

# Applied Causal Inference Powered by ML and AI

Victor Chernozhukov\*

Christian Hansen<sup>†</sup>

Nathan Kallus<sup>‡</sup>

Martin Spindler<sup>§</sup>

Vasilis Syrgkanis<sup>¶</sup>

February 28, 2024

Publisher: Online

Version 0.1.1

\* MIT

<sup>†</sup> Chicago Booth

<sup>‡</sup> Cornell University

<sup>§</sup> Hamburg University

<sup>¶</sup> Stanford University

# Causal Inference via Randomized Experiments

# 2

"Let us divide them in halves, let us cast lots, that one half of them may fall to my share, and the other to yours; I will cure them without bloodletting and sensible evacuation; but do you do as ye know [...] we shall see how many Funerals both of us shall have."

– Jan Baptist van Helmont [1].

In this chapter we begin discussion of causal inference by focusing on Randomized Control Trials (RCTs). In a randomized control trial, units are randomly divided into those that receive a treatment and those that receive no treatment. Under randomization and other assumptions, the difference in average outcomes between the treated and untreated groups is an average treatment (causal) effect (ATE). By considering pre-treatment covariates, we can improve the precision of the ATE estimate, explore heterogeneity across subgroups, or both. We describe methods for doing so and apply them to several RCTs. We introduce causal diagrams as a means of visualizing RCTs and their underlying causal assumptions. We conclude by outlining some limitations of RCTs.

2.1 Potential Outcomes Framework and Average Treatment Effects . . . . .	42
Random Assignment/Randomized Controlled Trials	46
Statistical Inference with Two Sample Means . . . . .	47
Pfizer/BioNTech Covid Vaccine RCT . . . . .	48
2.2 Pre-treatment Covariates and Heterogeneity . . . . .	50
Regression and Statistical Inference for ATEs . . . . .	52
Classical Additive Approach: Improving Precision Under Linearity . . . . .	52
The Interactive Approach: Always Improves Precision and Discovers Heterogeneity .	55
Reemployment Bonus RCT . . . . .	56
2.3 Drawing RCTs via Causal Diagrams . . . . .	57
2.4 The Limitations of RCTs	58
Externalities, Stability, and Equilibrium Effects . . . . .	58
Ethical, Practical, and Generalizability Concerns . . . . .	59
2.A Approximate Distribution of the Two Sample Means . . . . .	62
2.B Statistical Properties of the Classical Additive Approach* . . . . .	63
2.C Statistical Properties of the Interactive Regression Approach* . . . . .	64

## 2.1 Potential Outcomes Framework and Average Treatment Effects

In this section, we discuss the potential outcomes framework for analyzing causality and treatment effects. It offers an elegant way to formalize counterfactuals as a mathematical concept.

We begin by introducing the two *latent* (unobserved) variables

$$Y(1) \text{ and } Y(0).$$

They represent the potential or counterfactual random outcomes for an observational unit when the unit is subject to treatment (treatment state  $d = 1$ ) or no treatment (control or untreated state  $d = 0$ ) [2]. In an economic context, the treatment might be a training program or a policy intervention, and the outcome might be an individual's wage or employment status. In what follows, it is also useful to introduce the potential response or structural function:

$$d \mapsto Y(d),$$

which maps the potential treatment state  $d \in \{0, 1\}$  to the random potential outcome  $Y(d)$ .

In this formulation, we have dependence of the potential outcome  $Y(d)(\omega)$  on the underlying state of the world  $\omega$ . In our formalization,  $\omega$  will represent randomness across observational units and from any other sources.<sup>1</sup>

The quantities  $Y(1)$  and  $Y(0)$  are "counterfactual" because they can't be simultaneously observed. That is, we generally do not have identical replicas of the observational units that are simultaneously subject to both treatment and control. [3] calls the inability to observe an individual simultaneously under treatment and control "the fundamental problem of causal inference". The inability to observe each individual's treatment and control outcome means that causal inference shares many features with "missing data" problems, see, e.g. [4].

The individual treatment effect is

$$Y(1) - Y(0).$$

This effect will vary across individuals as well as with other sources of randomness encoded in  $\omega$ . As mentioned above, only one of the two terms is actually observed, and hence it is generally infeasible to uncover the individual treatment effect.<sup>2</sup> However, we can hope to estimate averages and the distribution

For simplicity, we do not consider multivalued or continuous treatments.

1: Recall that a random variable  $V$  is a mapping  $\omega \mapsto V(\omega)$  from the underlying state of the world  $\omega \in \Omega$  to the real line (or other metric space) such that we can assign a probability law to it.

2: As an example, we could uncover individual treatment effects if we had identical twins that could be put in treatment and control groups, and we believed that the only difference in outcomes between these twins is induced by treatment – that is,  $\omega$  only depends on genetic makeup. Such an example seems unrealistic at best.

of  $Y(d)$  at the population level to compute quantities such as the average treatment effect (ATE):

$$\delta = E[Y(1) - Y(0)] = E[Y(1)] - E[Y(0)].$$

Let  $D$  denote the actual *assigned treatment*, a random variable, which takes a value of 1 if the observational unit participated in the treatment and 0 otherwise.

**Assumption 2.1.1** (Consistency) *We observe*

$$Y := Y(D).$$

For example, if treatment ( $D = 1$ ) corresponded to completion of a job training program and control ( $D = 0$ ) corresponded to not completing the program, Assumption 2.1.1 says that the observed wage outcome is equal to  $Y(1)$  for a given person if she has completed the program (has  $D = 1$ ) and is equal to  $Y(0)$  if this person has not completed the training program (has  $D = 0$ ). Assumption 2.1.1 seems almost tautological, but it importantly rules out hidden variation in treatment. That is, it requires that the treatment and control states are well-defined and clearly aligned with the observed treatment status,  $D$ .

**Assumption 2.1.2** (No Interference) *Potential outcomes for any observational unit depend only on the treatment status of that unit and not on the treatment status of any other unit.*

Assumption 2.1.2 has implicitly been captured in our definition of potential outcomes,  $Y(d)$ , which give the outcome of each unit *when the unit* is subject to treatment state  $d$ . This formulation rules out scenarios where the treatment given to one unit may impact the outcome of a different unit. Such spillovers could occur, for example, on social networks where treating an individual could impact all of that individual's friends. Some forms of spillovers are readily accommodated by expanding the definition of treatment and correspondingly adjusting definition of potential outcomes,<sup>3</sup> but treating these extensions is beyond the scope of this book.<sup>4</sup>

Assumptions 2.1.1 and 2.1.2 encapsulate what is often referred to as the Stable Unit-Treatment Value Assumption (SUTVA); see, e.g. Imbens and Rubin [10].

The following analytical example may help gain better understanding of the potential outcomes framework.

3: For example, consider a case where each individual has two friends. We could define potential outcomes allowing for spillovers as  $Y(d_0, d_1, d_2)$  where  $d_0$  denotes the treatment state of an individual,  $d_1$  denotes the treatment state of the individual's friend 1, and  $d_2$  denotes the treatment state of the individual's friend 2.

4: For further reading we refer, among many others, to [5], [6], [7], [8] and [9].

**Example 2.1.1** [Analytical Example] Consider the following model

$$\begin{aligned} Y(1) &:= \theta_1 + \epsilon_1 \\ Y(0) &:= \theta_0 + \epsilon_0 \\ D &:= 1(\nu > 0), \\ Y &:= Y(D), \end{aligned}$$

where  $\theta_0$  and  $\theta_1$  are constants, and  $(\epsilon_0, \epsilon_1, \nu)$  are jointly normal random stochastic disturbances with mean 0 and covariance matrix  $\Sigma$ . Here,  $\nu$  represents factors that influence selection into the treatment state. In this example  $E[Y(1)] = \theta_1$ ,  $E[Y(0)] = \theta_0$ , and the ATE is  $\delta = \theta_1 - \theta_0$ . Importantly, only  $D$  and  $Y$  are observed.

Under Assumption 2.1.1, population data directly provide the conditional averages

$$E[Y | D = d] = E[Y(d) | D = d], \text{ for } d \in \{0, 1\}.$$

The difference of the two averages gives us the average predictive effect (APE) of treatment status on the outcome:

$$\pi = E[Y | D = 1] - E[Y | D = 0].$$

It measures the association of the treatment status with the outcome.

While the APE is identified – meaning computable from the population data – it may seem surprising (or not at all) that the APE in general does not agree with the ATE  $\delta$ :

$$\delta \neq \pi. \tag{2.1.1}$$

The difference between the APE and ATE is generally said to be due to *selection bias*. The meaning of selection bias is clarified through the following example, and clarified theoretically below.

**Example 2.1.2** (Selection Bias in Observational Data) Suppose we want to study the impact of smoking marijuana on life longevity. Suppose that smoking marijuana has no causal effect on life longevity:

$$Y = Y(0) = Y(1),$$

so that

$$\delta = E[Y(1)] - E[Y(0)] = 0.$$

However, the observed smoking behavior,  $D$ , is not assigned in an experimental study. Suppose that the behavior determining  $D$  is associated with poor health choices such as drinking alcohol, which are known to cause shorter life expectancy, so that  $E[Y | D = 1] < E[Y | D = 0]$ . In this case, we have negative a predictive effect:

$$\pi = E[Y | D = 1] - E[Y | D = 0] < 0 = \delta,$$

which differs from the true causal effect  $\delta = 0$ .

To sum up, in the smoking example, the chosen "treatment" variable  $D$  is potentially negatively associated with the potential health outcome, inducing the selection bias – the difference between the predictive effect and the causal effect.

**Example 2.1.3** (Analytical Version of the Smoking Example)

To capture dependence between  $Y(d)$  and  $v$  in the smoking context analytically, we can go back to Example 2.1.1, and make variables  $\epsilon_d$  and  $v$  be negatively associated:

$$E[\epsilon_d v] < 0.$$

The negative association between the  $\epsilon_d$  and  $v$  then results in the observed smoking status,  $D$ , being negatively associated with the potential outcomes  $Y(d)$ . Specifically, we have

$$E[Y|D = 1] < E[Y|D = 0],$$

which can be verified through additional analytical calculations or via simulation experiments (a homework).

It is useful to emphasize the main reason for having selection bias is that

$$E[Y(d)|D = 1] \neq E[Y(d)]$$

whenever  $D$  is not independent of  $Y(d)$ . If  $D$  and  $Y(d)$  were independent,

$$E[Y(d)|D = 1] = E[Y(d)]$$

would hold since in this case  $D$  is uninformative about the potential outcome and drops out from the conditional expectation.

To sum up, the problem with observational studies like our contrived Example 2.1.2 is that the "treatment" variable  $D$

is determined by individual behaviors which may be linked to potential outcomes. This linkage generates selection bias - the disagreement between APE and ATE. There are many ways of addressing selection bias, one of which is through an experiment, where we randomly assign the treatment to the units.

## Random Assignment/Randomized Controlled Trials

A way to clearly remove selection bias is through random assignment of treatment.

**Assumption 2.1.3** (Random Assignment/Exogeneity) *Suppose that treatment status is randomly assigned. Namely,  $D$  is statistically independent of each potential outcome  $Y(d)$  for  $d \in \{0, 1\}$ , which is denoted as*

$$D \perp\!\!\!\perp Y(d)$$

and  $0 < P(D = 1) < 1$ .

This assumption states that the treatment assignment mechanism is purely random, and ensures that there are units in treatment and in control.

**Example 2.1.4** (Analytical Example Continued) In the analytical example 2.1.1, Assumption 2.1.3 is satisfied if the stochastic shock  $v$  determining  $D$  is independent of stochastic shocks  $\epsilon_0$  and  $\epsilon_1$  determining  $Y(1)$  and  $Y(0)$ , i.e.

$$v \perp\!\!\!\perp (\epsilon_0, \epsilon_1).$$

A key result is that selection bias is removed under Assumption 2.1.3 which allows us to learn summaries of causal effects.

**Theorem 2.1.1** (Randomization Removes Selection Bias) *Under Assumption 2.1.3, the average outcome in treatment group  $d$  recovers the average potential outcome under the treatment status  $d$ :*

$$E[Y \mid D = d] = E[Y(d) \mid D = d] = E[Y(d)],$$

for each  $d \in \{0, 1\}$ . Hence the average predictive effect and average treatment effect coincide:

$$\begin{aligned} \pi &:= E[Y \mid D = 1] - E[Y \mid D = 0] \\ &= E[Y(1)] - E[Y(0)] =: \delta. \end{aligned}$$

Assumption 2.1.3 is often not plausible for observational data. In a *randomized controlled trial* (RCT)<sup>5</sup>, the aim is to ensure the plausibility of Assumption 2.1.3 by direct random assignment of treatment  $D$ . That is, subjects are randomly assigned a treatment state  $D$  by the experimenter without regard to any of their characteristics. Because the random assignment of the treatment is unrelated to all subject characteristics by construction, well-executed RCTs guarantee that Assumption 2.1.3 is satisfied. Because of this property, many consider RCTs as the gold standard in causal inference, and RCTs are routinely employed in a variety of important settings.<sup>6</sup> Examples include evaluating the efficacy of medical treatment, vaccinations, training programs, marketing campaigns, and other kinds of interventions.

**Example 2.1.5** (No Selection Bias in Experimental Data) Suppose that in the smoking example (Example 2.1.2), we worked with data where smoking or non-smoking was generated by perfectly enforced random assignment. In this case, we would have agreement between average predictive and treatment effects:  $\pi = \delta$ . While it is difficult to imagine a long-run RCT where study participants could be forced to smoke or not smoke marijuana (we discuss such limitations as well as ethical considerations in Section 2.4), RCTs are routinely employed in a variety of other important settings.

5: Synonyms are experiments and A/B tests.

6: Of course, RCTs must be correctly done to guarantee Assumption 2.1.3. For example, RCTs where experimental protocols are not followed continue to suffer from selection bias. There are also examples, *quasi-experiments*, where we may believe that Assumption 2.1.3 is plausible that do not correspond to explicit designed experiments.

## Statistical Inference with Two Sample Means

Inference is based on the independent sample  $\{(Y_i, D_i)\}_{i=1}^n$  obtained from an RCT, where index  $i$  denotes the observational unit. We assume that each  $(Y_i, D_i)$  has the same distribution as  $(Y, D)$ . Estimation of the two means  $\theta_d = E[Y | D = d]$  for  $d = 0$  and  $d = 1$  can be done by considering two group means

$$\hat{\theta}_d = \frac{E_n[Y1(D = d)]}{E_n[1(D = d)]}.$$

The two means example can also be treated as a special case of linear regression,<sup>7</sup> but we find it instructive to work out the details directly for the two group means. We provide these details in Section 2.A.

Under mild regularity conditions, we have that

$$\sqrt{n} \begin{pmatrix} \hat{\theta}_0 - \theta_0 \\ \hat{\theta}_1 - \theta_1 \end{pmatrix} \overset{a}{\sim} N(0, V),$$

7: Indeed, we can regress  $Y$  on  $D$  and  $1 - D$ ; that is, estimate the model  $Y = \theta_1 D + \theta_0(1 - D) + U$ . We can then apply the inferential machinery developed in the previous chapter.



where

$$V = \begin{pmatrix} \frac{\text{Var}(Y|1(D=0))}{P(D=0)} & 0 \\ 0 & \frac{\text{Var}(Y|1(D=1))}{P(D=1)} \end{pmatrix}$$

so that  $\hat{\delta} = \hat{\theta}_1 - \hat{\theta}_0$  obeys

$$\sqrt{n}(\hat{\delta} - \delta) \stackrel{a}{\sim} N(0, V_{11} + V_{22}).$$

To use this result in practice, variance components are usually estimated using the *plug-in principle*, which amounts to using the sample analogues of the expressions above.

Sometimes we are interested in relative effectiveness of treatment effects (for example, vaccine efficiency):

$$f(\theta) = (\theta_1 - \theta_0)/\theta_0 = \delta/\theta_0.$$

Relative effectiveness can be estimated by  $\hat{\delta}/\hat{\theta}_0 = f(\hat{\theta})$ , where  $\hat{\theta} = \{\hat{\theta}_d\}_{d \in \{0,1\}}$  and  $\theta = \{\theta_d\}_{d \in \{0,1\}}$ , with approximate distribution obtained using the *delta method*:

$$\sqrt{n}(f(\hat{\theta}) - f(\theta)) \approx G' \sqrt{n}(\hat{\theta} - \theta) \stackrel{a}{\sim} N(0, G'VG),$$

where  $G = \nabla f(\theta)$ ,  $\hat{\theta} = (\hat{\theta}_0, \hat{\theta}_1)'$ ,  $\theta = (\theta_0, \theta_1)$ .<sup>8</sup>

## Pfizer/BioNTech Covid Vaccine RCT

Pfizer/BNTX was the first vaccine approved for emergency use in the EU and US to reduce the risk of Covid-19 disease. See the Food and Drug Administration (FDA) [briefing](#) for details about the RCT and the summary data. Volunteers were randomly assigned to receive either a treatment (2-dose vaccination) or a placebo, without knowing which they received, and the doctors making the diagnoses did not know whether a given volunteer received a vaccination or not. In other words, the trial was a double-blind randomized control trial. The results of the study are presented in the following table.

8: The approximation follows from application of the first order Taylor expansion and continuity of the derivative  $\nabla f$  at  $\theta$ .



**Figure 2.1:** Tozinameran (Pfizer-BioNTech Covid-19 vaccine); Image Source: Wikipedia / Arne Müsseler

[Vaccination RCT R Notebook](#) and [Vaccination RCT Python Notebook](#) contain the analysis of the Pfizer-BioNTech Covid-19 Vaccine RCTs.

Efficacy Endpoint Subgroup	BNT162b2 N <sup>a</sup> =19965 Cases n1 <sup>b</sup> Surveillance Time <sup>c</sup> (n2 <sup>d</sup> )	Placebo N <sup>a</sup> =20172 Cases n1 <sup>b</sup> Surveillance Time <sup>c</sup> (n2 <sup>d</sup> )	Vaccine Efficacy % (95% CI) <sup>e</sup>
Overall	9 2.332 (18559)	169 2.345 (18708)	94.6 (89.6, 97.6)
<b>Age group (years)</b>			
16 to 17	0 0.003 (58)	1 0.003 (61)	100.0 (-3969.9, 100.0)
18 to 64	8 1.799 (14443)	149 1.811 (14566)	94.6 (89.1, 97.7)
65 to 74	1 0.424 (3239)	14 0.423 (3255)	92.9 (53.2, 99.8)
≥75	0 0.106 (805)	5 0.109 (812)	100.0 (-12.1, 100.0)

**Figure 2.2:** The aggregate data from the Pfizer RCT; source: FDA [briefing](#).

We see that the rate of Covid-19 infection was relatively low at the time. Specifically, the treatment group saw 9 Covid-19 cases per 19,965, while the control group saw 169 cases per 20,172.

The estimated average treatment effect is about

$$-792.7 \text{ cases per } 100,000,$$

and the 95% confidence band is<sup>9</sup>

$$[-922, -664].$$

Under Assumptions 2.1.3 and 2.2.1 the confidence band suggests that the Covid-19 vaccine caused a reduction in the risk of contracting Covid-19.

We also compute the Vaccine Efficacy metric, which according to [11], refers to the following measure:

$$VE = \frac{\text{Risk for Unvaccinated} - \text{Risk for Vaccinated}}{\text{Risk for Unvaccinated}}.$$

It describes the relative reduction in risk caused by vaccination. Estimating the VE is simple as we can plug-in the estimated group means. We can compute standard errors using the delta method or by simulation. We obtain that the overall vaccine efficacy is 94.6%, replicating the results shown in Figure 2.2. Our 95% confidence interval for VE, based on the normal approximation, is

$$[90.9\%, 98.2\%],$$

which differs only slightly from the FDA briefing table.<sup>10</sup>

**Remark 2.1.1** We notice that the confidence intervals for the VE for the two age groups of seniors are very wide, so to increase precision we pool them together and calculate the effectiveness of the vaccine for the two groups that are 65 or older. The resulting VE estimate is 95% and the two-sided

9: In this example, we don't need the underlying individual data to evaluate the effectiveness of the vaccine because the potential outcomes are Bernoulli random variables with mean  $E[Y(d)]$  and variance  $\text{Var}(Y(d)) = E[Y(d)(1 - EY(d))]$ .

10: The analysis in the FDA table is based on the inversion of exact binomial tests, the Cornfield procedure.

confidence interval based on the normal approximation is

$$[82\%, 106\%]$$

A more refined approach is possible, based on the inversion of exact binomial ratio Cornfield tests [12], which we report in [Vaccination RCT R Notebook](#) and [Vaccination RCT Python Notebook](#). This approach, using [Vaccination RCT R Notebook](#), yields a confidence interval of

$$[69\%, 99\%].$$

The reason is that the accumulated counts of binomials are too few for the Gaussian approximations to provide a high-quality approximation, so the exact binomial ratio test inversion delivers a more accurate confidence interval.

## 2.2 Pre-treatment Covariates and Heterogeneity

Sometimes we also have additional *pre-treatment* or *pre-determined* covariates  $W$ . We might be interested in either using these covariates to estimate average effects more precisely or to describe heterogeneity of the treatment effects. For example, we might be interested in the impact of a treatment across age or income groups.

For this purpose, we consider conditional average treatment effects (CATE):

$$\delta(W) = E[Y(1) | W] - E[Y(0) | W],$$

which compare the average potential outcomes conditional on a set of covariates  $W$ .

We can directly learn the conditional predictive effects (CAPE),

$$\pi(W) = E[Y | D = 1, W] - E[Y | D = 0, W],$$

from population data. However, these CAPE will generally not agree with the CATE. One assumption that will be sufficient for the CAPE and CATE to agree is having treatment assigned randomly and independently of covariates. As before, the use of RCTs help ensure the plausibility of this assumption.

**Assumption 2.2.1** (Random Assignment Independent of Co-

variates) Suppose that treatment status is randomly assigned. Namely,  $D$  is statistically independent of both the potential outcomes and a set of pre-determined covariates:

$$D \perp\!\!\!\perp (Y(0), Y(1), W),$$

and  $0 < P(D = 1) < 1$ .

This assumption spells out that, if we plan to use covariates in the analysis, randomization has to be made with respect to these covariates as well. In practice, it is often tempting to use post-treatment covariates, but the use of such variables runs the danger of violating Assumption 2.2.1. In the extreme case, conditioning on the post-treatment observed outcome  $Y$ , we find that  $\pi(Y) = 0$ , even when there is a treatment effect. In a less extreme case, conditioning on post-treatment variables related to the outcome can "control-away" part of the effect, diminishing estimates.

A common scenario where accidentally using a post-treatment covariate may occur is when researchers encounter missing data from imperfect data collection in following-up with control and treated units to collect demographic information. When we drop observations with missing data, we implicitly condition on a post-treatment variable (missingness) which can cause violations of Assumption 2.2.1.

The desire to assess randomization with respect to covariates motivates the following diagnostic procedure.

**Testing Covariance Balance.** The random assignment assumption induces covariate balance. Namely, the distribution of covariates should be the same under both treatment and control:

$$W|D = 1 \sim W|D = 0,$$

and, equivalently,

$$D|W \sim D.$$

A useful implication is that  $D$  is not predictable by  $W$ :

$$E[D | W] = E[D].$$

This latter condition is testable using regression tools. It amounts to saying that the  $R^2$  of a regression of  $D$  on  $W$  is 0.

For random variables  $A$  and  $B$ ,  $A \sim B$  denotes that  $A$  and  $B$  have the same distribution.

Under Assumption 2.2.1, Theorem 2.1.1 continues to hold, but we now have a stronger result.

**Theorem 2.2.1** (Randomization with Covariates) *Under Assumption 2.2.1, the expected value of  $Y$  conditional on treatment status  $D = d$  and covariates  $W$  coincides with the expected value of potential outcome  $Y(d)$  conditional on covariates  $W$ :*

$$E[Y \mid D = d, W] = E[Y(d) \mid D = d, W] = E[Y(d) \mid W],$$

for each  $d$ . Hence the conditional predictive and average treatment effects agree:

$$\pi(W) = \delta(W).$$

## Regression and Statistical Inference for ATEs

Empirical researchers often base statistical inference on the ATE using the classical additive linear regression model, where covariates enter additively in the model. This approach has some good practical properties and often empirically leads to improvements in precision over the simple two-means approach, though this precision improvement is not guaranteed. Another approach that we will emphasize is the interactive regression approach, where de-meaned covariates are also interacted with the base treatment. Including interactions of de-meaned covariates with the treatment always improves precision, and it also allows us to discover treatment effect heterogeneity.

### Classical Additive Approach: Improving Precision Under Linearity

We begin explaining the classical additive approach. Here, to simplify the exposition, we make the strong assumption that the conditional expectation function is exactly linear:

$$E[Y \mid D, W] = D\alpha + \beta'X, \quad (2.2.1)$$

where  $X = (1, W)$  contains an intercept and pre-treatment covariates  $W$ . This setup is clearly restrictive, but the statistical inference result will be valid without this assumption.<sup>11</sup> Later in the book, we will consider fully nonlinear models.

We assume that covariates are centered.<sup>12</sup>

$$E[W] = 0.$$

By Assumption 2.2.1, there is covariate balance:

$$E[W \mid D = 1] = E[W \mid D = 0].$$

11: See Section 2.B for details.

12: Theoretically, this is implemented by redefining  $W := W - E[W]$ . In estimation, this is implemented by redefining  $W_i := W_i - \mathbb{E}_n[W]$ .

Using centered covariates implies that

$$E[Y(0)] = E[E[Y | D = 0, X]] = \beta_1$$

$$E[Y(1)] = E[E[Y | D = 1, X]] = \beta_1 + \alpha.$$

That is, the average outcome in the untreated state is  $\beta_1$ , and the average treatment effect  $\delta = E[Y(1)] - E[Y(0)]$  equals  $\alpha$ .

Equation (2.2.1) implies that

$$Y = D\alpha + \beta'X + \epsilon, \quad \epsilon \perp (D, X), \quad (2.2.2)$$

implying that  $\alpha$  coincides with the coefficient in the BLP of  $Y$  on  $D$  and  $X$ . In fact, even if we don't assume the model (2.2.1), we still have that  $\alpha = \delta$ . That is, the projection coefficient  $\alpha$  recovers the ATE  $\delta$  without the linearity assumption as we detail in Section 2.B. Furthermore the statistical inference result stated below will hold without requiring linear conditional expectation functions as it is simply a statement about inference on the BLP.

We are interested in statistical inference on the ATE and Relative ATE<sup>13</sup>

$$\alpha \quad \text{and} \quad \alpha/\beta_1.$$

13: Relative ATE is often called *lift* in business applications.

Under regularity conditions, application of the OLS theory from Chapter 1 gives us

$$\begin{pmatrix} \sqrt{n}(\hat{\alpha} - \alpha) \\ \sqrt{n}(\hat{\beta}_1 - \beta_1) \end{pmatrix} \stackrel{a}{\sim} N(0, \mathbf{V}),$$

where covariance matrix  $\mathbf{V}$  has components:

$$V_{11} = \frac{E[\epsilon^2 \tilde{D}^2]}{(E[\tilde{D}^2])^2}, \quad V_{22} = \frac{E[\epsilon^2 \tilde{1}^2]}{(E[\tilde{1}^2])^2}, \quad V_{12} = V_{21} = \frac{E[\epsilon^2 \tilde{D} \tilde{1}]}{E[\tilde{1}^2]E[\tilde{D}^2]},$$

where  $\tilde{D} = D - E[D]$  is the residual after partialling out  $X$  from  $D$  linearly and  $\tilde{1} := (1 - D)$  is the residual after partialling out  $D$  and  $W$  from 1.

We also obtain the approximate normality for the Relative ATE using the delta method:

$$\sqrt{n}(\hat{\alpha}/\hat{\beta}_1 - \alpha/\beta_1) \stackrel{a}{\sim} N(0, G'VG),$$

where

$$G = [1/\beta_1, -\alpha/\beta_1^2]'$$

## Improvement in Precision under Linearity

Now we explain the role of covariates in potentially delivering improvements in precision of estimating the ATE. The underlying idea is that of "denoising." This improvement, however, hinges on the linear model (2.2.1). In the next section, we will obtain improvement without linearity assumptions.

We consider what happens when we do not include covariates in the regression. In this case, the OLS estimator  $\bar{\alpha}$  estimates the projection coefficient  $\alpha$  in the BLP using  $(1, D)$  alone:<sup>14</sup>

$$Y = \alpha D + \beta_1 + U, \quad E[U] = E[UD] = 0,$$

where the noise

$$U = \beta'(X - E[X]) + \epsilon$$

contains the part of  $Y$  that is linearly predicted by  $X$ ,  $\beta'(X - E[X]) = \beta'X - \beta_1$ . We then have that  $\bar{\alpha}$  obeys

$$\sqrt{n}(\bar{\alpha} - \alpha) \stackrel{a}{\sim} N(0, \bar{V}_{11}), \quad \bar{V}_{11} = \frac{E[U^2 \tilde{D}^2]}{(E[\tilde{D}^2])^2}.$$

Under the linear model (2.2.1), it follows that

$$V_{11} \leq \bar{V}_{11},$$

with the inequality being strict (" $<$ ") if  $\text{Var}(\beta'X) > 0$ .<sup>15</sup> That is, under (2.2.1), using pre-determined covariates improves the precision of estimating the ATE  $\alpha$ .

However, this improvement theoretically hinges on the correctness of the additive linear model. Statistical inference on the ATE based on the the normal approximation provided above remains valid without this assumption as long as robust standard errors are used.<sup>16</sup> However, the precision can be either higher or lower than that of the classical two-sample approach without covariates. That is, without (2.2.1),  $V_{11}$  and  $\bar{V}_{11}$  are not generally comparable.

**Remark 2.2.1** While the inferential result we derived is robust with respect to the linearity assumption on the CEF, the improvement in precision itself is **not** guaranteed in general and hinges on the validity of the linearity assumption. We provide simulation examples where controlling for pre-determined covariates linearly lowers the precision (increases robust standard errors) in [Covariates in RCT R Notebook](#) and

14: Here  $U = Y - \alpha D - \beta_1$  obeys

$$\begin{aligned} E[U | D = d] &= E[Y(d) - \alpha d - \beta_1 | D = d] \\ &= E[Y(d) - \alpha d - \beta_1] = 0, \end{aligned}$$

invoking random assignment and the definition of  $\alpha$  and  $\beta_1$ .

15: Verify this as a reading exercise.

16: We always use robust variance formulas throughout the book. However, the default inferential algorithms in R and Python often report the classical Student's formulas as variances, which critically rely on the linearity assumption.

Covariates in RCT Python Notebook.

## The Interactive Approach: Always Improves Precision and Discovers Heterogeneity

We can also consider estimation of CATE through the lens of an interactive linear regression model, which interacts treatment indicator  $D$  with regressors  $X$  constructed from original raw regressors  $W$ . Including these interactions respects the logic of approximating the conditional expectation of  $Y$  given  $D$  and raw regressors using linear functional forms. To simplify exposition, we first assume that the interactive model is exactly correct for the CEF:

$$E[Y | D, W] = \alpha'XD + \beta'X. \quad (2.2.3)$$

In Section 2.C, we explain how this approach works without this assumption.

As before, we assume

$$X = (1, W)', \quad E[W] = 0,$$

which can be achieved in practice by recentering. Here, we recover CATE via

$$\begin{aligned} \delta(W) &= E[Y(1) | W] - E[Y(0) | W] \\ &= E[Y | D = 1, W] - E[Y | D = 0, W] = \alpha'X. \end{aligned}$$

Using that  $E[W] = 0$ , the ATE is then

$$\delta = E[\delta(W)] = E[\alpha'X] = \alpha_1,$$

where  $\alpha_1$  is the first component of  $\alpha$ . The function  $\alpha_2'W$ , where  $\alpha_2$  is the vector all elements of  $\alpha$  excluding  $\alpha_1$ , therefore describes the deviation of CATE away from the ATE.

We can verify that  $\alpha$  is the coefficient of the linear projection equation:

$$Y = \alpha'DX + \beta'X + \epsilon, \quad \epsilon \perp (X, DX).$$

Therefore, we can treat

$$\bar{D} := DX$$

as a vector of technical treatments<sup>17</sup> and invoke the "partialling

Covariates in RCT R Notebook and Covariates in RCT Python Notebook explore the use of covariates to both improve precision and learn about heterogeneity via a simulation experiment.

<sup>17</sup>: A technical treatment refers to any variable obtained as a transformation of the original treatment variable.



out" approach for inference on components of  $\alpha$ . The variance formulas are given in Section 2.C.

**Remark 2.2.2** (Improvement in Precision Guarantee) Unlike the previous approach, the "interactive" approach always delivers improvements in precision for estimating  $\delta$ , even if the linearity in (2.2.3) does not hold; this was demonstrated by Lin [13]. Section 2.C explains this point in detail and provides a deeper dive into the properties of the interactive approach without assuming correct linear specification of the CEF.

## Reemployment Bonus RCT

Here we re-analyze the Pennsylvania re-employment bonus experiment [14], which was conducted in the 1980s by the U.S. Department of Labor to test the incentive effects of alternative compensation schemes for unemployment insurance (UI). In these experiments, UI claimants were randomly assigned either to a control group or one of five treatment groups. We focus our discussion on treatment group 4. In the control group the current rules of the UI applied. Individuals in the treatment groups were offered a cash bonus if they found a job within some pre-specified period of time (qualification period), provided that the job was retained for a specified duration; see the [Penn Data Codebook](#) for further details on the data.

We consider the

- ▶ classical 2-sample approach, no adjustment (CL)
- ▶ classical linear regression adjustment (CRA)
- ▶ interactive regression adjustment (IRA)
- ▶ interactive regression adjustment with double lasso (partially out by lasso) (IRA-DL)

We use the last approach in the spirit of exploration and experimentation. We describe the last approach and establish its validity in Chapter 4.

Estimates of the ATE on (log) unemployment duration and corresponding estimated standard errors are given in Table 2.1.

	CL	CRA	IRA	IRA-DL
Estimate	-0.0855	-0.0797	-0.0755	-0.0789
Std. Error	0.0359	0.0356	0.0356	0.0356

[Reemployment Bonus RCT R Notebook](#) and [Reemployment Bonus RCT Python Notebook](#) explore the use of covariates to improve precision and learn about heterogeneity in a Reemployment Bonus RCT.

**Table 2.1:** Estimates of the ATE of the reemployment bonus on log unemployment duration..

The different estimators deliver fairly similar point estimates suggesting that treatment group 4 experiences an average decrease in unemployment duration of around 8%. The three regression estimators deliver estimates that are slightly more precise (have lower standard errors) than the simple difference in means estimator.

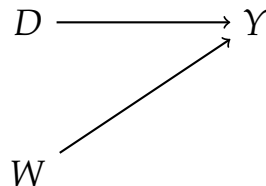
We also see that the regression estimators offer slightly lower estimates of the ATE than the difference in means estimator. These differences likely occur due to minor imbalances in the treatment allocation: People older than 54 tended to receive the treatment more than other groups of qualified UI claimants during the later period of the experiment. Loosely speaking, the regression estimators try to correct for this imbalance by "partialling out" the effect of this oversampling and averaging over differences net of these "imbalancing" effects. We will explain how regression adjustment corrects for imbalances in Chapter 5.

See [Reemployment Bonus RCT R Notebook](#) and [Reemployment Bonus RCT Python Notebook](#) for the results from the balance check.

## 2.3 Drawing RCTs via Causal Diagrams

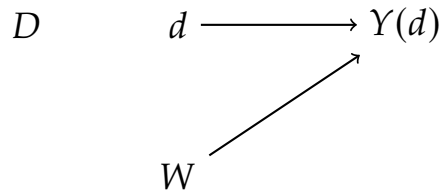
RCTs can be visualized using causal diagrams. These enable us to simply and clearly show the causal assumptions that underpin our model for retrieving treatment effects. Causal diagrams were introduced as early as 1920s by Sewall and Philip Wright ([15],[16]) and emerged as a fully formal tool due to the work of Judea Pearl and James H. Robins ([17], [18]).

In causal diagrams, random variables are denoted by nodes; and arrows between nodes represent causal effects. In our RCT set-up, we have that the assigned treatment variable causes outcome variable  $Y$ , and the pre-treatment variables  $W$  also cause the outcome variable  $Y$ , but they don't cause the treatment assignment  $D$ . This causal diagram is illustrated in Figure 2.3 below.



**Figure 2.3:** Causal Diagram for a RCT

Figure 2.4 depicts a version of the diagram that also includes potential outcomes as nodes.



**Figure 2.4:** A Causal Diagram for the RCT Research Design

In Figure 2.4, we show the potential outcomes  $Y(d)$  as a single node. The pre-treatment covariates affect this node, which is represented by the arrow from  $W$  to the  $Y(d)$  node. The assigned treatment variable  $D$  is independent of the node  $Y(d)$ , which is shown by the absence of an arrow connecting the two nodes. The arrow from  $d$  to  $Y(d)$  shows the causal dependency of  $Y(d)$  on the deterministic node  $d$ . The assigned treatment  $D$  is also shown to be independent of the node  $W$ . The potential outcome process  $d \mapsto Y(d)$  and treatment assignment jointly determine the realized outcome variable  $Y$  via the assignment  $Y := Y(D)$ .

We further develop the use of these concepts and the use of causal diagrams as a formal tool in Chapter 7 and Chapter 8.

## 2.4 The Limitations of RCTs

Here, we briefly outline some of the primary limitations of RCTs. We first consider threats to identification, outlining settings in which the stable unit treatment value assumption (SUTVA), an important assumption that underpins causal inference in an RCT setting, is unlikely to hold, and the implications for inference. We then address ethical and practical concerns in RCT implementation and generalizability.

### Externalities, Stability, and Equilibrium Effects

The traditional formulation of Rubin's causal model relies on SUTVA as described in Section 2.1. Part of SUTVA is the requirement that the potential outcomes of one unit should be unaffected by the assignment of treatments to other units [19]. In the following, we consider some cases where this assumption might not hold.

In a vaccine example, this assumption holds if treatment and control populations are "small" (infinitesimal) subpopulations of the entire general population. Our methods measure the average vaccine effects in these settings. However, if we vaccinate a

sufficiently large percentage of people, reaching herd immunity, the outcomes for the control group would be essentially the same as outcomes for the treated. SUTVA therefore would not hold.<sup>18</sup>

In economics, we refer to such spillover effects as externalities or, in some contexts, as general equilibrium effects. For example, there is a positive externality created by people who take the vaccine (and people that don't take vaccine "free ride," once the vaccination level is high enough). Consider another example. We might want to study the earning effect of getting a college degree versus not having a college degree. If treatment will target a relatively small subpopulation of people, there likely won't be any large general equilibrium wage effects. On the other hand, if the treatment will target a large subpopulation, the equilibrium wage will likely adjust (the college wage premium might decrease, for example). In another example, the outcomes for one individual in large-scale training programs may be affected by the number of people trained to perform the same job.

18: Because SUTVA does not hold in the vaccination context, it is customary to use relative measures of impact like "vaccine efficiency" because they may be a somewhat more stable measure when generalizing from "small" treated subpopulations to a "large" treated population.

## Ethical, Practical, and Generalizability Concerns

Many RCTs are infeasible because implementing them would be unethical. The general ethical principles and guidelines for research involving human subjects are set out in the 1978 Belmont report ([20]). The key ethical principles are "Respect for persons," "Beneficence," and "Justice." Human subject trials are subject to regulation by an institutional review board, which determines whether the trial is ethical with reference to these guiding principles, or whether it should be prevented from registering.

For example, we previously considered a hypothetical RCT where individuals are assigned to a smoking treatment group. The trial would violate the principle of "beneficence" as the researcher might be causing physical harm to study participants by assigning them to smoking. Thus, RCTs are rarely a feasible means of retrieving the causal effects of harmful interventions as they tend to be unethical.

RCTs may also face practical issues. They can be prohibitively expensive when the treatment is costly, data collection costs are high, or the sample size required for adequate power is high. These issues make it difficult to implement long-term RCTs and find evidence on the long-term effects of interventions, particularly because they are more likely to suffer from attrition.

It may also be politically infeasible for policymakers to enforce randomization of receipt of a desirable treatment.

Even in the best case, where an RCT is successfully implemented and we are confident in our retrieved average treatment effect, it may be difficult to generalize (or extrapolate) the result of an RCT in a specific context to a general finding. This difficulty might be because local conditions or implementation capacity materially differ between where interventions are staged or because the scale of the intervention is important.

## Notebooks

- ▶ [Vaccination RCT R Notebook](#) and [Vaccination RCT Python Notebook](#) contain the analysis of vaccination examples.
- ▶ [Covariates in RCT R Notebook](#) and [Covariates in RCT Python Notebook](#) explore the use of covariates to improve precision and learn about heterogeneity via a simulation experiment.
- ▶ [Reemployment Bonus RCT R Notebook](#) and [Reemployment Bonus RCT Python Notebook](#) explore the use of covariates to improve precision and learn about heterogeneity in a Reemployment Bonus RCT.

## Notes

RCTs have a profound influence on business, economics and science more generally. For example, RCTs are routinely used to study the efficacy of drugs and efficacy of various programs in labor and development economics, among other subfields of economics. The FDA moved to RCTs as the gold standard of proving that treatments work in 1970s-80s. In the tech industry and marketing, RCTs are also called "A/B Tests" and are now widely used. Many major tech companies have their own experimental platforms to carry out thousands of experiments.<sup>19</sup>

The expansion of the use of experimentation in economics is associated with the work of Richard Thaler, the recipient of the 2017 Alfred Nobel Memorial Prize in Economics;<sup>20</sup> Abhijit Banerjee, Esther Duflo, and Michael Kremer, the recipients of the 2019 Alfred Nobel Memorial Prize in Economics;<sup>21</sup> and John List, among many others.

19: See, for example, [ExP platform at Microsoft](#) and the [WebLab platform](#) at Amazon.

20: "for his contributions to behavioural economics." Source: [NobelPrize.org](#)

21: "for their experimental approach to alleviating global poverty." Source: [NobelPrize.org](#)

We touched upon very basic ideas here. The basic random design is just one of many possible randomized designs that allow us to uncover causal effects. For an in-depth analysis of design of experiments, please see lecture notes by Art Owen ([21]). For standard RCTs and causal analysis more generally, see the book by Imbens and Rubin [10]. Duflo et al. [22] is another good overview of the use of RCTs with a focus on development economics applications. For real examples of how RCTs are done and designed in practice, see, for example, the FDA registry of RCTs, the American Economic Association for a registry of RCTs in economics, or the [The Poverty Action Lab](#).

## Study Questions

1. Set-up a simulation experiment that illustrates the contrived smoking example, following the analytical example we've presented in the text. Illustrate the difference between estimates obtained via an RCT (smoking generated independently of potential outcomes) and an observational study (smoking choice is correlated with potential outcomes).
2. Sketch out the proof of the large sample properties of the two means estimator.
3. Study the notebook on vaccinations RCTs. Try to replicate the results in the FDA briefing table for each age 18-64 (exact replication is not required). Explain your calculations.
4. Study the notebook on the reemployment example. Experiment with putting even more flexible controls (e.g. use extra interactions of some controls). Report your findings.
5. Work and experiment with the Covariates in RCT notebook. Explain the main points being made.
6. Skim over the information on the Pfizer RCT design [briefing](#). Write down one paragraph summarizing the study design.
7. Skim over one of the RCTs registered with [AEA RCT Registry](#). Write down one paragraph summarizing the study design.

8. Think of some RCTs where stability (SUTVA) is likely to hold and some RCTs where it likely does not.
9. Explain why we can't learn individual treatment effects by first putting a unit in treatment and then putting the individual in control second (or the other way around). A hint is to think of all sources of randomness represented by  $\omega$ . Would the situation be different if you had a time machine?

## 2.A Approximate Distribution of the Two Sample Means

To demonstrate the result in the text, we note that

$$\hat{\theta}_d - \theta_d = \frac{\mathbb{E}_n[(Y(d) - EY(d))1(D = d)]}{\mathbb{E}_n[1(D = d)]}$$

for  $d \in \{0, 1\}$  because we can re-write the population group average as

$$\theta_d = E[Y(d)] = E[Y(d)] \frac{\mathbb{E}_n[1(D = d)]}{\mathbb{E}_n[1(D = d)]}.$$

Hence, for each  $d \in \{0, 1\}$ ,

$$\sqrt{n}(\hat{\theta}_d - \theta_d) = \sqrt{n} \frac{\mathbb{E}_n[(Y(d) - EY(d))1(D = d)]}{\mathbb{E}_n[1(D = d)]}.$$

By the law of large numbers,  $\mathbb{E}_n[1(D = d)] \approx P(D = d)$ ; so we have the approximation

$$\sqrt{n}\{\hat{\theta}_d - \theta_d\}_{d \in \{0,1\}} \approx \sqrt{n} \frac{\mathbb{E}_n[(Y(d) - EY(d))1(D = d)]}{P(D = d)}.$$

Note that the terms being averaged are

$$\frac{(Y_i(d) - E[Y(d)])1(D_i = d)}{P(D = d)}.$$

These terms have zero mean<sup>22</sup> and variance

$$\frac{E[(Y(d) - E[Y(d)])^2 1(D = d)^2]}{P(D = d)^2} = \frac{\text{Var}(Y | 1(D = d) = 1)}{P(D = d)}.$$

22: Why? Hint: Use the law of iterated expectations.

Also note the zero covariance:

$$E \left[ \frac{(Y(1) - E[Y(1)])1(D = 1)}{P(D = 1)} \frac{(Y(0) - E[Y(0)])1(D = 0)}{P(D = 0)} \right] = 0.$$

The application of the central limit theorem then yields the claimed result.

## 2.B Statistical Properties of the Classical Additive Approach\*

Here we analyze statistical inference on ATE using OLS and adjusting for  $X = (1, W)$ , without making the linearity assumptions we made in Section 2.2.

We consider the linear projection equation in the population:

$$Y = D\alpha + X'\beta + \epsilon, \quad \epsilon \perp (D, X).$$

Here, we have that  $D$  and  $X = (1, W)$  with  $E[W] = 0$ , so that  $\beta'X = \beta_1 + \beta_2'W$ . Moreover, we have that  $D \perp W$  in the RCT setting.

First, we'd like to verify that  $\alpha = E[Y(1)] - E[Y(0)]$  and  $\beta_1 = E[Y(0)]$ . For  $U := \beta_2'W + \epsilon$ , we can write

$$Y = D\alpha + \beta_1 + U, \quad U \perp (1, D).$$

$U \perp (1, D)$  holds because  $(1, D) \perp (W, \epsilon)$  using that  $E[W] = 0$  and that  $D \perp (W, \epsilon)$ . Therefore,  $D\alpha + \beta_1$  coincides with the population projection of  $Y$  onto  $(1, D)$ . Hence, the projection coefficients are the same as those obtained by the 2-sample approach in the population. Therefore,  $\beta_1 = E[Y(0)]$  and  $\alpha = E[Y(1)] - E[Y(0)]$ .

Second, we'd like to explain the details of the approximate normality for the estimators of sample OLS coefficients  $\hat{\beta}_1$ . The OLS theory of the first chapter implies that the OLS estimator  $\hat{\alpha}$  obeys

$$\sqrt{n}(\hat{\alpha} - \alpha) \approx \sqrt{n} \frac{E_n[\epsilon \tilde{D}]}{E_n[\tilde{D}^2]} \stackrel{a}{\sim} N(0, V_{11}),$$

where  $\tilde{D} = D - E[D]$  is the residual after partialling out  $X$  from  $D$  linearly,<sup>23</sup> and

$$V_{11} = \frac{E[\epsilon^2 \tilde{D}^2]}{(E[\tilde{D}^2])^2}.$$

23: Derive that  $\tilde{D} = D - E[D]$  from Assumption 2.2.1.



Applying the same theory for  $\beta_1$  (the intercept coefficient), yields<sup>24</sup>

$$\sqrt{n}(\hat{\beta}_1 - \beta_1) \approx \sqrt{n} \frac{\mathbb{E}_n[\epsilon \tilde{1}]}{\mathbb{E}_n[\tilde{1}^2]} \stackrel{a}{\sim} N(0, \mathbf{V}_{22}),$$

where  $\tilde{1} := (1 - D)$  is the residual after partialling out  $D$  and  $X$  from 1 and

$$\mathbf{V}_{22} = \frac{\mathbb{E}[\epsilon^2 \tilde{1}^2]}{(\mathbb{E}[\tilde{1}^2])^2}.$$

We can also establish that the estimators are jointly approximately normal with covariance

$$\mathbf{V}_{12} = \frac{\mathbb{E}[\epsilon^2 \tilde{D} \tilde{1}]}{\mathbb{E}[\tilde{1}^2] \mathbb{E}[\tilde{D}^2]}.$$

24: To explain the derivation, note that by partialling out  $D$  and  $W$  (recall that  $X = (1, W)$ ) from 1 and  $Y$ , we obtain

$$\tilde{Y} = \beta_1 \tilde{1} + \epsilon; \quad \tilde{1} := (1 - D).$$

The projection of 1 on  $D$  and  $W$  is given by  $D$  since  $D$  is binary and we've assumed  $\mathbb{E}[W] = 0$ .

## 2.C Statistical Properties of the Interactive Regression Approach<sup>★</sup>

Here we analyze the estimation of the ATE using OLS and adjusting for  $(W, DW)$  without making any linearity assumptions on the potential outcomes as we did in Section 2.2. We essentially show that the interactive model can be viewed as estimating the BLP of each of the two potential outcomes  $Y(0)$  and  $Y(1)$ . Using this fact one can then easily argue that the variance of the OLS estimate of the effect using the interactive model can only be lower than the variance of the unadjusted OLS estimate.

Letting  $X = (1, W)$  be an intercept and the pre-treatment covariates  $W$ , let us write the BLP of each of  $Y(0)$  and  $Y(1)$  using  $X$  as

$$Y(d) = \beta'_d X + \epsilon_d, \quad \epsilon_d \perp X, \quad d = 0, 1. \quad (2.C.1)$$

Under Assumption 2.2.1, (2.C.1) coincides with the BLP of  $Y$  using  $X$  in the  $D = d$  population. Letting  $\epsilon = D\epsilon_1 + (1 - D)\epsilon_0$ , we thus have

$$Y = \beta'_d X + \epsilon, \quad \mathbb{E}[\epsilon X \mid D = d] = 0, \quad d = 0, 1. \quad (2.C.2)$$

The BLPs in each of the two populations,  $D = 0$  and  $D = 1$ , can be combined across the populations to state the BLP of  $Y$  using  $(X, DX)$  marginally:

$$Y = \beta'_0 X + \beta'_\delta X D + \epsilon, \quad \epsilon \perp (X, DX), \quad (2.C.3)$$

where  $\beta_\delta = \beta_1 - \beta_0$ .<sup>25</sup> Such a linear rule is called *interactive* because it includes the interaction (meaning, product) of  $D$  and  $W$  as a regressor, in addition to  $D$  and  $W$ .

We assume that covariates are centered:

$$E[W] = 0.$$

Since  $X$  contains an intercept,  $\varepsilon_d \perp X$  implies  $E[\varepsilon_d] = 0$ . Together with centered covariates, we find that

$$E[Y(d)] = E[\beta'_d X + \varepsilon_d] = \beta_{d,1}.$$

This means that the ATE coincides with the coefficient on  $D$  in the BLP of  $Y$  using  $(X, DX)$ . That is,  $\beta_{\delta,1} = \delta$ .

We are often interested in the ATE and Relative ATE

$$\delta \quad \text{and} \quad \delta/E[Y(0)].$$

If we use OLS to estimate the BLP of  $Y$  using  $(X, DX)$ , then an application of the OLS theory in the previous chapter gives us that, under regularity conditions,

$$\begin{pmatrix} \sqrt{n}(\hat{\beta}_{\delta,1} - \delta) \\ \sqrt{n}(\hat{\beta}_{0,1} - E[Y(0)]) \end{pmatrix} \stackrel{a}{\sim} N(0, \mathbf{V}),$$

where covariance matrix  $\mathbf{V}$  has components:

$$V_{11} = \frac{E[\varepsilon^2 \tilde{D}^2]}{(E[\tilde{D}^2])^2}, \quad V_{22} = \frac{E[\varepsilon^2 \tilde{1}^2]}{(E[\tilde{1}^2])^2}, \quad V_{12} = V_{21} = \frac{E[\varepsilon^2 \tilde{D} \tilde{1}]}{E[\tilde{1}^2]E[\tilde{D}^2]},$$

where  $\tilde{D} = D - E[D]$  is the residual after partialling out linearly  $(1, W, DW)$  from  $D$  and  $\tilde{1} := (1 - D)$  is the residual after partialling out  $(D, W, DW)$  from 1.<sup>26</sup>

We can then obtain the approximate normality for the Relative ATE using the delta method:

$$\sqrt{n}(\hat{\beta}_{\delta,1}/\hat{\beta}_{0,1} - \delta/E[Y(0)]) \stackrel{a}{\sim} N(0, G'VG),$$

where

$$G = [1/E[Y(0)], -\delta/(E[Y(0)]^2)]'.$$

We can rewrite (2.C.3) as

$$Y = \beta_{0,1} + D\beta_{\delta,1} + U, \quad U = \beta'_{0,2}W + \beta'_{\delta,2}WD + \varepsilon.$$

From  $\varepsilon \perp (X, D, DX)$ ,  $E[W] = 0$ , and Assumption 2.2.1, we obtain that  $U \perp (1, D)$ , meaning that  $\beta_{0,1} + D\beta_{\delta,1}$  is the BLP

25: Note that (2.C.1) and (2.C.2) imply  $E[\varepsilon DX] = 0$  and  $E[\varepsilon X] = 0$  and thus that  $\varepsilon \perp (X, DX)$ .

26: The derivation follows identical steps as that in Section 2.B with the only exception that when defining  $\tilde{D}$  we need to partial out  $(1, W, DW)$  from  $D$  and when defining  $\tilde{1}$  we need to partial out  $(D, W, DW)$  from 1. However, since  $E[W] = E[DW] = 0$ , the two residuals take the same form of  $D - E[D]$  and  $1 - D$  correspondingly.

of  $Y$  using  $(1, D)$ . We can therefore estimate the ATE as the coefficient on  $D$  either in the OLS of  $Y$  on  $(1, D)$  or in the OLS of  $Y$  on  $(X, DX)$ . The former exactly coincides with the unadjusted estimator  $\hat{\delta}$  from Section 2.1, which obeys

$$\sqrt{n}(\hat{\delta} - \delta) \stackrel{a}{\sim} N(0, \bar{V}_{11}), \quad \bar{V}_{11} = \frac{E[U^2 \tilde{D}^2]}{(E[\tilde{D}^2])^2}.$$

Since  $\epsilon$  satisfies the BLP conditions for each of the treatment populations, i.e.  $E[\epsilon W \mid D = d] = 0$ , it then follows that

$$V_{11} \leq \bar{V}_{11}.$$

Moreover, the inequality is strict if  $\text{Var}(\beta'_{0,2} W) > 0$  or  $\text{Var}(\beta'_{1,2} W) > 0$ .<sup>27</sup> That is, pre-determined covariates improve the precision of estimating the ATE  $\delta$ , when using the interactive model, without any linearity assumptions on the CEF.

27: Verify this as a reading exercise.

# Bibliography

- [1] Jan Baptist van Helmont. *Oriatrike, Or, Physick Refined, the Common Errors Therein Refuted, and the Whole Art Reformed and Rectified*. Loyd, London, 1662 (cited on page 41).
- [2] Donald B. Rubin. 'Estimating causal effects of treatments in randomized and nonrandomized studies.' In: *Journal of Educational Psychology* 66.5 (1974), pp. 688–701 (cited on page 42).
- [3] P. Holland. 'Causal Inference, Path Analysis, and Recursive Structural Equations Models'. In: *Sociological Methodology*. Washington, DC: American Sociological Association, 1986, pp. 449–493 (cited on page 42).
- [4] Peng Ding and Fan Li. 'Causal Inference: A Missing Data Perspective'. In: *Statistical Science* 33.2 (2018), pp. 214–237. DOI: [10.1214/18-STS645](https://doi.org/10.1214/18-STS645) (cited on page 42).
- [5] Tyler J. VanderWeele, Guanglei Hong, Stephanie M. Jones, and Joshua L. Brown. 'Mediation and Spillover Effects in Group-Randomized Trials: A Case Study of the 4Rs Educational Intervention'. In: *Journal of the American Statistical Association* 108.502 (2013), pp. 469–482. (Visited on 02/17/2024) (cited on page 43).
- [6] Peter M. Aronow and Cyrus Samii. 'Estimating average causal effects under general interference, with application to a social network experiment'. In: *The Annals of Applied Statistics* 11.4 (2017), pp. 1912–1947. DOI: [10.1214/16-AOAS1005](https://doi.org/10.1214/16-AOAS1005) (cited on page 43).
- [7] Michael P. Leung. 'Treatment and Spillover Effects Under Network Interference'. In: *The Review of Economics and Statistics* 102.2 (2020), pp. 368–380 (cited on page 43).
- [8] Francis J. DiTraglia, Camilo García-Jimeno, Rossa O'Keefe-O'Donovan, and Alejandro Sánchez-Becerra. 'Identifying causal effects in experiments with spillovers and non-compliance'. In: *Journal of Econometrics* 235.2 (2023), pp. 1589–1624. DOI: <https://doi.org/10.1016/j.jeconom.2023.01.008> (cited on page 43).
- [9] Gonzalo Vazquez-Bare. 'Identification and estimation of spillover effects in randomized experiments'. In: *Journal of Econometrics* 237.1 (2023), p. 105237. DOI: <https://doi.org/10.1016/j.jeconom.2021.10.014> (cited on page 43).

- [10] Guido W. Imbens and Donald B. Rubin. *Causal Inference in Statistics, Social, and Biomedical Sciences*. Cambridge University Press, 2015 (cited on pages 43, 61).
- [11] Walter A Orenstein, Roger H Bernier, Timothy J Dondero, Alan R Hinman, James S Marks, Kenneth J Bart, and Barry Sirotkin. *Field evaluation of vaccine efficacy / Walter A. Orenstein ... [et al.]* 1984 (cited on page 49).
- [12] Jerome Cornfield. 'A statistical problem arising from retrospective studies'. In: *Proceedings of the Third Berkeley Symposium on Mathematical Statistics and Probability*. Vol. 4. University of California Press Berkeley, CA. 1956, pp. 135–148 (cited on page 50).
- [13] Winston Lin. 'Agnostic notes on regression adjustments to experimental data: Reexamining Freedman's critique'. In: *Annals of Applied Statistics* 7.1 (2013), pp. 295–318 (cited on page 56).
- [14] Yannis Biliias. 'Sequential testing of duration data: The case of the Pennsylvania 'reemployment bonus' experiment'. In: *Journal of Applied Econometrics* 15.6 (2000), pp. 575–594 (cited on page 56).
- [15] Philip G. Wright. *The Tariff on Animal and Vegetable Oils*. New York: The Macmillan company, 1928 (cited on page 57).
- [16] Sewall Wright. 'Correlation and Causation'. In: *Journal of Agricultural Research* 20.7 (Jan. 1921), pp. 557–585 (cited on page 57).
- [17] Judea Pearl. 'Causal diagrams for empirical research'. In: *Biometrika* 82.4 (1995), pp. 669–688 (cited on page 57).
- [18] Sander Greenland, Judea Pearl, and James M. Robins. 'Causal diagrams for epidemiologic research'. In: *Epidemiology* 10.1 (1999), pp. 37–48 (cited on page 57).
- [19] David R. Cox. *Planning of experiments*. Wiley, 1958 (cited on page 58).
- [20] *The Belmont report: Ethical principles and guidelines for the protection of human subjects of research*. Tech. rep. National Commission for the Protection of Human Subjects of Biomedical and Behavioral Research, 1978 (cited on page 59).
- [21] Art Owen. 'A First Course in Experimental Design: Notes from Stat 263/363'. Lecture notes. Accessed 1/17/2024. 2020 (cited on page 61).
- [22] Esther Duflo, Rachel Glennerster, and Michael Kremer. 'Using randomization in development economics research: A toolkit'. In: *Handbook of Development Economics* 4 (2007), pp. 3895–3962 (cited on page 61).