

35 Causality

This chapter was written by Victor Veitch and Alex D’Amour.

35.1 Introduction

35.1.1 Why is causality different than other forms of ML?

The bulk of machine learning considers relationships between observed variables with the goal of summarizing these relationships in a manner that allows predictions on similar data. However, for many problems, our main interest is to predict how system would change if it were observed under different conditions. For instance, in healthcare, we are interested in whether a patient will recover if given a certain treatment (as opposed to whether treatment and recovery are associated in the observed data). **Causal inference** addresses how to formalize such problems, determine whether they can be solved, and, if so, how to solve them. This chapter covers the fundamentals of this subject. Code examples for the discussed methods are available at <https://github.com/vveitch/causality-tutorials>.

To make the gap between observed data modeling and causal inference concrete, consider the relationships depicted in Figure 35.1 and Figure 35.2. Figure 35.1 shows the relationship between deaths by drowning and ice cream production in the United States in 1931 (the pattern holds across most years). Figure 35.2 shows the relationship between smoking and lung cancer across various countries. In each case, there is a strong positive association. Faced with this association, we might ask: could we reduce drowning deaths by banning ice cream? Could we reduce lung cancer by banning cigarettes? We intuitively understand that these interventional questions have different answers, despite the fact that the observed associations are similar. Determining the causal effect of some intervention in the world requires some such causal hypothesis about the world.

For concreteness, consider three possible explanations for the association between ice cream and drowning. Perhaps eating ice cream does cause people to drown—due to stomach cramps or similar. Or, perhaps, drownings increase demand for ice cream—the survivors eat huge quantities of ice cream to handle their grief. Or, the association may be due (at least in part) to a common cause: warm weather makes people more likely to eat ice cream and more likely to go swimming (and, hence, to drown). Under all three scenarios, we can observe exactly the same data, but the implications for an ice cream ban are very different. Hence, answering questions about what will happen under an intervention requires us to incorporate some causal knowledge of the world—e.g., which of these scenarios is plausible?

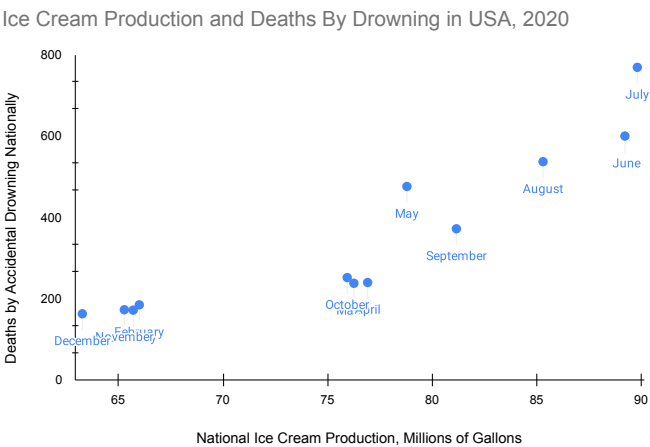


Figure 35.1: Ice cream production is strongly associated with deaths by drowning. Ice cream production data from the US Department of Agriculture National Agricultural Statistics Service. Drowning data from the National Center for Health Statistics at the United States Centers for Disease Control.

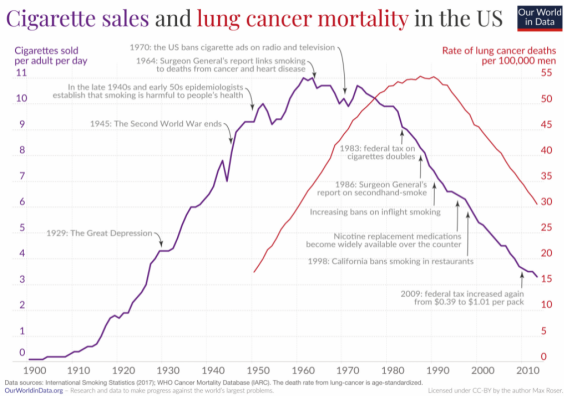


Figure 35.2: Smoking is strongly associated with lung cancer. Figure by Max Roser, ourworldindata.org/smoking-big-problem-in-brief

Our goal in this chapter is to introduce the essentials of estimating causal effects. The high-level approach has three steps.

- **Causal Estimands:** The first step is to formally define the quantities we want to estimate. These are summaries of how the world would change under intervention, rather than summaries of the world as it has already been observed. E.g., we want to formalize “The expected number of drownings in the United States if we ban ice cream”.
- **Identification:** The next step is to identify the causal estimands with quantities that can, in principle, be estimated from observational data. This step involves codifying our causal knowledge of the world and translating this into a statement such as, “The causal effect is equal to the expected number of drownings after adjusting for month”. This step tells us what causal questions we could answer with perfect knowledge of the observed data distribution.
- **Estimation:** Finally, we must estimate the observable quantity using a finite data sample. The form of causal estimands favors certain efficient estimation procedures that allow us to exploit non-parametric (e.g., machine learning) predictive models.

In this chapter, we’ll mainly focus on the estimation of the causal effect of an intervention averaged over all members of a population, known as the **Average Treatment Effect** or **ATE**. This is the most common problem in applied causal inference work. It is in some sense the simplest problem, and will allow us to concretely explain the use and importance of the fundamental causal concepts. These causal concepts include structural causal models, causal graphical models, the do-calculus, and efficient estimation using influence function techniques. This problem is also useful for understanding the role that standard predictive modeling and machine learning play in estimating causal quantities.

35.2 Causal Formalism

In causal inference, the goal is to use data to learn about how the outcome in the world would change under intervention. In order to make such inferences, we must also make use of our causal knowledge of the world. This requires a formalism that lets us make the notion of intervention precise and lets us encode our causal knowledge as assumptions.

35.2.1 Structural Causal Models

Consider a setting in which we observe four variables from a population of people: A_i , an indicator of whether or not person i smoked at a particular age, Y_i , an indicator of whether or not person i developed lung cancer at a later age, H_i , a “health consciousness” index that measures a person’s health-consciousness (perhaps constructed from a set of survey responses about attitudes toward health), and G_i , an indicator for whether the person has a genetic predisposition towards cancer. Suppose we observe a dataset of these variables drawn independently and identically from a population, $(A_i, Y_i, H_i) \stackrel{\text{iid}}{\sim} P^{\text{obs}}$, where “obs” stands for “observed”.

In standard practice, we model data like these using probabilistic models. Notably, there are many different ways to specify a probabilistic model for the same observed distribution. For example, we

could write a probabilistic model for P^{obs} as

$$A \sim P^{\text{obs}}(A) \quad (35.1)$$

$$H|A \sim P^{\text{obs}}(H|A) \quad (35.2)$$

$$Y|A, H \sim P^{\text{obs}}(Y|H, A) \quad (35.3)$$

$$G|A, H, Y \sim P^{\text{obs}}(G|A, H, Y) \quad (35.4)$$

This is a valid factorization, and sampling variables in this order would yield valid samples from the joint distribution P^{obs} . However, this factorization does not map well to a mechanistic understanding of how these variables are causally related in the world. In particular, it is perhaps more plausible that health-consciousness H causally precedes smoking status A , since a person's health-consciousness would influence their decision to smoke.

These intuitions about causal ordering are intimately tied to the notion of intervention. Here, we will focus on a notion of intervention that can be represented in terms of “structural” models that describe mechanistic relationships between variables. The fundamental objects that we will reason about are **structural causal models**, or SCM's. SCM's resemble probabilistic models, but they encode additional assumptions. Specifically, SCM's serve two purposes: they describe a probabilistic model *and* they provide semantics for transforming the data-generating process through intervention.

Formally, SCM's describe a mechanistic data generating process with an ordered sequence of equations that resemble assignment operations in a program. Each variable in a system is determined by combining other modeled variables (the causes) with exogenous “noise” according to some (unknown) deterministic function. For instance, a plausible SCM for P^{obs} might be

$$G \leftarrow f_G(\xi_0) \quad (35.5)$$

$$H \leftarrow f_H(\xi_1) \quad (35.6)$$

$$A \leftarrow f_A(H, \xi_2) \quad (35.7)$$

$$Y \leftarrow f_Y(G, H, A, \xi_3) \quad (35.8)$$

where the (unknown) functions f are fixed, and the variables ξ are unmeasured causes, modeled as independent random “noise” variables. Conceptually, the functions f_G, f_H, f_A, f_Y describe deterministic physical relationships in the real world, while the variables ξ are hidden causes that are sufficient to distinguish each unit i in the population. Because we assume that each observed unit i is drawn at random from a the population, we model ξ as random noise.

SCM's imply probabilistic models, but not the other way around. For example, our example SCM implies probabilistic model for the observed data based on the factorization $P^{\text{obs}}(G, H, A, Y) = P^{\text{obs}}(G)P^{\text{obs}}(H)P^{\text{obs}}(A|H)P^{\text{obs}}(Y|A, H)$. Thus, we could sample from the SCM in the same way we would from a probabilistic model: draw a set of noise variables ξ and evaluate each assignment operation in the SCM in order.

Beyond the probabilistic model, an SCM encodes additional assumptions about the effects of interventions. In an SCM, interventions are represented by replacing assignment statements. For example, if we were interested in the distribution of Y in the hypothetical scenario that smoking were eliminated, we could set the second line of the SCM to be $A \leftarrow 0$. Because the f functions in the SCM are assumed to be invariant mechanistic relationships, the SCM encodes the assumption that this edited SCM generates data that we would see if we really applied this intervention in

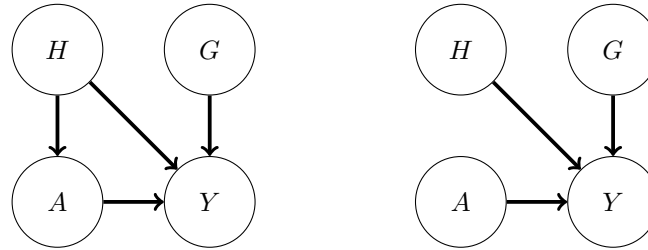


Figure 35.3: (Left) Causal graph illustrating relationships between smoking A , cancer Y , health consciousness H , and genetic cancer pre-disposition G . (Right) “Mutilated” causal graph illustrating relationships under an intervention on smoking A .

the world. Thus, the ordering of statements in an SCM are load-bearing: they imply substantive assumptions about how the world changes in response to interventions. This is in contrast to more standard probabilistic models where variables can be rearranged by applications of Bayes Rule without changing the substantive implications of the model.

We note that structural causal model may not incorporate all possible notions of causality. For example, laws based on conserved quantities or equilibria—e.g., the ideal gas law—do not trivially map to SCMs, though these are fundamental in disciplines such as physics and economics. Nonetheless, we will confine our discussion to SCMs.

35.2.2 Causal DAGs

SCMs encode many details about the assumed generative process of a system, but often it is useful to reason about causal problems at a higher level of abstraction. In particular, it is often useful to separate the causal structure of a problem from the particular functional form of those causal relationships. Causal graphs provide this level of abstraction. A causal graph specifies which variables causally affect other variables, but leaves the parametric form of the structural equations f unspecified. Given an SCM, the corresponding causal graph can be drawn as follows: for each line of the SCM, draw arrows from the variables on the right hand side to variables on the left hand side. The causal DAG for our smoking-cancer example is shown in Figure 35.3. In this way, causal DAGs are related to SCMs in the same way that probabilistic graphical models (PGMs) are related to probabilistic models.

In fact, in the same way that SCMs imply a probabilistic model, causal DAGs imply a PGM. Functionally, causal graphs behave as probabilistic graphical models (Chapter 4). They imply conditional independence relationships between the variables in the observed data in same way. They obey the Markov property: If $X \leftarrow Y \rightarrow Z$ then $X \perp\!\!\!\perp Z|Y$; recall d-separation (Section 4.2.4.1). Additionally, if $X \rightarrow Y \leftarrow Z$ then, usually, $X \not\perp\!\!\!\perp Z|Y$ (even if X and Z are marginally independent). In this case, Y is called a **collider** for X and Z .

Conceptually, the difference between causal DAGs and PGMs is that probabilistic graphical models encode our assumptions about statistical relationships, whereas causal graphs encode our (stronger) assumptions about causal relationships. Such causal relationships can be used to derive how statistical relationships would change under intervention.

Causal graphs also allow us to reason about the causal and non-causal origins of statistical dependencies in observed data without specifying a full SCM. In a causal graph, two variables—say, A and D —can be statistically associated in different ways. First, there can be a directed path from (ancestor) A to (descendant) D . In this case, A is a causal ancestor of D and interventions on A will propagate through to change D ; $P(D|\text{do}(A = a)) \neq P(D|\text{do}(A = a'))$. For example, smoking is a causal ancestor of cancer in our example. Alternatively, A and D could share a common cause—there is some variable C such that there is a directed path from C to A and from C to D . If A and D are associated only through such a path then interventions on A will not change the distribution of D . However, it is still the case that $P(D|A = a) \neq P(D|A = a')$ —observing different values of A changes our guess for the value of D . The reason is that A carries information about C , which carries information about D . For example, suppose we lived in a world where there was no effect of smoking on developing cancer (e.g., everybody vapes), there would nevertheless be an association between smoking and cancer because of the path $A \leftarrow H \rightarrow Y$. The existence of such “backdoor paths” is one core reason that statistical and causal association are not the same. Of course, more complicated variants of these associations are possible—e.g., C is itself only associated with A through a backdoor path—but this already captures the key distinction between causal and non-causal paths.

Recall that our aim in introducing SCMs and causal graphs is to enable us to formalize our causal knowledge of the world and to make precise what interventional quantities we’d like to estimate. Writing down a causal graph gives a simple formal way to encode our knowledge of the causal structure of a problem. Usefully, this causal structure is sufficient to directly reason about the implications of interventions without fully specifying the underlying SCM. The key observation is that if a variable A is intervened on then, after intervention, none of the other variables are causes of A . That is, when we replace a line of an SCM with a statement directly assigning a variable a particular value, we cut off all dependencies that variable had on its causal parents. Accordingly, in the causal graph, the intervened on variable has no parents. This leads us to the **graph surgery** notion of intervention: an intervention that sets A to a is the operation that deletes all incoming edges to A in the graph, and then conditions on $A = a$ in the resulting probability distribution (which is defined by the conditional independence structure of the post-surgery graph). We’ll use Pearl’s do notation to denote this operation. $P(\mathbf{X}|\text{do}(A = a))$ is the distribution of \mathbf{X} given $A = a$ under the mutilated graph that results from deleting all edges going into A . Similarly, $\mathbb{E}[\mathbf{X}|\text{do}(A = a)] \triangleq \mathbb{E}_{P(\mathbf{X}|\text{do}(A=a))}[\mathbf{X}]$. Thus, we can formalize statements such as “The average effect of receiving drug A ” as

$$\text{ATE} = E[Y|\text{do}(A = 1)] - E[Y|\text{do}(A = 0)], \quad (35.9)$$

where ATE stands for Average Treatment Effect.

For concreteness, consider our running example. We contrast the distribution that results by conditioning on A with the distribution that results from intervening on A :

$$P(Y, H, G|A = a) = P(Y|H, G, A = a)P(G)P(H|A = a) \quad (35.10)$$

$$P(Y, H, G|\text{do}(A = a)) = P(Y|H, G, A = a)P(G)P(H) \quad (35.11)$$

The key difference between these two distributions is that the standard conditional distribution describes a population where health consciousness H has the distribution that we observe among individuals with smoking status $A = a$, while the interventional distribution described a population where health consciousness H follows the marginal distribution among all individuals. For example, we would expect $P(H | A = \text{smoker})$ to put more mass on lower values of H than the marginal

health consciousness distribution than the marginal distribution $P(H)$, which would also include non-smokers. The intervention distribution thus incorporates a hypothesis of how smoking would affect the subpopulation individuals who tend to be too health conscious to smoke in the observed data.

35.2.3 Identification

A central challenge in causal inference is that many different SCM's can produce identical distributions of observed data. This means that, on the basis of observed data alone, we cannot uniquely identify the SCM that generated it. This is true no matter how large of a data sample is available to us.

For example, consider the setting where there is a treatment A that may or may not have an effect on outcome Y , and where both the treatment and outcome are known to be affected by some *unobserved* common binary cause U . Now, we might be interested in the causal estimand $E[Y|\text{do}(A = 1)]$. In general, we can't learn this quantity from the observed data. The problem is that, we can't tell apart the case where the treatment has a strong effect from the case where the treatment has no effect, but $U = 1$ both causes people to tend to be treated and causes increases the probability of a positive outcome. The same observation shows we can't learn the (more complicated) interventional distribution $P(Y|\text{do}(A = 1))$ (if we could learn this, then we'd get the average effect automatically).

Thus, an important part of causal inference is to augment the observed data with knowledge about the underlying causal structure of the process under consideration. Often, these assumptions can narrow the space of SCM's sufficiently so that there is only one value of the causal estimand that is compatible with the observed data. We say that the causal estimand is **identified** or **identifiable** under a given set of assumptions if those assumptions are sufficient to provide a unique answer. There are many different sets of sufficient conditions that yield identifiable causal effects; we call each set of sufficient conditions an **identification strategy**.

Given a set of assumptions about the underlying SCM, the most common way to show that a causal estimand is identified is by construction. Specifically, if the causal estimand can be written entirely in terms of observable probability distributions, then it is identified. We call such a function of observed distributions a **statistical estimand**. Once such a statistical estimand has been recovered, we can then construct and analyze an estimator for that quantity using standard statistical tools. As an example of a statistical estimand, in the SCM above, it can be shown the ATE as defined in Equation (35.9), is equal to the following statistical estimand

$$\text{ATE} \stackrel{(*)}{=} \tau^{\text{ATE}} \triangleq E[E[Y|H, A = 1] - E[Y|H, A = 0]], \quad (35.12)$$

where the equality $(*)$ only holds because of some specific properties of the SCM. Note that the RHS above only involves conditional expectations between observed variables (there are no *do* operators), so τ^{ATE} is only a function of observable probability distributions.

There are many kinds of assumptions we might make about the SCM governing the process under consideration. For example, the following are assertions we might make about the system in our running example:

1. The probability of developing cancer is additive on the logit scale in A , G , and H (i.e., logistic regression is a well-specified model).
2. For each individual, smoking can never decrease the probability of developing cancer.

3. Whether someone smokes is influenced by their health consciousness H , but not by their genetic predisposition to cancer G .

These assumptions range from strong parametric assumptions fully specifying the form of the SCM equations, to non-parametric assumptions that only specify what the inputs to each equation are, leaving the form fully unspecified. Typically, assumptions that fully specify the parametric form are very strong, and would require far more detailed knowledge of the system under consideration than we actually have. The goal in identification arguments is to find a set of assumptions that are weak enough that they might be plausibly true for the system under consideration, but which are also strong enough to allow for identification of the causal effect.

If we are not willing to make any assumptions about the functional form of the SCM, then our assumptions are just about which variables affect (and do not affect) the other variables. In this sense, such which-affects-which assumptions are minimal. These assumptions are exactly the assumptions captured by writing down a (possibly incomplete) causal DAG, showing which variables are parents of each other variable. The graph may be incomplete because we may not know whether each possible edge is present in the physical system. For example, we might be unsure whether the gene G actually has a causal effect on health consciousness H . It is natural to ask to what extent we can identify causal effects only on the basis of partially specified causal DAGs. It turns out much progress can be made based on such non-parametric assumptions; we discuss this in detail in Section 35.8.

We will also discuss certain assumptions that cannot be encoded in a causal graph, but that are still weaker than assuming that full functional forms are known. For example, we might assume that the outcome is affected additively by the treatment and any confounders, with no interaction terms between them. These weaker assumptions can enable causal identification even when assuming the causal graph alone does not.

It is worth emphasizing that every causal identification strategy relies on assumptions that have some content that cannot be validated in the observed data. This follows directly from the ill-posedness of causal problems: if the assumptions used to identify causal quantities could be validated, that would imply that the causal estimand was identifiable from the observed data alone. However, since we know that there are many values of the causal estimand that are compatible with observed data, it follows that the assumptions in our identification strategy must have unobservable implications.

35.2.4 Counterfactuals and the Causal Hierarchy

Structural causal models let us formalize and study a hierarchy of different kinds of query about the system under consideration. The most familiar is observational queries: questions that are purely about statistical associations (e.g., “Are smoking and lung cancer associated in the population this sample was drawn from?”). Next is interventional queries: questions about causal relationships at the population level (e.g., “How much does smoking increase the probability of cancer in a given population?”). The rest of this chapter is focused on the definition, identification, and estimation of interventional queries. Finally, there are counterfactual queries: questions about causal relationships at the level of specific individuals, had something been different (e.g., “Would Alice have developed cancer had she not smoked?”). This causal hierarchy was popularized by [Pea09a, Ch. 1].

Interventional queries concern the prospective effect of an intervention on an outcome; for example, if we intervene and prevent a randomly sampled individual from smoking, what is the probability they develop lung cancer? Ultimately, the probability statement here is about our uncertainty about the “noise” variables ξ in the SCM. These are the unmeasured factors specific to

the randomly selected individual. The distribution is determined by the population from which that individual is sampled. Thus, interventional queries are statements about populations. Interventional queries can be written in terms of conditional distributions using **do-notation**, e.g.

$$P(Y|\text{do}(A = 0)) \quad (35.13)$$

where conditioning on the “do” clause means that the line defining A in the underlying SCM was changed to an intervention setting $A \leftarrow 0$. In our example, this represents the distribution of lung cancer outcomes for an individual selected at random and prevented from smoking.

Counterfactual queries concern how an observed outcome might have been different had an intervention been applied in the past. Counterfactual queries are often framed in terms of attributing a given outcome to a particular cause. For example, would Alice have developed cancer had she not smoked? Did most smokers with lung cancer develop cancer because they smoked? Counterfactual queries are so called because they require a comparison of counterfactual outcomes within individuals. In the formalism of SCM’s, counterfactual outcomes for an individual i are generated by running the same values of ξ_i through differently intervened SCM’s. Counterfactual outcomes are often written in terms of *potential outcomes* notation. In our running smoking example, this would look like:

$$Y_i(a) \triangleq f_Y(G_i, H_i, a, \xi_{3,i}). \quad (35.14)$$

That is, $Y_i(a)$ is the outcome we would have seen had A been set to a while all of $G_i, H_i, \xi_{3,i}$ were kept fixed.

It is important to understand what distinguishes interventional and fundamentally counterfactual queries. Just because a query can be written in terms of potential outcomes does not make it a counterfactual query. For example, the average treatment effect, which is the canonical interventional query, is easy to write in potential outcomes notation:

$$\text{ATE} = \mathbb{E}[Y_i(1) - Y_i(0)]. \quad (35.15)$$

Instead, the key dividing line between counterfactual and interventional queries is whether the query requires knowing the joint distribution of potential outcomes within individuals, or whether marginal distributions of potential outcomes across individuals will suffice. An important signature of a counterfactual query is conditioning on the value of one potential outcome. For example, “the lung cancer rate among smokers who developed cancer, had they not smoked” is a counterfactual query, and can be written as:

$$\mathbb{E}[Y_i(0) \mid Y_i(1) = 1, A_i = 1] \quad (35.16)$$

Answering this query requires knowing how individual-level cancer outcomes are related (through $\xi_{3,i}$) across the worlds where the each individual i did and did not smoke. Notably, this query cannot be rewritten using do-notation, because it requires a distinction between $Y(0)$ and $Y(1)$ while the ATE can: $\mathbb{E}[Y \mid \text{do}(A = 1)] - \mathbb{E}[Y \mid \text{do}(A = 0)]$.

Counterfactual queries require categorically more assumptions for identification than interventional ones. For identifying interventional queries, knowing the DAG structure of an SCM is often sufficient, while for counterfactual queries, some assumptions about the functional forms in the SCM are necessary. This is because only one potential outcome is ever observed for each individual, so the dependence between potential outcomes within individuals is not observable. For example, the data

in our running example provide no information on how individual-level smoking and non-smoking cancer risk are related. Thus, answering a question like “Did smokers who developed cancer have lower non-smoking cancer risk than smokers who did not develop cancer?”, requires additional assumptions about how characteristics encoded in ξ_i are translated to cancer outcomes. To answer this question without such assumptions, we would need to observe smokers who developed cancer in the alternate world where they did not smoke. Because they compare how individuals would have turned out under different generating processes, counterfactual queries are often referred to as “cross-world” quantities. On the other hand, interventional queries only require understanding the marginal distributions of potential outcomes $Y_i(0)$ and $Y_i(1)$ across individuals; thus, no cross-world information is necessary at the individual level.

We conclude this section by noting that counterfactual outcomes and potential outcomes notation are often conceptually useful, even if they are not used to explicitly answer counterfactual queries. Many causal queries are more intuitive to formalize in terms of potential outcomes. E.g., “Would I have smoked if I was more health conscious?” may be more intuitive than “Would a randomly sampled individual from the same population have smoked had they been subject to an intervention that made them more health conscious?”. In fact, some schools of causal inference use potential outcomes, rather than DAGs, as their primary conceptual building block [See IR15]. Causal graphs and potential outcomes both provide ways to formalize interventional queries and causal assumptions. Ultimately, these are mathematically equivalent. Nevertheless, practically, they have different strengths. The main advantage of potential outcomes is that counterfactual statements often map more directly to our mechanistic understanding of the world. This can make it easier to articulate causal desiderata and causal assumptions we may wish to use. On the other hand, the potential outcomes notation does not automatically distinguish between interventional and counterfactual queries. Additionally, causal graphs often give an intuitive and easy way of articulating assumptions about structural causal models involving many variables—potential outcomes get quickly unwieldy. In short: both formalizations have distinct advantages, and those advantages are simply about how easy it is to translate our causal understanding of the world into crisp mathematical assumptions.

35.3 Randomized Control Trials

We now turn to the business of estimating causal effects from data. We begin with **randomized control trials**, which are experiments designed to make the causal concerns as simple as possible.

The simplest situation for causal estimation is when there are no common causes of A and Y . The world is rarely so obliging as to make this the case. However, sometimes we can design an experiment to enforce the no-common-causes structure. In randomized control trials we assign each participant to either the treatment or control group at random. Because random assignment does not depend on any property of the units in the study, there are no causes of treatment assignment, and hence also no common causes of Y and A .

In this case, it’s straightforward to see that $P(Y|\text{do}(A = a)) = P(Y|a)$. This is essentially by definition of the graph surgery: since A has no parents, the mutilated graph is the same as the original graph. Indeed, the graph surgery definition is chosen to make this true: any sensible formalization of causality should have this identification result.

It is common to use RCTs to study the average treatment effect,

$$\text{ATE} = E[Y|\text{do}(A = 1)] - E[Y|\text{do}(A = 0)]. \quad (35.17)$$

This is the expected difference between being assigned treatment and assigned no treatment for a randomly chosen member of the population. It's easy to see that in an RCT this causal quantity is identified as a parameter τ^{RCT} of the observational distribution:

$$\tau^{\text{RCT}} = \mathbb{E}[Y|A = 1] - \mathbb{E}[Y|A = 0].$$

Then, a natural estimator is:

$$\hat{\tau}^{\text{RCT}} \triangleq \frac{1}{n_A} \sum_{i:A_i=1} Y_i - \frac{1}{n - n_A} \sum_{i:A_i=0} Y_i, \quad (35.18)$$

where n_A is the number of units who received treatment. That is, we estimate the average treatment effect as the difference between the average outcome of the treated group and the average outcome of the untreated (control) group.¹

Randomized control trials are the gold standard for estimating causal effects. This is because we know *by design* that there are no confounders that can produce alternative causal explanations of the data. In particular, the assumption of the triangle DAG—there are no unobserved confounders—is enforced by design. However, there are limitations. Most obviously, randomized control trials are sometimes infeasible to conduct. This could be due to expense, regulatory restrictions, or more fundamental difficulties (e.g., in developmental economics, the response of interest is sometimes collected decades after treatment). Additionally, it may be difficult to ensure that the participants in an RCT are representative of the population where the treatment will be deployed. For instance, participants in drug trials may skew younger and poorer than the population of patients who will ultimately take the drug.

35.4 Confounder Adjustment

We now turn to the problem of estimating causal effects using observational (i.e., not experimental) data. The most common application of causal inference is estimating the average treatment effect (ATE) of an intervention. The ATE is also commonly called the **average causal effect**, or ACE. Here, we focus on the important special case where the treatment A is binary, and we observe the outcome Y as well as a set of common causes X that influence both A and Y .

35.4.1 Causal Estimand, Statistical Estimand, and Identification

Consider a problem where we observe treatment A , outcome Y , and covariates X , which are drawn i.i.d. from some unknown distribution P . We wish to learn the average treatment effect: the expected difference between being assigned treatment and assigned no treatment for a randomly chosen member of the population. Following the discussion in the introduction, there are three steps to learning this quantity: mathematically formalize the causal estimand, give conditions for the causal estimand to be identified as a statistical estimand, and, finally, estimate this statistical estimand from data. We now turn to the first two steps.

1. There is a literature on efficient estimation of causal effects in RCT's going back to Fisher [Fis25] that employ more sophisticated estimators. See also Lin [Lin13] and Bloniarz et al. [Blo+16] for more modern treatments.

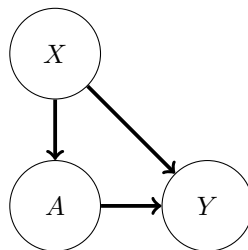


Figure 35.4: A causal DAG illustrating a situation where treatment A and outcome Y are both influenced by observed confounders X .

The average treatment effect is defined to be the difference between the average outcome if we intervened and set A to be 0, versus the average outcome if we intervened and set A to be 1. Using the do notation, we can write this formally as

$$\text{ATE} = \mathbb{E}[Y|\text{do}(A = 1)] - \mathbb{E}[Y|\text{do}(A = 0)]. \quad (35.19)$$

The next step is to articulate sufficient conditions for the ATE to be identified as a statistical estimand (a parameter of distribution P). The key issue is the possible presence of **confounders**. Confounders are “common cause” variables that affect both the treatment and outcome. When there are confounding variables in observed data, the sub-population of people who are observed to have received one level of the treatment A will differ from the rest of the population in ways that are relevant to their observed Y . For example, there is a strong positive association between horseback riding in childhood (treatment) and healthiness as an adult (outcome) [RB16]. However, both of these quantities are influenced by wealth X . The population of people who rode horses as children ($A = 1$) is wealthier than the population of people who did not. Accordingly, horseback-riding population will have better health outcomes even if there is no actual causal benefit of horseback riding for adult health.

We’ll express the assumptions required for causal identification in the form of a causal DAG. Namely, we consider the simple triangle DAG in Figure 35.4, where the treatment and outcome are influenced by *observed* confounders X . It turns out that the assumption encoded by this DAG suffices for identification. To understand why this is so, recall that the target causal effect is defined according to the distribution we would see if the edge from X to A was removed (that’s the meaning of do). The key insight is that because the intervention only modifies the relationship between X and A , the structural equation that generates outcomes Y given X and A , illustrated in Figure 35.4 as the $A \rightarrow Y \leftarrow X$, is the same even after the $X \rightarrow Y$ edge is removed. For example, we might believe that the physiological processes by which smoking status A and confounders X produce lung cancer Y remain the same, regardless of how the decision to smoke or not smoke was made. Secondly, because the intervention does not change the composition of the population, we would also expect the distribution of background characteristics X to be the same between the observational and intervened processes.

With these insights about invariances between observed and interventional data, we can derive a statistical estimand for the ATE as follows.

Theorem 2 (Adjustment with No Unobserved Confounders). *We observe $A, Y, X \sim P$. Suppose that*

1. (Confounders observed) *The data obeys the causal structure in Figure 35.4. In particular, X contains all common causes of A and Y and no variable in X is caused by A or Y .*
2. (Overlap) $0 < P(A = 1|X = x) < 1$ for all values of x . *That is, there are no individuals for whom treatment is always or never assigned.*

Then, the average treatment effect is identified as $ATE = \tau$, where

$$\tau = \mathbb{E}[\mathbb{E}[Y|A = 1, X]] - \mathbb{E}[\mathbb{E}[Y|A = 0, X]]. \quad (35.20)$$

Proof. First, we expand the ATE using the tower property of expectation, conditioning on X . Then, we apply the invariances discussed above:

$$ATE = \mathbb{E}[Y|\text{do}(A = 1)] - \mathbb{E}[Y|\text{do}(A = 0)] \quad (35.21)$$

$$= \mathbb{E}[\mathbb{E}[Y|\text{do}(A = 1), X]] - \mathbb{E}[\mathbb{E}[Y|\text{do}(A = 0), X]] \quad (35.22)$$

$$= \mathbb{E}[\mathbb{E}[Y|A = 1, X]] - \mathbb{E}[\mathbb{E}[Y|A = 0, X]] \quad (35.23)$$

The final equality is the key to passing from a causal to observational quantity. This follows because, from the causal graph, the conditional distribution of Y given A, X is the same in both the original graph, and in the mutilated graph created by removing the edge from X to A . This mutilated graph defines $P(Y|\text{do}(A = 1), X)$, so the equality holds.

The condition that $0 < P(A = 1|X = x) < 1$ is required for the first equality (the tower property) to be well defined. \square

Note that Equation (35.20) is a function of only conditional expectations and distributions that appear in the observed data distribution (in particular, it contains no “do” operators). Thus, if we can fully characterize the observed data distribution P , we can map that distribution to a unique ATE.

It is useful to note how τ differs from the naive estimand $\mathbb{E}[Y|A = 1] - \mathbb{E}[Y|A = 0]$ that just reports the treatment-outcome association without adjusting for confounding. The comparison is especially clear when we write out the outer expectation in τ explicitly as an integral over X :

$$\tau = \int \mathbb{E}[Y | A = 1, X]P(X)dX - \int \mathbb{E}[Y | A = 0, X]P(X)dX \quad (35.24)$$

We can write the naive estimand in a similar form by applying the tower property of expectation:

$$\mathbb{E}[Y | A = 1] - \mathbb{E}[Y | A = 0] = \int \mathbb{E}[Y | A = 1, X]P(X | A = 1)dX - \int \mathbb{E}[Y | A = 0, X]P(X | A = 0)dX \quad (35.25)$$

The key difference is the probability distribution over X that is being integrated over. The observational difference in means integrates over the over distinct conditional distributions of confounders X , depending on the value of A . On the other hand, in the ATE estimand τ , we integrate over the same distribution $P(X)$ for both levels of the treatment.

Overlap In addition to the assumption on the causal structure, identification requires that there is sufficient random variation in how treatments are assigned.

Definition 1. A distribution P on A, X satisfies **overlap** if $0 < P(A = 1|x) < 1$ for all x . It satisfies **strict overlap** if $\epsilon < P(A = 1|x) < 1 - \epsilon$ for all x and some $\epsilon > 0$.

Overlap is the requirement that any unit could have either recieved the treatment or not.

To see the necessity of overlap, consider estimating the effectiveness of a drug in a study where patient sex is a confounder, but the drug was only ever prescribed to male patients. Then, conditional on a patient being female, we would know that patient was assigned to control. Without further assumptions, it's impossible to know the effect of the drug on a population with female patients, because there would be no data to inform the expected outcome for treated female patients, that is, $\mathbb{E}[Y | A = 1, X = \text{female}]$. In this case, the statistical estimand (35.20) would not be identifiable. In the same vein, strict overlap ensures that the conditional distributions at each stratum of X can be estimated in finite samples.

Overlap can be particularly limiting in settings where we are adjusting for a large number of covariates (in an effort to satisfy no unobserved confounding). Then, certain combinations of traits may be very highly predictive of treatment assignment, even if individual traits are not. E.g., male patients over age 70 with BMI greater than 25 are very rarely assigned the drug. If such groups represent a significant fraction of the target population, or have significantly different treatment effects, then this issue can be problematic. In this case, the strict overlap assumption puts very strong restrictions on observational studies: for an observational study to satisfy overlap, most dimensions of the confounders X would need to closely mimic the balance we would expect in an RCT [D'A+21].

35.4.2 ATE Estimation with Observed Confounders

We now return to estimating the ATE using observed—i.e., not experimental—data. We've shown that in the case where we observe all common causes of the treatment and outcome, the ATE is causally identified with a statistical estimand τ . We now consider several strategies for estimating this quantity using a finite data sample. Broadly, these techniques are known as backdoor adjustment.²

Recall that the defining characteristic of a confounding variable is that it affects both treatment and outcome. Thus, an adjustment strategy may aim to account for the influence of confounders on the observed outcome, the influence of confounders on treatment, or both. We discuss each of these strategies in turn.

35.4.2.1 Outcome Model Adjustment

We begin with an approach to covariate adjustment that relies on modeling the conditional expectation of the outcome Y given treatment A and confounders X . This strategy is often referred to as g-computation or outcome adjustment.³ To begin, we define

Definition 2. The conditional expected outcome is the function Q given by

$$Q(a, x) = \mathbb{E}[Y | A = a, X = x]. \quad (35.26)$$

² As we discuss in Section 35.8, this backdoor adjustment references the estimand returned by the do-calculus to eliminate confounding from a backdoor path. This also generalizes the approaches discussed here to some cases where we do not observe all common causes.

³ The “g” stands for generalized, for now-inscrutable historical reasons [Rob86].

Substituting this definition into the definition of our estimand τ , Equation (35.20), we have $\tau = \mathbb{E}[Q(1, x) - Q(0, x)]$. This suggests a procedure for estimating τ : fit a model \hat{Q} for Q and then report

$$\hat{\tau}^Q \triangleq \frac{1}{n} \sum_i \hat{Q}(1, x_i) - \hat{Q}(0, x_i). \quad (35.27)$$

To fit \hat{Q} , recall that $E[Y|a, x] = \operatorname{argmin}_Q \mathbb{E}[(Y - Q(A, X))^2]$. That is, the minimizer (among all functions) of the squared loss risk is the conditional expected outcome.⁴ So, to approximate Q , we simply use mean squared error to fit a predictor that predicts Y from A and X .

The estimation procedure takes several steps. We first fit a model \hat{Q} to predict Y . Then, for each unit i , we predict that unit's outcome had they received treatment $\hat{Q}(1, x_i)$ and we predict their outcome had they not received treatment $\hat{Q}(0, x_i)$.⁵ If the unit actually did receive treatment ($a_i = 1$) then $\hat{Q}(0, x_i)$ is our guess about what would have happened in the counterfactual case that they did not. The estimated expected gain from treatment for this individual is $\hat{Q}(1, x_i) - \hat{Q}(0, x_i)$ —the difference in expected outcome between being treated and not treated. Finally, we estimate the outer expectation with respect to $P(X)$ —the true population distribution of the confounders—using the empirical distribution $\hat{P}(X) = 1/n \sum_i \delta_{x_i}$. In effect, this means we substitute the expectation (over an unknown distribution) by an average over the observed data.

Linear regression It's worth saying something more about the special case where Q is modeled as a linear function of both the treatment and all the covariates. That is, the case where we assume the identification conditions of Theorem 2 and we additionally assume that the true, causal law (the SCM) governing Y yields: $Q(A, X) = \mathbb{E}[Y|A, X] = \mathbb{E}[f_Y(A, X, \xi)|A, X] = \beta_0 + \beta_A A + \beta_X X$. Plugging in, we see that $Q(1, X) - Q(0, X) = \beta_A$ (and so also $\tau = \beta_A$). Then, the estimator for the average treatment effect reduces to the estimator for the regression coefficient β_A . This “fit linear regression and report the regression coefficient” remains a common way of estimating the association between two variables in practice. The expected-outcome-adjustment procedure here may be viewed as a generalization of this procedure that removes the linear parametric assumption.

35.4.2.2 Propensity Score Adjustment

Outcome model adjustment relies on modeling the relationship between the confounders and the outcome. A popular alternative is to model the relationship between the confounders and the treatment. This strategy adjusts for confounding by directly addressing sampling bias in the treated and control groups. This bias arises from the relationship between the confounders and the treatment. Intuitively, the effect of confounding may be viewed as due to the difference between $P(X|A = 1)$ and $P(X|A = 0)$ —e.g., the population of people who rode horses as children is wealthier than the population of people who did not. When we observe all confounding variables X , this degree of over- or under-representation can be adjusted away by reweighting samples such that the confounders X have the same distribution in the treated and control groups. When the confounders are balanced between the two groups, then any differences between them must be attributable to the treatment.

4. To be precise, this definition applies when X and Y are square-integrable, and the minimization taken over measurable functions.

5. this interpretation is justified by the same conditions as Theorem 2

A key quantity for balancing treatment and control groups is the **propensity score**, which summarises the relationship between confounders and treatment.

Definition 3. The **propensity score** is the function g given by $g(x) = P(A = 1|X = x)$.

To make use of the propensity score in adjustment, we first rewrite the estimand τ in a suggestive form:

$$\tau = \mathbb{E}\left[\frac{YA}{g(X)} - \frac{Y(1-A)}{1-g(X)}\right]. \quad (35.28)$$

This identity can be verified by noting that $\mathbb{E}[YA|X] = \mathbb{E}[Y|A = 1, X]P(A = 1|X) + 0$, rearranging for $\mathbb{E}[Y|A = 1, X]$, doing the same for $\mathbb{E}[Y|A = 0, X]$, and substituting in to Equation (35.20). Note that the identity is just a mathematical fact about the statistical estimand—it does not rely on any causal assumptions, and holds whether or not τ can be interpreted as a causal effect.

This expression suggests the **inverse probability of treatment weighted estimator**, or IPTW estimator:

$$\hat{\tau}^{\text{IPTW}} \triangleq \frac{1}{n} \sum_i \frac{Y_i A_i}{\hat{g}(X_i)} - \frac{Y_i(1 - A_i)}{1 - \hat{g}(X_i)}. \quad (35.29)$$

Here, \hat{g} is a estimate of the propensity score function. Recall from Section 13.2.1 that if a model is well-specified and the loss function is a proper scoring rule then risk minimizer $g^* = \operatorname{argmin}_g \mathbb{E}[L(A, g(X))]$ will be $g^*(X) = P(A = 1|X)$. That is, we can estimate the propensity score by fitting a model that predicts A from X . Cross-entropy and squared loss are both proper scoring rules, so we may use standard pipelines.

In summary, the procedure is to estimate the propensity score function (with machine learning), and then to plug the estimated propensity scores $\hat{g}(x_i)$ into Equation (35.29). The IPTW estimator computes a difference of weighted averages between the treated and untreated group. The effect is to upweight the outcomes of units who were unlikely to be treated but who nevertheless actually, by chance, recieved treatment (and similarly for untreated). Intuitively, such units are typical for the untreated population. So, their outcomes under treatment are informative about what would have happened had a typical untreated unit received treatment.

A word of warning is in order. Although the IPTW is asymptotically valid and popular in practice, it can be very unstable in finite samples. If estimated propensity scores are extreme for some values of x (that is, very close to 0 or 1), then the corresponding IPTW weights can be very large, resulting in a high-variance estimator. In some cases, this instability can be mitigated by instead using the Hajek version of the estimator.

$$\hat{\tau}^{\text{h-IPTW}} \triangleq \sum_i Y_i A_i \frac{1/\hat{g}(X_i)}{\sum_i A_i / \hat{g}(X_i)} - \sum_i Y_i (1 - A_i) \frac{1/(1-\hat{g}(X_i))}{\sum_i (1 - A_i) / (1-\hat{g}(X_i))}. \quad (35.30)$$

However, extreme weights can persist even after self-normalization, either because there are truly strata of X where treatment assignment is highly imbalanced, or because the propensity score estimation method has overfit. In such cases, it is common to apply heuristics such as weight clipping. See Khan and Ugander [KU21] for a longer discussion of inverse-propensity type estimators, including some practical improvements.

35.4.2.3 Double Machine Learning

We have seen how to estimate the average treatment effect using either the relationship between confounders and outcome, or the relationship between confounders and treatment. In each case, we follow a two step estimation procedure. First, we fit models for the expected outcome or the propensity score. Second, we plug these fitted models into a downstream estimator of the effect.

Unsurprisingly, the quality of the estimate of τ depends on the quality of the estimates \hat{Q} or \hat{g} . This is problematic because Q and g may be complex functions that require large numbers of samples to estimate. Even though we're only interested in the 1-dimensional parameter τ , the naive estimators described thus far can have very slow rates of convergence. This leads to unreliable inference or very large confidence intervals.

Remarkably, there are strategies for combining Q and g in estimators that, in principle, do better than using either Q or g alone. The **Augmented Inverse Probability of Treatment Weighted Estimator (AIPW)** is one such estimator. It is defined as

$$\hat{\tau}^{\text{AIPW}} \triangleq \frac{1}{n} \sum_i \hat{Q}(1, X_i) - \hat{Q}(0, X_i) + A_i \frac{Y_i - \hat{Q}(1, X_i)}{\hat{g}(x_i)} - (1 - A_i) \frac{Y_i - \hat{Q}(0, X_i)}{1 - \hat{g}(X_i)}. \quad (35.31)$$

That is, $\hat{\tau}^{\text{AIPW}}$ is the outcome adjustment estimator plus a stabilization term that depends on the propensity score. This estimator is a particular case of a broader class of estimators that are referred to as **semi-parametrically efficient** or **double machine-learning** estimators [Che+17e; Che+17d]. We'll use the later terminology here.

We now turn to understanding the sense in which double machine learning estimators are robust to misestimation of the nuisance functions Q and g . To this end, we define the influence curve of τ to be the function ϕ defined by⁶

$$\phi(X_i, A_i, Y_i; Q, g, \tau) \triangleq Q(1, X_i) - Q(0, X_i) + A_i \frac{Y_i - Q(1, X_i)}{g(x_i)} - (1 - A_i) \frac{Y_i - Q(0, X_i)}{1 - g(X_i)} - \tau. \quad (35.32)$$

By design, $\hat{\tau}^{\text{AIPW}} - \tau = \frac{1}{n} \sum_i \phi(\mathbf{X}_i; \hat{Q}, \hat{g}, \tau)$. We begin by considering what would happen if we simply knew Q and g , and didn't have to estimate them. In this case, the estimator would be $\hat{\tau}^{\text{ideal}} = \frac{1}{n} \sum_i \phi(\mathbf{X}_i; Q, g, \tau)$ and, by the central limit theorem, we would have:

$$\sqrt{n}(\hat{\tau}^{\text{ideal}} - \tau) \xrightarrow{d} \text{Normal}(0, \mathbb{E}[\phi(\mathbf{X}_i; Q, g, \tau)^2]). \quad (35.33)$$

This result characterizes the estimation uncertainty in the best possible case. If we knew Q and g , we could rely on this result for, e.g., finding confidence intervals for our estimate.

The question is: what happens when Q and g need to be estimated? For general estimators and nuisance function models, we don't expect the \sqrt{n} -rate of Equation (35.33) to hold. For instance, $\sqrt{n}(\hat{\tau}^Q - \tau)$ only converges if $\sqrt{n}\mathbb{E}[(\hat{Q} - Q)^2]^{\frac{1}{2}} \rightarrow 0$. That is, for the naive estimator we only get the \sqrt{n} rate for estimating τ if we can also estimate Q at the \sqrt{n} rate—a much harder task! This is the issue that the double machine learning estimator helps with.

6. Influence curves are the foundation of what follows, and the key to generalizing the analysis beyond the ATE. Unfortunately, going into the general mathematics would require a major digression, so we omit it. However, see references at the end of the chapter for some pointers to the relevant literature.

1 To understand how, we decompose the error in estimating τ as follows:

$$2 \sqrt{n}(\hat{\tau}^{\text{AIPTW}} - \tau) \quad (35.34)$$

$$3 = \frac{1}{\sqrt{n}} \sum_i \phi(\mathbf{X}_i; Q, g, \tau) \quad (35.35)$$

$$4 + \frac{1}{\sqrt{n}} \sum_i \phi(\mathbf{X}_i; \hat{Q}, \hat{g}, \tau) - \phi(\mathbf{X}_i; Q, g, \tau) - \mathbb{E}[\phi(\mathbf{X}; \hat{Q}, \hat{g}, \tau) - \phi(\mathbf{X}; Q, g, \tau)] \quad (35.36)$$

$$5 + \sqrt{n} \mathbb{E}[\phi(\mathbf{X}; \hat{Q}, \hat{g}, \tau) - \phi(\mathbf{X}; Q, g, \tau)] \quad (35.37)$$

6 We recognize the first term, Equation (35.35), as $\sqrt{n}(\hat{\tau}^{\text{ideal}} - \tau)$, the estimation error in the optimal
7 case where we know Q and g . Ideally, we'd like the error of $\hat{\tau}^{\text{AIPTW}}$ to be asymptotically equal to
8 this ideal case—which will happen if the other two terms go to 0.

9 The second term, Equation (35.36), is a penalty we pay for using the same data to estimate Q, g
10 and to compute τ . For many model classes, it can be shown that such “empirical process” terms go
11 to 0. This can also be guaranteed in general by using different data for fitting the nuisance functions
12 and for computing the estimator (see the next section).

13 The third term, Equation (35.37), captures the penalty we pay for misestimating the nuisance
14 functions. This is where the particular form of the AIPTW is key. With a little algebra, we can show
15 that

$$16 \mathbb{E}[\phi(\mathbf{X}; \hat{Q}, \hat{g}) - \phi(\mathbf{X}; Q, g)] = \mathbb{E}\left[\frac{1}{g(X)}(\hat{g}(X) - g(X))(\hat{Q}(1, X) - Q(1, X))\right] \quad (35.38)$$

$$17 + \frac{1}{1 - g(X)}(\hat{g}(X) - g(X))(\hat{Q}(0, X) - Q(0, X)). \quad (35.39)$$

18 The important point is that estimation errors of Q and g are multiplied together. Using the Cauchy-
19 Schwarz inequality, we find that $\sqrt{n} \mathbb{E}[\phi(\mathbf{X}; \hat{Q}, \hat{g}) - \phi(\mathbf{X}; Q, g)] \rightarrow 0$ as long as $\sqrt{n} \max_a \mathbb{E}[(\hat{Q}(a, X) -$
20 $Q(a, X))^2]^{\frac{1}{2}} \mathbb{E}[(\hat{g}(X) - g(X))^2]^{\frac{1}{2}} \rightarrow 0$. That is, the misestimation penalty will vanish so long as the
21 product of the misestimation errors is $o(\sqrt{n})$. For example, this means that τ can be estimated at
22 the (optimal) \sqrt{n} rate even when the estimation error of each of Q and g only decreases as $o(n^{-1/4})$.

23 The upshot here is that the double machine learning estimator has the special property that the
24 weak condition $\sqrt{n} \mathbb{E}(\hat{Q}(T, X) - Q(T, X))^2 \mathbb{E}(\hat{g}(X) - g(X))^2 \rightarrow 0$ suffices to imply that

$$25 \sqrt{n}(\hat{\tau}^{\text{AIPTW}} - \tau) \xrightarrow{d} \text{Normal}(0, \mathbb{E}[\phi(\mathbf{X}_i; Q, g, \tau)^2]) \quad (35.40)$$

26 (though strictly speaking this requires some additional technical conditions we haven't discussed).
27 This is *not* true for the earlier estimators we discussed, which require a much faster rate of convergence
28 for the nuisance function estimation.

29 The AIPTW estimator has two further nice properties that are worth mentioning. First, it is
30 **non-parametrically efficient**. This means that this estimator has the smallest possible variance
31 of any estimator that does not make parametric assumptions; namely, $\mathbb{E}[\phi(\mathbf{X}_i; Q, g, \tau)^2]$. This means,
32 for example, that this estimator yields the smallest confidence intervals of any approach that does not
33 rely on parametric assumptions. Second, it is **double robust**: the estimator is consistent (converges
34 to the true τ as $n \rightarrow \infty$) as long as at least one of either \hat{Q} or \hat{g} is consistent.

35.4.2.4 Cross Fitting

The term Equation (35.36) in the error decomposition above is the penalty we pay for reusing the same data to both fit Q, g and to compute the estimator. For many choices of model for Q, g , this term goes to 0 quickly as n gets large and we achieve the (best case) \sqrt{n} error rate. However, this property doesn't always hold.

As an alternative, we can always randomly split the available data and use one part for model fitting, and the other to compute the estimator. Effectively, this means the nuisance function estimation and estimator computation are done using independent samples. It can then be shown that the reuse penalty will vanish. However, this comes at the price of reducing the amount of data available for each of nuisance function estimation and estimator computation.

This strategy can be improved upon by a **cross fitting** approach. We divide the data into K folds. For each fold j we use the other $K - 1$ folds to fit the nuisance function models $\hat{Q}^{-j}, \hat{g}^{-j}$. Then, for each datapoint i in fold j , we take $\hat{Q}(a_i, x_i) = \hat{Q}^{-j}(a_i, x_i)$ and $\hat{g}(x_i) = \hat{g}^{-j}(x_i)$. That is, the estimated conditional outcomes and propensity score for each datapoint are predictions from a model that was not trained on that datapoint. Then, we estimate τ by plugging $\{\hat{Q}(a_i, x_i), \hat{g}(x_i)\}_i$ into Equation (35.31). It can be shown that this cross fitting procedure has the same asymptotic guarantee—the central limit theorem at the \sqrt{n} rate—as described above.

35.4.3 Uncertainty Quantification

In addition to the point estimate $\hat{\tau}$ of the average treatment effect, we'd also like to report a measure of the uncertainty in our estimate. For example, in the form of a confidence interval. The asymptotic normality of $\sqrt{n}\hat{\tau}$ (Equation (35.40)) provides a means for this quantification. Namely, we could base confidence intervals and similar on the limiting variance $\mathbb{E}[\phi(\mathbf{X}_i; Q, g, \tau)^2]$. Of course, we don't actually know any of Q, g , or τ . However, it turns out that it suffices to estimate the asymptotic variance with $\frac{1}{n} \sum_i \phi(\mathbf{X}_i; \hat{Q}, \hat{g}, \hat{\tau})^2$ [Che+17e]. That is, we can estimate the uncertainty by simply plugging in our fitted nuisance models and our point estimate of τ into

$$\hat{V}[\hat{\tau}] = 1/n \sum_i \phi(\mathbf{X}_i; \hat{Q}, \hat{g}, \hat{\tau})^2. \quad (35.41)$$

This estimated variance can then be used to compute confidence intervals in the usual manner. E.g., we'd report a 95% confidence interval for τ as $\hat{\tau} \pm 1.96\sqrt{\hat{V}[\hat{\tau}]/n}$.

Alternatively, we could quantify the uncertainty by bootstrapping. Note, however, that this would require refitting the nuisance functions with each bootstrap model. Depending on the model and data, this can be prohibitively computationally expensive.

35.4.4 Matching

One particularly popular approach to adjustment-based causal estimation is **matching**. Intuitively, the idea is to match each treated unit to an untreated unit that has the same (or at least similar) values of the confounding variables and then compare the observed outcomes of the treated unit and its matched control. If we match on the full set of common causes, then the difference in outcomes is, intuitively, a noisy estimate of the effect the treatment had on that treated unit. We'll now build this up a bit more carefully. In the process we'll see that matching can be understood as, essentially, a particular kind of outcome model adjustment.

For simplicity, consider the case where X is a discrete random variable. Define \mathcal{A}_x to be the set of treated units with covariate value x , and \mathcal{C}_x to be the set of untreated units with covariate value x . In this case, the matching estimator is:

$$\hat{\tau}^{\text{matching}} = \sum_x \hat{P}(x) \left(\frac{1}{|\mathcal{A}_x|} \sum_{i \in \mathcal{A}_x} Y_i - \frac{1}{|\mathcal{C}_x|} \sum_{j \in \mathcal{C}_x} Y_j \right), \quad (35.42)$$

where $\hat{P}(x)$ is an estimator of $P(X = x)$ —e.g., the fraction of units with $X = x$. Now, we can rewrite $Y_i = Q(A_i, X_i) + \xi_i$ where ξ_i is a unit-specific noise term defined by the equation. In particular, we have that $\mathbb{E}[\xi_i | A_i, X_i] = 0$. Subbing this in, we have:

$$\hat{\tau}^{\text{matching}} = \sum_x \hat{P}(x) (Q(1, x) - Q(0, x)) + \sum_x \frac{1}{|\mathcal{A}_x|} \sum_{i \in \mathcal{A}_x} \xi_i - \frac{1}{|\mathcal{C}_x|} \sum_{j \in \mathcal{C}_x} \xi_j. \quad (35.43)$$

We can recognize the first term as an estimator of usual target parameter τ (it will be equal to τ if $\hat{P}(x) = P(x)$). The second term is a difference of averages of random variables with expectation 0, and so each term will converge to 0 as long as $|\mathcal{A}_x|$ and $|\mathcal{C}_x|$ each go to infinity as we see more and more data. Thus, we see that the matching estimator is a particular way of estimating the parameter τ . The procedure can be extended to continuous covariates by introducing some notion of values of X being close, and then matching close treatment and control variables.

There are two points we should emphasize here. First, notice that the argument here has nothing to do with causal identification. Matching is a particular technique for estimating the observational parameter τ . Whether or not τ can be interpreted as an average treatment effect is determined by the conditions of Theorem 2—the particular estimation strategy doesn’t say anything about this. Second, notice that in essence matching amounts to a particular choice of model for \hat{Q} . Namely, $\hat{Q}(1, x) = \frac{1}{|\mathcal{A}_x|} \sum_{i \in \mathcal{A}_x} Y_i$ and similarly for $\hat{Q}(0, x)$. That is, we estimate the conditional expected outcome as a sample mean over units with the same covariate value. Whether this is a good idea depends on how good of a model for Q this is. In situations where better models are possible (e.g., a machine-learning model fits the data well), we might expect to get a more accurate estimate by using the conditional expected outcome predictor directly.

There is another important case we mention in passing. In general, when using adjustment based identification, it suffices to adjust for any function $\phi(X)$ of X such that $A \perp\!\!\!\perp X | \phi(X)$. To see that adjusting for only $\phi(X)$ suffices, first notice that $g(X) = P(A = 1 | X) = P(A = 1 | \phi(X))$ only depends on $\phi(X)$, and then recall that can write the target parameter as $\tau = \mathbb{E}[\frac{YA}{g(X)} - \frac{Y(1-A)}{1-g(X)}]$, whence τ only depends on X through $g(X)$. That is: replacing X by a reduced version $\phi(X)$ such that $g(X) = P(A = 1 | \phi(X))$ can’t make any difference to τ . Indeed, the most popular choice of $\phi(X)$ is the propensity score itself, $\phi(X) = g(X)$. This leads to **propensity score matching**, a two step procedure where we first fit a model for the propensity score, and then run matching based on the estimated propensity score values for each unit. Again, this is just a particular estimation procedure for the observational parameter τ , and says nothing about whether it’s valid to interpret τ as a causal effect.

35.4.5 Practical Considerations and Procedures

Lots of issues can arise in practice.

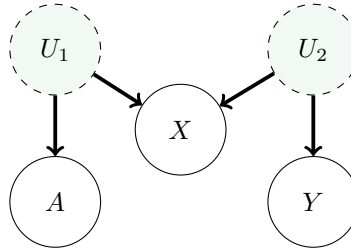


Figure 35.5: The M-bias causal graph. Here, A and Y are not confounded. However, conditioning on the covariate X opens a backdoor path, passing through U_1 and U_2 (because X is a collider). Thus, adjusting for X creates bias. This is true even though X need not be a pre-treatment variable.

35.4.5.1 What to adjust for

Choosing which variables to adjust for is a key detail in estimating causal effects using covariate adjustment. The criterion is clear when one has a full causal graph relating A , Y , and all covariates X to each other. Namely, adjust for all variables that are actually causal parents of A and Y . In fact, with access to the full graph, this criteria can be generalized somewhat—see Section 35.8.

In practice, we often don’t actually know the full causal graph relating all of our variables. As a result, it is common to apply simple heuristics to determine which variables to adjust for. Unfortunately, these heuristics have serious limitations. However, exploring these are instructive.

A key condition in Theorem 2 is that the covariates X that we adjust for must include all the common causes. In the absence of a full causal graph, it is tempting to condition on as many observed variables as possible to try to ensure this condition holds. However, this can be problematic. For instance, suppose that M is a mediator of the effect of A on Y —i.e., M lies on one of the directed paths between A and Y . Then, conditioning on M will block this path, removing some of the causal effect. Note that this does not always result in an attenuated, or smaller-magnitude, effect estimate. The effect through a given mediator may run in the opposite direction of other causal pathways from the treatment; thus conditioning on a mediator can inflate or even flip the sign of a treatment effect. Alternatively, if C is a collider between A and Y —a variable that is caused by both—then conditioning on C will induce an extra statistical dependency between A and Y .

Both pitfalls of the “condition on everything” heuristic discussed above both involve conditioning on variables that are downstream of the treatment A . A natural response to this is to limit conditioning to all pre-treatment variables, or those that are causally upstream of the treatment. Importantly, if there is a valid adjustment set in the observed covariates X , then there will also be a valid adjustment set among the pre-treatment covariates. This is because any open backdoor path between A and Y must include a parent of A , and the set of pre-treatment covariates includes these parents. However, it is still possible that conditioning on the full set of pre-treatment variables can induce new backdoor paths between A and Y through colliders. In particular, if there is a covariate D that is separately confounded with the treatment A and the outcome Y then D is a collider, and conditioning on D opens a new backdoor path. This phenomenon is known as m-bias because of the shape of the graph [Pea09c], see Figure 35.5.

A practical refinement of the pre-treatment variable heuristic is given in VanderWeele TJ [VT11]. Their heuristic suggests conditioning on all pre-treatment variables that are causes of the treatment, outcome, or both. The essential qualifier in this heuristic is that the variable is causally upstream of treatment and/or outcome. This eliminates the possibility of conditioning on covariates that are only confounded with treatment and outcome, avoiding m-bias. Notably, this heuristic requires more causal knowledge than the above heuristics, but does not require detailed knowledge of how different covariates are causally related to each other.

The VanderWeele TJ [VT11] criterion is a useful rule of thumb, but other practical considerations often arise. For example, if one has more knowledge about the causal structure among covariates, it is possible to optimize adjustment sets to minimize the variance of the resulting estimator [RS20]. One important example of reducing variance by pruning adjustment sets is the exclusion of variables that are known to only be a parent of the treatment, and not of the outcome (so called instruments, as discussed in Section 35.5).

Finally, adjustment set selection criteria operate under the assumption that there actually exists a valid adjustment set among observed covariates. When there is no set of observed covariates in X that block all backdoor paths, then any adjusted estimate will be biased. Importantly, in this case, the bias does not necessarily decrease as one conditions on more variables. For example, conditioning on an instrumental variable often results in an estimate that has higher bias, in addition to the higher variance discussed above. This phenomenon is known as bias amplification or z-bias; see Section 35.7.2. A general rule of thumb is that variables that explain away much more variation in the treatment than in the outcome can potentially amplify bias, and should be treated with caution.

35.4.5.2 Overlap

Recall that in addition to no-unobserved-confounders, identification of the average treatment effect requires overlap: the condition that $0 < P(A = 1|x) < 1$ for the population distribution P . With infinite data, any amount of overlap will suffice for estimating the causal effect. In realistic settings, even near failures can be problematic. Equation (35.40) gives an expression for the (asymptotic) variance of our estimate: $\mathbb{E}[\phi(\mathbf{X}_i; \hat{Q}, \hat{g}, \hat{\tau})^2]/n$. Notice that $\phi(\mathbf{X}_i; \hat{Q}, \hat{g}, \hat{\tau})^2$ involves terms that are proportional to $1/g(X)$ and $1/(1 - g(X))$. Accordingly, the variance of our estimator will balloon if there are units where $g(x) \approx 0$ or $g(x) \approx 1$ (unless such units are rare enough that they don't contribute much to the expectation).

In practice, a simple way to deal with potential overlap violation is to fit a model \hat{g} for the treatment assignment probability—which we need to do anyways—and check that the values $\hat{g}(x)$ are not too extreme. In the case that some values are too extreme, the simplest resolution is to cheat. We can simply exclude all the data with extreme values of $\hat{g}(x)$. This is equivalent to considering the average treatment effect over only the subpopulation where overlap is satisfied. This changes the interpretation of the estimand. The restricted subpopulation ATE may or may not provide a satisfactory answer to the real-world problem at hand, and this needs to be justified based on knowledge of the real-world problem.

35.4.5.3 Choice of Estimand and Average Treatment Effect on the Treated

Usually, our goal in estimating a causal effect is qualitative. We want to know what the sign of the effect is, and whether it's large or small. The utility of the ATE is that it provides a concrete query

we can use to get a handle on the qualitative question. However, it is not sacrosanct; sometimes we're better off choosing an alternative causal estimand that still answers the qualitative question but which is easier to estimate statistically. The **average treatment effect on the treated** or **ATT**,

$$\text{ATT} \triangleq \mathbb{E}_{X|A=1}[\mathbb{E}[Y|X, \text{do}(A=1)] - E[Y|X, \text{do}(A=0)]], \quad (35.44)$$

is one such an estimand that is frequently useful.

The ATT is useful when many members of the population are very unlikely to receive treatment, but the treated units had a reasonably high probability of receiving the control. This can happen if, e.g., we sample control units from the general population, but the treatment units all self-selected into treatment from a smaller subpopulation. In this case, it's not possible to (non-parametrically) determine the treatment effect for the control units where no similar unit took treatment. The ATT solves this obstacle by simply omitting such units from the average.

If we have the causal structure Figure 35.4, and the overlap condition $P(A=1|X=x) < 1$ for all $X=x$ then the ATT is causally identified as

$$\tau^{\text{ATT}} = \mathbb{E}_{X|A=1}[\mathbb{E}[Y|A=1, X] - E[Y|A=0, X]]. \quad (35.45)$$

Note that the required overlap condition here is weaker than for identifying the ATE. (The proof is the same as Theorem 2.)

The estimation strategies for the ATE translate readily to estimation strategies for the ATT. Namely, estimate the nuisance functions the same way and then simply replace averages over all data points by averages over the treated datapoints only. In principle, it's possible to do a little better than this by making use of the untreated datapoints as well. A corresponding double machine learning estimator is

$$\hat{\tau}^{\text{ATT-AIPTW}} \triangleq \frac{1}{n} \sum_i \frac{A_i}{P(A=1)} (Y - \hat{Q}(0, X_i)) - \frac{(1-A_i)g(X)}{P(A=1)(1-g(X))} (Y - \hat{Q}(0, X_i)). \quad (35.46)$$

. The variance of this estimator can be estimated by

$$\phi^{\text{ATT}}(\mathbf{X}_i; Q, g, \tau) \triangleq \frac{1}{n} \sum_i \frac{A_i}{P(A=1)} (Y - \hat{Q}(0, X_i)) - \frac{(1-A_i)g(X)}{P(A=1)(1-g(X))} (Y - \hat{Q}(0, X_i) - \frac{A\tau}{P(A=1)}) \quad (35.47)$$

$$\hat{\mathbb{V}}[\hat{\tau}^{\text{ATT-AIPTW}}] \triangleq \frac{1}{n} \sum_i \phi^{\text{ATT}}(\mathbf{X}_i; \hat{Q}, \hat{g}, \hat{\tau}^{\text{ATT-AIPTW}}). \quad (35.48)$$

Notice that the estimator for the ATT doesn't require estimating $Q(1, X)$. This can be a considerable advantage when the treated units are rare.

See Chernozhukov et al. [Che+17e] for details.

35.4.6 Summary and Practical Advice

We have seen a number of estimators that follow the general procedure:

1. fit statistical or machine-learning models $\hat{Q}(a, x)$ as a predictor for Y , and/or $\hat{g}(x)$ as a predictor for A

- 1
- 2 2. compute the predictions $\hat{Q}(0, x_i), \hat{Q}(1, x_i), \hat{g}(x_i)$ for each data point, and
- 3 3. combine these predictions into an estimate of the average treatment effect.

4

5 Importantly, no single estimation approach is a silver bullet. For example, the double machine-
6 learning estimator has appealing theoretical properties, such as asymptotic efficiency guarantees and
7 a recipe for estimating uncertainty without needing to bootstrap the model fitting. However, in
8 terms of the quality of point estimates, the double ML estimators can sometimes underperform their
9 more naive counterparts [KS07]. In fact, there are cases where each of outcome regression, propensity
10 weighting, or doubly robust methods will outperform the others.

11 One difficulty in choosing an estimator in practice is that there are fewer guardrails in causal
12 inference than there are in standard predictive modeling. In predictive modeling, we construct a
13 train-test split and validate our prediction models using the true labels or outcomes in the held-out
14 dataset. However, for causal problems, the causal estimands are functionals of a different data-
15 generating process from the one that we actually observed. As a result, it is impossible to empirically
16 validate many aspects of causal estimation using standard techniques.

17 The effectiveness of a given approach is often determined by how much we trust the specification of
18 our propensity score or outcome regression models $\hat{g}(x)$ and $\hat{Q}(a, x)$, and how well the treatment and
19 control groups overlap in the dataset. Using flexible models for the nuisance functions g and Q can
20 alleviate some of the concerns about model misspecification, but our freedom to use such models is
21 often constrained by dataset size. When we have the luxury of large data, we can use flexible models;
22 on the other hand, when the dataset is relatively small, we may need to use a smaller parametric
23 family or stringent regularization to obtain stable estimates of Q and g . Similarly, if overlap is poor
24 in some regions of the covariate space, then flexible models for Q may be highly variable, and inverse
25 propensity score weights may be large. In these cases, IPTW or AIPTW estimates may fluctuate
26 wildly as a function of large weights. Meanwhile, outcome regression estimates will be sensitive to
27 the specification of the Q model and its regularization, and can incur bias that is difficult to measure
28 if the specification or regularization does not match the true outcome process.

29 There are a number of practical steps that we can take to sanity-check causal estimates. The
30 simplest check is to compute many different ATE estimators (e.g., outcome regression, IPTW, doubly
31 robust) using several comparably complex estimators of Q and g . We can then check whether they
32 agree, at least qualitatively. If they do agree then this can provide some peace of mind (although it
33 is not a guarantee of accuracy). If they disagree, caution is warranted, particularly in choosing the
34 specification of the Q and g models.

35 It is also important to check for failures of overlap. Often, issues such as disagreement between
36 alternative estimators can be traced back to poor overlap. A common way to do this, particularly
37 with high-dimensional data, is to examine the estimated (ideally cross-fitted) propensity scores $\hat{g}(x_i)$.
38 This is a useful diagnostic, even if the intention is to use an outcome regression approach that only
39 incorporates and estimated outcome regression function $\hat{Q}(a, x_i)$. If overlap issues are relevant, it
40 may be better to instead estimate either the average treatment effect on the treated, or the “trimmed”
41 estimand given by discarding units with extreme propensities.

42 Uncertainty quantification is also an essential part of most causal analyses. This frequently take
43 the form of an estimate of the estimator’s variance, or a confidence interval. This may be important
44 for downstream decision-making, and can also be a useful diagnostic. We can calculate variance either
45 by bootstrapping the entire procedure (including refitting the models in each bootstrap replicate),
46 or computing analytical variance estimates from the AIPTW estimator. Generally, large variance

47

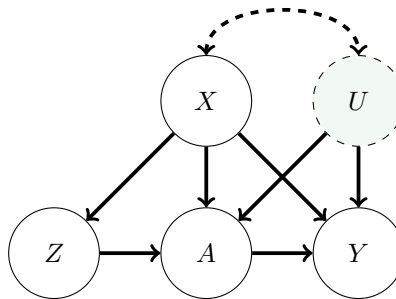


Figure 35.6: Causal graph illustrating the Instrumental Variable setup. The treatment A and outcome Y are both influenced by unobserved confounder U . Nevertheless, identification is sometimes possible due to the presence of the instrument Z . We also allow for observed covariates X that we may need to adjust for. The dashed arrow between U and X indicates a statistical dependency where we remain agnostic to the particular causal relationship.

estimates may indicate issues with the analysis. For example, poor overlap will often (although not always) manifest as extremely large variances under either of these methods. Small variance estimates should be treated with caution, unless other checks, such as overlap checks, or stability across different Q and g models, also pass.

The previous advice only addresses the statistical problem of estimating τ from a data sample. It does not speak to whether or not τ can reasonably be interpreted as an average treatment effect. Considerable care should be devoted to whether or not the assumption that there are no unobserved confounders is reasonable. There are several methods for assessing the sensitivity of the ATE estimate to violations of this assumption. See Section 35.7. Bias due to unobserved confounding can be substantial in practice—often overwhelming bias due to estimation error—so it is wise to conduct such an analysis.

35.5 Instrumental Variable Strategies

Adjustment-based methods rely on observing all confounders affecting the treatment and outcome. In some situations, it is possible to identify interesting causal effects even when there are unobserved confounders. We now consider strategies based on **instrumental variables**. The instrumental variable graph is shown in Figure 35.6. The key ingredient is the instrumental variable Z , a variable that has a causal effect on Y only through its causal effect on A . Informally, the identification strategy is to determine the causal effect of Z on Y , the causal effect of Z on A , and then combine these into an estimate of the causal effect of A on Y .

For this identification to strategy to work the instrument must satisfy three conditions. There are observed variables (confounders) X such that:

1. **Instrument Relevance** $Z \not\perp A|X$: the instrument must actually affect the treatment assignment.
2. **Instrument Unconfoundedness** Any backdoor path between Z and Y is blocked by X , even conditional on A .

3. Exclusion Restriction All directed paths from Z to Y pass through A . That is, the instrument affects the outcome *only* through its effect on A .

(It may help conceptually to first think through the case where X is the empty set—i.e., where the only confounder is the unobserved U). These assumptions are necessary for using instrumental variables for causal identification, but they are not quite sufficient. In practice, they must be supplemented by an additional assumption that depends more closely on the details of the problem at hand. Historically, this additional assumption was usually that both the instrument-treatment and treatment-outcome relationship are linear. We'll examine some less restrictive alternatives below.

Before moving on to how to use instrumental variables for identification, let's consider how we might encounter instruments in practice. The key is that it's often possible to find, and measure, variables that affect treatment and that are assigned (as if) at random. For example, suppose we are interested in measuring the effect of taking a drug A on some health outcome Y . The challenge is that whether a study participant actually takes the drug can be confounded with Y —e.g., sicker people may be more likely to take their medication, but have worse outcomes. However, the assignment of treatments to patients can be randomized and this random assignment can be viewed as an instrument. This **random assignment with non-compliance** scenario is common in practice. The random assignment—the instrument—satisfies relevance (so long as assigning the drug affects the probability of the patient taking the drug). It also satisfies unconfoundedness (because the instrument is randomized). And, it plausibly satisfies exclusion restriction: telling (or not telling) a patient to take a drug has no effect on their health outcome except through influencing whether or not they actually take the drug. As a second example, the **judge fixed effects** research design uses the identity of the judge assigned to each criminal case to infer the effect of incarceration on some life outcome of interest (e.g., total lifetime earnings). Relevance will be satisfied so long as different judges have different propensities to hand out severe sentences. The assignment of trial judges to cases is randomized, so unconfoundedness will also be satisfied. And, exclusion restriction is also plausible: the particular identity of the judge assigned to your case has no bearing on your years-later life outcomes, except through the particular sentence that you're subjected to.

It's important to note that these assumptions require some care, particularly exclusion restriction. Relevance can be checked directly from the data, by fitting a model to predict the treatment from the instrument (or vice versa). Unconfoundedness is often satisfied by design: the instrument is randomly assigned. Even when literal random assignment doesn't hold, we often restrict to instruments where unconfoundedness is "obviously" satisfied—e.g., using number of rainy days in a month as an instrument for sun exposure. Exclusion restriction is trickier. For example, it might fail in the drug assignment case if patients who are not told to take a drug respond by seeking out alternative treatment. Or, it might fail in the judge fixed effects case if judges hand out additional, unrecorded, punishments in addition to incarceration. Assessing the plausibility of exclusion restriction requires careful consideration based on domain expertise.

We now return to the question of how to make use of an instrument once we have it in hand. As previously mentioned, getting causal identification using instrumental variables requires supplementing the IV assumptions with some additional assumption about the causal process.

35.5.1 Additive Unobserved Confounding

We first consider **additive unobserved confounding**. That is, we assume that the structural causal model for the outcome has the form:⁷

$$Y \leftarrow f(A, X) + f_U(U). \quad (35.49)$$

In words, we assume that there are no interaction effects between the treatment and the unobserved confounder—everyone responds to treatment in the same way. With this additional assumption, we see that $\mathbb{E}[Y|X, \text{do}(A = a)] - \mathbb{E}[Y|X, \text{do}(A = a')] = f(a, X) - f(a', X)$. In this setting, our goal is to learn this contrast.

Theorem 3 (Additive Confounding Identification). *If the instrumental variables assumptions hold and also additive unobserved confounding holds, then there is a function $\tilde{f}(a, x)$ where*

$$\mathbb{E}[Y|x, \text{do}(A = a)] - \mathbb{E}[Y|x, \text{do}(A = a')] = \tilde{f}(a, x) - \tilde{f}(a', x), \quad (35.50)$$

for all x, a, a' and such that \tilde{f} satisfies

$$\mathbb{E}[Y|z, x] = \int \tilde{f}(a, x)p(a|z, x)da. \quad (35.51)$$

Here, $p(a|z, x)$ is the conditional probability density of treatment.

In particular, if there is a unique function g that satisfies

$$\mathbb{E}[Y|z, x] = \int g(a, x)p(a|z, x)da, \quad (35.52)$$

then $g = \tilde{f}$ and this relation identifies the target causal effect.

Before giving the proof, let's understand the point of this identification result. The key insight is that both the left hand side of Equation (35.52) and $p(a|z, x)$ (appearing in the integrand) are identified by the data, since they involve only observational relationships between observed variables. So, \tilde{f} is identified implicitly as one of the functions that makes Equation (35.52) true. If there is a unique such function, then this fully identifies the causal effect.

Proof. With the additive unobserved confounding assumption, the instrument unconfoundedness implies that $U \perp\!\!\!\perp Z|X$. Then, we have that:

$$\mathbb{E}[Y|Z, X] = \mathbb{E}[f(A, X)|Z, X] + \mathbb{E}[f_U(U)|Z, X] \quad (35.53)$$

$$= \mathbb{E}[f(A, X)|Z, X] + \mathbb{E}[f_U(U)|X] \quad (35.54)$$

$$= \mathbb{E}[\tilde{f}(A, X)|Z, X], \quad (35.55)$$

where $\tilde{f} = f(A, X) + \mathbb{E}[f_U(U)|X]$. Now, identifying just \tilde{f} would suffice for us, because we could then identify contrasts between treatments: $f(a, x) - f(a', x) = \tilde{f}(a, x) - \tilde{f}(a', x)$. (The term $\mathbb{E}[f_U(U)|x]$ cancels out). Accordingly, we rewrite Equation (35.55) as:

$$\mathbb{E}[Y|z, x] = \int \tilde{f}(a, x)p(a|z, x)da. \quad (35.56)$$

□

⁷. We roll the unit-specific variables ξ into U to avoid notational overload.

It's worth dwelling briefly on how the IV assumptions come into play here. The exclusion restriction is implied by the additive unobserved confounding assumption, which we use explicitly. We also use the unconfoundedness assumption to conclude $U \perp\!\!\!\perp Z|X$. However, we do not use relevance. The role of relevance here is in ensuring that few functions solve the relation Equation (35.52). Informally, the solution g is constrained by the requirement that it hold for all values of Z . However, different values of Z only add non-trivial constraints if $p(a|z, x)$ differ depending on the value of z —this is exactly the relevance condition.

Estimation The basic estimation strategy is to fit models for $\mathbb{E}[Y|z, x]$ and $p(a|z, x)$ from the data, and then solve the implicit equation Equation (35.52) to find g consistent with the fitted models. The procedures for doing this can vary considerably depending on the particulars of the data (e.g., if Z is discrete or continuous) and the choice of modeling strategy. We omit a detailed discussion, but [see e.g. NP03; Dar+11; Har+17; SSG19; BKS19; Mua+20; Dik+20] for various concrete approaches. It's also worth mentioning an additional nuance to the general procedure. Even if relevance holds, there will often be more than one function that satisfies Equation (35.52). So, we have only identified \tilde{f} as a member of this set of functions. In practice, this ambiguity is defeated by making some additional structural assumption about \tilde{f} . For example, we model \tilde{f} with a neural network, and then choose the network satisfying Equation (35.52) that has minimum l_2 -norm on the parameters (i.e., we pick the l_2 -regularized solution).

35.5.2 Instrument Monotonicity and Local Average Treatment Effect

We now consider an alternative assumption to additive unobserved confounding that is applicable when both the instrument and treatment are binary. It will be convenient to conceptualize the instrument as assignment-to-treatment. Then, the population divides into four subpopulations:

1. Compliers, who take the treatment if assigned to it, and who don't take the treatment otherwise.
2. Always takers, who take the treatment no matter their assignment
3. Never takers, who refuse the treatment no matter their assignment
4. Defiers, who refuse the treatment if assigned to it, and who take the treatment if not assigned.

Our goal in this setting will be to identify the average treatment effect among the compliers. The **local average treatment effect** (or **complier average treatment effect**) is defined to be⁸

$$\text{LATE} = \mathbb{E}[Y|\text{do}(A = 1), \text{complier}] - \mathbb{E}[Y|\text{do}(A = 0), \text{complier}]. \quad (35.57)$$

The LATE requires an additional assumption for identification. Namely, **instrument monotonicity**: being assigned (not assigned) the treatment only increases (decreases) the probability that each unit will take the treatment. Equivalently, $P(\text{defier}) = 0$.

We can then write down the identification result.

⁸. We follow the econometrics literature in using “LATE” because “CATE” is already commonly used for conditional average treatment effect.

Theorem 4. *Given the instrumental variable assumptions and instrument monotonicity, the local average treatment is identified as a parameter τ^{LATE} of the observational distribution; that is, $\text{LATE} = \tau^{\text{LATE}}$. Namely,*

$$\tau^{\text{LATE}} = \frac{\mathbb{E}[\mathbb{E}[Y|X, Z = 1] - \mathbb{E}[Y|X, Z = 0]]}{\mathbb{E}[P(A = 1|X, Z = 1) - P(A = 1|X, Z = 0)]}. \quad (35.58)$$

Proof. We now show that, given the IV assumptions and monotonicity, $\text{LATE} = \tau^{\text{LATE}}$. First, notice that

$$\tau^{\text{LATE}} = \frac{\mathbb{E}[Y|\text{do}(Z = 1)] - \mathbb{E}[Y|\text{do}(Z = 0)]}{P(A = 1|\text{do}(Z = 1)) - P(A = 1|\text{do}(Z = 0))}. \quad (35.59)$$

This follows from backdoor adjustment, Theorem 2, applied to the numerator and denominator separately. Our strategy will be to decompose $\mathbb{E}[Y|\text{do}(Z = z)]$ into the contributions from the compliers, the units that ignore the instrument (the always/never takers), and the defiers. To that end, note that $P(\text{complier}|\text{do}(Z = z)) = P(\text{complier})$ and similarly for always/never takers and defiers—interventions on the instrument don't change the composition of the population. Then,

$$\mathbb{E}[Y|\text{do}(Z = 1)] - \mathbb{E}[Y|\text{do}(Z = 0)] \quad (35.60)$$

$$= (\mathbb{E}[Y|\text{complier}, \text{do}(Z = 1)] - \mathbb{E}[Y|\text{complier}, \text{do}(Z = 0)])P(\text{complier}) \quad (35.61)$$

$$+ (\mathbb{E}[Y|\text{always/never}, \text{do}(Z = 1)] - \mathbb{E}[Y|\text{always/never}, \text{do}(Z = 0)])P(\text{always/never}) \quad (35.62)$$

$$+ (\mathbb{E}[Y|\text{defier}, \text{do}(Z = 1)] - \mathbb{E}[Y|\text{defier}, \text{do}(Z = 0)])P(\text{defier}). \quad (35.63)$$

The key is the effect on the complier subpopulation, Equation (35.61). First, by definition of the complier population, we have that:

$$\mathbb{E}[Y|\text{complier}, \text{do}(Z = z)] = \mathbb{E}[Y|\text{complier}, \text{do}(A = z)]. \quad (35.64)$$

That is, the causal effect of the treatment is the same as the causal effect of the instrument in this subpopulation—this is the core reason why access to an instrument allows identification of the local average treatment effect. This means that

$$\text{LATE} = \mathbb{E}[Y|\text{complier}, \text{do}(Z = 1)] - \mathbb{E}[Y|\text{complier}, \text{do}(Z = 0)]. \quad (35.65)$$

Further, we have that $P(\text{complier}) = P(A = 1|\text{do}(Z = 1)) - P(A = 1|\text{do}(Z = 0))$. The reason is simply that, by definition of the subpopulations,

$$P(A = 1|\text{do}(Z = 1)) = P(\text{complier}) + P(\text{always taker}) \quad (35.66)$$

$$P(A = 1|\text{do}(Z = 0)) = P(\text{always taker}). \quad (35.67)$$

Now, plugging the expression for $P(\text{complier})$ and Equation (35.65) into Equation (35.61) we have that:

$$(\mathbb{E}[Y|\text{complier}, \text{do}(Z = 1)] - \mathbb{E}[Y|\text{complier}, \text{do}(Z = 0)])P(\text{complier}) \quad (35.68)$$

$$= \text{LATE} \times (P(A = 1|\text{do}(Z = 1)) - P(A = 1|\text{do}(Z = 0))) \quad (35.69)$$

This gives us an expression for the local average treatment effect in terms of the effect of the instrument on the compliers and the probability that a unit takes the treatment when assigned/not-assigned.

The next step is to show that the remaining instrument effect decomposition terms, Equations (35.62) and (35.63), are both 0. Equation (35.62) is the causal effect of the instrument on the always/never takers. It's equal to 0 because, by definition of this subpopulation, the instrument has no causal effect in the subpopulation—they ignore the instrument! Mathematically, this is just $\mathbb{E}[Y|\text{always/never, do}(Z = 1)] = \mathbb{E}[Y|\text{always/never, do}(Z = 0)]$. Finally, Equation (35.63) is 0 by the instrument monotonicity assumption: we assumed that $P(\text{defier}) = 0$.

In totality, we now have that Equations (35.61) to (35.63) reduces to:

$$\mathbb{E}[Y|\text{do}(Z = 1)] - \mathbb{E}[Y|\text{do}(Z = 0)] \quad (35.70)$$

$$= \text{LATE} \times (P(A = 1|\text{do}(Z = 1)) - P(A = 1|\text{do}(Z = 0))) + 0 + 0 \quad (35.71)$$

Rearranging for LATE and plugging in to Equation (35.59) gives claimed identification result. \square

35.5.2.1 Estimation

For estimating the local average treatment effect under the monotone instrument assumption, there is a double-machine learning approach that works with generic supervised learning approaches. Here, we want an estimator $\hat{\tau}^{\text{LATE}}$ for the parameter

$$\tau^{\text{LATE}} = \frac{\mathbb{E}[\mathbb{E}[Y|X, Z = 1] - \mathbb{E}[Y|X, Z = 0]]}{\mathbb{E}[P(A = 1|X, Z = 1) - P(A = 1|X, Z = 0)]}. \quad (35.72)$$

To define the estimator, it's convenient to introduce some additional notation. First, we define the nuisance functions:

$$\mu(z, x) = \mathbb{E}[Y|z, x] \quad (35.73)$$

$$m(z, x) = P(A = 1|x, z) \quad (35.74)$$

$$p(x) = P(Z = 1|x). \quad (35.75)$$

We also define the score ϕ by:

$$\phi_{Z \rightarrow Y}(\mathbf{X}; \mu, p) \triangleq \mu(1, X) - \mu(0, X) + \frac{Z(Y - \mu(1, X))}{p(X)} - \frac{(1 - Z)(Y - \mu(0, X))}{1 - p(X)} \quad (35.76)$$

$$\phi_{Z \rightarrow A}(\mathbf{X}; m, p) \triangleq m(1, X) - m(0, X) + \frac{Z(A - m(1, X))}{p(X)} - \frac{(1 - Z)(A - m(0, X))}{1 - p(X)} \quad (35.77)$$

$$\phi(\mathbf{X}; \mu, m, p, \tau) \triangleq \phi_{Z \rightarrow Y}(\mathbf{X}; \mu, p) - \phi_{Z \rightarrow A}(\mathbf{X}; m, p) \times \tau \quad (35.78)$$

Then, the estimator is defined by a two stage procedure:

1. Fit models $\hat{\mu}, \hat{m}, \hat{p}$ for each of μ, m, p (using supervised machine learning).

2. Define $\hat{\tau}^{\text{LATE}}$ as the solution to $\frac{1}{n} \sum_i \phi(\mathbf{X}_i; \hat{\mu}, \hat{m}, \hat{p}, \hat{\tau}^{\text{LATE}}) = 0$. That is,

$$\hat{\tau}^{\text{LATE}} = \frac{\frac{1}{n} \sum_i \phi_{Z \rightarrow Y}(\mathbf{X}_i; \hat{\mu}, \hat{p})}{\frac{1}{n} \sum_i \phi_{Z \rightarrow A}(\mathbf{X}_i; \hat{m}, \hat{p})} \quad (35.79)$$

It may help intuitions to notice that the double machine learning estimator of the LATE is effectively the double machine learning estimator of the average treatment effect of Z on Y divided by the double machine learning estimator of the average treatment effect of Z on A .

Similarly to Section 35.4, the nuisance functions can be estimated by:

1. fit a model $\hat{\mu}$ that predicts Y from Z, X by minimizing mean square error
2. fit a model \hat{m} that predicts A from Z, X by minimizing mean cross-entropy
3. fit a model \hat{p} that predicts Z from X by minimizing mean cross-entropy.

As in Section 35.4, reusing the same data for model fitting and computing the estimator can potentially cause problems. This can be avoided with use a cross-fitting procedure as described in Section 35.4.2.4. In this case, we split the data into K folds and, for each fold k , use all the but the k th fold to compute estimates $\hat{\mu}_{-k}, \hat{m}_{-k}, \hat{p}_{-k}$ of the nuisance parameters. Then we compute the nuisance estimates for each datapoint i in fold k by predicting the required quantity using the nuisance model fit on the other folds. That is, if unit i is in fold k , we compute $\hat{\mu}(z_i, x_i) \triangleq \hat{\mu}_{-k}(z_i, x_i)$ and so forth.

The key result is that if we use the cross-fit version of the estimator and the estimators for the nuisance functions converge to their true values in the sense that

1. $\mathbb{E}(\hat{\mu}(Z, X) - \mu(Z, X))^2 \rightarrow 0$, $\mathbb{E}(\hat{m}(Z, X) - m(Z, X))^2 \rightarrow 0$, and $\mathbb{E}(\hat{p}(X) - p(X))^2 \rightarrow 0$
2. $\sqrt{\mathbb{E}[(\hat{p}(X) - p(X))^2]} \times (\sqrt{\mathbb{E}[(\hat{\mu}(Z, X) - \mu(Z, X))^2]} + \sqrt{\mathbb{E}[(\hat{m}(Z, X) - m(Z, X))^2]}) = o(\sqrt{n})$

then (with some omitted technical conditions) we have asymptotic normality at the \sqrt{n} -rate:

$$\sqrt{n}(\hat{\tau}^{\text{LATE-cf}} - \tau^{\text{LATE}}) \xrightarrow{d} \text{Normal}\left(0, \frac{\mathbb{E}[\phi(\mathbf{X}; \mu, m, p, \tau^{\text{LATE}})^2]}{\mathbb{E}[m(1, X) - m(0, X)]^2}\right). \quad (35.80)$$

As with double machine learning for the confounder adjustment strategy, the key point here is that we can achieve the (optimal) \sqrt{n} rate for estimating the LATE under a relatively weak condition on how well we estimate the nuisance functions—what matters is the *product* of the error in p and the errors in μ, m . So, for example, a very good model for how the instrument is assigned (p) can make up for errors in the estimation of the treatment-assignment (m) and outcome (μ) models.

The double machine learning estimator also gives a recipe for quantifying uncertainty. To that end, define

$$\hat{\tau}_{Z \rightarrow A} \triangleq \frac{1}{n} \sum_i \phi_{Z \rightarrow A}(\mathbf{X}_i; \hat{m}, \hat{p}) \quad (35.81)$$

$$\hat{\mathbb{V}}[\hat{\tau}^{\text{LATE}}] \triangleq \frac{1}{\hat{\tau}_{Z \rightarrow A}^2} \frac{1}{n} \sum_i \phi(\mathbf{X}_i; \hat{\mu}, \hat{m}, \hat{p}, \hat{\tau}^{\text{LATE}})^2. \quad (35.82)$$

Then, subject to suitable technical conditions, $\hat{\mathbb{V}}[\hat{\tau}^{\text{LATE-cf}}]$ can be used as an estimate of the variance of the estimator. More precisely,

$$\sqrt{n}(\hat{\tau}^{\text{LATE}} - \tau^{\text{LATE}}) \xrightarrow{d} \text{Normal}(0, \hat{\mathbb{V}}[\hat{\tau}^{\text{LATE}}]). \quad (35.83)$$

Then, confidence intervals or p -values can be computed using this variance in the usual way. The main extra condition required for the variance estimator to be valid is that the nuisance parameters must all converge at rate $O(n^{-1/4})$ (so an excellent estimator for one can't fully compensate for terrible estimators of the others). In fact, even this condition is unnecessary in certain special cases—e.g., when p is known exactly, which occurs when the instrument is randomly assigned. See Chernozhukov et al. [Che+17e] for technical details.

35.5.3 Two Stage Least Squares

Commonly, the IV assumptions are supplemented with the following linear model assumptions:

$$A_i \leftarrow \alpha_0 + \alpha Z_i + \delta_A X_i + \gamma_A X_i + \xi_i^A \quad (35.84)$$

$$Y_i \leftarrow \beta_0 + \beta A_i + \delta_Y X_i + \gamma_Y X_i + \xi_i^Y \quad (35.85)$$

That is, we assume that the real-world process for treatment assignment and the outcome are both linear. In this case, plugging Equation (35.84) into Equation (35.85) yields

$$Y_i \leftarrow \tilde{\beta}_0 + \beta \alpha Z_i + \tilde{\delta} X_i + \tilde{\gamma} X_i + \tilde{\xi}_i. \quad (35.86)$$

The point is that β , the average treatment effect of A on Y , is equal to the coefficient $\beta\alpha$ of the instrument in the outcome-instrument model divided by the coefficient α of the instrument in the treatment-instrument model. So, to estimate the treatment effect, we simply fit both linear models and divide the estimated coefficients. This procedure is called **two stage least squares**.

The simplicity of this procedure is seductive. However, the required linearity assumptions are hard to satisfy in practice and frequently lead to severe issues. A particularly pernicious version of this is that linear-model misspecification together with weak relevance can yield standard errors for the estimate that are far too small. In practice, this can lead us to find large, significant estimates from two stage least squares when the truth is actually a weak or null effect. See [Rei16; You19; ASS19; Lal+21] for critical evaluations of two stage least squares in practice.

35.6 Difference in Differences

Unsurprisingly, time plays an important role in causality. Causes precede effects, and we should be able to incorporate this knowledge into causal identification. We now turn to a particular strategy for causal identification that relies on observing each unit at multiple time points. Data of this kind is sometimes called **panel data**. We'll consider the simplest case. There are two time periods. In the first period, none of the units are treated, and we observe an outcome Y_{0i} for each unit. Then, a subset of the units are treated, denoted by $A_i = 1$. In the second time period, we again observe the outcomes Y_{1i} for each unit, where now the outcomes of the treated units are affected by the treatment. Our goal is to determine the average effect receiving the treatment had on the treated units. That is, we want to know the average difference between the outcomes we actually observed for the treated units, and the outcomes we would have observed on those same units if they had not been treated. The general strategy we look at is called **difference in differences**.

As a concrete motivating example, consider trying to determine the effect raising minimum wage on employment. The concern here is that, in an efficient labor market, increasing the price of workers

will reduce the demand for them, thereby driving down employment. As such, it seems increasing minimum wage may hurt the people the policy is nominally intended to help. The question is: how strong is this effect in practice? Card and Krueger [CK94a] studied this effect using difference in differences. The Philadelphia metropolitan area includes regions in both Pennsylvania and New Jersey (different US states). On April 1st 1992, New Jersey raised its minimum wage from \$4.25 to \$5.05. In Pennsylvania, the wage remained constant at \$4.25. The strategy is to collect employment data from fast food restaurants (which pay many employees minimum wage) in each state before and after the change in minimum wage. In this case, for restaurant i , we have Y_{0i} , the number of full time employees in February 1992, and Y_{1i} , the number of full time employees in November 1992. The treatment is simply $A_i = 1$ if the restaurant was located in New Jersey, and $A_i = 0$ if located in Pennsylvania. Our goal is to estimate the average effect of the minimum wage hike on employment in the restaurants affected by it (i.e., the ones in New Jersey).

The assumption in classical difference-in-differences is the following structural equation:

$$Y_{ti} \leftarrow W_i + S_t + \tau A_i 1[t = 1] + \xi_{ti}, \quad (35.87)$$

with $\mathbb{E}[\xi_{ti}|W_i, S_t, A_i] = 0$. Here, W_i is a unit specific effect that is constant across time (e.g., the location of the restaurant or competence of the management) and S_t is a time-specific effect that applies to all units (e.g., the state of the US economy at each time). Both of these quantities are treated as unobserved, and not explicitly accounted for. The parameter τ captures the target causal effect. The (strong) assumption here is that unit, time, and treatment effects are all additive. This assumption is called **parallel trends**, because it is equivalent to assuming that, in the absence of treatment, the trend over time would be the same in both groups. It's easy to see that under this assumption, we have:

$$\tau = \mathbb{E}[Y_{1i} - Y_{0i}|A = 1] - \mathbb{E}[Y_{1i} - Y_{0i}|A = 0]. \quad (35.88)$$

That is, the estimand first computes the difference across time for both the treated and untreated group, and then computes the difference between these differences across the groups. The obvious estimator is then

$$\hat{\tau} = \frac{1}{n_A} \sum_{i:A_i=1} Y_{1i} - Y_{0i} - \frac{1}{n - n_A} \sum_{i:A_i=0} Y_{1i} - Y_{0i}, \quad (35.89)$$

where n_A is the number of treated units.

The root identification problem addressed by difference-in-differences is that $\mathbb{E}[W_i|A_i = 1] \neq \mathbb{E}[W_i|A_i = 0]$. That is, restaurants in New Jersey may be systematically different from restaurants in Pennsylvania in unobserved ways that affect employment.⁹ This is why we can't simply compare average outcomes for the treated and untreated. The identification assumption is that this unit-specific effect is the only source of statistical association with treatment; in particular we assume the time-specific effect has no such issue: $\mathbb{E}[S_{1i} - S_{0i}|A_i = 1] = \mathbb{E}[S_{1i} - S_{0i}|A_i = 0]$. Unfortunately, this assumption can be too strong. For instance, administrative data shows employment in Pennsylvania falling relative to employment in New Jersey between 1993 and 1996 [AP08, §5.2]. Although this

9. This is similar to the issue that arises from unobserved confounding, except W_i need not be a cause of the treatment assignment.

doesn't directly contradict the parallel trends assumption used for identification, which needs to hold only in 1992, it does make it seem less credible.

To weaken the assumption, we'll look at a version that requires parallel trends to hold only after adjusting for covariates. To motivate this, we note that there were several different types of fast food restaurant included in the employment data. These vary, e.g., in the type of food they serve, and in cost per meal. Now, it seems reasonable the trend in employment may depend on the type of restaurant. For example, more expensive chains (such as Kentucky Fried Chicken) might be more affected by recessions than cheaper chains (such as McDonald's). If expensive chains are more common in New Jersey than in Pennsylvania, this effect can create a violation of parallel trends—if there's recession affecting both states, we'd expect employment to go down more in New Jersey than in Pennsylvania. However, we may find it credible that McDonald's restaurants in New Jersey have the same trend as McDonald's in Pennsylvania, and similarly for Kentucky Fried Chicken.

The next step is to give a definition of the target causal effect that doesn't depend on a parametric model, and a non-parametric statement of the identification assumption to go with it. In words, the causal estimand will be the average treatment effect on the units that received the treatment. To make sense of this mathematically, we'll introduce a new piece of notation:

$$P^{A=1}(Y|\text{do}(A=a)) \triangleq \int P(Y|A=a, \text{parents of } Y) dP(\text{parents of } Y|A=1) \quad (35.90)$$

$$\mathbb{E}^{A=1}[Y|\text{do}(A=a)] \triangleq \mathbb{E}_{P^{A=1}(Y|\text{do}(A=a))}[Y]. \quad (35.91)$$

In words: recall that the ordinary do operator works by replacing $P(\text{parents}|A=a)$ by the marginal distribution $P(\text{parents})$, thereby breaking the backdoor associations. Now, we're replacing the distribution $P(\text{parents}|A=a)$ by $P(\text{parents}|A=1)$, irrespective of the actual treatment value. This still breaks all backdoor associations, but is a better match for our target of estimating the treatment effect only among the treated units.

To formalize a causal estimand using the do calculus, we need to assume some partial causal structure. We'll use the graph in Figure 35.7. With this in hand, our causal estimand is the average treatment effect on the units that received the treatment, namely:

$$\text{ATT}^{\text{DiD}} = \mathbb{E}^{A=1}[Y_1 - Y_0|\text{do}(A=1)] - \mathbb{E}^{A=1}[Y_1 - Y_0|\text{do}(A=0)] \quad (35.92)$$

In the minimum wage example, this is the average effect of the minimum wage hike on employment in the restaurants affected by it (i.e., the ones in New Jersey).

Finally, we formalize the identification assumption that, conditional on X , the trends in the treated and untreated groups are the same. The **conditional parallel trends** assumption is:

$$\mathbb{E}^{A=1}[Y_1 - Y_0|X, \text{do}(A=0)] = \mathbb{E}[Y_1 - Y_0|X, A=0]. \quad (35.93)$$

In words, this says that for treated units with covariates X , the trend we would have seen had we not assigned treatment is the same as the trend we actually saw for the untreated units with covariates X . That is, if New Jersey had not raised its minimum wage, then McDonald's in New Jersey would have the same expected change in employment as McDonald's in Pennsylvania.

With this in hand, we can give the main identification result:

Theorem 5 (Difference in Differences Identification). *We observe $A, Y_0, Y_1, X \sim P$. Suppose that*

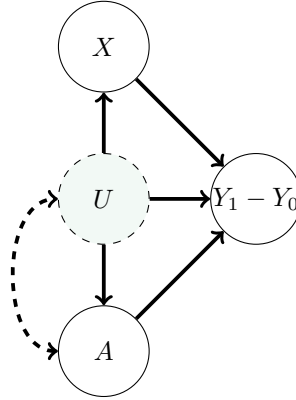


Figure 35.7: Causal graph assumed for the difference-in-differences setting. Here, the outcome of interest is the difference between the pre- and post-treatment period, $Y_1 - Y_0$. This difference is influenced by the treatment, unobserved factors U , and observed covariates X . The dashed arrow between U and A indicates a statistical dependency between the variables, but where we remain agnostic to the precise causal mechanism. For example, in the minimum wage example, U might be the average income in restaurant's neighbourhood, which is dependent on the state, and hence also the treatment.

1. (Causal Structure) The data follows the causal graph in Figure 35.7.
2. (Conditional Parallel Trends) $\mathbb{E}^{A=1}[Y_1 - Y_0 | X, \text{do}(A = 0)] = \mathbb{E}[Y_1 - Y_0 | X, A = 0]$.
3. (Overlap) $P(A = 1) > 0$ and $P(A = 1 | X = x) < 1$ for all values of x in the sample space. That is, there are no covariate values that only exist in the treated group.

Then, the average treatment effect on the treated is identified as $\text{ATT}^{\text{DiD}} = \tau^{\text{DiD}}$, where

$$\tau^{\text{DiD}} = \mathbb{E}[\mathbb{E}[Y_1 - Y_0 | A = 1, X] - \mathbb{E}[Y_1 - Y_0 | A = 0, X] | A = 1]. \quad (35.94)$$

Proof. First, by unrolling definitions, we have that

$$\mathbb{E}^{A=1}[Y_1 - Y_0 | \text{do}(A = 1), X] = \mathbb{E}[Y_1 - Y_0 | A = 1, X]. \quad (35.95)$$

The interpretation is the near-tautology that the average effect among the treated under treatment is equal to the actually observed average effect among the treated. Next,

$$\mathbb{E}^{A=1}[Y_1 - Y_0 | \text{do}(A = 0), X] = \mathbb{E}[Y_1 - Y_0 | A = 0, X]. \quad (35.96)$$

is just the conditional parallel trends assumption. The result follows immediately.

(The overlap assumption is required to make sure all the conditional expectations are well defined). \square

35.6.1 Estimation

With the identification result in hand, the next task is to estimate the observational estimand Equation (35.94). To that end, we define $\tilde{Y} \triangleq Y_1 - Y_0$. Then, we've assumed that $\tilde{Y}, X, A \stackrel{\text{iid}}{\sim} P$ for some unknown distribution P , and our target estimand is $\mathbb{E}[\mathbb{E}[\tilde{Y}|A = 1, X] - \mathbb{E}[\tilde{Y}|A = 0, X]|A = 1]$. We can immediately recognize this as the observational estimand that occurs in estimating the average treatment effect through adjustment, described in Section 35.4.5.3. That is, even though the causal situation and the identification argument are different between the adjustment setting and the difference in differences setting, the statistical estimation task we end up with is the same. Accordingly, we can use all of the estimation tools we developed for adjustment. That is, all of the techniques there—expected outcome modeling, propensity score methods, double machine learning, and so forth—were purely about the *statistical* task, which is the same between the two scenarios.

So, we're left with the same general recipe for estimation we saw in Section 35.4.6. Namely,

1. fit statistical or machine-learning models $\hat{Q}(a, x)$ as a predictor for $\tilde{Y} = Y_1 - Y_0$, and/or $\hat{g}(x)$ as a predictor for A
2. compute the predictions $\hat{Q}(0, x_i), \hat{Q}(1, x_i), \hat{g}(x_i)$ for each data point, and
3. combine these predictions into an estimate of the average treatment effect on the treated.

The estimator in the third step can be the expected outcome model estimator, the propensity weighted estimator, the double machine learning estimator, or any other strategy that's valid in the adjustment setting.

35.7 Credibility Checks

Once we've chosen an identification strategy, fit our models, and produced an estimate, we're faced with a basic question: should we believe it? Whether the reported estimate succeeds in capturing the true causal effect depends on whether the assumptions required for causal identification hold, the quality of the machine learning models, and the variability in the estimate due to only having access to a finite data sample. The latter two problems are already familiar from machine learning and statistical practice. We should, e.g., assess our models by checking performance on held out data, examining feature importance, and so forth. Similarly, we should report measures of the uncertainty due to finite sample (e.g., in the form of confidence intervals). Because these procedures are already familiar practice, we will not dwell on them further. However, model evaluation and uncertainty quantification are key parts of any credible causal analysis.

Assessing the validity of identification assumptions is trickier. First, there are assumptions that can in fact be checked from data. For example, overlap should be checked in analysis using backdoor adjustment or difference in differences, and relevance should be checked in the instrumental variable setting. Again, checking these conditions is absolutely necessary for a credible causal analysis. But, again, this involves only familiar data analysis, so we will not discuss it further. Next, there are the causal assumptions that cannot be verified from data; e.g., no unobserved confounding in backdoor adjustment, the exclusion restriction in IV, and conditional parallel trends in DiD. Ultimately, the validity of these assumptions must be assessed using substantive causal knowledge of the particular problem under consideration. However, it is possible to conduct some supplementary analyses that make the required judgement easier. We now discuss two such techniques.

35.7.1 Placebo Checks

In many situations we may be able to find a variable that can be interpreted as a “treatment” that is known to have no effect on the outcome, but which we expect to be confounded with the outcome in a very similar fashion to the true treatment of interest. For example, if we’re trying to estimate the efficacy of a COVID vaccine in preventing symptomatic COVID, we might take our placebo treatment to be vaccination against HPV. We do not expect that there’s any causal effect here. However, it seems plausible that latent factors that cause an individual to seek (or avoid) HPV vaccination and COVID vaccination are similar; e.g., health conscientiousness, fear of needles, and so forth. Then, if our identification strategy is valid for the COVID vaccine, we’d also expect it to be to be valid for HPV vaccination. Accordingly, our estimation procedure we use for estimating the COVID effect should, when applied to HPV, yield $\hat{\tau} \approx 0$. Or, more precisely, the confidence interval should contain 0. If this does not happen, then we may suspect that there are still some confounding factors lurking that are not adequately handled by the identification procedure.

A similar procedure works when there is a variable that can be interpreted as an outcome which is known to not be affected by the treatment, but that shares confounders with the outcome we’re actually interested in. For example, in the COVID vaccination case, we might take the null outcome to be symptomatic COVID within 7 days of vaccination [Dag+21]. Our knowledge of both the biological mechanism of vaccination and the amount of time it takes to develop symptoms after COVID infection (at least 2 days) lead us to conclude that it’s unlikely that the treatment has a causal effect on the outcome. However, the properties of the treated people that affect how likely they are to develop symptomatic COVID are largely the same in the 7 day and, e.g., 6 month window. That includes factors such as risk aversion, baseline health, and so forth. Again, we can apply our identification strategy to estimate the causal effect of the treatment on the null outcome. If the confidence interval does not include 0, then we should doubt the credibility of the analysis.

35.7.2 Sensitivity Analysis to Unobserved Confounding

We now specialize to the case of estimating the average causal effect of a binary treatment by adjusting for confounding variables, as described in Section 35.4. In this case, causal identification is based on the assumption of ‘no unobserved confounding’; i.e., the assumption that the observed covariates include all common causes of the treatment assignment and outcome. This assumption is fundamentally untestable from observed data, but its violation can induce bias in the estimation of the treatment effect—the unobserved confounding may completely or in part explain the observed association. Our aim in this part is to develop a sensitivity analysis tool to aid in reasoning about potential bias induced by unobserved confounding.

Intuitively, if we estimate a large positive effect then we might expect the real effect is also positive, even in the presence of mild unobserved confounding. For example, consider the association between smoking and lung cancer. One could argue that this association arises from a hormone that both predisposes carriers to both an increased desire to smoke and to a greater risk of lung cancer. However, the association between smoking and lung cancer is large—is it plausible that some unknown hormonal association could have a strong enough influence to explain the association? Cornfield et al. [Cor+59] showed that, for a particular observational dataset, such an unmeasured hormone would need to increase the probability of smoking by at least a factor of nine. This is an unreasonable effect size for a hormone, so they conclude it’s unlikely the causal effect can be

explained away.

We would like a general procedure to allow domain experts to make judgments about whether plausible confounding is “mild” relative to the “large” effect. In particular, the domain expert must translate judgments about the strength of the unobserved confounding into judgments about the bias induced in the estimate of the effect. Accordingly, we must formalize what is meant by strength of unobserved confounding, and to show how to translate judgments about confounding strength into judgments about bias.

A prototypical example, due to Imbens [Imb03] (building on [RR83]), illustrates the broad approach. As above, the observed data consists of a treatment A , an outcome Y , and covariates X that may causally affect the treatment and outcome. Imbens [Imb03] then posits an additional unobserved binary confounder U for each patient, and supposes that the observed data and unobserved confounder were generated according to the following assumption, known as **Imbens’ Sensitivity Model**:

$$U_i \stackrel{\text{iid}}{\sim} \text{Bern}(1/2) \quad (35.97)$$

$$A_i | X_i, U_i \stackrel{\text{ind}}{\sim} \text{Bern}(\text{sig}(\gamma X_i + \alpha U_i)) \quad (35.98)$$

$$Y_i | X_i, A_i, U_i \stackrel{\text{ind}}{\sim} \mathcal{N}(\tau A_i + \beta X_i + \delta U_i, \sigma^2). \quad (35.99)$$

where sig is the sigmoid function.

If we had observed U_i , we could estimate $(\hat{\tau}, \hat{\gamma}, \hat{\beta}, \hat{\alpha}, \hat{\delta}, \hat{\sigma}^2)$ from the data and report $\hat{\tau}$ as the estimate of the average treatment effect. Since U_i is not observed, it is not possible to identify the parameters from the data. Instead, we make (subjective) judgments about plausible values of α —how strongly U_i affects the treatment assignment—and δ —how strongly U_i affects the outcome. Contingent on plausible $\alpha = \alpha^*$ and $\delta = \delta^*$, the other parameters can be estimated. This yields an estimate of the treatment effect $\hat{\tau}(\alpha^*, \delta^*)$ under the presumed values of the sensitivity parameters.

The approach just outlined has a major drawback: it relies on a parametric model for the full data generating process. The assumed model is equivalent to assuming that, had U been observed, it would have been appropriate to use logistic regression to model treatment assignment, and linear regression to model the outcome. This assumption also implies a simple, parametric model for the relationships governing the observed data. This restriction is out of step with modern practice, where we use flexible machine-learning methods to model these relationships. For example, the assumption forbids the use of neural networks or random forests, though such methods are often state-of-the-art for causal effect estimation.

Austen plots We now turn to developing an alternative an adaptation of Imbens’ approach that fully decouples sensitivity analysis and modeling of the observed data. Namely, the **Austen plots** of [VZ20]. An example Austen plot is shown in Figure 35.8. The high-level idea is to posit a generative model that uses a simple, interpretable parametric form for the influence of the unobserved confounder, but that *puts no constraints on the model for the observed data*. We then use the parametric part of the model to formalize “confounding strength” and to compute the induced bias as a function of the confounding.

Austen plots further adapt two strategies pioneered by Imbens [Imb03]. First, we find a parameterization of the model so that the sensitivity parameters, measuring strength of confounding, are on a standardized, unitless scale. This allows us to compare the strength of hypothetical unobserved confounding to the strength of observed covariates, measured from data. Second, we plot the curve

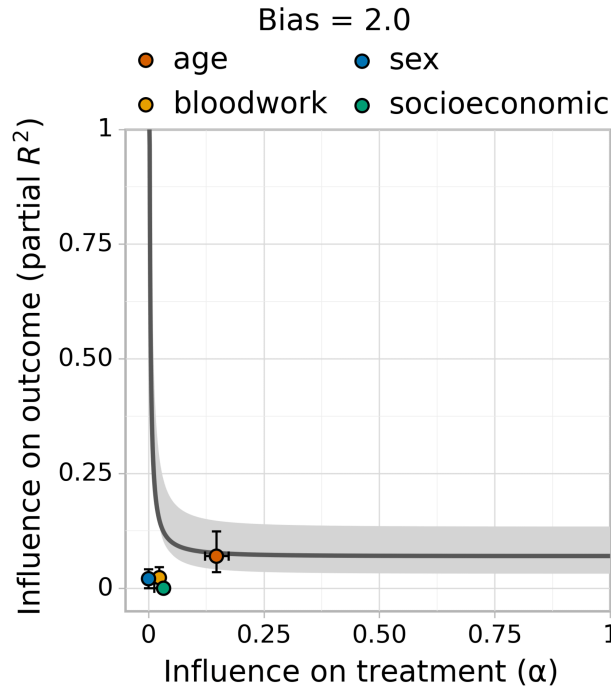


Figure 35.8: Austen plot showing how strong an unobserved confounder would need to be to induce a bias of 2 in an observational study of the effect of combination blood pressure medications on diastolic blood pressure [Dor+16]. We chose this bias to equal the nominal average treatment effect estimated from the data. We model the outcome with Bayesian Additive Regression Trees and the treatment assignment with logistic regression. The curve shows all values treatment and outcome influence that would induce a bias of 2. The colored dots show the influence strength of (groups of) observed covariates, given all other covariates. For example, an unobserved confounder with as much influence as the patient’s age might induce a bias of about 2.

of all values of the sensitivity parameter that would yield given level of bias. This moves the analyst judgment from “what are plausible values of the sensitivity parameters?” to “are sensitivity parameters this extreme plausible?”

Figure 35.8, an Austen plot for an observational study of the effect of combination medications on diastolic blood pressure, illustrates the idea. A bias of 2 would suffice to undermine the qualitative conclusion that the blood-pressure treatment is effective. Examining the plot, an unobserved confounder as strong as age could induce this amount of confounding, but no other (group of) observed confounders has so much influence. Accordingly, if a domain expert thinks an unobserved confounder as strong as age is unlikely then they may conclude that the treatment is likely effective. Or, if such a confounder is plausible, they may conclude that the study fails to establish efficacy.

Setup The data are generated independently and identically $(Y_i, A_i, X_i, U_i) \stackrel{\text{iid}}{\sim} P$, where U_i is not observed and P is some unknown probability distribution. The approach in Section 35.4 assumes that the observed covariates X contain all common causes of Y and A . If this ‘no unobserved confounding’ assumption holds, then the ATE is equal to parameter, τ , of the observed data distribution, where

$$\tau = \mathbb{E}[\mathbb{E}[Y|X, A = 1] - \mathbb{E}[Y|X, A = 0]]. \quad (35.100)$$

This observational parameter is then estimated from a finite data sample. Recall from Section 35.4 that this involves estimating the conditional expected outcome $Q(A, X) = \mathbb{E}[Y|A, X]$ and the propensity score $g(X) = P(A = 1|X)$, then plugging these into an estimator $\hat{\tau}$.

We are now concerned with the case of possible unobserved confounding. That is, where U causally affects Y and A . If there is unobserved confounding then the parameter τ is not equal to the ATE, so $\hat{\tau}$ is a biased estimate. Inference about the ATE then divides into two tasks. First, the statistical task: estimating τ as accurately as possible from the observed data. And, second, the causal (domain-specific) problem of assessing $\text{bias} = \text{ATE} - \tau$. We emphasize that our focus here is bias due to causal misidentification, not the statistical bias of the estimator. Our aim is to reason about the bias induced by unobserved confounding—the second task—in a way that imposes no constraints on the modeling choices for \hat{Q} , \hat{g} and $\hat{\tau}$ used in the statistical analysis.

Sensitivity Model Our sensitivity analysis should impose no constraints on how the *observed* data is modeled. However, sensitivity analysis demands some assumption on the relationship between the observed data and the *unobserved* confounder. It is convenient to formalize such assumptions by specifying a probabilistic model for how the data is generated. The strength of confounding is then formalized in terms of the parameters of the model (the sensitivity parameters). Then, the bias induced by the confounding can be derived from the assumed model. Our task is to posit a generative model that both yields a useful and easily interpretable sensitivity analysis, and that avoids imposing any assumptions about the observed data.

To begin, consider the functional form of the sensitivity model used by Imbens [Imb03].

$$\text{logit}P(A = 1|x, u) = h(x) + \alpha u \quad (35.101)$$

$$\mathbb{E}[Y|a, x, u] = l(a, x) + \delta u, \quad (35.102)$$

for some functions h and l . That is, the propensity score is logit-linear in the unobserved confounder, and the conditional expected outcome is linear.

By rearranging Equation (35.101) to solve for u and plugging in to Equation (35.102), we see that it’s equivalent to assume $\mathbb{E}[Y|t, x, u] = \tilde{l}(t, x) + \tilde{\delta} \text{logit}P(A = 1|x, u)$. That is, the unobserved confounder u only influences the outcome through the propensity score. Accordingly, by positing a distribution on $P(A = 1|x, u)$ directly, we can circumvent the need to explicitly articulate U (and h).

Definition 35.7.1. Let $\tilde{g}(x, u) = P(A = 1|x, u)$ denote the propensity score given observed covariates x and the unobserved confounder u .

The insight is that we can posit a sensitivity model by defining a distribution on \tilde{g} directly. We choose:

$$\tilde{g}(X, U)|X \sim \text{Beta}(g(X)(1/\alpha - 1), (1 - g(X))(1/\alpha - 1)).$$

That is, the full propensity score $\tilde{g}(X, U)$ for each unit is assumed to be sampled from a Beta distribution centered at the observed propensity score $g(X)$. The sensitivity parameter α plays the same role as in Imbens' model: it controls the influence of the unobserved confounder U on treatment assignment. When α is close to 0 then $\tilde{g}(X, U)|X$ is tightly concentrated around $g(X)$, and the unobserved confounder has little influence. That is, U minimally affects our belief about who is likely to receive treatment. Conversely, when α is close to 1 then \tilde{g} concentrates near 0 and 1; i.e., knowing U would let us accurately predict treatment assignment. Indeed, it can be shown that α is the change in our belief about how likely a unit was to have gotten the treatment, given that they were actually observed to be treated (or not):

$$\alpha = \mathbb{E}[\tilde{g}(X, U)|A = 1] - \mathbb{E}[\tilde{g}(X, U)|A = 0]. \quad (35.103)$$

With the \tilde{g} model in hand, we define the **Austen Sensitivity Model** as follows:

$$\tilde{g}(X, U)|X \sim \text{Beta}(g(X)(1/\alpha - 1), (1 - g(X))(1/\alpha - 1)) \quad (35.104)$$

$$A|X, U \sim \text{Bern}(\tilde{g}(X, U)) \quad (35.105)$$

$$\mathbb{E}[Y|A, X, U] = Q(A, X) + \delta(\text{logit}\tilde{g}(X, U) - \mathbb{E}[\text{logit}\tilde{g}(X, U)|A, X]). \quad (35.106)$$

This model has been constructed to satisfy the requirement that the propensity score and conditional expected outcome are the g and Q actually present in the observed data:

$$\mathbb{P}(A = 1|X) = \mathbb{E}[\mathbb{E}[T|X, U]|X] = \mathbb{E}[\tilde{g}(X, U)|X] = g(X)$$

$$\mathbb{E}[Y|A, X] = \mathbb{E}[\mathbb{E}[Y|A, X, U]|A, X] = Q(A, X).$$

The sensitivity parameters are α , controlling the dependence between the unobserved confounder the treatment assignment, and δ , controlling the relationship with the outcome.

Bias We now turn to calculating the bias induced by unobserved confounding. By assumption, X and U together suffice to render the average treatment effect identifiable as:

$$\text{ATE} = \mathbb{E}[\mathbb{E}[Y|A = 1, X, U] - \mathbb{E}[Y|A = 0, X, U]].$$

Plugging in our sensitivity model yields,

$$\text{ATE} = \mathbb{E}[Q(1, X) - Q(0, X)] + \delta(\mathbb{E}[\text{logit}\tilde{g}(X, U)|X, A = 1] - \mathbb{E}[\text{logit}\tilde{g}(X, U)|X, A = 0]).$$

The first term is the observed-data estimate τ , so

$$\text{bias} = \delta(\mathbb{E}[\text{logit}\tilde{g}(X, U)|X, A = 1] - \mathbb{E}[\text{logit}\tilde{g}(X, U)|X, A = 0]).$$

Then, by invoking Beta-Bernoulli conjugacy and standard Beta identities,¹⁰ we arrive at,

Theorem 6. *Under the Austen sensitivity model, Equation (35.106), an unobserved confounder with influence α and δ induces bias in the estimated treatment effect equal to*

$$\text{bias} = \frac{\delta}{1/\alpha - 1} \mathbb{E}\left[\frac{1}{g(X)} + \frac{1}{1 - g(X)}\right].$$

10. We also use the recurrence relation $\psi(x + 1) - \psi(x) = 1/x$, where ψ is the digamma function.

That is, the amount of bias is determined by the sensitivity parameters and by the *realized* propensity score. Notice that more extreme propensity scores lead to more extreme bias in response to unobserved confounding. This means, in particular, that conditioning on a covariate that affects the treatment but that does not directly affect the outcome (an instrument) will increase any bias due to unobserved confounding. This general phenomena is known as **z-bias**.

Sensitivity Parameters The Austen model provides a formalization of confounding strength in terms of the parameters α and δ and tells us how much bias is induced by a given strength of confounding. This lets us translate judgments about confounding strength to judgments about bias. However, it is not immediately obvious how to translate qualitative judgements such as “I think any unobserved confounder would be much less important than Age” to judgements about the possible values of the sensitivity parameters.

First, because the scale of δ is not fixed, it may be difficult to compare the influence of potential unobserved confounders to the influence of reference variables. To resolve this, we reexpress the outcome-confounder strength in terms of the (non-parametric) partial coefficient of determination:

$$R_{Y,\text{par}}^2(\alpha, \delta) = 1 - \frac{\mathbb{E}(Y - \mathbb{E}[Y|A, X, U])^2}{\mathbb{E}(Y - Q(A, X))^2}.$$

The key to computing the reparameterization is the following result

Theorem 7. *Under the Austen sensitivity model, Equation (35.106), the outcome influence is*

$$R_{Y,\text{par}}^2(\alpha, \delta) = \delta^2 \sum_{a=0}^1 \frac{\mathbb{E}[\psi_1(g(X)^a(1-g(X))^{1-a}(1/\alpha - 1) + 1[A=a])]}{\mathbb{E}[(Y - Q(A, X))^2]},$$

where ψ_1 is the trigamma function.

See Veitch and Zaveri [VZ20] for the proof.

By design, α —the strength of confounding influence on on treatment assignment—is already on a fixed, unitless scale. However, because the measure is tied to the model it may be difficult to interpret, and it is not obvious how to compute reference confounding strength values from the observed data. The next result clarifies these issues.

Theorem 8. *Under the Austen sensitivity model, Equation (35.106),*

$$\alpha = 1 - \frac{\mathbb{E}[\tilde{g}(X, U)(1 - \tilde{g}(X, U))]}{\mathbb{E}[g(X)(1 - g(X))]}.$$

See Veitch and Zaveri [VZ20] for the proof. That is, the sensitivity parameter α is measures how much more extreme the propensity scores become when we condition on U . That is, α is a measure of the extra predictive power U adds for A , above and beyond the predictive power in X . It may also be insightful to notice that

$$\alpha = R_{A,\text{par}}^2 = 1 - \frac{\mathbb{E}[(A - \tilde{g}(X, U))^2]}{\mathbb{E}[(A - g(X))^2]}. \quad (35.107)$$

That is, α is just the (non-parametric) partial coefficient of determination of U on A —the same measure used for the outcome influence. (To see this, just expand the expectations conditional on $A = 1$ and $A = 0$).

Estimating bias In combination, Theorems 6 and 7 yield an expression for the bias in terms of α and $R_{Y,\text{par}}^2$. In practice, we can estimate the bias induced by confounding by fitting models for \hat{Q} and \hat{g} and replacing the expectations by means over the data.

35.7.2.1 Calibration using observed data

The analyst must make judgments about the influence a hypothetical unobserved confounder might have on treatment assignment and outcome. To calibrate such judgments, we'd like to have a reference point for how much the observed covariates influence the treatment assignment and outcome. In the sensitivity model, the degree of influence is measured by partial R_Y^2 and α . We want to measure the degree of influence of an observed covariate Z given the other observed covariates $X \setminus Z$.

For the outcome, this can be measured as:

$$R_{Y \cdot Z|T, X \setminus Z}^2 \triangleq 1 - \frac{\mathbb{E}(Y - Q(A, X))^2}{\mathbb{E}(Y - \mathbb{E}[Y|A, X \setminus Z])^2}.$$

In practice, we can estimate the quantity by fitting a new regression model \hat{Q}_Z that predicts Y from A and $X \setminus Z$. Then we compute

$$R_{Y \cdot Z|T, X \setminus Z}^2 = 1 - \frac{\frac{1}{n} \sum_i (y_i - \hat{Q}(t_i, x_i))^2}{\frac{1}{n} \sum_i (y_i - \hat{Q}_Z(t_i, x_i \setminus z_i))^2}.$$

Using Theorem 8, we can measure influence of observed covariate Z on treatment assignment given $X \setminus Z$ in an analogous fashion to the outcome. We define $g_{X \setminus Z}(X \setminus Z) = P(A = 1|X \setminus Z)$, then fit a model for $g_{X \setminus Z}$ by predicting A from $X \setminus Z$, and estimate

$$\hat{\alpha}_{Z|X \setminus Z} = 1 - \frac{\frac{1}{n} \sum_i \hat{g}(x_i)(1 - \hat{g}(x_i))}{\frac{1}{n} \sum_i \hat{g}_{X \setminus Z}(x_i \setminus z_i)(1 - \hat{g}_{X \setminus Z}(x_i \setminus z_i))}.$$

Grouping covariates The estimated values $\hat{\alpha}_{X \setminus Z}$ and $\hat{R}_{Y, X \setminus Z}^2$ measure the influence of Z conditioned on all the other confounders. In some cases, this can be misleading. For example, if some piece of information is important but there are multiple covariates providing redundant measurements, then the estimated influence of each covariate will be small. To avoid this, group together related or strongly dependent covariates and compute the influence of the entire group in aggregate. For example, grouping income, location, and race as ‘socioeconomic variables’.

35.7.2.2 Practical Use

We now have sufficient results to produce Austen plots such as Figure 35.8. At a high level, the procedure is:

1. Produce an estimate $\hat{\tau}$ using any modeling tools. As a component of this, estimate the propensity score \hat{g} and conditional outcome model \hat{Q}
2. Pick a level of bias that would suffice to change the qualitative interpretation of the estimate (e.g., the lower bound of a 95% confidence interval).

3. Plot the values of α and $R_{Y,\text{par}}^2$ that would suffice to induce that much bias. This is the black curve on the plot. To calculate these values, use Theorems 6 and 7 together with the estimated \hat{g} and \hat{Q} .

4. Finally, compute reference influence level for (groups of) observed covariates. In particular, this requires fitting reduced models for the conditional expected outcome and propensity that do not use the reference covariate as a feature.

In practice, an analyst only needs to do the model fitting parts themselves. The bias calculations, reference value calculations, and plotting can be done automatically with standard libraries.¹¹

Austen plots are predicated on Equation (35.106). This assumption replaces the purely parametric Equation (35.99) with a version that eliminates any parametric requirements on the observed data. However, we emphasize that Equation (35.106) does, implicitly, impose some parametric assumption on the structural causal relationship between U and A, Y . Ultimately, any conclusion drawn from the sensitivity analysis depends on this assumption, which is not justified on any substantive grounds. Accordingly, such sensitivity analyses can only be used to informally guide domain experts. They do not circumvent the need to thoroughly adjust for confounding. This reliance on a structural assumption is a generic property of sensitivity analysis.¹² Indeed, there are now many sensitivity analysis models that allow the use of any machine learning model in the data analysis [e.g., RRS00; FDF19; She+11; HS13; BK19; Ros10; Yad+18; ZSB19; Sch+21a]. However, none of these are yet in routine use in practice. We have presented Austen plots here not because they make an especially virtuous modeling assumption, but because they are (relatively) easy to understand and interpret.

Austen plots are most useful in situations where the conclusion from the plot would be ‘obvious’ to a domain expert. For instance, in Figure 35.8, we can be confident that an unobserved confounder similar to socioeconomic status would not induce enough bias to change the qualitative conclusion. By contrast, Austen plots should not be used to draw conclusions such as, “I think a latent confounder could only be 90% as strong as ‘age’, so there is evidence of a small non-zero effect”. Such nuanced conclusions might depend on issues such as the particular sensitivity model we use, or finite-sample variation of our bias and influence estimates, or on incautious interpretation of the calibration dots. These issues are subtle, and it would be difficult resolve them to a sufficient degree that a sensitivity analysis would make an analysis credible.

Calibration using observed data The interpretation of the observed-data calibration requires some care. The sensitivity analysis requires the analyst to make judgements about the strength of influence of the unobserved confounder U , *conditional on the observed covariates* X . However, we report the strength of influence of observed covariate(s) Z , *conditional on the other observed covariates* $X \setminus Z$. The difference in conditioning sets can have subtle effects.

Cinelli and Hazlett [CH20] give an example where Z and U are identical variables in the true model, but where influence of U given A, X is larger than the influence of Z given $A, X \setminus Z$. (The influence of Z given $X \setminus Z, U$ would be the same as the influence of U given X). Accordingly, an analyst is *not* justified in a judgement such as, “I know that U and Z are very similar. I see Z has substantial influence, but the dot is below the line. Thus, U will not undo the study conclusions.” In

11. See github.com/vveitch/causality-tutorials/blob/main/Sensitivity_Analysis.ipynb

12. In extreme cases, there can be so little unexplained variation in A or Y that only a very weak confounder could be compatible with the data. In this case, essentially assumption free sensitivity analysis is possible [Man90].

essence, if the domain expert suspects a strong interaction between U and Z then naively eyeballing the dot-vs-line position may be misleading. A particular subtle case is when U and Z are independent variables that both strongly influence A and Y . The joint influence on A creates an interaction effect between them when A is conditioned on (the treatment is a collider). This affects the interpretation of $R_{Y,U|X,A}^2$. Indeed, we should generally be skeptical of sensitivity analysis interpretation when it is expected that a strong confounder has been omitted. In such cases, our conclusions may depend substantively on the particular form of our sensitivity model, or other unjustifiable assumptions.

Although the interaction problem is conceptually important, its practical significance is unclear. We often expect the opposite effect: if U and Z are dependent (e.g., race and wealth) then omitting U should increase the apparent importance of Z —leading to a conservative judgement (a dot artificially towards the top right part of the plot).

35.8 The Do Calculus

We have seen several strategies for identifying causal effects as parameters of observational distributions. Confounder adjustment (Section 35.4) relied only on the assumed causal graph (and overlap), which specified that we observe all common causes of A and Y . On the other hand, instrumental variable methods and difference-in-differences each relied on both an assumed causal graph and partial functional form assumptions about the underlying structural causal model. Because functional form assumptions can be quite difficult to justify on substantive grounds, it's natural to ask when causal identification is possible from the causal graph alone. That is, when can we be agnostic to the particular functional form of the structural causal models?

There is a general “**calculus of intervention**”, known as the **do-calculus**, that gives a general recipe for determining when the causal assumptions expressed in a causal graph can be used to identify causal effects [Pea09c]. The do-calculus is a set of three rewrite rules that allows us to replace statements where we condition on variables being set by intervention, e.g. $P(Y|\text{do}(A = a))$, with statements involving only observational quantities, e.g. $\mathbb{E}_X[P(Y|A = a, X)]$. When causal identification is possible, we can repeatedly apply the three rules to boil down our target causal parameter into an expression involving only the observational distribution.

35.8.1 The three rules

To express the rules, let X , Y , Z , and W be arbitrary disjoint sets of variables in a causal DAG G .

Rule 1 The first rule allows us to insert or delete observations z :

$$p(y|\text{do}(x), z, w) = p(y|\text{do}(x), w) \text{ if } (Y \perp Z|X, W)_{G_{\overline{X}}} \quad (35.108)$$

where $G_{\overline{X}}$ denotes cutting edges going into X , and $(Y \perp Z|X, W)_{G_{\overline{X}}}$ denotes conditional independence in the mutilated graph. The rule follows from d-separation in the mutilated graph. This rule just says that conditioning on irrelevant variables leaves the distribution invariant (as we would expect).

Rule 2 The second rule allows us to replace $\text{do}(z)$ with conditioning on (seeing) z . The simplest case where can do this is: if Z is a root of the causal graph (i.e., it has no causal parents) then $p(y|\text{do}(z)) = p(y|z)$. The reason is that the do operator is equivalent to conditioning in the mutilated

causal graph where all the edges into Z are removed, but, because Z is a root, the mutilated graph is just the original causal graph. The general form of this rule is:

$$p(y|\text{do}(x), \text{do}(z), w) = p(y|\text{do}(x), z, w) \text{ if } (Y \perp\!\!\!\perp Z|X, W)_{G_{\overline{XZ}}} \quad (35.109)$$

where $G_{\overline{XZ}}$ cuts edges going into X and out of Z . Intuitively, we can replace $\text{do}(z)$ by z as long as there are no backdoor (non-directed) paths between z and y . If there are in fact no such paths, then cutting all the edges going out of Z will mean there are no paths connecting Z and Y , so that $Y \perp\!\!\!\perp Z$. The rule just generalizes this line of reasoning to allow for extra observed and intervened variables.

Rule 3 The third rule allows us to insert or delete actions $\text{do}(z)$:

$$p(y|\text{do}(x), \text{do}(z), w) = p(y|\text{do}(x), w) \text{ if } (Y \perp\!\!\!\perp Z|X, W)_{G_{\overline{XZ^*}}} \quad (35.110)$$

where $G_{\overline{XZ^*}}$ cuts edges going into X and Z^* , and where Z^* is the set of Z -nodes that are not ancestors of any W -node in $G_{\overline{X}}$. Intuitively, this condition corresponds to intervening on X , and checking whether the distribution of Y is invariant to *any* intervention that we could apply on Z .

35.8.2 Revisiting Backdoor Adjustment

We begin with a more general form of the adjustment formula we used in Section 35.4.

First, suppose we observe all of A 's parents, call them X . For notational simplicity, we'll assume for the moment that X is discrete. Then,

$$p(Y = y|\text{do}(A = a)) = \sum_x p(Y = y|x, \text{do}(A = a))p(x|\text{do}(A = a)) \quad (35.111)$$

$$= \sum_x p(Y = y|x, A = a)p(x). \quad (35.112)$$

The first line is just a standard probability relation (marginalizing over z). We are using causal assumptions in two ways in the second line. First, $p(x|\text{do}(A = a)) = p(x)$: the treatment has no causal effect on Z , so interventions on A don't change the distribution of Z . This is rule 3, Equation (35.110). Second, $p(Y = y|z, \text{do}(A = a)) = p(Y = y|z, A = a)$. This equality holds because conditioning on the parents blocks all non-directed paths from A to Y , reducing the causal effect to be the same as the observational effect. The equality is an application of rule 2, Equation (35.109).

Now, what if we don't observe all the parents of A ? The key issue is **backdoor paths**: paths between A and Y that contain an arrow into A . These paths are the general form of the problem that occurs when A and Y share a common cause. Suppose that we can find a set of variables S such that (1) no node in S is a descendant of A ; and (2) S blocks every backdoor path between A and Y . Such a set is said to satisfy the **backdoor criterion**. In this case, we can use S instead of the parents of X in the adjustment formula, Equation (35.112). That is,

$$p(Y = y|\text{do}(A = a)) = \mathbb{E}_S[p(Y = y|S, A = a)]. \quad (35.113)$$

The proof follows the invocation of rules 3 and 2, in the same way as for the case where S is just the parents of A . Notice that requiring S to not contain any descendants of A means that we don't risk

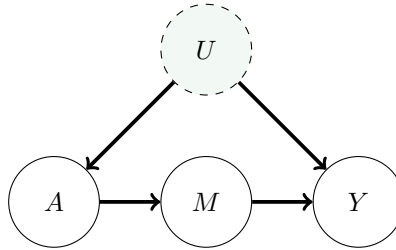


Figure 35.9: Causal graph illustrating the frontdoor criterion setup. The effect of the treatment A on outcome Y is entirely mediated by mediator M . This allows us infer the causal effect even if the treatment and outcome are confounded by U .

conditioning on any variables that mediate the effect, nor any variables that might be colliders—either would undermine the estimate.

The backdoor adjustment formula generalizes the adjust-for-parents approach and adjust-for-all-common-causes approach of Section 35.4. That’s because both the parents of A and the common causes satisfy the backdoor criterion.

In practice, the full distribution $p(Y = y | \text{do}(A = a))$ is rarely used as the causal target. Instead, we try to estimate a low-dimensional parameter of this distribution, such as the average treatment effect. The adjustment formula immediately translates in the obvious way. If we define

$$\tau = \mathbb{E}_S[\mathbb{E}[Y | A = 1, S] - \mathbb{E}[Y | A = 0, S]],$$

then we have that $\text{ATE} = \tau$ whenever S satisfies the backdoor criteria. The parameter τ can then be estimated from finite data using the methods described in Section 35.4, using S in place of the common causes X .

35.8.3 Frontdoor Adjustment

Backdoor adjustment is applicable if there’s at least one observed variable on every backdoor path between A and Y . As we have seen, identification is sometimes still possible even when this condition doesn’t hold. Frontdoor adjustment is another strategy of this kind. Figure 35.9 shows the causal structure that allows this kind of adjustment strategy. Suppose we’re interested in the effect of smoking A on developing cancer Y , but we’re concerned about some latent genetic confounder U .

Suppose that all of the directed paths from A to Y pass through some set of variables M . Such variables are called **mediators**. For example, the effect of smoking on lung cancer might be entirely mediated by the amount of tar in the lungs and measured tissue damage. It turns out that if all such mediators are observed, and the mediators do not have an unobserved common cause with A or Y , then causal identification is possible. To understand why this is true, first notice that we can identify the causal effect of A on M and the causal effect of M on A , both by backdoor adjustment. Further, the mechanism of action of A on Y is: A changes M which in turn changes Y . Then, we

1
2 can combine these as:

$$3 \quad p(Y|\text{do}(A = a)) = \sum_m p(Y|\text{do}(M = m))p(M = m|\text{do}(A = a)) \quad (35.114)$$

$$4 \quad = \sum_m \sum_{a'} p(Y|a', m)p(a')p(m|a) \quad (35.115)$$

5
6
7
8
9 The second line is just backdoor adjustment applied to identify each of the do expressions (note that
10 A blocks the M - Y backdoor path through U).

11 Equation (35.115) is called the **front-door formula** [Pea09b, §3.3.2]. To state the result in more
12 general terms, let us introduce a definition. We say a set of variables M satisfies the **front-door**
13 **criterion** relative to an ordered pair of variables (A, Y) if (1) M intercepts all directed paths from
14 A to Y ; (2) there is no unblocked backdoor path from A to M ; and (3) all backdoor paths from M
15 to Y are blocked by A . If M satisfies this criterion, and if $p(A, M) > 0$ for all values of A and M ,
16 then the causal effect of A on Y is identifiable and is given by Equation (35.115).

17 Let us interpret this theorem in terms of our smoking example. Condition 1 means that smoking
18 A should have no effect on cancer Y except via tar and tissue damage M . Conditions 2 and 3 mean
19 that the genotype U cannot have any effect on M except via smoking A . Finally, the requirement
20 that $p(A, M) > 0$ for all values implies that high levels of tar in the lungs must arise not only due to
21 smoking, but also other factors (e.g., pollutants). In other words, we require $p(A = 0, M = 1) > 0$ so
22 we can assess the impact of the mediator in the untreated setting.

23 We can now use the do-calculus to derive the frontdoor criterion; following [PM18b, p236]. Assuming
24 the causal graph G shown in Figure 35.9:

$$\begin{aligned} 25 \quad p(y|\text{do}(a)) &= \sum_m p(y|\text{do}(a), m)p(m|\text{do}(a)) && \text{(probability axioms)} \\ 26 \quad &= \sum_m p(y|\text{do}(a), \text{do}(m))p(m|\text{do}(a)) && \text{(rule 2 using } G_{\overline{ST}}) \\ 27 \quad &= \sum_m p(y|\text{do}(a), \text{do}(m))p(m|a) && \text{(rule 2 using } G_{\underline{S}}) \\ 28 \quad &= \sum_m p(y|\text{do}(m))p(m|a) && \text{(rule 3 using } G_{\overline{ST}^*}) \\ 29 \quad &= \sum_{a'} \sum_m p(y|\text{do}(m), a')p(a'|\text{do}(m))p(m|a) && \text{(probability axioms)} \\ 30 \quad &= \sum_{a'} \sum_m p(y|m, a')p(a'|\text{do}(m))p(m|a) && \text{(rule 2 using } G_{\underline{T}}) \\ 31 \quad &= \sum_{a'} \sum_m p(y|m, a')p(a')p(m|a) && \text{(rule 3 using } G_{\overline{T}^*}) \end{aligned}$$

32
33
34
35
36
37
38
39
40
41
42
43
44 **Estimation** To estimate the causal distribution from data using the frontdoor criterion we need to
45 estimate each of $p(y|m, a)$, $p(a)$, and $p(m|a)$. In practice, we can fit models $\hat{p}(y|m, a)$ by predicting
46 Y from M and A , and $\hat{p}(m|a)$ by predicting M from A . Then, using the empirical distribution to
47

estimate $p(a)$, the final estimate is:

$$\frac{1}{|A|} \sum_{a'} \sum_m \hat{p}(y|m, a') \hat{p}(m|a), \quad (35.116)$$

where $|A|$ is the number of treatments.

We usually have more modest targets than the full distribution $p(y|\text{do}(a))$. For instance, we may be content with just estimating the average treatment effect. It's straightforward to derive a formula for this using the frontdoor adjustment. Similarly to backdoor adjustment, more advanced estimators of the ATE through frontdoor effect are possible in principle. For example, we might combine fitted models for $\mathbb{E}[Y|m, a]$ and $P(M|a)$. See Fulcher et al. [Ful+20] for an approach to robust estimation via front door adjustment, as well as a generalization of the front door approach to more general settings.

35.9 Further Reading

There is an enormous and growing literature on the intersection of causality and machine learning.

First, there are many textbooks on theoretical and practical elements of causal inference. These include Pearl [Pea09c], focused on causal graphs, Angrist and Pischke [AP08], focused on econometrics, Hernán and Robins [HR20b], with roots in epidemiology, Imbens and Rubin [IR15], with origin in statistics, and Morgan and Winship [MW15], for a social sciences perspective. The introduction to causality in Shalizi [Sha22, §7] is also recommended, particularly the treatment of matching.

Double machine-learning has featured prominently in this chapter. This is a particular instantiation of non-parametric estimation. This topic has substantial theoretical and practical importance in modern causal inference. The double machine learning work includes estimators for many commonly encountered scenarios [Che+17e; Che+17d]. Good references for a lucid explanation of how and why non-parametric estimation works include [Ken16; Ken17; FK21]. Usually, the key guarantees of non-parametric estimator are asymptotic. Generally, there are many estimators that share optimal asymptotic guarantees (e.g. the AIPTW estimator given in Equation (35.31)). Although these are asymptotically equivalent, in finite samples their behavior can be very different. There are estimators that preserve asymptotic guarantees but aim to improve performance in practical finite sample regimes [e.g., vR11].

There is also considerable interest in the estimation of heterogeneous treatment effects. The question here is: what effect would this treatment have when applied to a unit with such-and-such specific characteristics? E.g., what is the effect of this drug on women over the age of 50? The causal identification arguments used here are more-or-less the same as for the estimation of average case effects. However, the estimation problems can be substantially more involved. Some reading includes [Kün+19; NW20; Ken20; Yad+21].

There are several commonly applicable causal identification and estimation strategies beyond the ones we've covered in this chapter. **Regression discontinuity designs** rely on the presence of some sharp, arbitrary non-linearity in treatment assignment. For example, eligibility for some aid programs is determined by whether an individual has income below or above a fixed amount. The effect of the treatment can be studied by comparing units just below and just above this threshold. **Synthetic controls** are a class of methods that try to study the effect of a treatment on a given unit by constructing a synthetic version of that unit that acts as a control. For example, to study the

1
 2 effect of legislation banning smoking indoors in California, we can construct a synthetic California
 3 as a weighted average of other states, with weights chosen to balance demographic characteristics.
 4 Then, we can compare the observed outcome of California with the outcome of the synthetic control,
 5 constructed as the weighted average of the outcomes of the donor states. See Angrist and Pischke
 6 [AP08] for a textbook treatment of both strategies. Closely related are methods that use time series
 7 modeling to create synthetic outcomes. For example, to study the effect of an advertising campaign
 8 beginning at time T on product sales Y_t , we might build a time series model for Y_t using data in the
 9 $t < T$ period, and then use this model to predict the values of $(\hat{Y}_t)_{t>T}$ we would have seen had the
 10 campaign not been run. We can estimate the causal effect by comparing the factual, realized Y_t to
 11 the predicted, counterfactual, \hat{Y}_t . See Brodersen et al. [Bro+15] for an instantiation of this idea.
 12 In this chapter, our focus has been on using machine learning tools to estimate causal effects.
 13 There is also a growing interest in using the ideas of causality to improve machine learning tools.
 14 This is mainly aimed at building predictors that are robust against when deployed in new domains
 15 [SS18c; SCS19; Arj+20; Mei18b; PBM16a; RC+18; Zha+13a; Sch+12b; Vei+21] or that do not rely
 16 on particular ‘spurious’ correlations in the training data [RPH21; Wu+21; Gar+19; Mit+20; WZ19;
 17 KCC20; KHL20; TAH20; Vei+21].