz}\lambda_{xy} \quad (59)$$

$$\sigma_{yw} = \lambda_{uy}\lambda_{uw} + \lambda_{xy}\lambda_{xz}\lambda_{zw} \quad (60)$$

$$\sigma_{zw} = \lambda_{zw} \quad (61)$$

The average causal effect is $\text{ACE} = \lambda_{xy}$. By Equation 2, the regression coefficient without adjusting for Z , $\beta_{yx.w}$, equals

$$\beta_{yx.w} = \lambda_{xy} - \frac{\lambda_{xz}\lambda_{zw}\lambda_{uw}\lambda_{uy}}{1 - (\lambda_{xz}\lambda_{zw})^2} \quad (62)$$

Now we must compute the regression coefficient adjusting for Z , $\beta_{yx.wz}$. Following Cramér (1946) $\beta_{yx.wz}$ can be written as

$$\beta_{yx.wz} = \frac{\rho_{yx.z} - \rho_{xw.z}\rho_{yw.z}}{(1 - \rho_{xw.z}^2)^{1/2}} \times \frac{\sigma_{y.z}}{\sigma_{x.z}} \quad (63)$$

Where

$$\rho_{yx.z} = \frac{\sigma_{yx} - \sigma_{xz}\sigma_{yz}}{(1 - \sigma_{xz}^2)^{1/2}(1 - \sigma_{yz}^2)^{1/2}} \quad (64)$$

denotes the partial correlation of Y with X after adjusting for Z and

$$\sigma_{y.z} = (1 - \sigma_{yz}^2)^{1/2}$$

denotes the partial standard deviation of Y after adjusting for Z . Since X is d -separated from W given Z , we know that $\rho_{xw.z} = 0$ (we can also verify this by checking that $\beta_{xw.z} = 0$). Also note that $\sigma_{yz} = \lambda_{xz}\lambda_{xy} = \sigma_{xz}\sigma_{yx}$. We thus obtain

$$\beta_{yx.wz} = \rho_{yx.z} \times \frac{(1 - \sigma_{yz}^2)^{1/2}}{(1 - \sigma_{xz}^2)^{1/2}} \quad (65)$$

$$= \frac{\sigma_{yx} - \sigma_{xz}\sigma_{yz}}{(1 - \sigma_{xz}^2)^{1/2}(1 - \sigma_{yz}^2)^{1/2}} \times \frac{1 - \sigma_{yz}^2}{1 - \sigma_{xz}^2} \quad (66)$$

$$= \frac{\sigma_{yx} - \sigma_{xz}^2\sigma_{yx}}{1 - \sigma_{xz}^2} \quad (67)$$

$$= \sigma_{yx} \frac{1 - \sigma_{xz}^2}{1 - \sigma_{xz}^2} = \sigma_{yx} = \lambda_{xy} \quad (68)$$

Model 16. Path-tracing leads to the following covariances,

$$\sigma_{yx} = \lambda_{xy} \quad (69)$$

$$\sigma_{xz} = \lambda_{xz} \quad (70)$$

$$\sigma_{yz} = \lambda_{xz}\lambda_{xy} + \lambda_{uz}\lambda_{uy} \quad (71)$$

We have: $\text{ACE} = \lambda_{xy}$, $\beta_{yx} = \lambda_{xy}$, and $\beta_{yx.z} = \lambda_{xy} - \frac{\lambda_{xz}\lambda_{uz}\lambda_{uy}}{1 - \lambda_{xz}^2}$.

Model 17. Path-tracing leads to the following covariances,

$$\sigma_{yx} = \lambda_{xy} \quad (72)$$

$$\sigma_{xz} = \lambda_{xz} + \lambda_{xy}\lambda_{yz} \quad (73)$$

$$\sigma_{yz} = \lambda_{yz} + \lambda_{xz}\lambda_{xy} \quad (74)$$

We have: $\text{ACE} = \lambda_{xy}$, $\beta_{yx} = \lambda_{xy}$, and $\beta_{yx.z} = \lambda_{xy} \times \frac{(1 - \lambda_{xz}\sigma_{xz})}{1 - \sigma_{xz}^2} - \frac{\sigma_{xz}\lambda_{yz}}{1 - \sigma_{xz}^2}$.

Model 18. Path-tracing leads to the following covariances,

$$\sigma_{yx} = \lambda_{xy} \quad (75)$$

$$\sigma_{xz} = \lambda_{xy}\lambda_{yz} \quad (76)$$

$$\sigma_{yz} = \lambda_{yz} \quad (77)$$

We have: $\text{ACE} = \lambda_{xy}$, $\beta_{yx} = \lambda_{xy}$, and $\beta_{yx.z} = \lambda_{xy} \times \frac{1 - \lambda_{yz}^2}{1 - \lambda_{xy}^2\lambda_{yz}^2}$.