

# Causality & Design-Based Inference

Jake Bowers and Thomas Leavitt

July 26, 2019

## **Abstract**

Counterfactual causal quantities cannot be observed, but researchers can use statistical procedures — namely, estimators and hypothesis tests — to draw inferences from data that can be observed. In this chapter, we present a unified account of estimation and testing for causal inference, showing how a study’s research design can provide a foundation for both estimators and tests. We show how certain characteristics of research designs can justify claims that a given estimator or test has “good” properties (e.g., unbiasedness, consistency, controlled error rates). We develop ideas first in the context of a randomized controlled experiment. In that context, we juxtapose estimation of and tests about causal effects and then provide an explicit comparison of Fisherian and Neymanian hypothesis tests. We then extend our analysis to research designs that are either partially controlled (e.g., experiments with non-compliance and/or attrition) or uncontrolled (e.g., observational studies). We show the ways in which knowledge and assumptions about research design — as well as assessments of how inferences would change were these assumptions to be false — constitute a reliable basis for causal inference. We conclude by discussing the value of design-based causal inference in light of recent debates on its role in social-scientific inquiry more broadly.

# 1 Design-Based Causal Inference

No one knows the true causal effect of an intervention. In an experiment, a researcher can assign some units to treatment and others to control; yet, one cannot see how treated units would have acted were they assigned to control nor how the control units would have acted were they assigned to treatment.<sup>1</sup> In the face of this fundamental ignorance, statisticians have developed two prominent approaches to inferring unobservable causal effects using data that can be observed. An analyst can either (1) generate a guess about (usually average) treatment effects or (2) posit a hypothesis about the effects of a treatment (such as the hypothesis that a treatment had no effects) and then assess the consistency of observable data with that null hypothesis relative to a class of alternative hypotheses (such as the hypothesis that a treatment had a positive effect).

In what follows, we will define criteria by which a procedure qualifies as “good” in the context of both estimation and testing, and subsequently explain the role research design plays in whether estimators and tests satisfy these criteria. We consider estimators and tests about causal effects, first, in the context of a randomized study design under full control of the researcher and, second, in cases in which the researcher does not fully control the study design. We show the ways in which either complete knowledge or assumptions (and the ways in which they could be violated) about the study design constitute what [Fisher \(1935\)](#) referred to as a “reasoned basis for inference.”

## 2 Causality and Research Design

### 2.1 Defining Causal Effects

Consider a study in which there is a finite population of  $1, \dots, N$  units and the index  $i \in \{1, \dots, N\}$  runs over these units. Each individual,  $i$ , can be in either the treatment condition,  $z_i = 1$ , or the control condition,  $z_i = 0$ . Under the Stable Unit Treatment Value Assumption (SUTVA) ([Cox, 1958](#); [Rubin, 1980, 1986](#)), each individual has a treated potential outcome,  $y_{t,i}$  (unit  $i$ ’s outcome if given the intervention), and a control potential outcome,  $y_{c,i}$  (unit  $i$ ’s outcome if not given the intervention).<sup>2</sup> An individual causal effect,  $\tau_i$ , for each of the  $i \in \{1, \dots, N\}$  units is a function of

---

<sup>1</sup>[Holland \(1986, 947\)](#) refers to the inability to observe both potential outcomes for a single unit at the same moment in time as the “fundamental problem of causal inference.”

<sup>2</sup>For simplicity, we consider studies in which there are two conditions—treatment and control—although the same general principles apply to studies with multiple conditions (see [Dasgupta et al., 2015](#)).

each unit's two potential outcomes,  $\tau_i \equiv f(y_{c,i}, y_{t,i})$ , such as,  $\tau_i \equiv \frac{y_{t,i}}{y_{c,i}}$  or  $\tau_i \equiv y_{t,i} - y_{c,i}$ . For example, for this paper, we focus specifically on the additive, individual causal effect defined as  $\tau_i \equiv y_{t,i} - y_{c,i}$ . Researchers, however, can never observe both potential outcomes for each unit; instead, one can observe only  $y_i$ , which can be equal to either  $y_{c,i}$  or  $y_{t,i}$  depending on whether unit  $i$  is assigned to treatment ( $z_i = 1$ ) or control ( $z_i = 0$ ). We therefore represent observed outcomes by the function  $y_i = z_i y_{t,i} + (1 - z_i) y_{c,i}$ . For example, researchers may want to make an inference about the  $1, \dots, N$  individual causal effects, which we collect into the vector  $\boldsymbol{\tau}' = [\tau_1 \quad \tau_2 \quad \cdots \quad \tau_N]$ , or about a function of these individual causal effects, such as the average causal effect,  $\bar{\tau} = \left(\frac{1}{N}\right) \sum_{i=1}^N \tau_i$ . We say that researchers want to “make an inference” because neither  $\boldsymbol{\tau}$  nor  $\bar{\tau}$  can be directly observed. Researchers, of course, don't simply want to “make an inference”; they want to make inferences that can reliably track the true causal quantity of interest. This chapter shows how inferential procedures based on the research design can have such a reliable relationship with true causal quantities and explains what it means for a procedure to be “based on research design.”

## 2.2 Defining a Research Design

At its most basic level, a research design refers to the process by which units come to be in one study condition instead of another, i.e., each  $z_i$  comes to equal 1 or 0. More formally, we denote the collection of the values of  $z_i$  for all  $i \in \{1, \dots, N\}$  units by the vector  $\mathbf{z}' = [z_1 \quad \cdots \quad z_N]$  and define a research design as (1) a set of possible ways (events) in which the whole vector  $\mathbf{z}$  could occur and (2) a probability distribution on this set of possible events. In a controlled study design (i.e., an experiment), we think of a researcher as “assigning” conditions to all units in the study. When a researcher does not control how a unit  $i$  takes on a value of  $z_i$ , we think of that unit as “selecting” into its own condition. As we lay the groundwork of concepts and notation, we write “assignment” and assume control by the researcher, but we will apply the general framework later to uncontrolled research designs in which units “select” into study conditions.

The set of possible ways in which  $\mathbf{z}$  occur depends on the process by which units are assigned to study conditions. If individuals can be in either the treatment or control condition irrespective of any other individual in the population, then we call this process individual assignment. As a simple example consider the “coin flip” assignment process: In this case, the proportion of  $N$  individuals

in either the treatment or control condition can vary across different assignments. We refer to this process as *simple* individual assignment. A researcher can implement simple, individual assignment via an actual physical, stochastic process, such as  $N$  flips of a (potentially biased) coin, although in practice researchers will typically use random number generators (RNGs).

Under completely unconstrained simple, individual assignment, the number of units in the treatment condition can range from 0 to  $N$  and the number of units in the control condition can likewise range from  $N - 0$  to  $N - N$ . For example, in an experiment with 6 units, simple, individual assignment can allow 0, 1, 2, 3, 4, 5 or 6 units to be in the treatment and control conditions, respectively. More formally, we write the set,  $\Omega$ , of possible ways that a researcher can assign all individuals to study conditions as follows:

$$(1) \quad \Omega = \{0, 1\}^N = \left\{ \begin{bmatrix} 0 \\ 0 \\ \vdots \\ 0 \\ 0 \end{bmatrix}, \begin{bmatrix} 1 \\ 0 \\ \vdots \\ 0 \\ 0 \end{bmatrix}, \begin{bmatrix} 0 \\ 1 \\ \vdots \\ 0 \\ 0 \end{bmatrix}, \dots, \begin{bmatrix} 1 \\ 0 \\ \vdots \\ 1 \\ 1 \end{bmatrix}, \begin{bmatrix} 0 \\ 1 \\ \vdots \\ 1 \\ 1 \end{bmatrix}, \begin{bmatrix} 1 \\ 1 \\ \vdots \\ 1 \\ 1 \end{bmatrix} \right\}.$$

We can write the number of possible assignments in the set  $\Omega$  by  $|\Omega|$  (the “cardinality of Omega”), under simple assignment as follows:

$$\begin{aligned} |\Omega| &= \binom{N}{0} + \binom{N}{1} + \dots + \binom{N}{N-1} + \binom{N}{N} \\ &= \sum_{n_t=0}^N \binom{N}{n_t}, \end{aligned}$$

where  $n_t = \sum_{i=1}^N z_i$  is the number of units in the treatment condition, which can range from 0 to  $N$ , and  $\binom{N}{n_t} = \frac{N!}{(N-n_t)!n_t!}$  is the number of ways to choose  $n_t$  units from a total of  $N$  units. Conversely,  $n_c = \sum_{i=1}^N (1 - z_i)$  is the number of units in the control condition, which can range from  $N - 0$  to  $N - N$ . In practice, researchers who control the assignment process will typically forbid assignments in which all units are in either condition (the first and last assignments in Equation (1)), in which case  $|\Omega| = \sum_{n_t=1}^{N-1} \binom{N}{n_t}$ .

The “coin flipping” design helps us introduce the formal elements of a research design. In

practice, individual coin flips can lead to lopsided designs in which many units are in one condition or another. An alternative design that enables the researcher to control the numbers of units in each condition is *complete* assignment.

Complete, individual assignment differs from simple, individual assignment only in that the value of  $n_t$  is fixed across all possible assignments. We have described simple, individual assignment via the example of coin flips. Complete, individual assignment can be thought of as draws from an urn. Imagine, for example, that an urn contains  $N$  balls, of which  $n_t$  are blue balls and  $N - n_t = n_c$  are red balls. The researcher could draw the first ball from the urn and assign the first unit in the study to the treatment condition if the ball is blue and to the control condition if the ball is red. The second draw could follow the same rule for the assignment of the second unit, and so and so forth until no more balls remain in the urn. This form of assignment ensures that exactly  $n_t$  units are in the treatment condition and  $N - n_t = n_c$  units are in the control condition. More formally, complete, individual assignment excludes any assignment,  $\mathbf{z}$ , with more or less treatment units,  $n_t$ , than that which is predetermined by the researcher. Therefore, under complete assignment, the number of possible assignments is simply  $|\Omega| = \binom{N}{n_t}$ .

Simple and complete assignment can also happen at the *cluster* (as opposed to the *individual*) level. In this setup, we have not only a set of  $1, \dots, N$  individuals, but also a set of  $1, \dots, K$  clusters, where each cluster,  $k \in \{1, \dots, K\}$ , contains  $n_k \geq 1$  individual units and  $N = \sum_{k=1}^K n_k$ . In cluster assignment designs, all of the  $i \in \{1, \dots, n_k\}$  units in the  $k$ th cluster are assigned to either the treatment condition,  $z_{ik} = 1$ , or the control condition,  $z_{ik} = 0$ . In simple, cluster assignment, the number of possible assignments is given by  $|\Omega| = \sum_{k_t=0}^K \binom{K}{k_t}$ , where  $k_t$  denotes the number of treatment clusters, although (just as in simple, individual assignment) researchers will typically ensure that  $k_t \notin \{0, K\}$ . Under complete, cluster assignment, the number of treatment clusters is fixed, such that the number of assignments is  $|\Omega| = \binom{K}{k_t}$ .

Lastly, *blocked assignment* is when individuals or clusters are assigned (either simply or completely) to study conditions within blocks, which we index from  $b \in \{1, \dots, B\}$ . Blocks are typically constructed on the basis of individuals' or clusters' values of baseline covariates. Baseline covariates are measured prior to assignment and, hence, their values are fixed regardless of the condition

to which a unit or cluster is assigned. Under simple, individual, block assignment, the number of possible assignments is  $|\Omega| = \prod_{b=1}^B \left( \sum_{n_{t,b}=1}^{N_b-1} \binom{N_b}{n_{t,b}} \right)$ , where  $N_b$  is the number of units in block  $b$ ,  $n_{t,b}$  is the number of units in the treatment condition in block  $b$  and  $n_{t,b} \notin \{0, N_b\}$  for all  $b$ . Under complete, individual, block assignment  $|\Omega| = \prod_{b=1}^B \binom{N_b}{n_{t,b}}$ . One can analogously deduce the number of possible assignments under either simple or complete cluster, block assignment. As we will explain in subsequent sections, block assignment carries important implications for properties of both estimators and hypothesis tests.

Given a set of possible assignments,  $\Omega$ , arising from an assignment mechanism, the remaining component of a research design is a probability distribution on this set of assignments. In a *uniform, randomized experiment*, the probability of each assignment is simply  $\frac{1}{|\Omega|}$  for all  $\mathbf{z}$ , whereby each assignment has an identical probability of realization. Yet the probability distribution on  $\Omega$  need not be uniform, even in a randomized experiment. Design-based inference means only that the stochastic properties of estimators and tests be based on this probability distribution on  $\Omega$ , regardless of whether that distribution is uniform or not. As we now move to discussions of both estimation and testing, notice throughout that whenever we refer to random quantities, the randomness of those quantities stems solely from the probability distribution on  $\Omega$ .

### 2.3 An Illustrative Example

In the sections to follow, we demonstrate our arguments via a simple hypothetical example that consists of  $N = 6$  units and an individual assignment process (complete, individual assignment) in which 3 units are assigned to treatment ( $n_t = 3$ ) and to control ( $n_c = 3$ ). Let's further imagine that (unbeknownst to the researcher) the 6 units' potential outcomes and individual causal effects are as follows in Table 1:

$\mathbf{y}_c$	$\mathbf{y}_t$	$\boldsymbol{\tau}$
20	22	2
8	12	4
11	11	0
10	15	5
14	18	4
1	4	3

**Table 1:** True values of  $\mathbf{y}_c$ ,  $\mathbf{y}_t$  and  $\boldsymbol{\tau}$ , where  $\tau_i = y_{t,i} - y_{c,i}$ , for the study population

Based on the complete, individual assignment process in which there are  $N = 6$  units and  $n_t = 3$  of which are assigned to treatment, the set of  $|\Omega| = \binom{6}{3} = 20$  possible assignments is given by Equation (2):

$$(2) \quad \Omega = \left\{ \begin{bmatrix} 1 \\ 1 \\ 1 \\ 0 \\ 0 \\ 0 \end{bmatrix}, \begin{bmatrix} 1 \\ 1 \\ 0 \\ 1 \\ 0 \\ 0 \end{bmatrix}, \dots, \begin{bmatrix} 0 \\ 0 \\ 1 \\ 0 \\ 1 \\ 1 \end{bmatrix}, \begin{bmatrix} 0 \\ 0 \\ 0 \\ 1 \\ 1 \\ 1 \end{bmatrix} \right\}.$$

The assignment one draws from the set  $\Omega$  determines which potential outcomes one observes. One can observe treatment potential outcomes only for units assigned to treatment and control potential outcomes only for units assigned to control. (Recall the function  $y_i = z_i y_{t,i} + (1 - z_i) y_{c,i}$ , which determines the potential outcome one observes for each individual  $i$ .) As Table 2 shows, for each of the  $\binom{6}{3} = 20$  possible assignments, there are  $\binom{6}{3} = 20$  corresponding possible realizations of observed data  $y$ :

$z_1$	$y_c$	$y_t$	$y_1$	$z_2$	$y_c$	$y_t$	$y_2$	$z_{19}$	$y_c$	$y_t$	$y_{19}$	$z_{20}$	$y_c$	$y_t$	$y_{20}$
1	?	22	22	1	?	22	22	0	20	?	20	0	20	?	20
1	?	12	12	1	?	12	12	0	8	?	8	0	8	?	8
1	?	11	11	0	11	?	11	1	?	11	11	0	11	?	11
0	10	?	10	1	?	15	15	0	10	?	10	1	?	15	15
0	14	?	14	0	14	?	14	1	?	18	18	1	?	18	18
0	1	?	1	0	1	?	1	1	?	4	4	1	?	4	4

**Table 2:** All possible realizations of experimental data from a completely randomized study with 6 units and 3 treated units.

Only one such possible realization of data in Table 2 can be observed, but, knowing that there are 20 possible realizations, allows the researcher to use procedures—e.g., estimators or hypothesis tests—to make inferences about unobservable causal quantities based on the single observed realization. We want procedures for drawing causal inferences to have properties that are “good” (a notion that we will define more precisely in later sections). These properties describe or measure a procedure’s performance in two contexts: (1) studies with a fixed, finite population and (2) a hypothetical scenario in which the size of a given study increases towards  $\infty$ , while all other relevant factors remain constant. We refer to the latter context as one of *asymptotic growth*, which we conceptualize below.

If we have an experimental pool of 6 units, that are not a sample using a known procedure from a well-defined population, what does “asymptotic growth” mean? We follow [Brewer \(1979\)](#) and [Middleton and Aronow \(2015\)](#) in using the idea of “copies” as a way to talk about how estimators and tests behave as study sizes increase. In short, this conception of asymptotic growth states that (1) the original population of  $N$  units is copied  $h - 1$  times; (2) within each of the  $h$  copies, exactly  $n_t$  units are assigned to the treatment condition and the remaining  $n_c = N - n_t$  units are assigned to the control condition; (3) the  $h$  copies are then collected into a single population with  $hN$  total units,  $hn_t$  treated units and  $hn_c$  control units.

In the context of our working example, this conception of growth stipulates that the study population of  $N = 6$  units is embedded in a sequence of populations of increasing sizes in which the initial population is simply copied  $h - 1$  times:

$y_c$	$y_t$	$\tau$
20	22	2
8	12	4
11	11	0
10	15	5
14	18	4
1	4	3

$y_c$	$y_t$	$\tau$
20	22	2
8	12	4
11	11	0
10	15	5
14	18	4
1	4	3

$y_c$	$y_t$	$\tau$
20	22	2
8	12	4
11	11	0
10	15	5
14	18	4
1	4	3

$y_c$	$y_t$	$\tau$
20	22	2
8	12	4
11	11	0
10	15	5
14	18	4
1	4	3

$y_c$	$y_t$	$\tau$
20	22	2
8	12	4
11	11	0
10	15	5
14	18	4
1	4	3

$y_c$	$y_t$	$\tau$
20	22	2
8	12	4
11	11	0
10	15	5
14	18	4
1	4	3

$y_c$	$y_t$	$\tau$
20	22	2
8	12	4
11	11	0
10	15	5
14	18	4
1	4	3

$y_c$	$y_t$	$\tau$
20	22	2
8	12	4
11	11	0
10	15	5
14	18	4
1	4	3

$y_c$	$y_t$	$\tau$
20	22	2
8	12	4
11	11	0
10	15	5
14	18	4
1	4	3

**Table 3:** Finite populations under asymptotic growth in which  $h \in \{1, 2, 3, 4, \dots\}$

Notice that over this sequence of increasing finite populations in [Table 3](#), all relevant factors other than  $N$  remain constant: The proportions of treatment and control units remain fixed, the mean of control and treatment potential outcomes remain fixed, as do their variances and their covariance. Notice, however, that the number of possible assignments increases over this sequence of increasing finite populations from  $\binom{6}{3} = 20$  to  $\binom{12}{6} = 924$  to  $\binom{18}{9} = 48620$  to  $\binom{24}{12} = 2704156$  and



so forth.

We will show that, in either the finite context—given in Table 1—or in the asymptotic context—given in Table 3—whether a procedure is “good” depends on whether it maintains fidelity to the research design—i.e., the probability distribution on the set  $\Omega$ . In other words, we will show that, in a randomized experiment, a “good” procedure is one that heeds the dictum of Senn (2004, 3729) who, in the voice of R. A. Fisher, states that “[a]s ye randomise so shall ye analyse.”

### 3 Estimation

As we have mentioned above, no one can observe both potential outcomes for any given unit in a given study population. One can, however, generate a guess about some function of the study population’s individual causal effects (e.g., the mean causal effect) using observed outcomes. We call this unobservable, causal quantity the *estimand*. The *estimator*, by contrast, refers to the procedure that generates a guess about the estimand. An *estimate* is the actual output of the estimator once it is applied to a given data set.

One estimand is the mean causal effect,  $\bar{\tau} = \left(\frac{1}{N}\right) \sum_{i=1}^N \tau_i$ , which, to return to the example from Table 1, is  $\bar{\tau} = \frac{2+4+0+5+4+3}{6} = 3$ . A procedure for generating a guess about  $\bar{\tau}$  is the Difference-in-Means estimator, which we can define in terms of observable quantities as follows:

$$(3) \quad \hat{\tau}(\mathbf{Z}, \mathbf{Y}) = \frac{\mathbf{Z}'\mathbf{Y}}{\mathbf{Z}'\mathbf{Z}} - \frac{(\mathbf{1} - \mathbf{Z})'\mathbf{Y}}{(\mathbf{1} - \mathbf{Z})'(\mathbf{1} - \mathbf{Z})} \\ = \left( \frac{1}{\sum_{i=1}^N Z_i} \right) \sum_{i=1}^N Z_i Y_i - \left( \frac{1}{\sum_{i=1}^N (1 - Z_i)} \right) \sum_{i=1}^N (1 - Z_i) Y_i.$$

In the example from Table 1, the random variables<sup>3</sup> of  $\mathbf{Z}$  and  $\mathbf{Y}$  can take on any of the possible values,  $(\mathbf{z}_1, \mathbf{y}_1), \dots, (\mathbf{z}_{20}, \mathbf{y}_{20})$ , given in Table 2. If we apply the estimator in Equation 3 to the possible realizations of data in Table 2, then there are 20 possible estimates that correspond to

---

<sup>3</sup>To distinguish between fixed quantities and quantities that can take on different values with certain probabilities (i.e., random quantities), we now use uppercase letters for random quantities and lowercase letters for fixed or realized quantities.

each of the 20 possible realization of data:

$$\hat{\tau}(\mathbf{z}_1, \mathbf{y}_1) = 6.6667, \hat{\tau}(\mathbf{z}_2, \mathbf{y}_2) = 7.6667, \dots, \hat{\tau}(\mathbf{z}_{19}, \mathbf{y}_{19}) = -1.6667, \hat{\tau}(\mathbf{z}_{20}, \mathbf{y}_{20}) = -0.6667.$$

The researcher can observe only one of these 20 possible estimates and this single estimate should be generated by an estimator that is “good.” More specifically, three “good” properties of an estimator are *unbiasedness*, *consistency* and *precision*. An unbiased estimator is one in which, although any single estimate may be close to or far from the true value of the estimand, the *expected value of the estimator*—i.e., the probability-weighted mean of all possible estimates—is equal to the value of the estimand. Consistency states that, as the number of units in the study increases asymptotically, holding all other factors constant, the probability distribution of an estimator concentrates increasingly around the truth. (For any fixed  $\varepsilon > 0$ , the probability that the estimate and its target differ by no more than  $\varepsilon$  tends to 1.) Lastly, a precise estimator is one in which the expected distance of an estimate from the true causal effect is small.

In the following discussion, we show the role that research design plays in whether the Difference-in-Means estimator is unbiased, consistent and/or precise with respect to the estimand  $\bar{\tau} = \left(\frac{1}{N}\right) \sum_{i=1}^N \tau_i$ . We also show how designs can yield estimators that are more or less precise. Researchers may want to estimate quantities other than  $\bar{\tau}$ . Although we do not discuss such cases, the general principles for determining whether an estimator is unbiased, consistent and/or precise with respect to the causal quantity of interest is the same: Researchers can define an estimand that they seek to infer, define an estimator by which they would estimate this quantity under all possible realizations of data, and subsequently assess whether this estimator is unbiased, consistent and/or precise based only on the probabilities with which possible data are realized.

### 3.1 Unbiasedness

We now show that a “good” estimator of the unknown estimand,  $\bar{\tau}$ , is the estimator given in Equation (3),  $\hat{\tau}(\mathbf{Z}, \mathbf{Y})$ . In particular, we will show that this estimator satisfies the criterion of

---

<sup>4</sup>When the numbers of treatment and control units are *not* fixed, such as in simple, individual assignment, the Difference-in-Means estimator remains unbiased in a uniform randomized experiment so long as at least one unit is always in the treatment and control conditions, respectively. In general, when the numbers of treatment and control are *not* fixed, the Difference-in-Means estimator is not necessarily unbiased, such as in cluster uniform random assignment when clusters are of unequal sizes (see [Middleton and Aronow, 2015](#)).

no systematic error—i.e., unbiasedness—in a *uniform randomized experiment* when the numbers of treatment and control units are both fixed.<sup>4</sup> Whether the Difference-in-Means estimator is unbiased with respect to  $\bar{\tau}$  depends solely on the known research design, i.e., whether or not there is a uniform probability distribution on the set of assignments.

Returning to the example from Table 1, recall that if we apply the estimator in Equation 3 to the possible realizations of data in Table 2, then there are 20 possible estimates that correspond to each of the 20 possible realizations of data:

$$\hat{\tau}(\mathbf{z}_1, \mathbf{y}_1) = 6.6667, \hat{\tau}(\mathbf{z}_2, \mathbf{y}_2) = 7.6667, \dots, \hat{\tau}(\mathbf{z}_{19}, \mathbf{y}_{19}) = -1.6667, \hat{\tau}(\mathbf{z}_{20}, \mathbf{y}_{20}) = -0.6667.$$

Informally, an unbiased estimator produces a guess about the estimand with no systematic error. Slightly more formally, an estimator is unbiased if the average of all possible estimates is equal to the true value of the estimand. This average of estimates, however, must be weighted by the probabilities of observing each possible estimate; we call this average the “expected value” and denote the expected value of the Difference-in-Means estimator by  $\mathbb{E}[\hat{\tau}(\mathbf{Z}, \mathbf{Y})]$ .

To calculate the expected value of the Difference-in-Means estimator to assess properties like bias and consistency, we need to know the probability associated with each of these 20 possible estimates. We know that the estimator is a function of two random quantities,  $\mathbf{Z}$  and  $\mathbf{Y}$ , but  $\mathbf{Y}$  inherits randomness only from  $\mathbf{Z}$  since  $Y_i = Z_i y_{ti} + (1 - Z_i) y_{c,i}$  for all  $i \in \{1, \dots, N\}$  units. Therefore, the probability that one would see each estimate depends only on  $\mathbf{Z}$ . So, we calculate the expected value of the estimator in general as follows:

$$\mathbb{E}[\hat{\tau}(\mathbf{Z}, \mathbf{Y})] = \hat{\tau}(\mathbf{z}_1, \mathbf{y}_1) \Pr(\mathbf{Z} = \mathbf{z}_1) + \dots + \hat{\tau}(\mathbf{z}_{|\Omega|}, \mathbf{y}_{|\Omega|}) \Pr(\mathbf{Z} = \mathbf{z}_{|\Omega|}).$$

In the context of the running example, there are 20 possible estimates corresponding to each of the  $\mathbf{z}_1, \dots, \mathbf{z}_{20}$  possible assignments, and the probability of each of those possible assignments is

$\frac{1}{20}$ . Therefore, the expected value of the Difference-in-Means estimator is:

$$\begin{aligned}\mathbb{E} \left[ \hat{\tau}(\mathbf{Z}, \mathbf{Y}) \right] &= \hat{\tau}(\mathbf{z}_1, \mathbf{y}_1) \Pr(\mathbf{Z} = \mathbf{z}_1) + \cdots + \hat{\tau}(\mathbf{z}_{20}, \mathbf{y}_{20}) \Pr(\mathbf{Z} = \mathbf{z}_{20}) \\ &= 6.6667 \left( \frac{1}{20} \right) + \cdots + -0.6667 \left( \frac{1}{20} \right) \\ &= 3.\end{aligned}$$

In this example, the expected value of the estimator,  $\mathbb{E} \left[ \hat{\tau}(\mathbf{Z}, \mathbf{Y}) \right]$ , is exactly equal to the true mean causal effect,  $\bar{\tau}$ . The estimator is unbiased given the design. If  $\Pr(\mathbf{Z} = \mathbf{z})$  did not equal  $\frac{1}{20}$  for all  $\mathbf{z}$ —i.e., if some assignments were more or less probable than others—then the Difference-in-Means estimator might not be unbiased. In general, the equality between  $\mathbb{E} \left[ \hat{\tau}(\mathbf{Z}, \mathbf{Y}) \right]$  and  $\bar{\tau}$  holds when (1) units are assigned to study conditions as individuals (not as groups, i.e., *clusters*), (2) there is always at least one unit in the treatment condition and one in the control condition and (3) each possible assignment has an identical probability of realization. In other words, the Difference-in-Means estimator is unbiased with respect to the mean additive causal effect in a uniform randomized experiment under either complete, individual assignment or simple, individual assignment so long as there is always at least one unit in each of the study conditions.

Notice that the Difference-in-Means estimator did not require large numbers of units or assumptions about the distributions of potential outcomes to be unbiased. The potential outcomes could be any set of values and the property of unbiasedness would still hold. However, the unbiasedness property did require that there be a uniform probability distribution on the set of possible assignments. But in a controlled research design, like a randomized experiment, the researcher knows whether or not this condition is true.

### 3.2 Consistency

“No systematic error” is not the same as “close to the truth.” <sup>36</sup> explains the need for another conception of “good” estimator when he writes, “[u]nbiasedness is too weak a property, since it says nothing about approximating the truth.” While we know that an unbiased estimator yields estimates that are, on average, equal to the true value of the estimand, any single estimate might be far from the truth. For example, in our running example, not one of the 20 possible estimates

is actually equal to the true mean causal effect of  $\bar{\tau} = 3$ , even though the probability-weighted average of those 20 estimates is equal to 3. Another “good” characteristic of an estimator is to produce values close to the truth, especially as the size of the experiment grows larger, while other factors remain constant, and more information is supplied to the estimator from the design.

In contrast to unbiasedness, consistency states that as the number of units in the study grows asymptotically, holding all other factors constant, the probability of an estimate within an arbitrarily small distance,  $\varepsilon$ , from the truth is equal to 1. More formally, we can define consistency as follows:

$$(4) \quad \lim_{h \rightarrow \infty} \Pr \left( \left| \hat{\tau}(\mathbf{Z}, \mathbf{Y}) - \bar{\tau} \right| < \varepsilon \right) = 1 \text{ for all } \varepsilon > 0$$

or equivalently as

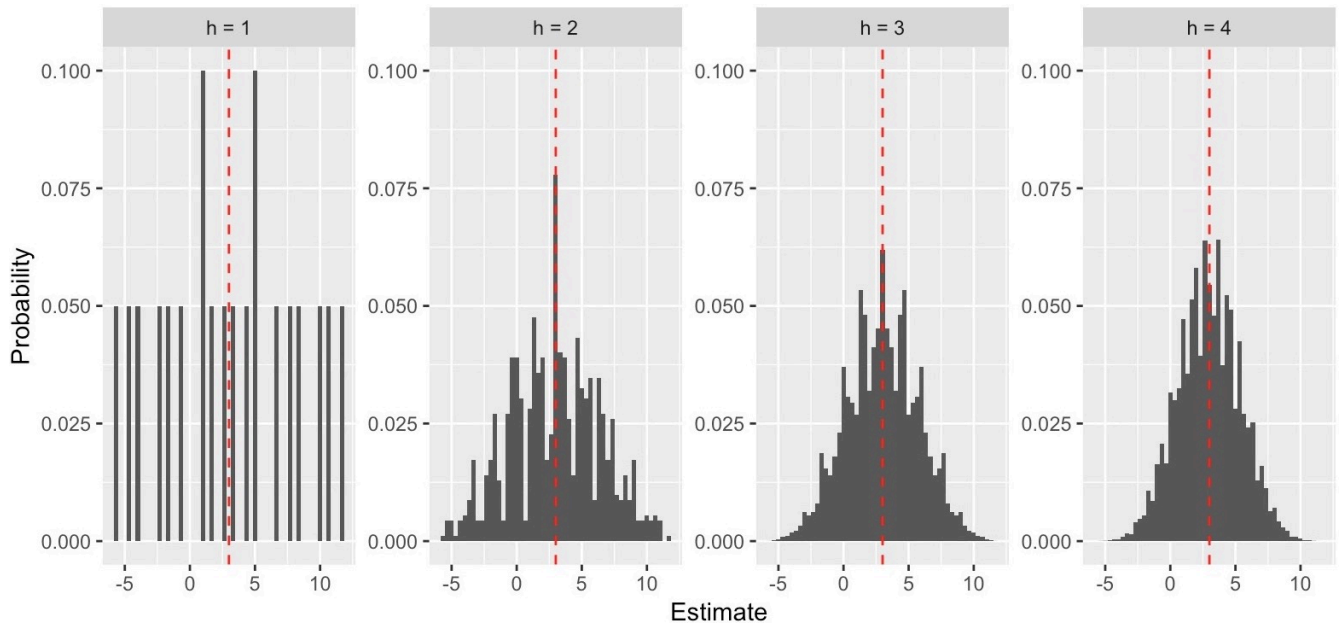
$$(5) \quad \lim_{h \rightarrow \infty} \Pr \left( \hat{\tau}(\mathbf{Z}, \mathbf{Y}) \in (\bar{\tau} - \varepsilon, \bar{\tau} + \varepsilon) \right) = 1 \text{ for all } \varepsilon > 0,$$

where, referring back to the conception of asymptotic growth in Table 3,  $h$  is the number of copies of the original finite population from Table 1.

To unpack Equations 4 and 5, Figure 1 shows what happens to the distribution of the Difference-in-Means estimator under complete random assignment as  $h \rightarrow \infty$ :

The general trend is that the probability of estimates close to the true mean causal effect,  $\bar{\tau} = 3$ , grows larger and larger and ultimately converge in probability (over the sequence of increasing finite populations) to 1. For example, following Equation (5), consider the probability that an estimate lies on the interval  $(3 - \varepsilon, 3 + \varepsilon)$  and let  $\varepsilon = 1$ . The respective probabilities of an estimate on this interval for  $h \in \{1, 2, 3, 4\}$  are 0.1, 0.2359, 0.2559 and 0.2994, and as  $h \rightarrow \infty$  that probability tends to 1. This property holds for  $\varepsilon = 1$  as well as for any positive value of  $\varepsilon$  that one could choose. For example, we could have let  $\varepsilon = 0.5$ , in which case the corresponding probabilities for  $h \in \{1, 2, 3, 4\}$  are 0.1, 0.1234, 0.1522 and 0.1643, and the limiting probability as  $h \rightarrow \infty$  is also 1.

In this particular case, consistency follows (in part) from unbiasedness. As the size of  $N$  increases towards  $\infty$ , while all other factors remain constant, the probability of an estimate arbi-



**Figure 1:** Distribution of Difference-in-Means estimator as  $h \rightarrow \infty$

trarily close to the estimator’s expected value is equal to 1. Unbiasedness ensures that the expected value of the estimator is equal to the true value of the estimand. Therefore, as the size of the study population grows towards  $\infty$ , the estimator produces a value arbitrarily close to the truth (not just to the estimator’s expected value) with probability 1. Both the unbiasedness and consistency of the Difference-in-Means estimator, moreover, arise solely from the research design. Even though the distribution of the estimator starts to look more and more Normal as the size of the study population increases to  $\infty$ , we made no such distributional assumptions to show the estimator’s unbiasedness and consistency.

### 3.3 Precision

While unbiased and consistent estimators are desirable, such estimators may yield estimates far from the truth with high probability in actual experiments with fixed study populations. One estimator is more precise than another estimator for a given study design if it produces estimates that are closer to the truth on average. In other words, a “good” estimator also has a low variance; in our case, a more precise estimator than the Difference-in-Means would make the 20 possible estimates in Table 2 closer to the true mean causal effect, on average. We now consider, first, the factors that make the Difference-in-Means estimator produce guesses with lower expected distance from the true mean causal effect and, second, the procedure one can use to conservatively estimate

the variance of the Difference-in-Means estimator.

Neyman (1923) derived an exact analytic expression for the variance of the Difference-in-Means estimator based solely on the research design of a randomized experiment as follows:

$$(6) \quad \sigma_{\hat{\tau}}^2 = \frac{1}{N-1} \left( \frac{n_t \sigma_{y_c}^2}{n_c} + \frac{(N-n_t) \sigma_{y_t}^2}{n_t} + 2\sigma_{y_c, y_t} \right),$$

where  $\sigma_{y_c}^2$  is the variance of control potential outcomes,  $\sigma_{y_t}^2$  is the variance of treated potential outcomes and  $\sigma_{y_c, y_t}$  is the covariance of control and treated potential outcomes.

Equation 6 suggests that one can increase the precision of the Difference-in-Means estimator (i.e., reduce the estimator’s variance) by increasing the number of treatment units and/or the number of control units. Precision can also increase by decreasing the variances of treatment,  $\sigma_{y_t}^2$ , and control,  $\sigma_{y_c}^2$ , potential outcomes. For a simple and clear account of the factors that increase precision, see Gerber and Green (2012, Section 3.2).

A standard design choice that researchers can make to increase precision is by blocking. That is, a researcher can, first, construct blocks that are similar in terms of covariates related to potential outcomes and, second, assign units to study conditions within blocks. Blocked assignment works by excluding assignments that yield estimates far from the true mean effect, on average. To see this point, we return to the example in Table 1 and introduce  $\mathbf{x}$ , which is a vector of baseline covariates,  $x_i$ , for all  $i \in \{1, \dots, N\}$  units. The vector  $\mathbf{x}$  is a fixed quantity that is measured for all units prior to assignment; hence,  $\mathbf{x}$  cannot change as a function of whichever assignment is realized.

$\mathbf{y}_c$	$\mathbf{y}_t$	$\tau$	$\mathbf{x}$
20	22	2	1
8	12	4	1
11	11	0	0
10	15	5	1
14	18	4	1
1	4	3	0

**Table 4:** True values of  $\mathbf{y}_c$ ,  $\mathbf{y}_t$ ,  $\tau$  and the baseline covariate  $\mathbf{x}$

From Table 4, we can see that  $\mathbf{x}$  is related to both  $\mathbf{y}_c$  and  $\mathbf{y}_t$ . The treatment potential outcomes are greater, on average, among units whose baseline covariate values are equal to 1 compared to units whose baseline covariate values are equal to 0. The same is true for control potential outcomes.

If the researcher puts units whose covariates values are equal to 1 in one block and units whose covariate values are equal to 0 in another, and then assigns half of units to treatment and control within blocks, the set of possible assignments,  $\Omega_b$ , would be as follows:

$$(7) \quad \Omega_b = \left\{ \begin{bmatrix} 1 \\ 1 \\ 1 \\ 0 \\ 0 \\ 0 \end{bmatrix}, \begin{bmatrix} 1 \\ 1 \\ 0 \\ 0 \\ 0 \\ 1 \end{bmatrix}, \dots, \begin{bmatrix} 0 \\ 0 \\ 1 \\ 1 \\ 1 \\ 0 \end{bmatrix}, \begin{bmatrix} 0 \\ 0 \\ 0 \\ 1 \\ 1 \\ 1 \end{bmatrix} \right\}.$$

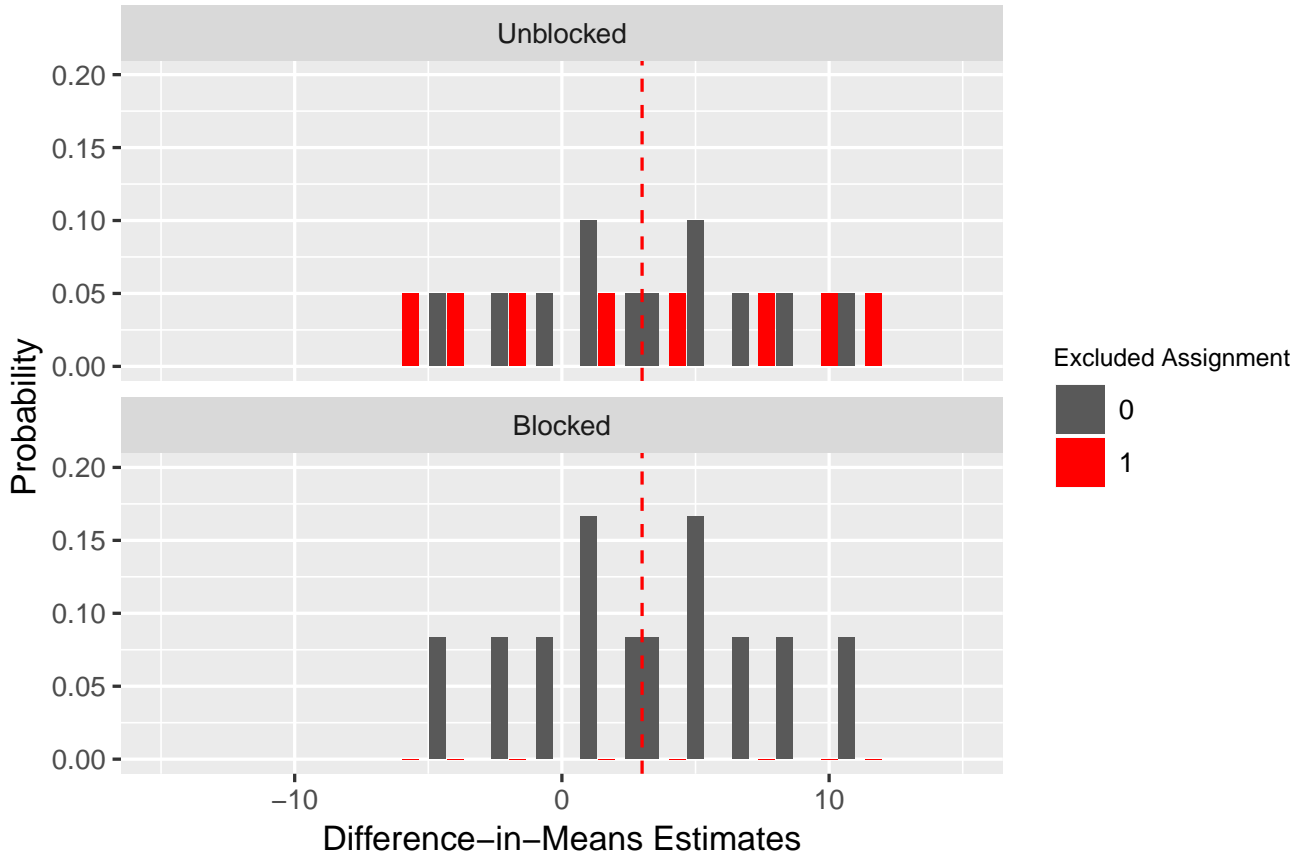
The set  $\Omega_b$  above has only 12 possible assignments as opposed to the 20 possible assignments in Equation 2 under complete random assignment without blocks. In particular, the assignments of  $\mathbf{z}_2, \mathbf{z}_3, \mathbf{z}_7, \mathbf{z}_8, \mathbf{z}_{13}, \mathbf{z}_{14}, \mathbf{z}_{18}, \mathbf{z}_{19} \in \Omega$  are excluded from  $\Omega_b$ .

Figure 2 shows the eight estimates corresponding to the eight assignments that were included in unblocked assignment, but excluded in blocked assignment. On average, these 8 estimates are farther from the true mean effect than are the other 12 estimates. More concretely, the average squared distance of the 8 excluded estimates from the truth is 36.91667 and the same average squared distance of the 12 included assignments is 18.12963. Hence, this blocked design increases precision by reducing the probability (to 0) of estimates that are, on average, far from the truth and increasing the probability of estimates that are, on average, closer to the truth.<sup>5</sup>

---

<sup>5</sup>Note that, if treatment assignment probabilities differ *across* blocks (but are uniform *within* blocks), then the standard Difference-in-Means estimator may be biased. In such cases, an unbiased estimator would be the Difference-in-Means estimator that generates an estimate within each block and subsequently weights each block-specific estimate by the proportion of units in that block.





**Figure 2:** Distribution of Difference-in-Means Estimates under (1) unblocked and (b) blocked Assignment

Thus far, we have focused on design choices that can decrease the variance of the Difference-in-Means estimator. But note that the variance of the Difference-in-Means estimator is, like the true mean causal effect, a fixed, unobservable quantity. As we can see from Equation (6), the variance of the Difference-in-Means estimator depends on the variance of treatment and control potential outcomes, as well as their covariance, none of which can be directly observed. We need a “good” procedure by which we can estimate the variance of the Difference-in-Means estimator if we are to reasonably infer the precision of the Difference-in-Means estimator. We will use such a variance estimator in the context of not only evaluating estimators, but also hypothesis testing.

One can unbiasedly estimate two of the three unknown quantities in Equation 6. Following Cochran (1977, Theorem 2.4), unbiased estimators of  $\sigma_{y_c}^2$  and  $\sigma_{y_t}^2$ , respectively, are:  $\hat{\sigma}_{y_c}^2 = \left(\frac{n-1}{n(n_c-1)}\right) \sum_{i:Z_i=0}^n (y_{ci} - \hat{\mu}_{y_c})^2$  and  $\hat{\sigma}_{y_t}^2 = \left(\frac{n-1}{n(n_t-1)}\right) \sum_{i:Z_i=1}^n (y_{ti} - \hat{\mu}_{y_t})^2$ , where  $n_c = n - n_t$ . We cannot write an unbiased estimator for  $\sigma_{y_c, y_t}$  since no two potential outcomes for any unit can be jointly

observed. [Neyman \(1923\)](#) noted, however, that one could use a conservative procedure for estimating the quantity in Equation 6 by assuming the largest possible value of  $2\sigma_{y_c, y_t}$ , which, by the Cauchy-Schwarz inequality and the AM-GM inequality (i.e., inequality of arithmetic and geometric means), is  $\sigma_{y_c}^2 + \sigma_{y_t}^2$ .

After substituting  $\sigma_{y_c}^2 + \sigma_{y_t}^2$  for  $2\sigma_{y_c, y_t}$ , the analytic expression for the variance of the Difference-in-Means estimator (assuming  $2\sigma_{y_c, y_t} = \sigma_{y_c}^2 + \sigma_{y_t}^2$ ) is

$$\frac{1}{N-1} \left( \frac{n_t \sigma_{y_c}^2}{(N-n_t)} + \frac{(N-n_t) \sigma_{y_t}^2}{n_t} + \sigma_{y_c}^2 + \sigma_{y_t}^2 \right),$$

which can be simplified to

$$(8) \quad \frac{N}{N-1} \left( \frac{\sigma_{y_c}^2}{n_c} + \frac{\sigma_{y_t}^2}{n_t} \right).$$

Now there are only two unknown quantities in Equation 8, each of which can be unbiasedly estimated. Hence, one can now unbiasedly estimate the quantity in 8 via the conservative variance estimator of

$$(9) \quad \hat{\sigma}_{\hat{\tau}}^2 = \frac{N}{N-1} \left( \frac{\hat{\sigma}_{y_c}^2}{n_c} + \frac{\hat{\sigma}_{y_t}^2}{n_t} \right).$$

This estimator is conservative because, since it unbiasedly estimates the quantity in (8), its expected value is equal to or greater than the true variance of the estimator given in (6).<sup>6</sup>

Thus far, we have explained the role that research design—i.e., the probability distribution on the set of assignments,  $\Omega$ —plays in determining whether estimators are unbiased, consistent and precise. We have also explained how one can infer the variance of an estimator via a conservative procedure. We have used a simple example of complete, uniform assignment to illustrate these points; yet an estimator that is unbiased, consistent and/or relatively precise in this design may not be so in another design. For example, the Difference-in-Means estimator is not necessarily unbiased

---

<sup>6</sup>Recent work has derived a consistent estimator for an upper bound on the term  $2\sigma_{y_c, y_t}$  that is always less than or equal to  $\sigma_{y_c}^2 + \sigma_{y_t}^2$  ([Aronow et al., 2014](#)). Such an estimator enables researchers to more precisely estimate the the variance of the Difference-in-Means estimator and, as we will discuss later, increase the power of hypothesis tests about the mean causal effect.

when there is a non-uniform probability distribution on  $\Omega$ ; however, the Horvitz-Thompson (i.e., inverse probability weighted) estimator (Horvitz and Thompson, 1952) is unbiased in such a design (see Aronow and Middleton, 2013). Such design-based inference differs from model-based inferences in that the former remains reliable without the need to impose a probability model on potential outcomes or to model the functional form (e.g., a linear model) that links the treatment variable to potential outcomes. The only probability model in design-based inference is the assignment process itself, which, in the case of a randomized experiment, is known to be the true model of the data generating process.

## 4 Hypothesis Testing

Focusing on  $\bar{\tau}$  engages with the fundamental problem of causal inference by aggregation and  $\hat{\bar{\tau}}$  can be shown to be an unbiased, consistent and relatively precise estimator. An alternative approach begins with claims about causal effects rather than guesses about them, stating hypotheses about  $\bar{\tau}$  or even about individual effects,  $\tau_i$ . This approach, then tests those claims, and so the procedures that we will assess in this section are properties of tests rather than of estimators. A researcher who performs a hypothesis test for causal inference first states a null hypothesis about a relationship between unobserved potential outcomes and a (usually composite) alternative hypothesis, which we will define more precisely below. The researcher can then assess the probability that the research design generates data more extreme than the observed data under the null hypothesis, where the data are summarized via a test statistic that maps the data — the observed outcome, perhaps adjusted to reflect the implications of the hypothesis, and treatment assignment — to a single number. For example,  $t(Z_i, Y_i)$  could be  $\left(\frac{1}{\sum_{i=1}^N Z_i}\right) \sum_{i=1}^N Z_i Y_i - \left(\frac{1}{\sum_{i=1}^N (1-Z_i)}\right) \sum_{i=1}^N (1-Z_i) Y_i$ , which is the same formula that we labeled  $\hat{\bar{\tau}}$  in Equation (3) but, in the context of hypothesis testing, is not an estimator, but a data summary. We call the probability of a test statistic more extreme than the observed test statistic a probability-value or  $p$ -value. Typically, when the  $p$ -value is lower than a pre-specified “significance level” of the test, which we denote by  $\alpha \in (0, 1)$ , a researcher *rejects* the null hypothesis — meaning that the researcher declares that the observed data are not consistent with the hypothesis. When the  $p$ -value is greater than or equal to the significance

level of the test, the researcher *fails to reject* the null hypothesis — meaning that the researcher declares that there is not enough information to state that the observed data are inconsistent with the hypothesized state of the world. Sometimes, we talk about hypothesis testing as an attempt to distinguish signal from noise—a high  $p$ -value tells us that we cannot distinguish signal from noise, and a low  $p$ -values tells us that we can do so.

Hypothesis tests are subject to at least two types of errors. One could, first, reject the null hypothesis when it is true (a type I error) or, second, fail to reject the null hypothesis when it is false (a type II error). Two features of hypothesis tests related to these two potential errors are the  $\alpha$  size of the test and the *power* of the test. We now define the  $\alpha$  size (as distinct from the  $\alpha$  level) and power of hypothesis tests.

A test's  $\alpha$  level is, in the words of [Rosenbaum \(2010, Glossary\)](#), that test's “promise” that the probability of a Type I error (i.e., the probability of a  $p$ -value less than  $\alpha$  when the null hypothesis is true) is less than or equal to the  $\alpha$  level. The test's  $\alpha$  size, on the other hand, is the test's true probability of a Type I error, which, in general, can be greater than, equal to or less than the  $\alpha$  level “promised” by the test. In contrast to the  $\alpha$  level and size of a test, a test's power is the probability of a  $p$ -value less than the  $\alpha$  level when the null hypothesis is false. In other words, power is 1 minus the Type II error probability; hence, as the power of a test increases, the Type II error probability decreases.

In the subsections to follow, we first define “good” properties of hypothesis tests. We then describe tests of causal hypotheses in two distinct traditions traceable to [Fisher \(1935\)](#) (and subsequently developed most extensively by [Rosenbaum, 2002, 2010](#)), and [Neyman and Pearson \(1933\)](#). We then explain the role that research design plays in justifying whether tests in either of these two traditions have good properties.

#### 4.1 What Makes a Hypothesis Test a “Good” Test?

We have already discussed three properties of good estimators—namely, unbiasedness, consistency and precision—but what makes a test of a null hypothesis relative to an alternative hypothesis a “good” test? Just as we did for estimation, we describe “good” features of hypothesis tests in the context of a fixed, finite population and in the context of a hypothetical scenario in which a

study’s size increases towards  $\infty$  by increasing the number of copies of the study. The first two “good” properties (a Type I error probability less than the  $\alpha$  level and an unbiased test) refer to the former context and the third property (a consistent test) refers to the latter context. Informally, a good hypothesis test should mislead us rarely — it should rarely encourage us to declare that we have discovered a signal in the noise when no signal exists, and it should often find such signals when they do exist.

Regardless of the size of a given study population, a hypothesis test ought to, first, control its  $\alpha$  size (true Type I error probability) such that the test’s  $\alpha$  size is less than or equal to the test’s  $\alpha$  level. Second, a hypothesis test ought to be an *unbiased test* (not to be confused with an *unbiased estimator*); that is, the probability of rejecting the null hypothesis when it is false and the alternative hypothesis is true should be at least as great as the probability of rejecting the null hypothesis when it is true and the alternative hypothesis is false (see [Lehmann and Romano, 2005](#), Chapter 4). Intuitively, we want to *reject* something that is *false* and we want to *not reject* something that is *true*. A test that leads us to reject true nulls with greater probability than we reject false nulls does not yield inferences that track the true causal effect. A good hypothesis test ought to be an unbiased test in this sense.

Turning now to the asymptotic context described in Section 2.3, a hypothesis test ought to be a *consistent test* (not the same as a *consistent estimator*); that is, as the size of the study population increases asymptotically while all other relevant factors remain constant, the probability of rejecting the null hypothesis when it is false and the alternative is true should tend to 1 (see [Lehmann and Romano, 2005](#), Chapter 11). We now show the role that research design plays in enabling “good” hypothesis tests in two different design-based traditions, Fisherian ([Fisher, 1935](#)) and Neymanian ([Neyman and Pearson, 1933](#)).

## 4.2 Fisherian Hypothesis Testing

Hypothesis testing in the tradition of [Fisher \(1935\)](#), and later developed most extensively by [Rosenbaum \(2002, 2010\)](#), assesses the consistency of the observed data with a null hypothesis vis-à-vis an alternative hypothesis. A strong null hypothesis,<sup>7</sup> which we denote by  $\tau_0$ , pos-

---

<sup>7</sup>We use the term “strong” instead of “sharp,” as used by [Fisher \(1935\)](#), to contrast *strong* null hypotheses with *weak* null hypotheses which we discuss below.

tulates an individual treatment effect for all  $i \in \{1, \dots, N\}$  units in a given study population. For example, one strong null hypothesis is  $\boldsymbol{\tau}'_0 = [5 \ 5 \ \dots \ 5 \ 5]$ , another might be  $\boldsymbol{\tau}'_0 = [0.5 \ -10 \ \dots \ 200 \ -74.25]$ . The *strong null hypothesis of no effect* (which we henceforth refer to as the strong null of no effect) specifically postulates that  $\tau_i = 0$  for all  $i \in \{1, \dots, N\}$  units—i.e., that  $\boldsymbol{\tau}'_0 = [0 \ 0 \ \dots \ 0 \ 0]$ .

The consistency of the observed data with a strong null hypothesis vis-à-vis an alternative is typically assessed via  $p$ -values. To reiterate, a  $p$ -value is the probability of a test statistic at least as extreme as the observed test statistic from the perspective of the null hypothesis: As we will show below, the hypothetical world of the null generates, along with the known research design, the probability distribution against which we compare our single observed test statistic. In the context of Fisherian hypothesis tests, we can formally represent upper ( $p_u$ ), lower ( $p_l$ ), and two-sided ( $p_t$ ),  $p$ -values as follows:

$$\begin{aligned}
 (10) \quad p_u &= \sum_{j=1}^{|\Omega|} \mathbb{1} \left[ t(\mathbf{z}_j, \mathbf{y}_{0_j}) \geq T \right] \Pr(\mathbf{Z} = \mathbf{z}_j) \\
 p_l &= \sum_{j=1}^{|\Omega|} \mathbb{1} \left[ t(\mathbf{z}_j, \mathbf{y}_{0_j}) \leq T \right] \Pr(\mathbf{Z} = \mathbf{z}_j) \\
 p_t &= \min \{1, 2 \min \{p_u, p_l\}\},
 \end{aligned}$$

where the index  $j \in \{1, \dots, |\Omega|\}$  runs over all possible assignments in the set of assignments  $\Omega$ ,  $\mathbb{1}$  is an indicator function that is 1 if the argument  $[\cdot]$  is true and 0 if false,  $t(\mathbf{z}_j, \mathbf{y}_{0_j})$  is the null test statistic (using  $\mathbf{y}_{0_j}$  to refer to the vector of observed outcomes for the  $j$ th assignment implied by the null hypothesis,  $H_0$ ), and  $T$  is the observed test statistic.<sup>8</sup>

To provide an illustration of Fisherian  $p$ -values, we return to the example in Table 1 and imagine that the assignment  $\mathbf{z}'_s = [1 \ 0 \ 0 \ 1 \ 1 \ 0]$  happened to be the one randomly selected. In this case, the realization of data would be as follows:

If the researcher uses the Difference-in-Means test statistic to provide a single, numerical sum-

---

<sup>8</sup>The expression for a two-sided  $p$ -value,  $p_t = \min \{1, 2 \min \{p_u, p_l\}\}$ , comes from Rosenbaum (2010, 33) who states that ‘[i]n general, if you want a two-sided P-value, compute both one-sided P-values, double the smaller one, and take the minimum of this value and 1.’ The rationale is that doubling a one-sided  $p$ -value compensated for, in essence, testing twice.

$\mathbf{z}_8$	$\mathbf{y}_c$	$\mathbf{y}_t$	$\mathbf{y}_8$
1	?	22	22
0	8	?	8
0	11	?	11
1	?	15	15
1	?	18	18
0	1	?	1

**Table 5:** Realization of Data if  $\mathbf{z}_8$  were the randomly drawn assignment

mary of the observed data in Table 5, then the observed test statistic is  $t(\mathbf{z}, \mathbf{y}) = 11.6667$ . Let’s assume that the researcher wants to assess the consistency of this observed test statistic with the strong null of no effect—i.e., that  $\tau_i = 0$  for all  $i$ —relative to the alternative hypothesis of a positive effect—i.e., that  $\tau_i$  is nonnegative for all  $i$  and positive for at least one  $i$ .

Potential outcomes are only partially observed, but the researcher can “fill in” the missing potential outcomes following the strong null hypothesis  $H_0 : y_{t,i} = y_{c,i}$  for all  $i$ . Below, we can also show what this hypothesis implies for the observed outcomes  $y_i$  recalling that  $y_i = Z_i y_{t,i} + (1 - Z_i) y_{c,i}$  and writing  $y_{c0,i}$  to mean “value of  $y_{c,i}$  under  $H_0$ ”:

$$(11) \quad \begin{aligned} y_{c0,i} &= y_i - z_i \tau_{0i} \\ y_{t0,i} &= y_i + (1 - z_i) \tau_{0i}, \end{aligned}$$

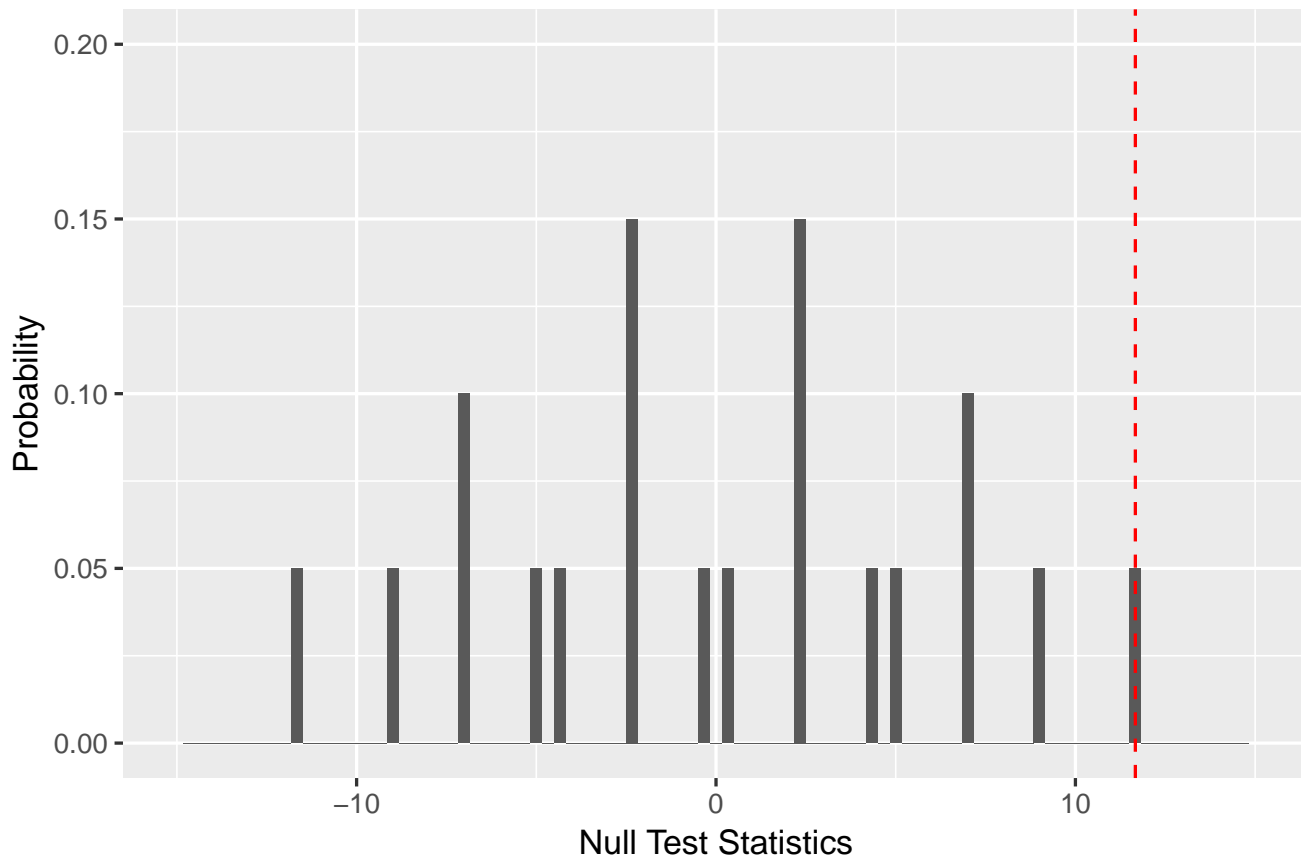
which in the case of the strong null of no effect implies that units’ null potential outcomes and observed outcomes are as in Table 6 below:

$\mathbf{z}_8$	$\mathbf{y}_{c08}$	$\mathbf{y}_{t08}$	$\mathbf{y}_8$
1	22	22	22
0	8	8	8
0	11	11	11
1	15	15	15
1	18	18	18
0	1	1	1

**Table 6:** Null potential outcomes if  $\mathbf{z}_8$  were the realized assignment and under the hypothesis that  $y_{t,i} = y_{c,i}$  for all  $i$

Considering the strong null of no effect as a model of unobserved causal relationships for the sake of argument, the researcher knows exactly what all other possible realizations of data would look like under each possible assignment in  $\Omega$ . Hence, the researcher can summarize all other possible realizations of data under the null with the same Difference-in-Means test statistic, generating a probability distribution of those null test statistics, and then calculate the probability of a null

test statistic greater than or equal to the observed test statistic of 11.6667. (See Figure 3 below.)



**Figure 3:** Distribution of Difference-in-Means test statistic under strong null of no effect when  $\mathbf{z}_8$  is the realized assignment

In this case, the upper  $p$ -value, ( $p_u$  from Equation 10), is 0.05. If the value of  $\alpha$  had been pre-set to a level greater than 0.05, then the researcher would reject the null. Table 1 shows that the strong null of no effect is false ( $y_{t,i} \neq y_{c,i}$  for all  $i$ ) and the alternative of a positive effect is true; hence, this particular choice to reject the strong null of no effect relative to the alternative of a positive effect happened to be a good one.

To assess whether the hypothesis testing procedure is a good one overall, we want to establish that the test (1) has a true Type I error probability less than the  $\alpha$  level, (2) is unbiased and (3) is consistent. We now illustrate these three properties in turn.

A key concept in evaluating tests is the idea of the size of the test. That is, if the level of the test is the hoped for amount of error, the size is the actual amount of error (Rosenbaum, 2010, Glossary). When the researcher tests a strong null hypothesis that is true, the  $\alpha$  size of the



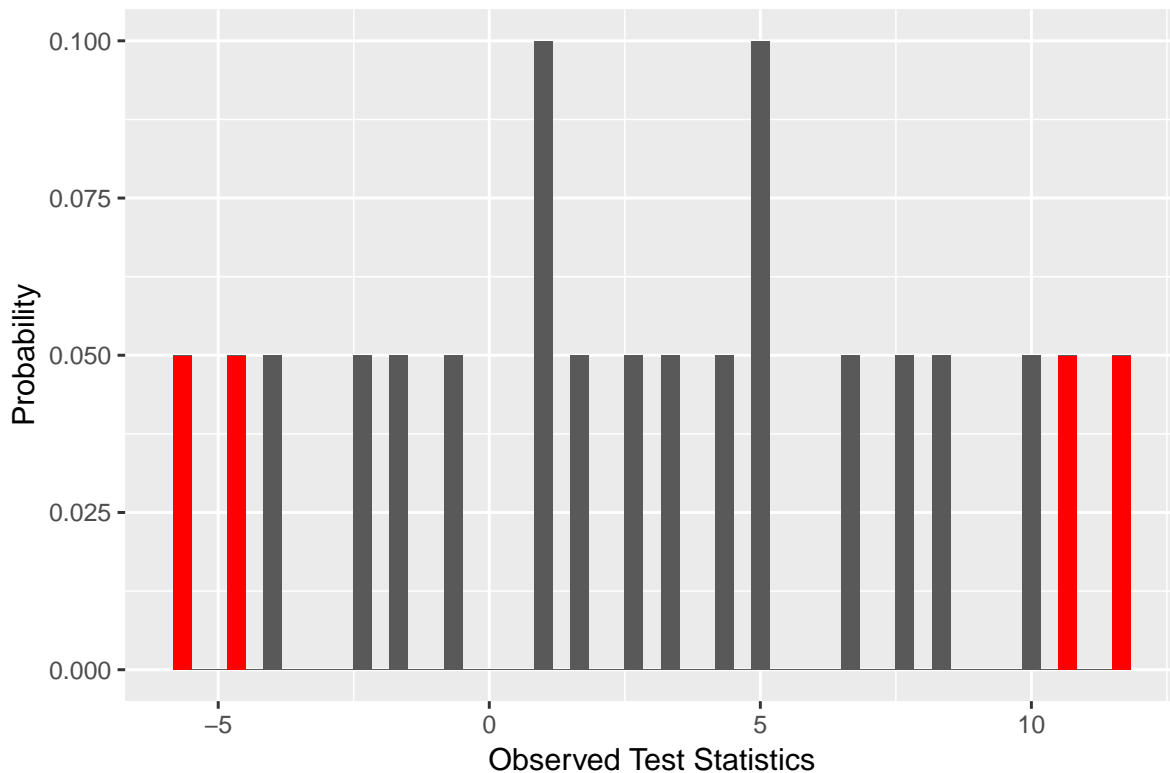
test should always be equal to or less than the test’s  $\alpha$  level. In other words, the Type I error probability (or size) should fulfill that which was “promised” by the test’s significance level.

To see this point, we return to the example given in Table 1, where the true vector of individual causal effects is  $\boldsymbol{\tau}' = [2 \ 4 \ 0 \ 5 \ 4 \ 3]$ , and Table 2 describes the 20 possible realizations of data. When the null hypothesis,  $\boldsymbol{\tau}_0$ , is false, the potential outcomes implied by the null hypothesis will vary depending on which data are realized. However, when the null hypothesis,  $\boldsymbol{\tau}_0$ , is true, i.e., when  $\boldsymbol{\tau}_0 = \boldsymbol{\tau}$ , the potential outcomes implied by the null hypothesis are fixed across all possible realizations of data as shown by Table 7:

$\mathbf{z}_1$	$\mathbf{y}_{c0}$	$\mathbf{y}_{t0}$	$\mathbf{y}_1$	$\mathbf{z}_2$	$\mathbf{y}_{c0}$	$\mathbf{y}_{t0}$	$\mathbf{y}_2$	...	$\mathbf{z}_{19}$	$\mathbf{y}_{c0}$	$\mathbf{y}_{t0}$	$\mathbf{y}_{19}$	$\mathbf{z}_{20}$	$\mathbf{y}_{c0}$	$\mathbf{y}_{t0}$	$\mathbf{y}_{20}$
1	20	22	22	1	20	22	22		0	20	22	20	0	20	22	20
1	8	12	12	1	8	12	12		0	8	12	8	0	8	12	8
1	11	11	11	0	11	11	11	...	1	11	11	11	0	11	11	11
0	10	15	10	1	10	15	15		0	10	15	10	1	10	15	15
0	14	18	14	0	14	18	14		1	14	18	18	1	14	18	18
0	1	4	1	0	1	4	1		1	1	4	4	1	1	4	4

**Table 7:** Null Potential Outcomes for all Possible Realizations of Data when the Null Hypothesis is True. The observed outcomes column,  $\mathbf{y}_j = \mathbf{y}_{c0} + \mathbf{z}_j\boldsymbol{\tau}'$ .

If we set the significance level of the test to  $\alpha = 0.10$ , then it is true by definition that the probability an observed test statistic lies in the lower tail of its distribution is less than or equal to 0.10, and the same is true for the probability that an observed test statistic lies in the upper tail of its distribution. Figure 4 shows the distribution of this test statistic in which both the lower and upper tails according to  $\alpha = 0.10$  are shaded in red.



**Figure 4:** Distribution of Observed Test Statistic. The upper and lower tails containing 10% of the area under the curve are shaded in red.

Notice that, no matter which of the 20 possible realizations of data actually is realized, the null potential outcomes are identical to the true potential outcomes in Table 1 and, hence, each of the 20 possible null distributions of the test statistic are identical to the distribution of the observed test statistic. Since all 20 possible null distributions are identical to the distribution of the observed test statistic, the probability that an observed test statistic lies in one of the tails of the null distribution must also be less than or equal to 0.10. This means that the probability of a null test statistic more extreme than the observed test statistic must be less than or equal to  $\alpha$ .

Fisherian hypothesis tests possess the property that the Type I error probability is less than or equal to the test's  $\alpha$  level. Yet such a property does *not* imply that the power of the test is greater than the test's Type I error probability (i.e., that the test is unbiased). To show that a test is unbiased relative to a specific class of alternative hypotheses, we now need to more precisely define the alternative to the null hypothesis.

In the sections above, we have referred to causal effects by the vector  $\boldsymbol{\tau}$ , but now we define the size of causal effects using the vectors of control and treatment potential outcomes,  $\mathbf{y}_c$  and  $\mathbf{y}_t$ .

More specifically, following Rosenbaum (2002), we define a treatment effect that is “larger” than another treatment effect as follows:

**Definition 1.** *One treatment effect,  $(\mathbf{y}_c^*, \mathbf{y}_t^*)$ , has a larger effect than another treatment effect,  $(\mathbf{y}_c, \mathbf{y}_t)$ , if and only if  $y_{t,i}^* \geq y_{t,i}$  and  $y_{c_i}^* \leq y_{c_i}$  for all  $i \in \{1, \dots, N\}$  units, where  $\mathbf{y}_t^* \neq \mathbf{y}_t$  or  $\mathbf{y}_c^* \neq \mathbf{y}_c$ .*

Such an ordering of causal effects is consistent with many models of treatment effects, such as additive, multiplicative, tobit, and dilated effects (Rosenbaum, 1999, 2002, 2010), not solely with the model of a constant, additive effect that we use here.

We now use this formal definition of “larger effect” in terms of potential outcomes to define a desirable property of a test statistic: “Larger” effects yield test statistic values greater than those produced by “smaller” effects. Rosenbaum (2002, Chapter 2, Lemma 6) shows that an “effect increasing” test statistic with respect to two possible realizations of data,  $(\mathbf{z}, \mathbf{y})$  and  $(\mathbf{z}, \mathbf{y}^*)$ , satisfies this property. Following Rosenbaum (2002), we define an effect increasing test statistic as follows:

**Definition 2.** *A test statistic  $t(\cdot, \cdot)$  is effect increasing when  $t(\mathbf{z}, \mathbf{y}) \leq t(\mathbf{z}, \mathbf{y}^*)$  whenever  $y_i \leq y_i^*$  for all  $i \in \{1, \dots, N\} : z_i = 1$  and  $y_i^* \leq y_i$  for all  $i \in \{1, \dots, N\} : z_i = 0$ .*

An effect increasing test statistic ensures that, when the null hypothesis is false and the alternative of a larger effect is true, each possible realization of data yields a test statistic value that is greater than or equal to the corresponding test statistic value when the null hypothesis is true and the alternative of a larger effect is false. To understand this property, consider the following example in Table 8 of two possible causal effects,  $(\mathbf{y}_c, \mathbf{y}_t)$  and  $(\mathbf{y}_c^*, \mathbf{y}_t^*)$ , in which the former is a null effect and the latter is a larger, positive causal effect (see Definition 1):

$\mathbf{y}_c$	$\mathbf{y}_t$	$\mathbf{y}_c^*$	$\mathbf{y}_t^*$
20	20	20	22
8	8	3	12
11	11	10	11
10	10	10	15
14	14	9	19
1	1	1	4

**Table 8:** No Causal Effect (left table,  $(\mathbf{y}_c, \mathbf{y}_t)$ ,) and a Positive Causal Effect (right table,  $(\mathbf{y}_c^*, \mathbf{y}_t^*)$ ).

Regardless of whichever 3 out of the 6 units are assigned to treatment in  $\mathbf{z}$ , the larger causal effect,  $(\mathbf{y}_c^*, \mathbf{y}_t^*)$ , will always yield a greater value of the observed outcome  $y_i$  for all treated units than the outcome we would see with the same  $\mathbf{z}$  in the no causal effect state; and we would also see a smaller value of  $y_i$  for all control units than what we would see in the no causal effect state. Table 9 shows that an effect increasing test statistic ensures that the observed test statistic of a larger effect is always greater than or equal to the observed test statistic of a smaller effect.

$\mathbf{Z}$	No Effects $t(\mathbf{Z}, \mathbf{Y})$	Positive Effects $t(\mathbf{Z}, \mathbf{Y}^*)$
$\mathbf{z}_1$	4.67	8.33
$\mathbf{z}_2$	4.00	9.67
$\mathbf{z}_3$	6.67	10.67
$\mathbf{z}_4$	-2.00	3.00
$\mathbf{z}_5$	6.00	11.67
$\mathbf{z}_6$	8.67	12.67
$\mathbf{z}_7$	0.00	5.00
$\mathbf{z}_8$	8.00	14.00
$\mathbf{z}_9$	-0.67	6.33
$\mathbf{z}_{10}$	2.00	7.33
$\mathbf{z}_{11}$	-2.00	2.67
$\mathbf{z}_{12}$	0.67	3.67
$\mathbf{z}_{13}$	-8.00	-4.00
$\mathbf{z}_{14}$	0.00	5.00
$\mathbf{z}_{15}$	-8.67	-2.67
$\mathbf{z}_{16}$	-6.00	-1.67
$\mathbf{z}_{17}$	2.00	7.00
$\mathbf{z}_{18}$	-6.67	-0.67
$\mathbf{z}_{19}$	-4.00	0.33
$\mathbf{z}_{20}$	-4.67	1.67

**Table 9:** Observed mean-difference test statistics under all possible assignments for both a no-effects and positive-effects true causal effect.

In addition to the property that increasing causal effects map monotonically onto increasing

test statistics, we also want the  $p$ -values for tests of the null hypothesis when the null is false and the alternative is true to be smaller compared to the  $p$ -values when the null is true and the alternative is false. An effect increasing test statistic also suffices for this property (For a formal proof, see [Rosenbaum, 2002](#), Chapter 2, Proposition 5.)

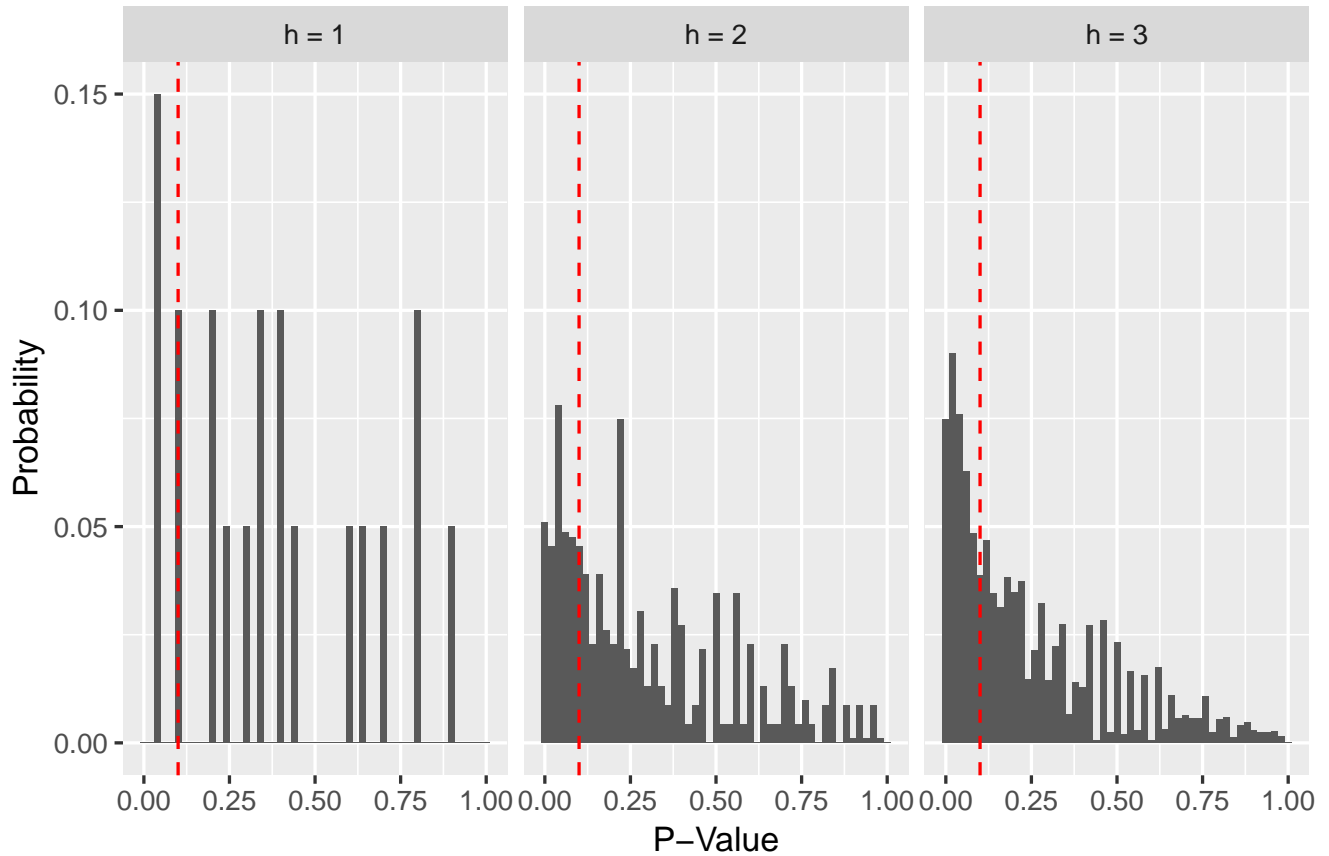
Table 10 shows that, for all  $\mathbf{z}_1, \dots, \mathbf{z}_{20}$ , the  $p$ -value of the strong null of no effect when the positive effect,  $(\mathbf{y}_c^*, \mathbf{y}_t^*)$ , is true is less than or equal to the strong null's  $p$ -value when no-effect,  $(\mathbf{y}_c, \mathbf{y}_t)$ , is true:

$\mathbf{Z}$	No-Effects P-Value for $(\mathbf{y}_c, \mathbf{y}_t)$	Positive Effects P-Value for $(\mathbf{y}_c^*, \mathbf{y}_t^*)$
$\mathbf{z}_1$	0.25	0.20
$\mathbf{z}_2$	0.30	0.10
$\mathbf{z}_3$	0.15	0.05
$\mathbf{z}_4$	0.70	0.40
$\mathbf{z}_5$	0.20	0.10
$\mathbf{z}_6$	0.05	0.05
$\mathbf{z}_7$	0.55	0.40
$\mathbf{z}_8$	0.10	0.05
$\mathbf{z}_9$	0.60	0.35
$\mathbf{z}_{10}$	0.40	0.20
$\mathbf{z}_{11}$	0.70	0.45
$\mathbf{z}_{12}$	0.45	0.30
$\mathbf{z}_{13}$	0.95	0.90
$\mathbf{z}_{14}$	0.55	0.30
$\mathbf{z}_{15}$	1.00	0.80
$\mathbf{z}_{16}$	0.85	0.65
$\mathbf{z}_{17}$	0.40	0.25
$\mathbf{z}_{18}$	0.90	0.80
$\mathbf{z}_{19}$	0.75	0.65
$\mathbf{z}_{20}$	0.80	0.60

**Table 10:** Comparing  $p$ -values for tests of the strong null of no effect when no effects are true  $((\mathbf{y}_c, \mathbf{y}_t))$  and false  $((\mathbf{y}_c^*, \mathbf{y}_t^*))$ .

We have just shown that Fisherian tests using effect increasing test statistics are unbiased: They provide more evidence against false claims than against true claims. In addition to being unbiased, Fisherian tests are consistent. That is, as the size of the study population grows toward  $\infty$ , while other factors remain constant, the power of the test tends to 1. Returning to the example in Table 1, notice that the probability of a  $p$ -value less than  $\alpha$  increases along the sequence of finite

populations of increasing size given in Table 3. For example, with an  $\alpha$  level of  $\alpha = 0.10$ , Figure 5 shows that the probability of a  $p$ -value less than  $\alpha = 0.10$  grows greater and greater:



**Figure 5:** Distribution of Fisherian  $p$ -values for test of strong null of no effect under all realizations of data as the size of the experimental pool grows from one copy of the study,  $h = 1$  ( $n = 6$ ); to two copies of the study,  $h = 2$  ( $n = 12$ ); to three copies of the study,  $h = 3$  ( $n = 18$ ).

Figure 5 shows that when the study population size grows from 6 to 12 to 18 units, while all other factors are held constant, the power increases from 0.15 to roughly 0.2933 and 0.3741, respectively. In general, as the population size increases along the sequence given in Table 3, the power of the test of the strong null of no effect will tend to 1. In other words, as we draw upon more and more data, we *reject a false* null with greater and greater probability.

We used the Difference-in-Means as a test statistic for our example data to demonstrate the properties of Fisherian hypothesis tests in randomized experiments. We could have shown the same properties using rank-based test statistics or standardized mean differences or many other functions of treatment and outcomes. All of those test statistics have the effect increasing property, and so would yield unbiased and consistent tests of strong null hypotheses. In fact, one of the

challenges of Fisherian testing is its flexibility in terms of test statistics. We do not engage with those decisions here, but encourage interested readers to see (Rosenbaum, 2010, Chapter 2) for some discussion, as well as ? and Bowers et al. (2013); ? for examples of using this flexibility to assess substantively interesting hypotheses about pareto optimal causal effects or the propagation of causal effects on networks, respectively.

### 4.3 Hypothesis Testing in the Neymanian Tradition

While Fisherian tests allow the use of many different kinds of test statistics, Neymanian hypothesis tests are much more closely related to the estimation of mean causal effects (see Section 3). In any given study, one can observe only a single *estimate*. However, a researcher can postulate (provisionally for the sake of argument) a weak null hypothesis,  $H_0 : \bar{\tau} = \bar{\tau}_0$ , relative to some alternative hypothesis, such as  $H_a : \bar{\tau} > \bar{\tau}_0$  and subsequently assess the probability of an estimate more extreme than the estimate the researcher actually observed if the weak null hypothesis were true. Such Neymanian hypothesis tests differ from Fisherian tests in several fundamental ways: Neymanian hypothesis tests (1) hypothesize about the mean causal effect, not the individual causal effect for each unit in the study; (2) require that the Difference-in-Means estimator be unbiased such that a hypothesis about the mean causal effect implies the same value for the mean (i.e., expected value) of the estimator; (3) require that the researcher estimate the variance of the Difference-in-Means estimator (recall that the distribution of the test statistic in the Fisherian test is known under random assignment); and (4) make the assumption that the Difference-in-Means estimator is well approximated by a Normal distribution and, hence, that the estimator rescaled via the  $z$ -score is well approximated by a standard Normal distribution.

We can see points (1) – (4) by looking at the common expressions for Neymanian  $p$ -values:

$$(12) \quad \begin{aligned} p_u &= 1 - \Phi \left( \frac{\hat{\tau} - \bar{\tau}_0}{\sqrt{\hat{\sigma}_{\hat{\tau}}^2}} \right) \\ p_l &= \Phi \left( \frac{\hat{\tau} - \bar{\tau}_0}{\sqrt{\hat{\sigma}_{\hat{\tau}}^2}} \right) \\ p_t &= 2 \left( 1 - \Phi \left( \frac{|\hat{\tau} - \bar{\tau}_0|}{\sqrt{\hat{\sigma}_{\hat{\tau}}^2}} \right) \right), \end{aligned}$$

where  $\hat{\tau}$  is the familiar Difference-in-Means estimator (from Section 3) now interpreted as a test statistic, not an estimator,  $\hat{\sigma}_{\hat{\tau}}^2$  is the conservative variance estimator (also from Section 3) now used to describe the reference distribution for a null hypothesis rather than the precision of an estimator,  $\bar{\tau}_0$  is a weak null hypothesis and  $\Phi(\cdot)$  is the standard Normal CDF.

The expressions in Equation 12 return the probability of an estimate at least as extreme as the observed estimate if the weak null hypothesis,  $\bar{\tau}_0$ , were true. Notice, though, that Equation 12 assumes that the Difference-in-Means test statistic is well approximated by a Normal distribution and, hence, that the  $z$ -score is well approximated by a standard Normal distribution. In addition, the weak null hypothesis,  $\bar{\tau}_0$ , is technically a claim about the mean of the Difference-in-Means test statistic, which we can denote by  $\mathbb{E}_0 \left[ \hat{\tau} \right]$ . However, because the Difference-in-Means estimator is an unbiased estimator, its expected value is always equal to the mean causal effect; this is why we use  $\bar{\tau}_0$  in Equation (12) rather than  $\mathbb{E}_0 \left[ \hat{\tau} \right]$ . Finally, note that the true variance of the Difference-in-Means test statistic is unknown, but the Normal CDF requires two arguments—a value for the mean and a value for the variance—to assign a probability to estimates as least as extreme as the observed estimate. Rather than postulate a hypothesis about the variance of the estimator, like one does for the mean of the estimator, Neymanian tests “plug in” an estimate from the conservative variance estimator to calculate a  $p$ -value via the standard Normal CDF.

In many situations, we can easily justify the assumption that the null hypothesis reference distribution using the Difference-in-Means test statistic is well approximated by a Normal distribution.



The finite population Central Limit Theorem (CLT) (see [Li and Ding, 2017](#)) and associated theory (see [Höglund, 1978](#)) implies that (under certain regularity conditions) the limiting distribution of the Difference-in-Means test statistic is Normal. [Figure 1](#) shows that the Difference-in-Means test statistic looks more and more Normal as the study size increases while other factors remain constant. The Difference-in-Means test statistic is indeed asymptotically Normal and the  $p$ -values in [Equation 12](#) are all asymptotically valid—i.e., as  $N$  grows toward  $\infty$ , the probability of a Type I error is less than or equal to the  $\alpha$  level of the test.

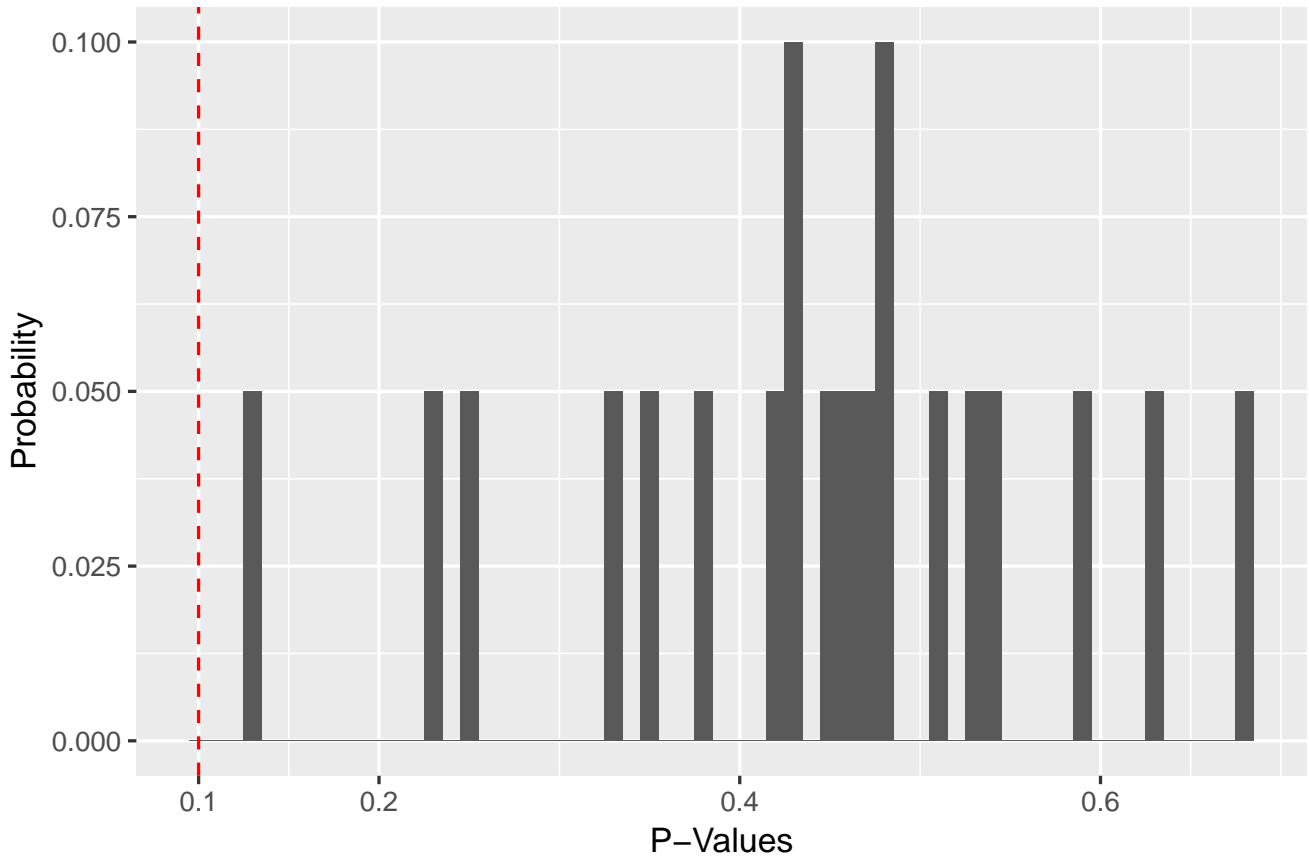
However, in small experiments in which the Difference-in-Means test statistic is poorly approximated by a Normal distribution, tests of a weak null hypothesis relative to an alternative may have either a Type I error probability greater than the test’s  $\alpha$  level (when the null is true) or low power (when the null is false). [Table 5](#) and [Figure 3](#) demonstrate this point — in which the assignment  $\mathbf{z}_8$  happened to be the one randomly drawn by the researcher. In this example, the observed Difference-in-Means test statistic is 11.6667 and the conservatively estimated variance is 10.3111. If we were to test the weak null hypothesis of no effect, i.e., that  $\bar{\tau}_0 = 0$ , against the alternative hypothesis that  $\bar{\tau}_a > 0$ , then the upper, one-tailed  $p$ -value would be as follows:

$$(13) \quad \left( 1 - \Phi \left( \frac{11.6667 - 0}{10.3111} \right) \right) = 0.1289,$$

which yields a higher  $p$ -value than the upper  $p$ -value we calculated via the Fisherian test of the strong null of no effects (which was  $p = .05$ ).

In general, if we were to calculate an upper  $p$ -value for a test of the weak null over all 20 possible realizations of data, [Figure 6](#) shows that the power of the test is low. In fact, with an  $\alpha$  level of 0.10, the power of the test of the weak null against the alternative of a positive mean causal effect is 0—i.e., the probability of calculating a  $p$ -value less than or equal to  $\alpha$  is 0.

In this particular case of a small experiment, we can see that the Neymanian hypothesis test has low power, lower than a Fisherian hypothesis test of the strong null of no effect — not shown here, but implied from the fact of at least one  $p$ -value less than 0.1 as shown by [Figure 3](#) (which implies the same weak null of  $\bar{\tau}_0 = 0$  that we used in the Neymanian test). For the same reason, even though Neymanian tests are asymptotically valid, they are not necessarily valid in small studies.



**Figure 6:** Distribution of Neymanian  $p$ -values under all realizations of data

An additional implication of the differences between Neymanian and Fisherian hypothesis tests is that the former is consistent<sup>9</sup>, but is not necessarily unbiased in finite contexts, even when the Difference-in-Means test statistic is well approximated by a Normal distribution. For intuition on this point, note that, when the alternative of a larger effect is true and the null of no effects is false, the Difference-in-Means test statistic will yield values that are systematically larger than the values it would produce if the null of no mean effect were true — larger causal effects leads to larger test statistic values. Yet the alternative of a larger effect could be such that, when it is true, the *variance estimates* are systematically larger than what the same variance estimates would be when the null is true — for example, if a positive mean causal effect is caused by a few outliers that react strongly to treatment. Since the  $z$ -score scales the Difference-in-Means estimates

<sup>9</sup>Although we know of no paper that provides a direct proof, we can see this property indirectly from [Wu and Ding \(2018\)](#) and [Lin \(2013\)](#). In their discussion of Theorem 1, [Wu and Ding \(2018, 8\)](#) state that “[a]n appealing consequence of these facts is that the FRT [Fisher randomization test] has power.” We thank Jason Wu and Peng Ding for pointing out that they could have said instead that the power of the test under a fixed alternative goes to 1 as the study size increases (i.e., that the test is consistent).

by the variance estimates, the systematically greater variance estimates when the alternative is true could yield *smaller*  $z$ -scores (and, hence, larger  $p$ -values) compared to when the null is true. Hence, Neymanian hypothesis tests of weak causal hypotheses are not necessarily unbiased in finite contexts, even when the assumption of Normality holds. For more on Fisherian versus Neymanian hypothesis tests, see [Ding \(2017\)](#), as well as the discussion from [Aronow and Offer-Westort \(2017\)](#), [Chung \(2017\)](#), [Bailey \(2017\)](#) and [Loh et al. \(2017\)](#).

Up to this point, we discussed the role that research design plays in the quality of procedures for causal inference—estimation and testing. To simplify exposition, we have been referring to situations where the researcher completely controls and, hence, knows the research design. We now consider situations in which the researcher does not control or only partially controls the research design.

## 5 Partially Controlled Research Designs: Noncompliance and Attrition

We refer to designs with imperfect compliance and/or attrition as partially controlled designs in that the researcher *does* control the probability distribution on the set of possible assignments, but does *not* control whether units actually comply with the assigned treatment or whether units report their outcomes. Causal inferences in such partially controlled designs often require more assumptions to make causal inferences, which are typically defined no longer on the whole study population, but a specific stratum of units in the study.

### 5.1 Noncompliance

Noncompliance occurs when units who are assigned to receive treatment or control do not actually receive it, where the random variable  $D_i \in \{0, 1\}$  is an indicator variable for whether unit  $i$  has received or not received the assigned treatment. We are often substantively interested in the causal effect of actual receipt of treatment, not its mere assignment. Under complete, uniform random assignment, we have shown above that the probability that  $Z_i = 1$  is identical for all  $i = 1, \dots, N$  units. Yet we're interested in the causal effect of  $D_i$  and since  $D_i$  is an outcome of  $Z_i$  (i.e., measured after  $Z_i$ ), the probability that  $D_i = 1$  is no longer identical for all units. Hence, a naive estimator of the difference in observed outcomes between those who did and did not receive the treatment (i.e., the per-protocol estimator) is not necessarily unbiased.

Since whether or not units actually receive (or comply with) the treatment is an outcome variable measured after assignment, we can define units' compliance status in terms of their unobservable potential outcomes:

$z_i = 0$	$z_i = 1$	Stratum
$d_{c,i} = 1$	$d_{t,i} = 1$	<i>Always-Taker</i>
$d_{c,i} = 0$	$d_{t,i} = 1$	<i>Complier</i>
$d_{c,i} = 1$	$d_{t,i} = 0$	<i>Defier</i>
$d_{c,i} = 0$	$d_{t,i} = 0$	<i>Never-Taker</i>

**Table 11:** Compliance Strata

Notice that the probability that  $D_i = 1$  for an *Always-Taker* is 1 and the probability that  $D_i = 1$  for a *Never-Taker* is 0. The only types of units for which the probability that  $D_i = 1$  remain identical and on the interval  $(0, 1)$  are *Compliers* and *Defiers*.

[Angrist et al. \(1996\)](#) show that scholars can *consistently estimate* the mean causal effect among *Compliers* under three further assumptions in addition to SUTVA and a uniform probability distribution on the set of assignments,  $\Omega$ , which is sometimes referred to as the exogeneity of the instrument,  $\mathbf{Z}$ . These three assumptions in addition to SUTVA and uniform random assignment are:

1. Exclusion restriction—i.e., the instrument affects the outcome only through the receipt of treatment;
2. No Defiers—i.e., that there are no units for which  $d_{c,i} = 1$  and  $d_{t,i} = 0$ ;
3. At least one Complier—i.e., there exists at least one unit in the study for which  $d_{c,i} = 0$  and  $d_{t,i} = 1$ .

To see this point, note that one can write the mean causal effect as a sum of the causal effects among the four strata above weighted by the proportion of units in each stratum:

$$\bar{\tau} = \delta_{AT}\pi_{AT} + \delta_C\pi_C + \delta_D\pi_D + \delta_{NT}\pi_{NT},$$

where  $\delta_s, \pi_s : s \in \{AT, C, D, NT\}$  represent the mean causal effect in each stratum and the proportion of units in that stratum, respectively.

By the exclusion restriction assumption, the causal effect of  $Z$  on  $Y$  must be 0 for Always-Takers and Never-Takers (i.e.,  $\delta_{AT} = 0$  and  $\delta_{NT} = 0$ ). And by the assumption of no Defiers, the proportion of Defiers,  $\pi_D$ , is 0. Hence, the mean causal effect among Compliers is

$$\begin{aligned}\bar{\tau} &= \delta_C \pi_C \\ \frac{\bar{\tau}}{\pi_C} &= \delta_C,\end{aligned}$$

i.e., that the mean causal effect on all units scaled by the proportion of Compliers is equal to the mean causal effect among Compliers.

We showed in Subsection 3.1 above that, under complete, uniform random assignment, the Difference-in-Means estimator is unbiased regardless of the distributions of potential outcomes. Both  $\mathbf{Y}$  and  $\mathbf{D}$  are outcomes of  $\mathbf{Z}$ ; hence, the Difference-in-Means estimators of  $\hat{\tau}(\mathbf{Z}, \mathbf{Y})$  and  $\hat{\tau}(\mathbf{Z}, \mathbf{D})$  unbiasedly estimate  $\bar{\tau}$  and  $\pi_C$ , respectively. The Wald (or IV) estimator (Wald, 1940), which is also sometimes referred to as the Bloom estimator (Bloom, 1984) or the CACE (Complier Average Causal Effect) estimator (Gerber and Green, 2012), is defined as  $\frac{\hat{\tau}(\mathbf{Z}, \mathbf{Y})}{\hat{\tau}(\mathbf{Z}, \mathbf{D})}$ . In short, it is the ratio of these two unbiased estimators. This ratio estimator consistently, though not necessarily unbiasedly, estimates  $\frac{\bar{\tau}}{\pi_C}$  (see Angrist and Pischke, 2008, Sections 4.6 – 4.7).

A design based test of the hypothesis of no mean causal effect among Compliers is difficult due to the absence of an analytic expression for the variance of the CACE ratio estimator. (For more on this topic, see Imbens and Rosenbaum, 2005; Kang et al., 2018). On the other hand, tests of the null hypothesis that the causal effect is 0 for all Compliers is relatively straightforward.

The null hypothesis of no causal effect among Compliers implies the strong null of no effect under the assumptions of the exclusion restriction and no Defiers. Since these two assumptions jointly imply that the individual causal effect of assignment is 0 for all Always-Takers and Never-Takers, and that there are no Defiers, then no causal effect among Compliers implies no causal effect among all units. For example, let’s imagine that (like in Table 5)  $\mathbf{z}_8$  were the assignment that the researcher happened to randomly draw, except that now we observe the following imperfect compliance as shown in Table 12:

We know that if  $d_{t,i} = 0$ , then, by the no Defiers assumption, that unit must be a Never-Taker.

$\mathbf{z}_8$	$\mathbf{y}_c$	$\mathbf{y}_t$	$\mathbf{y}_8$	$\mathbf{d}_8$	$\mathbf{d}_c$	$\mathbf{d}_t$
1	?	22	22	0	?	0
0	8	?	8	0	0	?
0	11	?	11	0	0	?
1	?	15	15	1	?	1
1	?	18	18	0	?	0
0	1	?	1	1	1	?

**Table 12:** Realization of Data if  $\mathbf{z}_8$  were the randomly drawn assignment

And if  $d_{c,i} = 1$ , then, by the same assumption, that unit must be an Always-Taker. We don't know, however, whether treated units for which  $d_t = 1$  are Compliers or Always-Takers and we don't know whether control units for which  $d_c = 0$  are Compliers or Never-Takers. Yet under the null hypothesis of no causal effect among Compliers, the individual causal effect is 0 regardless of whether a given unit is a Complier, Always-Taker or Never-Taker. Hence, we can fill in missing potential outcomes without knowing those units compliance strata as follows:

$\mathbf{z}_8$	$\mathbf{y}_c$	$\mathbf{y}_t$	$\mathbf{y}_8$	$\mathbf{d}_8$	$\mathbf{d}_c$	$\mathbf{d}_t$
1	22	22	22	0	0	0
0	8	8	8	0	0	?
0	11	11	11	0	0	?
1	15	15	15	1	?	1
1	18	18	18	0	0	0
0	1	1	1	1	1	1

**Table 13:** Potential outcome under strong null of no effect if  $\mathbf{z}_8$  were realized assignment

One can now use an effect increasing test statistic, such as the Difference-in-Means test statistic, to test the hypothesis of no Complier causal effect against either the alternative of a positive Complier causal effect or a negative Complier causal effect. [Hansen and Bowers \(2009\)](#) present an application of this idea in the context of the one-sided compliance in a cluster-randomized get-out-the-vote campaign with a binary outcome — replacing what would be a complex two-stage logistic multilevel model with a relatively simple Fisherian hypothesis test. And [Imbens and Rosenbaum \(2005\)](#) show how the Fisherian approach produces valid confidence intervals when the Neyman-style approach using two-stage least squares fails to control the false positive rate when the instrument is weak (i.e., there are few Compliers).

Under the assumptions of the exclusion restriction and no Defiers, we do not need to know which units are Compliers in order to assert the hypothesis of no Complier causal effect. However, if one were to posit a hypothesis other than no Complier causal effect, then one would also need to also posit some hypothesis about which units are Compliers and which are not. Hypothesis testing

with imperfect compliance is therefore more complicated when testing hypotheses other than that of no Complier causal effect. (For more on Fisherian approaches to instrumental variable analysis, see Kang et al., 2018; Rosenbaum, 1996, among others.)

## 5.2 Attrition, or Missing Outcomes

Our second step away from the ideal case of complete control over the research design is to allow for the possibility of missing outcomes. Let  $r_{t,i}$  be an indicator, i.e.,  $r_{t,i} \in \{0, 1\}$ , for whether subject  $i$  would respond to an attempt to measure an outcome and let  $r_{c,i} \in \{0, 1\}$  be an indicator for whether subject  $i$  would respond if assigned to control. We can represent whether an individual's outcomes are missing or not based on Equation 14:

$$(14) \quad Y_i = \begin{cases} y_{c,i} + [y_{t,i} - y_{c,i}]Z_i & \text{if } R_i = 1 \\ \text{NA} & \text{if } R_i = 0, \end{cases}$$

where  $R_i = Z_i r_{t,i} + (1 - Z_i) r_{c,i}$ .

From Equation 14 above, we can see that if  $R_i = 1$ , then the researcher will observe  $y_{c,i}$  for unit  $i$  if  $Z_i = 0$  and  $y_{t,i}$  for unit  $i$  if  $Z_i = 1$ . By contrast, if  $R_i = 0$ , then  $Y_i$  will be unobserved—i.e., NA.

We can define four distinct strata of subjects with regard to attrition in order to help us understand how attrition can affect the properties of estimators and hypothesis tests. Just as we focused on a subgroup of units in the case of uncontrolled compliance, we can only infer the causal effects on certain subgroups when outcomes are missing (even when assignments are randomized).

$z_i = 0$	$z_i = 1$	Stratum
$r_{ci} = 1$	$r_{ti} = 1$	<i>Always-Reporter</i>
$r_{ci} = 0$	$r_{ti} = 1$	<i>If-Treated-Reporter</i>
$r_{ci} = 1$	$r_{ti} = 0$	<i>If-Untreated-Reporter</i>
$r_{ci} = 0$	$r_{ti} = 0$	<i>Never-Reporter</i>

**Table 14:** Attrition Strata

In the context of attrition, can we define a set of assumptions, as we did under imperfect compliance, in which estimators and tests satisfy the properties they should on some subset of

the experimental data? It turns out that we can do so only if whether or not a unit attrits is independent of treatment assignment. Without further assumptions, we know only that the random variable  $R_i$  is independent of the random variable  $Z_i$  among only Always-Reports and Never-Reporters.<sup>10</sup> Therefore, if we assume that all experimental units belong to one of these two types, then our estimators and tests maintain the properties they should among the set of *Always-Reporters*.<sup>11</sup> (We cannot observe outcomes for *Never-Reporters* and, hence, cannot estimate or test hypotheses about causal effects on units who are *Never-Reporters*.) On this set of *Always-Reporters*, one can estimate causal effects, as well as test strong or weak causal hypotheses as we did above.

## 6 Uncontrolled Research Designs: Observational Studies

In this section, we finally relax the assumption that the researcher has control over how the treatment variable is assigned. The key distinction between experimental and observational studies is that in a randomized experiment, the researcher knows the probabilities with which units are *assigned* to treatment and control conditions; however, in an observational study, the researcher observes units only after they have *selected* into study conditions with unknown probabilities. How, then, is one to generate statistical inferences using estimators and tests focusing on causal effects when the probability distribution on the set of assignments,  $\Omega$ , is unknown? A common design-based approach to this problem is to define units' treatment assignment probabilities as an unknown function of a set of baseline covariates. In the ideal (and often unattainable case), by appropriate conditioning on these baseline covariates, the researcher can estimate and test hypotheses about causal effects via procedures that meet the desirable properties described at the outset of this chapter. Much of the work on observational studies emphasizes appropriate conditioning on baseline covariates, as well as methods to diagnosis the success of such conditioning strategies (e.g., [Hansen and Bowers, 2008](#); [Hartman and Hidalgo, 2018](#), among others). Realistically, the

---

<sup>10</sup>Note that  $R_i$  is independent of  $Z_i$  if and only if the probability that  $R_i$  takes on any value in its sample space does not vary conditional on any value that  $Z_i$  takes on in its sample space—i.e., that  $\Pr(R_i = 1 | Z_i = 1) = \Pr(R_i = 1 | Z_i = 0)$  and  $\Pr(R_i = 0 | Z_i = 1) = \Pr(R_i = 0 | Z_i = 0)$ . The only two types of subjects who satisfy such independence are Always-Reporters and Never-Reporters.

<sup>11</sup>Other approaches to estimation and testing relax the assumption that outcome missingness is independent of treatment assignment and devise procedures that bound inferences about treatment effects under best- or worst-case scenarios, e.g., “trimming bounds” and “extreme value bounds” ([Lee, 2009](#); [Manski, 1990](#)).



design-based approach in observational studies might be called an “as-if-randomized” approach: For example, a researcher might make choices about comparison groups such that, within group, treatment selection *appears* random.

The model of an observational study states that units are individually assigned to treatment or control by  $N$  *independent*, but not necessarily *identically distributed*, coin tosses. More specifically, for all  $i \in \{1, \dots, N\}$  units, we let  $\Pr(Z_i = 1)$  be equal to  $\lambda(\mathbf{x}_i)$  where  $\lambda(\cdot)$  is an unknown function and  $\mathbf{x}_i$  is unit  $i$ 's fixed vector of baseline covariate values. Even if we don't know  $\lambda(\cdot)$ , if  $\mathbf{x}_i = \mathbf{x}_j$  for any two units  $i$  and  $j \neq i$ , then it follows that  $\Pr(Z_i = 1) = \Pr(Z_j = 1)$ . Of course, we still don't know the function  $\lambda(\cdot)$  and, hence, don't know the actual values of  $\Pr(Z_i = 1)$  and  $\Pr(Z_j = 1)$ ; we know only that these two values are equal. Therefore, if we construct a block,  $b$ , that consists of units  $i$  and  $j \neq i$ , then each possible assignment within that block has an equal probability of realization. And if each possible assignment has an equal probability, then an observational study can be analyzed as if it is a uniform, block randomized experiment. (For more on the analysis of block, randomized experiments, see [Gerber and Green, 2012](#), Chapter 4.) This approach allows us to avoid directly estimating  $\lambda(\cdot)$ , although there are alternative approaches that do so and subsequently use these estimated values (known as estimated propensity scores) as a basis for inference (see [Robins et al., 2000](#)).

Scholars can therefore attempt to make an observational study as experiment-like as possible by creating matched blocks on the basis of observed covariates that determine units' treatment assignment probabilities. A range of matching algorithms exist to improve covariate balance and, hence, make observational studies like experiments as much as possible (at least on the basis of observed covariates). For more on this topic, see [Hansen \(2004\)](#), [Diamond and Sekhon \(2013\)](#), [Sävje et al. \(2017\)](#) and [Zubizarreta \(2012\)](#) among others. After matching (or perhaps weighting or non-parametric modeling) and favorable comparisons with actually randomized designs, researchers then must confront the fact that their observational studies are not randomized studies. This leads directly to the topic of sensitivity analyses.

## 7 Sensitivity to Assumptions in Uncontrolled Research Designs

Thus far, we have considered cases in which either the probability distribution on  $\Omega$  is known by random assignment or units' assignment probabilities are a function of only *observed* covariates. But in an observational study, we rarely know, let alone observe, all relevant covariates. We now consider deviations from the assumption that units' assignment probabilities are determined by only observed covariates and subsequently assess how one's inferences would change under violations of this assumption.

A powerful, design-based framework for such a sensitivity analysis is given by [Rosenbaum \(2002\)](#). Before explaining this framework, we need to define a few additional terms. First, the *treatment odds* for unit  $i \in \{1, \dots, N\}$  is  $\frac{\pi_i}{(1-\pi_i)}$ , which is simply the  $i$ th unit's probability of assignment to treatment divided by that unit's probability of assignment to control. The *treatment odds ratio* for any two units  $i$  and  $j \neq i$  is simply the ratio of the  $i$ th unit's treatment odds and the  $j$ th unit's treatment odds. If units' treatment odds are a function of only observed covariates *and* the researcher is able to obtain balance on all of these observed covariates, then the treatment odds for units  $i, j \neq i : \mathbf{x}_i = \mathbf{x}_j$  is identical and their treatment odds ratio is 1.

[Rosenbaum \(2002\)](#) considers what would happen when units' treatment odds are a function not only of observed covariates,  $\mathbf{x}$ , but also an unobserved covariate  $u$ . Under the assumption of a logistic functional form between all units' treatment odds and baseline covariates, as well as the constraint that  $0 \leq u \leq 1$ , one can write the treatment odds of the  $i$ th unit as follows:

$$\begin{aligned} \frac{\pi_i}{(1-\pi_i)} &= \exp \{ \kappa(\mathbf{x}_i) + \gamma u_i \} \\ \log \left( \frac{\pi_i}{(1-\pi_i)} \right) &= \kappa(\mathbf{x}_i) + \gamma u_i, \end{aligned}$$

where  $\kappa(\cdot)$  is an unknown function and  $\gamma$  is an unknown parameter, and the the treatment odds

ratio for units  $i$  and  $j$  is:

$$\begin{aligned} \frac{\left(\frac{\pi_i}{1-\pi_i}\right)}{\left(\frac{\pi_j}{1-\pi_j}\right)} &= \frac{\exp\{\kappa(\mathbf{x}_i) + \gamma u_i\}}{\exp\{\kappa(\mathbf{x}_j) + \gamma u_j\}} \\ &= \exp\left\{\left(\kappa(\mathbf{x}_i) + \gamma u_i\right) - \left(\kappa(\mathbf{x}_j) + \gamma u_j\right)\right\}. \end{aligned}$$

If  $\mathbf{x}_i = \mathbf{x}_j$ , then  $\kappa(\mathbf{x}_i) = \kappa(\mathbf{x}_j)$  and, hence, the treatment odds ratio is simply:

$$\exp\left\{\gamma(u_i - u_j)\right\}.$$

Since  $u_i, u_j \in [0, 1]$ , the minimum and maximum possible values of  $(u_i - u_j)$  are  $-1$  and  $1$ . Therefore, the minimum and maximum possible values of the treatment odds ratio are  $\exp\{-\gamma\}$  and  $\exp\{\gamma\}$ . After noting that  $\exp\{-\gamma\} = \frac{1}{\exp\{\gamma\}}$ , we can bound the treatment odds ratio between  $i$  and  $j$  as follows:

$$(15) \quad \frac{1}{\exp\{\gamma\}} \leq \frac{\left(\frac{\pi_i}{1-\pi_i}\right)}{\left(\frac{\pi_j}{1-\pi_j}\right)} \leq \exp\{\gamma\}.$$

We can denote  $\exp\{\gamma\}$  by  $\Gamma$  and subsequently consider how one's inferences would change for various values of  $\Gamma$ .

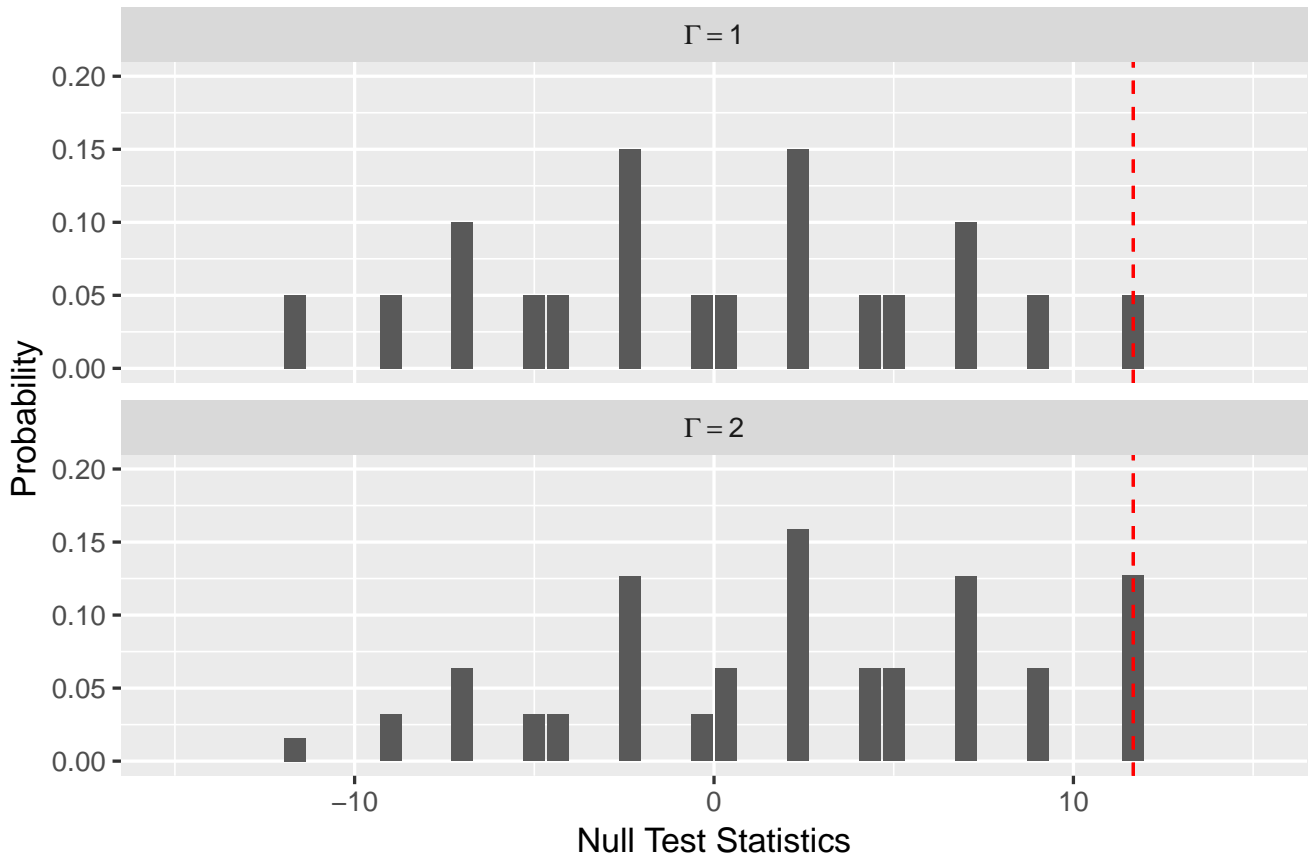
For example, let's say that a researcher obtains balance via stratification on all observed covariates—such that the design closely resembles a uniform, block randomized experiment—and subsequently tests a strong null hypothesis under the assumption that all units' treatment odds are identical. Now the researcher considers deviations from this assumption. Different assumptions about  $u$  and  $\gamma$  imply differing probabilities of possible assignments, which, as [Rosenbaum \(2002, Chapter 4\)](#) shows, can be represented by:

$$(16) \quad \Pr(\mathbf{Z} = \mathbf{z}) = \frac{\exp\{\gamma \mathbf{z}' \mathbf{u}\}}{\sum_{\mathbf{z} \in \Omega} \exp\{\gamma \mathbf{z}' \mathbf{u}\}} = \prod_{b=1}^B \frac{\exp\{\gamma \mathbf{z}' \mathbf{u}\}}{\sum_{\mathbf{z} \in \Omega} \exp\{\gamma \mathbf{z}' \mathbf{u}\}}.$$

One conservative approach for calculating such probabilities would be to assume that, within

each matched set, the value of  $u$  for all treated units is 1 and the value of  $u$  for all control units is 0. Hence, the treatment odds ratio of treated units to control units within each matched set is  $\exp\{\gamma\}$ . Now if the researcher were to assume that  $\Gamma = 2$ , then such an assumption would imply that  $\gamma = \log(2) \approx 0.6931$ . Now the researcher can calculate the probability of each possible assignment in  $\Omega$  using Equation 16. Notice that in this sensitivity analysis, the only random quantity is  $\mathbf{Z}$ . Potential outcomes are still fixed and the researcher can test any strong null hypothesis; however, under such a sensitivity analysis,  $p$ -values will be calculated based on the probabilities of each assignment given in (16) above.

To return to the working example from Table 1, let's imagine that the realized assignment was  $\mathbf{z}_8$ , which yielded an observed test statistic of 11.6667 and a  $p$ -value of 0.05. (See Figure 3 above.) We calculated that  $p$ -value under the assumption that all units had identical probabilities of assignment, i.e., that  $\Gamma = 1$  for all  $i$  and  $j$ . If we were to relax that assumption and assume, for example, that  $\Gamma = 2$ , then the  $p$ -value would increase. Figure 7 illustrates how the respective null distributions would differ under  $\Gamma = 1$  and  $\Gamma = 2$ , respectively:



**Figure 7:** Distributions of Difference-in-Means test statistic under strong null of no effect when  $\mathbf{z}_8$  is the realized assignment under (a) a model where an unobserved covariate has no effect on odds of treatment,  $\Gamma = 1$ , and (b) a model where an observed covariate doubles the odds of treatment,  $\Gamma = 2$ .

Notice that under  $\Gamma = 1$  and  $\Gamma = 2$ , the observed test statistic remains fixed at 11.6667. The set of 20 null test statistic also remains fixed for  $\Gamma = 1$  and  $\Gamma = 2$ . The value of  $\Gamma$  changes only the probability associated with each of the null test statistics. As we can see, when  $\Gamma = 1$ , the  $p$ -value is 0.05, but when  $\Gamma = 2$ , the  $p$ -value increases to approximately 0.1270. The researcher is no longer able to reject the strong null of no effect at a level of  $\alpha = 0.10$  when  $\Gamma = 2$  compared to when  $\Gamma = 1$ , and this is one way in which researchers can assess the sensitivity of their inferences to assumptions about the research design.

This approach is not the only way to formalize the impact of the assumptions underlying the "as-if randomized" research designs used for causal inference when researchers have little to no control over the selection process of the main explanatory variable. See also ? for an approach focusing on regression coefficients and regression based adjustment as well as an application by ?

of ideas like these to problems in international relations.

## 8 Conclusion

Statistics and research design help us learn about general and abstract social science theory using concrete and specific observations. Observation helps us learn about theory, but observation occurs using the tools of research design and observations are summarized, described and interpreted using the tools of statistics. In fact, we have shown that certain, counterfactual causal quantities can never be directly observed, and that our ability to report with confidence about such unobserved causal effects depends heavily on statistical tools, which, in turn depend on research design for their operation. The persuasiveness of claims about links from estimate or  $p$ -value to unobserved causal effect specific to a fixed pool of units to general theory depend on the clarity of each step. At the most nitty-gritty level, we want our tools to work well — our estimators and tests should err in known and controlled manners, and should err rarely.<sup>12</sup> We have asked: “How can we know that tools do the work that we want them to do?” And we have shown how to use the facts of research design to answer that question and to justify use of these common tools. This means that when we want to persuade an audience that our findings support a given theory (or urge modification of it), we do not need our audiences to believe that (a) we have a random sample from a well-defined population, (b) that the outcome arises from some known probability process (like a Normal or Zero-Inflated Poisson distribution), (c) that the treatment or selection process relates to background covariates in some known (often linear and additive) fashion. Instead, in a randomized experiment, we ask that a reader believe that our research design correctly describes the physical processes that occurred in the research (a request that is no hard to verify and assess). In an observational study, we ask readers to agree that the as-if randomized approach is reasonable (and we present direct evidence comparing observational designs to randomized experiments in order to make the provisional as-if randomized approach easier to grant).

In explaining about reliable procedures, we have used a very simple set of examples. Our simplifications include the use of only two experimental conditions (treatment and control); yet the general modes of inference described in this chapter can be straightforwardly applied to factorial

---

<sup>12</sup>We did not assess tools like randomizers or sampling schemes in this paper, but those tools are equally important in the effort to advance theory through observation.

experiments and other contexts with multiple treatments (see [Dasgupta et al., 2015](#)). A second such simplification has been the Stable Unit Treatment Value Assumption (SUTVA) ([Cox, 1958](#); [Rubin, 1980, 1986](#)), which, in the case of a binary treatment variable, implies that all units have only two potential outcomes. Yet both estimation (see [Aronow et al., 2017](#); [Hudgens and Halloran, 2008](#)) and testing (see [Bowers et al., 2013](#); [Rosenbaum, 2007](#); ?) are possible when units have more than two potential outcomes due to interference between unit, for example, in the context of experiments on networks. Third, we have focused on randomized experiments, and only briefly pointed to the strategy of “as-if randomized” approaches to estimation and testing in observational studies. We explained that such approaches, when paired with sensitivity analyses, can enable persuasive statistical inferences about causal effects when the researcher lacks control over the design. That is, once an observational research design compares favorably to the standard of an equivalent experiment (the way that a matched-design can be compared to a block-randomized experiment), statistical inference about causal effects can use the same procedures, provisionally justified in the same way, as in a randomized experiment if also followed by a sensitivity analysis.

This focus on the basics — on ensuring that our statistical tools do what they should — leaves larger questions unexplored. For example, some researchers would prefer to make inferences about not only counterfactual causal effects among units in a given study, but also about future units in data contexts that differ from the one under study. One might desire unbiased and consistent estimators, as well as valid and powerful tests, not only based on the research design generating the data collected here and now, but also for unknown future research designs guiding data collection elsewhere and at other times. Forecasting causal effects is an active research area. For only a few recent works on the topic, see (????). Whether or not a researcher or policy-maker would like a formal forecast of the causal effects of an intervention from one study to a new context (in time, space, and/or units), information provided by a single study to the motivating theoretical question still depends on the reliability of the tools used to conduct and analyze the study. We have focused on showing how the reliability of such procedures is based on the research design itself and leave questions about the properties of procedures for forecasting causal effects as a separate yet important topic.

Design-based causal inference emphasizes inferring a counterfactual quantity, not a quantity in a population or of an outcome generating model. Such an emphasis arises naturally from a wide range of social science research contexts, such as when (1) units are not a random sample — for example, when the units are administrative units like schools or countries or convenience samples arising without a known chance process; or (2) probability models of outcomes—and their structural relationship to explanatory variables—are difficult to write or justify — for example, when an outcome can be plausibly modelled by multiple different likelihood functions. In such contexts, simple comparisons based on the research design can advance social science theory and avoid debates about data models. When strong theory generates clear probability and structural models, then a model-based justification for statistical inference might be preferred, although we would want such modes of inferences to satisfy the same properties discussed in this chapter: Tests should not mislead and estimators should produce estimates close to the truth. The model-based approach to showing whether these characteristics hold is well described in most statistics textbooks, and we recommend ? for an overview.

A general question nevertheless remains: Are design-based procedures better or worse than model-based procedures for advancing social science theory or policy learning? Design-based inference is simple, easily interpretable and can ensure that estimators and tests are reliable based on few assumptions, where the assumptions tend to be easily defended in terms of the known features of the research design. But does design-based inference possess reliable properties only for narrowly defined research questions? To be sure, there is nothing intrinsic to design-based inference that requires scholars to use only specific estimators that focus only on specific estimands, such as mean causal effects, or to test only certain hypotheses about no causal effects instead of others. And, although we did not show it here, it is straightforward to assess the properties of non-standard estimators and tests by representing the research design and simulating from it (see ?, for an example of a framework for simulation based assessment of estimators and tests). For only one example of the flexibility of design-based approaches, imagine that we wondered whether a certain non-linear structural model well described a relationship between a causal driver (like an experimental treatment) and an outcome. In this case, [Bowers et al. \(2013\)](#); ? show how the



evidence against structural models of unobserved potential outcomes can be generated and hypothesis tests can be created. That is, there is nothing about design-based approaches which precludes the use of structural models. That said, often a clear difference of means can teach enough about a complicated structural theory such that there is no need to complicate the research design or statistical inference tasks. In the end, all else equal, reliable procedures advance scientific knowledge more than unreliable procedures do. For this reason, one of the benefits of engaging with design-based inference is that it brings clarity to the task of judging and choosing our statistical tools and provokes us to directly confront and grapple with the conditions under which evidence can be reliably interpreted as evidence for or against causal claims.

## References

- Angrist, J. D., G. W. Imbens, and D. B. Rubin (1996). Identification of causal effects using instrumental variables. *Journal of the American Statistical Association* 91(434), 444–455. 35
- Angrist, J. D. and J.-S. Pischke (2008). *Mostly Harmless Econometrics: An Empiricist’s Companion*. Princeton, NJ: Princeton University Press. 36
- Aronow, P. M., , and C. Samii (2017). Estimating average causal effects under general interference, with application to a social network experiment. *Annals of Applied Statistics* 11(4), 1912–1947. 46
- Aronow, P. M., D. P. Green, D. K. Lee, et al. (2014). Sharp bounds on the variance in randomized experiments. *The Annals of Statistics* 42(3), 850–871. 17
- Aronow, P. M. and J. A. Middleton (2013). A class of unbiased estimators of the average treatment effect in randomized experiments. *Journal of Causal Inference* 1(1), 135–154. 18
- Aronow, P. M. and M. R. Offer-Westort (2017). Understanding ding’s apparent paradox. *Statistical Science* 32(3), 346–348. 34
- Bailey, R. A. (2017). Inference from randomized (factorial) experiments. *Statistical Science* 32(3), 352–355. 34
- Bloom, H. S. (1984). Accounting for no-shows in experimental evaluation designs. *Evaluation Review* 8(2), 225–246. 36
- Bowers, J., M. Fredrickson, and C. Panagopoulos (2013). Reasoning about interference between units: A general framework. *Political Analysis* 21(1), 97–124. 30, 46, 47
- Brewer, K. (1979). A class of robust sampling designs for large-scale surveys. *Journal of the American Statistical Association* 74(368), 911–915. 7
- Chung, E. (2017). Randomization-based tests for “no treatment effects”. *Statistical Science* 32(3), 349–351. 34
- Cochran, W. G. (1977). *Sampling Techniques* (3rd ed.). Hoboken, NJ: John Wiley & Sons. 16
- Cox, D. R. (1958). *Planning of Experiments*. New York, NY: Wiley. 1, 46
- Dasgupta, T., N. S. Pillai, and D. B. Rubin (2015). Causal inference from 2k factorial designs by using potential outcomes. *Journal of the Royal Statistical Society: Series B (Statistical Methodology)* 77(4), 727–753. 1, 46
- Diamond, A. and J. S. Sekhon (2013). Genetic matching for estimating causal effects: A general multivariate matching method for achieving balance in observational studies. *Review of Economics and Statistics* 95(3), 932–945. 40
- Ding, P. (2017). A paradox from randomization-based causal inference. *Statistical Science* 32(3), 331–345. 34
- Fisher, R. A. (1935). *The Design of Experiments*. Edinburgh, SCT: Oliver and Boyd. 1, 19, 20

- Gerber, A. S. and D. P. Green (2012). *Field Experiments: Design, Analysis, and Interpretation*. New York, NY: W.W. Norton. 14, 36, 40
- Hansen, B. B. (2004). Full matching in an observational study of coaching for the sat. *Journal of the American Statistical Association* 99(467), 609–618. 40
- Hansen, B. B. and J. Bowers (2008). Covariate balance in simple, stratified and clustered comparative studies. *Statistical Science* 23(2), 219–236. 39
- Hansen, B. B. and J. Bowers (2009). Attributing effects to a cluster-randomized get-out-the-vote campaign. *Journal of the American Statistical Association* 104(487), 873–885. 37
- Hartman, E. and F. D. Hidalgo (2018). An equivalence approach to balance and placebo tests. *American Journal of Political Science* 62(4), 1000–1013. 39
- Höglund, T. (1978). Sampling from a finite population. a remainder term estimate. *Scandinavian Journal of Statistics* 5(1), 69–71. 32
- Holland, P. W. (1986). Statistics and causal inference. *Journal of the American Statistical Association* 81(396), 945–960. 1
- Horvitz, D. G. and D. J. Thompson (1952). A generalization of sampling without replacement from a finite universe. *Journal of the American Statistical Association* 47(260), 663–685. 18
- Hudgens, M. G. and M. E. Halloran (2008). Toward causal inference with interference. *Journal of the American Statistical Association* 103(482), 832–842. 46
- Imbens, G. W. and P. R. Rosenbaum (2005). Robust, accurate confidence intervals with a weak instrument: Quarter of birth and education. *Journal of the Royal Statistical Society: Series A (Statistics in Society)* 168(1), 109–126. 36, 37
- Kang, H., L. Peck, and L. Keele (2018). Inference for instrumental variables: A randomization inference approach. *Journal of the Royal Statistical Society. Series A: Statistics in Society* 181(4), 1231–1254. 36, 38
- Lee, D. S. (2009). Training, wages, and sample selection: Estimating sharp bounds on treatment effects. *The Review of Economic Studies* 76(3), 1071–1102. 39
- Lehmann, E. L. and J. P. Romano (2005). *Testing Statistical Hypotheses* (Third ed.). Springer Texts in Statistics. New York, NY: Springer-Verlag. 20
- Li, X. and P. Ding (2017). General forms of finite population central limit theorems with applications to causal inference. *Journal of the American Statistical Association* 112(520), 1759–1769. 32
- Lin, W. (2013). Agnostic notes on regression adjustments to experimental data: Reexamining freedman’s critique. *The Annals of Applied Statistics* 7(1), 295–318. 33
- Loh, W. W., T. S. Richardson, and J. M. Robins (2017). An apparent paradox explained. *Statistical Science* 32(3), 356–361. 34
- Manski, C. F. (1990). Nonparametric bounds on treatment effects. *The American Economic*

- Review 80(2)*, 319–323. 39
- Middleton, J. A. and P. M. Aronow (2015). Unbiased estimation of the average treatment effect in cluster-randomized experiments. *Statistics, Politics and Policy* 6(1-2), 39–75. 7, 9
- Neyman, J. (1923). Sur les applications de la théorie des probabilités aux expériences agricoles: Essai des principes. *Roczniki Nauk Rolniczych* 10, 1–51. 14, 17
- Neyman, J. and E. S. Pearson (1933). On the problem of the most efficient tests of statistical hypotheses. *Philosophical Transactions of the Royal Society of London* 231(694–706), 289–337. 19, 20
- Robins, J. M., M. A. Hernán, and B. Brumback (2000). Marginal structural models and causal inference in epidemiology. *Epidemiology* 11(5), 550–560. 40
- Rosenbaum, P. R. (1996). Identification of causal effects using instrumental variables: Comment. *Journal of the American Statistical Association* 91(434), 465–468. 38
- Rosenbaum, P. R. (1999). Reduced sensitivity to hidden bias at upper quantiles in observational studies with dilated treatment effects. *Biometrics* 55(2), 560–564. 26
- Rosenbaum, P. R. (2002). *Observational Studies* (Second ed.). New York, NY: Springer. 19, 20, 26, 28, 41, 42
- Rosenbaum, P. R. (2007). Interference between units in randomized experiments. *Journal of the American Statistical Association* 102(477), 191–200. 46
- Rosenbaum, P. R. (2010). *Design of Observational Studies*. New York, NY: Springer. 19, 20, 21, 23, 26, 30
- Rubin, D. B. (1980). Comment on ‘randomization analysis of experimental data in the fisher randomization test’ by basu, d. *Journal of the American Statistical Association* 75(371), 591–593. 1, 46
- Rubin, D. B. (1986). Which ifs have causal answers? (comment on ‘statistics and causal inference’ by paul w. holland). *Journal of the American Statistical Association* 81, 961–962. 1, 46
- Sävje, F., M. J. Higgins, and J. S. Sekhon (2017). Generalized full matching. Working Paper. 40
- Senn, S. (2004). Controversies concerning randomization and additivity in clinical trials. *Statistics in Medicine* 23(24), 3729–3753. 8
- Wald, A. (1940). The fitting of straight lines if both variables are subject to error. *The Annals of Mathematical Statistics* 11(3), 284–300. 36
- Wu, J. and P. Ding (2018, October). Randomization tests for weak null hypotheses. In *eprint arXiv:1809.07419v2*. 33
- Zubizarreta, J. R. (2012). Using mixed integer programming for matching in an observational study of kidney failure after surgery. *Journal of the American Statistical Association* 107(500), 1360–1371. 40