

Causal Inference

Ashley I Naimi

Spring 2022

Contents

1	Intro	2
2	Correlation and Causation	2
3	Introduction to Causal Inference	3
4	Potential Outcomes Notation	4
5	Estimand, Estimator, Estimate	6
5.1	Estimand	6
5.2	Estimator	9
5.3	Estimates	16
6	Identifiability: Average Treatment Effect	16
6.1	Counterfactual Consistency	17
6.2	Interference	18
6.3	Exchangeability	19
6.4	Conditional Exchangeability	20
6.5	Correct Model Specification	22
6.6	Positivity	23
7	Identifiability: Instrumental Variables and Local Average Treatment Effects	29
7.1	Randomization and ITT	30
7.2	The Instrumental Variable Estimator	31
7.3	Local Average Treatment Effects	32
7.4	IV Identifiability Assumptions	33
7.5	Weak Instruments	35
7.6	Homogeneity Makes the LATE \rightarrow ATE	35
8	Non-Identifiability and Partial Identifiability: Bounding Effects	35
9	Takeaway	38

1 Intro

In the last section, we discussed the concepts of a cohort, censoring and truncation, the relation between cumulative distribution functions (CDFs, a.k.a cumulative incidence functions, cumulative risk, other), and how to use the Kaplan-Meier, Aalen-Johansen, and Gray's CIF estimators to obtain cause-specific or sub-distribution estimates of the CDF. In this section, we'll build on the practice of estimating risks. In particular, we'll discuss what happens when we're interested in exposure effects, particularly in an observational setting. This requires accounting for problems (potential biases) that can arise. We'll start with a basic introduction to causal inference and potential outcomes. We'll discuss identifiability and what this means. Finally, we'll be introduced to some strategies that exist when our parameters of interest are not identified.

2 Correlation and Causation

In the *The Grammar of Science*, Karl [Pearson \(1911\)](#) wrote: “[b]eyond such discarded fundamentals as ‘matter’ and ‘force’ lies still another fetish amidst the inscrutable arcana of modern science, namely, the category of cause and effect.” He suggested that rather than pursue an understanding of cause-effect relations, scientists would be best served by measuring correlations through tables that classify individuals into specific categories. “Such a table is termed a contingency table, and the ultimate scientific statement of description of the relation between two things can always be thrown back upon such a contingency table.”

Over a century later, a majority of statistics courses tend to treat causal inference by simply stating that “correlation is not causation.” This treatment is hardly sufficient, for at least two reasons: 1) As scientists, our primary interest is (should be) in cause-effect relations; 2) People continue to conflate correlation with causation¹. For both of these reasons, we very much need to **clarify the conditions that would allow us to understand causality better**. This is what “causal inference” is all about.

Generally, I adopt the view that **the causal and statistical aspects of a scientific study should be kept as separate as possible**. The objective is to first define the effect and articulate the conditions under which causal

¹ Daniel Westreich and I reviewed a book whose authors were so caught up in the allure of “Big Data”, they thoroughly forgot that correlation \neq causation. See [Naimi and Westreich \(2014\)](#)

inference is possible for this effect, and then to understand what statistical tools will enable us to answer the causal question.² Causal inference tells us what we should estimate, and whether we can. Statistics tells us how to estimate it. By implication, we should avoid treating statistical models as if they were causal. For example, the practice of reading the risk ratio, odds ratio, or risk difference for an exposure of interest from a generalized linear (statistical) model³ will sometimes work under very specific conditions, but is not the best approach for quantifying exposure effects (Naimi and Whitcomb, 2020). In this section, we will cover exactly how the cumulative risk function can be integrated into this framework.

An additional element that we won't have an opportunity to get into in this course is how we can avoid making unnecessary parametric assumptions when estimating causal effects. Specifically, the objective is to avoid imposing unnecessary parametric forms⁴ on the causal models that we believe are generating the data. Machine learning in particular is central to this idea of not "imposing unnecessary parametric forms," which is one reason why it's becoming so popular. We'll be introduced to what all this means when we discuss the correct model specification assumption.

² Loosely speaking: Causal inference is the "what?" Statistics is the "how?" Epidemiology is the "why?"

³ or the hazard ratio from a Cox model, or the mean ratio from a Poisson model, or host of other types of regression models

⁴ e.g., additivity, linearity, distributional, and other.

3 Introduction to Causal Inference

"Causal inference" deals primarily with the **formal mechanisms by which we can combine data, assumptions, and models to interpret a correlation (or association) causally**.⁵ The framework in which we define what we mean by "causal relation" or "causal effect" is the **potential outcomes framework**.

A central notion in the potential outcomes framework is the counterfactual. This notion stems from the intuitive and informal practice of interpreting cause-effect relations as **circumstances (e.g., health outcomes) that would have arisen had things (e.g., exposures) been different**.

While this intuition serves an important purpose, it is not sufficient for doing rigorous science. Suppose we ask: "what is the effect of smoking on the 5-year cumulative CVD risk, irrespective of smoking's effect on body weight?" This question may seem clear and intuitive. To answer this question, we would do a study in which we collect data, enter these into a computer, perform some calculations, and obtain a number (we'd usually like to interpret as the "effect").

⁵ There are a number of introductory books and articles on causal inference in the empirical sciences. Here are some excellent options: [Hernán and Robins \(Forthcoming\)](#), [Pearl et al. \(2016\)](#), [Imbens and Rubin \(2015\)](#), [Cunningham \(2021\)](#)

But there is a problem.⁶ The calculations performed by the computer are **rigorously defined (i.e., unambiguous) mathematical objects**. On the other hand, **English language sentences about cause effect relations are ambiguous**. For example, the “effect of smoking” can mean many different things:

- All people smoke any tobacco ever versus no people smoke tobacco ever.
- All people smoke 3 cigarettes per day versus all people smoke 2 cigarettes per day.
- All people who have smoked any tobacco in the last 15 years cease to smoke any tobacco whatsoever.

Similarly, “irrespective of” can mean a number of things:

- The effect of smoking on CVD risk that would be observed in a hypothetical world where smoking did not affect body mass?
- The effect of smoking on CVD risk if everyone were set to “normal” body mass?
- The effect of smoking on CVD risk if everyone were held at the body mass they had in the month prior to study entry?

But the numerical strings of data and the computer algorithms applied to these data are well defined mathematical objects, which do not admit such ambiguity. Depending on several choices, including the data, how variables are coded, and the modeling strategy, the computer is being told which question to answer. There is a lot of potential uncertainty in the space between the English language sentences we use to ask causal questions, and the computer algorithms we use to answer those questions. Causal inference is about clarifying this uncertainty.

4 Potential Outcomes Notation

The building blocks for causal inference are **potential outcomes** (Rubin, 2005).

Importantly, these are conceptually distinct from **observed outcomes**. That is, the outcome that one might observe in a dataset is not the same as the potential outcome.

Potential outcomes are functions of exposures. For a given exposure x , we will write the potential outcome as Y^x .⁷ **This is interpreted as “the outcome**

⁶ This problem was articulated by Robins 1987, and I am using a version of the example from his paper.

⁷ Alternate notation includes: Y_x , $Y(x)$, $Y \mid \text{Set}(X = x)$, and $Y \mid \text{do}(X = x)$.

(Y) that would be observed if X were set to some value x ". For example, if X is binary [denoted $X \in (0, 1)$], then Y^x is the outcome that would be observed if $X = 0$ or $X = 1$. If we wanted to be specific about the value of x , we could write $Y^{x=0}$ or $Y^{x=1}$ (or, more succinctly, Y^0 or Y^1).

In the survival setting, the notation for potential outcomes is usually modified. This is because the standard notation in survival analysis does not typically include the outcome Y directly in the notation. Recall the cumulative risk and survival functions are denoted $F(t)$ and $S(t)$, respectively.

To denote potential outcomes in the survival setting, one would typically write $S_{T^x}(t)$, $F_{T^x}(t)$, or some variation thereof (e.g., [Hernán et al., 2005](#)). The important thing is that, once a potential outcome is written, it must be defined appropriately. In other words, one should always include a sentence such as: "where $S_{T^x}(t)$ denotes the probability of surviving past time $T > t$ that would be observed if the exposure were set to $X = x$."

Similarly, when the exposure and/or outcome are measured repeatedly over follow-up, notation must account for that. We thus use subscripts to denote when the variable was measured. For example, if the exposure is measured three times, we can denote the first measurement X_0 , the second X_1 , and the third X_2 . Additionally, we use **overbars** to denote the past history of a variable over follow-up time. For example, for the three time-point scenario, \overline{X}_1 denotes the set $\{X_0, X_1\}$ (i.e., the exposure values for the first time-point, the second time-point, but NOT the third time-point.) Finally, we sometimes need to index future exposures, which we do using **underbars**. For example, \underline{X}_1 denotes the set $\{X_1, X_2\}$ (i.e., the exposure values for the second time-point and the third time-point, but NOT the first time-point.)

More generally, for some arbitrary point over follow-up j , \overline{X}_j denotes $\{X_0, X_1, X_2, \dots, X_j\}$.⁸ In contrast, \underline{X}_j denotes $\{X_j, X_{j+1}, X_{j+2}, \dots, X_J\}$. Note here that we use capital " J " to denote the last follow-up time in a study.

We can then define potential outcomes as a function of these exposure histories and/or futures. For example:

$$\begin{aligned} Y^{\overline{x}_j=1} &= Y^{x_0=1, x_1=1, \dots, x_j=1} \\ Y^{\underline{x}_j=1} &= Y^{x_j=1, x_{j+1}=1, \dots, x_J=1} \\ Y^{\overline{x}_{j-1}=1, \underline{x}_j=0} &= Y^{x_0=1, x_1=1, \dots, x_{j-1}=1, x_j=0, x_{j+1}=0, \dots, x_J=0} \end{aligned}$$

⁸ Note that this notation presumes a discrete time framework, where study time is broken up into distinct chunks that can be indexed by $0, 1, \dots$. This is related to the issue of discrete versus continuous time survival outcomes mentioned in Section 1, Technical Note on page 6.

Note that these potential outcomes can also be reformulated in the context of survival outcomes. For instance, if we're interested in the effect of and exposure X on the cumulative risk $F(t)$ or survival function $S(t)$ of a time-to-event outcome T , we can write:

$$F_{T^x}(t), \quad S_{T^x}(t),$$

And we can make similar adaptations to the superscripted x indexing the T for time-dependent exposures.



Study Question:

Suppose you collect data from a single person and find that they are exposed. Can you interpret the outcome you observe to be the potential outcome that would have been observed had they been exposed? Why or why not?

5 Estimand, Estimator, Estimate

5.1 Estimand

Causal inference starts with a clear idea of the effect of interest (the target causal parameter, or **estimand**). We use potential outcomes to do this, but it is useful and important first to distinguish between estimands, estimators, and estimates.



Study Question:

You are familiar with the well known odds ratio equation for a 2×2 table: (ab/cd) . Is this an estimand, estimator, or estimate? Why?

The **estimand** is the (mathematical) object we want to quantify. It is, for example, the causal risk difference, risk ratio, or odds ratio for our exposure and outcome of interest. In our smoking CVD example, we might be interested in:

$$E(Y^1 - Y^0), \quad \frac{E(Y^1)}{E(Y^0)}, \quad \frac{Odds(Y^1 = 1)}{Odds(Y^0 = 1)},$$

where $Odds(Y^x = 1) = E(Y^x)/[1 - E(Y^x)]$, and where $E(\cdot)$ is the expectation operator taken with respect to the total population.⁹ There are

⁹ Throughout this course, if the outcome Y is binary, then $E(Y) \equiv P(Y = 1)$. Or, the expectation of Y is equivalent to the probability that $Y = 1$. This assumes that the binary outcome variable Y is coded as $\{0, 1\}$, and not, e.g., $\{1, 2\}$. For the more technically oriented,

$$E(Y) = \int y f(y) dy$$

where $f(y)$ is the probability density function of Y .

many other causal estimands besides these (effect of treatment on the treated, complier average causal effect, survivor average causal effect, stochastic effects, other).

Furthermore, the estimand need not always be causal (Casella and Berger, 2002). We may be interested in a statistical estimand, such as the conditional risk difference, risk ratio, or odds ratio:

$$E(Y | X = 1) - E(Y | X = 0), \quad \frac{E(Y | X = 1)}{E(Y | X = 0)}, \quad \frac{Odds(Y | X = 1)}{Odds(Y | X = 0)},$$

What's important is that one is clear about the objective. For example, in Naimi (2016) we defined counterfactual disparity measures as:

$$E(Y^m | X = 1) - E(Y^m | X = 0)$$

which is a mixed statistical and counterfactual estimand. It is a measure of disparity (statistical estimand) that would be observed if some variable M were set to a value m (counterfactual estimand).

The causal estimands presented above represent **average treatment effects** (on the risk difference, risk ratio, and odds ratio scale, respectively). This effect is sometimes referred to as a marginal treatment effect, because it averages (or marginalizes) the effect over the entire sample. For instance, if we consider the risk difference, it is easy to show that¹⁰

$$E(Y^1 - Y^0) = \frac{1}{N} \sum_{i=1}^N Y_i^1 - \frac{1}{N} \sum_{i=1}^N Y_i^0$$

However, we may want to estimate this effect in a subset of the population. For instance, $E(Y^1 - Y^0 | C = c)$ is the effect of $x = 1$ versus $x = 0$ among those with $C = c$. There are many different conditional treatment (in contrast to marginal) effects, this latter one being one of the simplest. Another common conditional treatment effect is the effect of treatment on the treated (ETT):

$$E(Y^1 - Y^0 | X = 1)$$

This effect compares the outcomes that would be observed if the exposure were set to 1 (Y^1) versus if the exposure were set to 0 (Y^0) among those who were observed to be or actually exposed in the sample ($X = 1$).

¹⁰ Recall that Y^x is not the observed (or sample) value of the outcome, so how do we actually get this average? When we discuss identifiability, we will see how we use observed data to quantify these contrasts.

To illustrate the relevance of this effect, consider the following (entirely fictional) scenario: Suppose that during gestation of a high-risk pregnancy, two clinical options are available to manage the risk of fetal death: premature delivery induction versus expectant management. Suppose further a researcher is interested in quantifying the effect of inducing delivery prematurely on fetal and infant death. This researcher collects data on a cohort of high-risk pregnant women, including whether delivery was induced prematurely, fetal/infant death, and a host of confounding variables. All parties involved agree the study is designed perfectly (no confounding, measurement error, loss to follow-up). They calculate the average treatment effect of premature delivery induction on fetal and infant death on the risk difference scale:

$$E(Y^1 - Y^0) = 0.15$$

This researcher concludes that, if all high-risk pregnancies were induced prematurely ($X = 1$), 15 more out of every 100 pregnancies would end in death, relative to what would happen if all high-risk pregnancies were left to expectant management ($X = 0$). In light of this incredibly high excess risk of death, this researcher advises abandoning the practice of premature delivery induction entirely.

Another researcher questions the relevance of the average treatment effect. They argue that physicians would never induce delivery prematurely in all versus no high-risk pregnancies. Rather, the more interesting question is: **for those women whose pregnancies were actually induced**, what would the risk of death have been had they not been induced? Underlying this more nuanced question is an understanding that physicians may be inducing pregnancy in women because of a number of reasons that make these women different from the rest. These differences, in turn, can lead to a different effect. This researcher thus calculates the effect of treatment on the treated:

$$E(Y^1 - Y^0 \mid X = 1) = -0.05$$

This other researcher concludes that, among those whose pregnancies were actually delivered prematurely, the risk of death would have been higher had they not been delivered prematurely.

This hypothetical example demonstrates a fundamental difference between

the ATE and the ETT: for those high-risk pregnancies that were not induced prematurely, the act of inducing premature delivery would not be beneficial. But for those high-risk pregnancies that were induced prematurely, the act of inducing premature delivery was beneficial. The ATE averages the beneficial and non-beneficial effects in the entire population, to yield an overall non-beneficial effect. The ETT isolates the beneficial effect among those who actually received the intervention. Thus, in this hypothetical example, premature delivery actually did benefit those who received it, even though it would not benefit everybody.

There are many other estimands that can be defined, including the local average treatment effect [Angrist et al. \(1996\)](#), the survivor average causal effect [Tchetgen Tchetgen \(2014\)](#), the complier average causal effect [Shrier et al. \(2014\)](#), principal strata effects [Frangakis and Rubin \(2002\)](#), stochastic effects [Munoz and van der Laan \(2012\)](#), incremental propensity score effects [Naimi et al. \(2021\)](#), and others. We will not discuss these in the context of this course, but it's good to be aware of their existence.

5.2 Estimator

The estimand is the object we want to estimate. The **estimator** is an equation that allows us to use our data to quantify the estimand. Suppose, for example, we were explicitly interested in quantifying the causal risk difference for the relation between smoking and 5 year CVD risk. To do this, we **have to** start by quantifying the associational risk difference, but there are many ways to do this (e.g., ordinary least squares, maximum likelihood, or many others).

To be specific, let's simulate some hypothetical data on the relation between smoking and CVD. Let's look at ordinary least squares, maximum likelihood, the generalized method of moments, and augmented inverse probability weighting (AIPW) as estimators:

```
remotes::install_github("yqzhong7/AIPW")
library(AIPW)

install.packages("SuperLearner", repos = "https://cloud.r-project.org/", dependencies=TRUE)

##
```

```
## The downloaded binary packages are in
## /var/folders/z_/cty0tpg97wz_x1d1zgdhwllr0000gs/T//Rtmp9FuIE6/downloaded_packages
```

```
library(SuperLearner)

# define the expit function
expit<-function(z){1/(1+exp(-(z)))}
set.seed(123)
n<-1e6
confounder<-rbinom(n,1,.5)
smoking<-rbinom(n,1,expit(-2+log(2)*confounder))
CVD<-rbinom(n,1,.1+.05*smoking+.05*confounder)

# the data
head(data.frame(CVD,smoking,confounder))
```

```
##   CVD smoking confounder
## 1    0      0          0
## 2    0      0          1
## 3    1      0          0
## 4    1      0          1
## 5    0      0          1
## 6    0      0          0
```

```
round(mean(confounder),3)
```

```
## [1] 0.499
```

```
round(mean(smoking),3)
```

```
## [1] 0.166
```

```
round(mean(CVD),3)
```

```
## [1] 0.133
```

#OLS

```
round(coef(lm(CVD~smoking+confounder)),4)
```

```
## (Intercept)      smoking  confounder
##          0.1000      0.0485      0.0501
```

#ML1

```
round(coef(glm(CVD~smoking+confounder,family=poisson("identity"))),4)
```

```
## (Intercept)      smoking  confounder
##          0.0999      0.0487      0.0502
```

#ML2

```
round(coef(glm(CVD~smoking+confounder,family=binomial("identity"))),4)
```

```
## (Intercept)      smoking  confounder
##          0.1000      0.0487      0.0501
```

#GMM

```
# round(gmm(CVD~smoking+confounder,x=cbind(smoking, confounder))$coefficients,4)
```

#AIPW

```
AIPW_SL <- AIPW$new(Y = CVD,
                    A = smoking,
                    W = confounder,
                    Q.SL.library = c("SL.mean","SL.glm"),
                    g.SL.library = c("SL.mean","SL.glm"),
                    k_split = 3,
                    verbose=FALSE)$
```

```
fit()$
summary()
```

```
round(AIPW_SL$result[3,1],4)
```

```
## [1] 0.0488
```

Let's try to fit the same estimators using python:

```
## Need to install reticulate package to use python in R
install.packages("reticulate", repos='http://lib.stat.cmu.edu/R/CRAN', dependencies=T)
```

```
##
## The downloaded binary packages are in
## /var/folders/z_/cty0tpg97wz_x1d1zgdhwllr0000gs/T//Rtmp9FuIE6/downloaded_packages
```

```
library(reticulate)

## need to tell R where python is
use_python("/opt/homebrew/Caskroom/miniforge/base/envs/test_env/bin/python3")
```

Now we can run python directly from within RMarkdown:

```
import statsmodels.api as sm
```

```
## /Library/Frameworks/R.framework/Versions/4.1-arm64/Resources/library/reticulate/python/rpytools/load
## module = _import(
## /opt/homebrew/Caskroom/miniforge/base/envs/test_env/lib/python3.9/site-packages/statsmodels/compat/p
## from pandas import Int64Index as NumericIndex
```

```
import pandas as pd
from zepid.causal.doublyrobust import AIPTW

# defining the variables and data frames
x = r.python_data_x
y = r.python_data_y

# adding the constant term for the ols, glm, and gmm functions
x = sm.add_constant(x)

# fitting the ols and glm models
result_OLS = sm.OLS(y, x).fit()
result_GLM1 = sm.GLM(y, x, family=sm.families.Poisson(sm.families.links.identity())).fit()
```

```
## /opt/homebrew/Caskroom/miniforge/base/envs/test_env/lib/python3.9/site-packages/statsmodels/genmod/g
## warnings.warn((f"The {type(family.link).__name__} link function "
```

```
result_GLM2 = sm.GLM(y, x, family=sm.families.Binomial(sm.families.links.identity())).fit()
```

```
# no gmm
```

```
# fitting the aipw estimator (NB: key difference in contrast to R code above is that we are
# not using super learner here, nor are we using cross fitting)
```

```
## /opt/homebrew/Caskroom/miniforge/base/envs/test_env/lib/python3.9/site-packages/statsmodels/genmod/g
## warnings.warn((f"The {type(family.link).__name__} link function "
```

```
df = r.python_data
aipw = AIPTW(df, exposure='smoking', outcome='CVD')
# Treatment model
aipw.exposure_model('confounder')
# Outcome model
```

```
## =====
## Propensity Score Model
##                               Generalized Linear Model Regression Results
## =====
## Dep. Variable:                smoking    No. Observations:                1000000
## Model:                      GLM         Df Residuals:                  999998
## Model Family:                Binomial    Df Model:                      1
## Link Function:               Logit       Scale:                        1.0000
## Method:                      IRLS       Log-Likelihood:                 -4.4115e+05
## Date:                        Wed, 09 Feb 2022    Deviance:                      8.8229e+05
## Time:                        23:14:25          Pearson chi2:                   1.00e+06
## No. Iterations:              5            Pseudo R-squ. (CS):             0.01657
## Covariance Type:             nonrobust
## =====
##               coef      std err          z      P>|z|      [0.025      0.975]
## -----
```

```
## Intercept      -2.0089      0.004    -459.096      0.000      -2.017      -2.000
## confounder      0.7073      0.006     126.906      0.000      0.696      0.718
## =====
## =====
```

```
aipw.outcome_model('smoking + confounder')
# Calculating estimate
```

```
## =====
## Outcome Model
##              Generalized Linear Model Regression Results
## =====
## Dep. Variable:          CVD    No. Observations:          1000000
## Model:                  GLM    Df Residuals:              999997
## Model Family:           Binomial    Df Model:              2
## Link Function:          Logit    Scale:                  1.0000
## Method:                  IRLS    Log-Likelihood:         -3.8764e+05
## Date:                   Wed, 09 Feb 2022    Deviance:          7.7529e+05
## Time:                   23:14:25    Pearson chi2:         9.99e+05
## No. Iterations:         5    Pseudo R-squ. (CS):      0.009028
## Covariance Type:        nonrobust
## =====
##              coef      std err          z      P>|z|      [0.025      0.975]
## -----
## Intercept      -2.1853      0.005    -462.114      0.000      -2.195      -2.176
## smoking         0.3756      0.007     51.626      0.000      0.361      0.390
## confounder      0.4420      0.006     72.992      0.000      0.430      0.454
## =====
## =====
```

```
aipw.fit()

# printing the results from each method
print(result_OLS.params)
```

```
## const          0.100020
```

```
## smoking      0.048529
## confounder   0.050066
## dtype: float64
```

```
print(result_GLM1.params)
```

```
## const      0.099946
## smoking    0.048708
## confounder  0.050153
## dtype: float64
```

```
print(result_GLM2.params)
```

```
## const      0.099956
## smoking    0.048672
## confounder  0.050140
## dtype: float64
```

```
aipw.summary()
```

```
## =====
##           Augmented Inverse Probability of Treatment Weights
## =====
## Treatment:      smoking      No. Observations:      1000000
## Outcome:        CVD          No. Missing Outcome:    0
## g-Model:        Logistic     Missing Model:        None
## Q-Model:        Logistic
## =====
## Risk Difference:      0.049
## 95.0% two-sided CI: (0.047 , 0.051)
## -----
## Risk Ratio:          1.391
## 95.0% two-sided CI: (1.374 , 1.408)
## =====
```

In our simple setting with 1 million observations, ordinary least squares, maximum likelihood, the generalized method of moments, and AIPW yield the same associational risk difference (as expected) even though they are (for some, completely) different **estimators**.¹¹

It is important to note that these estimates are not causal risk differences, but are associational. Even the results from the AIPW estimator are *associational*, even though this method is much more clearly motivated from within the causal inference framework (Robins and Greenland, 1994). To interpret them as causal effects, we have to evaluate whether we can **identify** the effect we want to estimate. We discuss this next.

5.3 Estimates

Finally, the values obtained from each estimation approach (~ 0.05) are our **estimates**.

6 Identifiability: Average Treatment Effect

In our simulation example, we estimated the associational (as opposed to causal) risk difference using four different estimators (ordinary least squares, two different maximum likelihood estimators, and AIPW). Estimating associations is all we can do with empirical data. Any time you use software to obtain a point estimate, you get an associational measure, irrespective of the method used.¹²

But our primary interest is often in causal quantities. In our simulated case, we want to estimate the causal risk difference for the effect of smoking on CVD. We can only do so if this causal risk difference is **identified**. Formally, *a parameter (e.g., causal risk difference) is identified if we can write it as a function of the observed data*.

The causal risk difference is defined as a contrast of potential outcomes. Referring back to our simulated example,¹³ we want to estimate the causal risk difference which is an example of an average treatment effect:

$$E(Y^1 - Y^0),$$

where Y^1 , Y^0 are the potential CVD outcomes that would be observed if smok-

¹¹ A slightly deeper dive into these concepts can be found in Naimi and Whitcomb (2020) Estimating Risk Ratios and Risk Differences Using Regression. Am J Epidemiol. 189(6):508-10

¹² This is true with ANY estimator, including IP-weighting, g computation, g estimation, or double robust approaches, such as AIPW (as demonstrated) or targeted maximum likelihood estimation.

¹³ To simplify the explanation here, I am ignoring the fact that we conditioned on (or adjusted for) confounders C . Of course, without adjusting for C , we get a confounded estimate. However, if we adjust for C , we no longer obtain the average treatment effect. Instead, we obtain the conditional treatment effect. There are important distinctions between average and conditional treatment effects that we will discuss in a subsequent section.

ing were set to 1 and 0, respectively. On the other hand, the associational risk difference is defined as a contrast of observed outcomes:

$$E(Y \mid X = 1) - E(Y \mid X = 0),$$

where each term in this equation is interpreted as the risk of CVD **among those who had** $X = x$.

The causal risk difference is identified if the following equation holds:

$$E(Y^x) = E(Y \mid X = x).$$

This equation says that the risk of CVD that would be observed if everyone were set to $X = x$ is equal to the risk of CVD that we observe among those with $X = x$. In this equation, the right hand side equation is written entirely in terms of observed data ($Y = 1$). The left hand side is a function of unobserved potential outcomes ($Y^x = 1$). Because potential outcomes are unobservable abstractions, this equivalence will only hold if we can make some assumptions.

6.1 Counterfactual Consistency

The first is **counterfactual consistency**, which states that the potential outcome that would be observed if we set the exposure to the observed value is the observed outcome (Hernán, 2005, Hernan and Taubman (2008), Hernán and VanderWeele (2011), VanderWeele and Hernán (2013)).¹⁴ Formally, counterfactual consistency states that:

$$\text{if } X = x \text{ then } Y^x = Y$$

The status of this assumption remains unaffected by the choice of analytic method (e.g., standard regression versus g methods). Rather, this assumption's validity depends on the nature of the exposure assignment mechanism.

One way to grasp what counterfactual consistency is about is to use the example of the “effect” of obesity on mortality (Hernan and Taubman, 2008). We know that obesity is associated with an increased risk of mortality, but interpreting this excess risk into a causal statement is tricky. In an observational study, the association between obesity and mortality is obtained by contrasting

¹⁴ While somewhat convoluted, this assumption is primarily about legitimizing the connection between our observational study, and future interventions in actual populations based on this study. In our observational study, we **see** people with with a certain value of the exposure. In a future intervention, we **set** people to a certain value of the exposure. The differences between seeing and setting can be profound.

the risk of mortality among, say, obese versus non-obese individuals. However, causally acting on this information would require us to find a way to make obese individuals non-obese. This might consist of getting obese individuals to diet, exercise, start smoking, or to undergo a single leg amputation (!). Each of these interventions could reduce BMI, and thus getting obese individuals to become non-obese. However, each intervention will likely have (dramatically) different effects on mortality.

The key here is that obesity is not a manipulable construct (on the other hand, dieting, exercise, smoking, and leg amputation, more or less, are). As a result, precisely translating what we mean by “the effect of obesity” is difficult. The same problem arises with other variables, such as the “effect of education,” the “effect of race/ethnicity,” and the “effect of socioeconomic status,” to name a few (Naimi and Kaufman, 2015).

6.2 Interference

We must also assume **no interference**, which states that the potential outcome for any given individual does not depend on the exposure status of another individual (Hudgens and Halloran, 2008, Naimi and Kaufman (2015)). If this assumption were not true, we would have to write the potential outcomes as a function of the exposure status of multiple individuals. For example, for two different people indexed by i and j , we might write: $Y_i^{x_i, x_j}$.¹⁵ Notation and methods that account for interference can become very complex very quickly (Tchetgen Tchetgen and VanderWeele, 2012, Halloran and Hudgens (2016), Hudgens and Halloran (2008)). As a result, we will not consider the impact of interference here, except only to say that different estimands and estimators should be used to properly account for them.

Together, counterfactual consistency and no interference allow us to make some progress in writing the potential risk $E(Y^x)$ as a function of the observed risk $E(Y \mid X = x)$. Specifically, by counterfactual consistency and no interference, we can do the following:

$$E(Y^x) = E(Y \mid X = x) \tag{1}$$

$$= E(Y^x \mid X = x) \tag{2}$$

¹⁵ Together, counterfactual consistency and no interference make up the stable-unit treatment value assumption (SUTVA), first articulated by Rubin (1980).

6.3 Exchangeability

A third assumption is **exchangeability**, which implies that the potential outcomes under a specific exposure (Y^x) are independent of the observed exposures X (Greenland and Robins, 1986, Greenland et al. (1999), Greenland and Robins (2009)). To explain the intuition behind exchangeability (Hernán and Robins, Forthcoming), consider a setting in which we are estimating the effect of aspirin on headache incidence in a cohort of individuals aged 18-40 years.¹⁶ To do this experiment, a researcher randomly assigns 50% of the cohort to aspirin, and the remaining 50% to placebo. However, to overcome some logistical complications, before actually giving them aspirin/placebo, this researcher hands out cards that indicate whether the participant was assigned to aspirin (red card) versus placebo (blue card).

¹⁶ Assume that our sample size is sufficiently large so as to avoid any sampling variability problems.

After the cards/aspirin/placebo are distributed and the follow-up period transpires, the researcher tallies up the number of headaches in each exposure group. He finds the following results:

$$\text{Aspirin (Red Card): } E(Y \mid X = 1) = 0.6$$

$$\text{Placebo (Blue Card): } E(Y \mid X = 0) = 0.1$$

However, after reviewing the study protocol, he realizes that he accidentally assigned placebo to those with the red card, and aspirin to those with the blue card, instead of the other way around. Fortunately, this has no actual impact on the study, with the exception of needing to switch the aspirin label with the placebo label. Why? Randomization (in a sufficiently large enough sample) creates independencies between outcome that would be observed under some exposure value (the potential outcome) and the observed exposure. In our case, $E(Y^{x=1}) = 0.1$, and this is the case whether the exposure received was placebo ($X = 0$) or aspirin ($X = 1$):

$$E(Y^{x=1}) = 0.1 \implies \begin{cases} E(Y^{x=1} \mid X = 1) = 0.1 \\ E(Y^{x=1} \mid X = 0) = 0.1 \end{cases}$$

Thus, because of randomization the following mathematical relation is implied:

$$E(Y^x \mid X) = E(Y^x) \quad (3)$$

which is exactly what we need to progress the identifiability statement above:

$$E(Y^x) = E(Y \mid X = x) \quad (4)$$

$$= E(Y^x \mid X = x) \text{ by consistency and no interference} \quad (5)$$

$$= E(Y^x) \text{ by exchangeability} \quad (6)$$



Study Question:

Why is the word “exchangeable” used to describe this concept? What, precisely, is being “exchanged”?

6.4 Conditional Exchangeability

With exchangeability, we are able to drop the observed exposure on the right side of the conditioning statement. However, we motivated this exchangeability assumption via simple randomization. What about when we have an observational study where the exposure is not randomized? It turns out that the validity of results from an observational study still rests upon the idea of randomization. For example, if we conduct an analysis in observational data where we adjust for 3 confounding variables, and we believe these three variables are sufficient to control for all confounding (and there are no other threats to validity, such as selection or information bias), then we can show that the same set of steps required to equate the average potential outcomes $E(Y^x)$ with the average observed outcome among those with $X = x$: $E(Y \mid X = x)$.

Consider our aspirin and headache example above, instead rather than randomly assign 50% of the individuals to aspirin and 50% to placebo, imagine that for people who in an average week sleep < 7 hours per night, we use a coin that chooses heads 75% if the time to assign aspirin, and 25% of the time to assign placebo. And for people who sleep ≥ 7 hours per night, we use a 50:50 coin to assign aspirin and placebo.

Using an aspirin:placebo assignment proportion of 75:25 for “non-sleepers”, and 50:50 for “sleepers” creates an association between sleeping quantity and aspirin assignment. If sleeping quantity also has an association with headache, what we’ve done is created a confounding relation between aspirin versus placebo and headache via sleeping quantity. Because of this

confounding relation, we can no longer re-write the conditional expectation $E(Y^x \mid X = x)$ as $E(Y^x)$.

However, if we adjust for sleeping quantity in our analysis, we can partly recover the procedure we need to equate these quantities:

$$E(Y^x) = \sum_c E(Y \mid X = x, C) \quad (7)$$

$$= \sum_c E(Y^x \mid X = x, C) \text{ by consistency and no interference} \quad (8)$$

$$= \sum_c E(Y^x \mid C) \text{ by conditional exchangeability} \quad (9)$$

$$= E(Y^x) \text{ by marginalization} \quad (10)$$

The only difference is that now we have to incorporate an additional step in which we “average” or marginalize over the distribution of C to obtain a weighted average of the $E(Y^x)$ in the sample or population.



Technical Note:

Consider the marginalization step in the identification equation above. This step involves transitioning from $\sum_c E(Y^x \mid C)$ to $E(Y^x)$. This simply denotes taking a weighted sum of $E(Y^x \mid C)$, where the weights are defined as a probability function of C . For example, if $C \in \{0, 1, \dots, k\}$, then this sum becomes:

$$E(Y^x \mid C = 0)P(C = 0) + E(Y^x \mid C = 1)P(C = 1) + \dots + E(Y^x \mid C = k)P(C = k)$$

More generally (i.e., for a more general case where C is not necessarily categorical), we can rewrite this as:

$$E(E(Y^x \mid C))$$

where the outer expectation is taken over C , and the inner expectation is taken over Y^x . This equation is sometimes referred to as the law of iterated expectations, the law of total expectation, or the tower rule. It plays an important role in causal inference, such as when we define (and sometimes implement) the g computation estimator. It is useful to understand, both when reading the technical literature, as well as when implementing variations of the technique in software.

6.5 Correct Model Specification

Although it seems that we have successfully written the potential risk as a function of the observed data, we are in need of two more assumptions. The first is **correct model specification**. This assumption is required when we rely on models to estimate effects, which is particularly relevant in an observational study when we have several confounders we have to adjust for.

Consider the example above, where we had to adjust for C to equate the potential and observed outcomes. In our simple example, we only considered one confounding variable (sleep quantity), but in a typical observational study, we'd adjust for quite a few variables. Consider further that we'd typically employ a statistical regression model (e.g., linear, logistic, Poisson, Cox, or other) to actually implement our adjustment, which might look something like¹⁷:

$$E(Y \mid X, C_1, C_2, C_3, C_4) = \beta_0 + \beta_1 X + \beta_2 C_1 + \beta_3 C_2 + \beta_4 C_3 + \beta_5 C_4$$

The problem with using the above model is that it makes fairly strong assumptions about exactly *how* Y is related to X and the confounders. Specifically, this equation states (or assumes) that the conditional mean of Y is related to all the variables additively such that a single unit increase in each variable results in a linear and independent increase in the mean of Y .

However, consider that for five variables there can be a total of¹⁸

$$\binom{5}{2} = \frac{5!}{2!(5-2)!} = 10$$

two-way interactions that we could potentially add to the model. Additionally, we could include higher-order interactions, for example, a three-way interaction between X , C_1 , and C_3 . In fact, if we considered higher order interactions, for this simple model would could have up to:

$$2^5 - 5 - 1 = 26$$

k -way interactions (including 2, 3, 4, and 5 way). If we exclude any of the relevant interactions from among this set, our model would be misspecified. This misspecification could result in bias, which could create an dependence between the observed exposure and the potential outcomes, and would thus not

¹⁷ Such a model would be what we'd use in SAS, Stata, R, or any other software when we use the regression function and include only main effects terms in the model

¹⁸ This equation is referred to as the binomial coefficient.

allow us to equate the potential and observed outcomes the way we needed to assume exchangeability.

There are other ways in which this correct model specification assumption can be violated, including making incorrect linearity (or nonlinearity) assumptions, choosing the wrong link function in a generalized linear model (see, e.g., [Weisberg and Welsh, 1994](#)), or making the wrong distributional assumption about the conditional mean of the outcome. It is for these reasons (among others) that machine learning methods are becoming so popular. They do not make such (parametric) assumptions about how the data were generated. However, they do come with some important trade-offs that should be considered before use.



Study Question:

Can you come up with a clearly articulated connection between correct model misspecification and exchangeability?

6.6 Positivity

There is another assumption known as **positivity**,¹⁹ and requires exposed and unexposed individuals within all confounding levels ([Mortimer et al., 2005](#), [Westreich and Cole \(2010\)](#)). There are two kinds of positivity violations (non-positivity): structural (or deterministic) and stochastic²⁰ (or random).

Structural non-positivity occurs when individuals with certain covariate values cannot be exposed. For example, in occupational epidemiology work-status (employed/unemployed in workplace under study) is a confounder, but individuals who leave the workplace can no longer be exposed to a work-based exposure. Alternatively, stochastic non-positivity arises when the sample size is not large enough to populate all confounder strata with observations.

Problems because of positivity arise for two reasons. The first is definitional. Consider the step in our equation above where we marginalize over C to equate the potential and observed outcomes. In the case where C is binary and we want to estimate the potential outcome if everyone were exposed to $X = 1$, this step could be re-written as:

$$E(Y^{x=1}) = E(Y \mid X = 1, C = 1)P(C = 1) + E(Y \mid X = 1, C = 0)P(C = 0)$$

¹⁹ Also known as the experimental treatment assignment assumption.

²⁰ The word **stochastic** is derived from the greek word "to aim," as in "to aim for a target."

Now imagine that for those with $C = 1$, it is either impossible to have $X = 1$ (structural nonpositivity) or we just don't have anyone in our sample with $X = 1$ (stochastic nonpositivity). Mathematically, it does not make sense to write $E(Y \mid X = 1, C = 1)$ because there are no individuals with $X = 1$ and $C = 1$. We thus cannot define this conditional average.

The second problem with positivity violations has to do with estimators. Consider, for example, a simple inverse probability weight that corresponds to the above scenario (i.e., if $C = 1$, there are no individuals with $X = 1$):

$$\frac{1}{P(X = 1 \mid C = 1)}$$

In this case, the probability in the denominator is zero. And because $1/0$ is undefined, we can't use IP-weighting to estimate the effect we're after with this estimator. The same type of problem arises even if there are only a very small number people in the sample with $X = 1$ if $C = 1$. In this latter case, imagine that the probability of being exposed is very small, say 0.0001. Then, the above weight would be equivalent to $1/0.0001 = 10,000$. The above weight means that one or more of these individuals will contribute 10,000 observations to the weighted analysis (usually well more than the original sample). These types of problems result in instability of the estimator (because the results end up being heavily dependent on only a few individuals in the sample with large weights).

When faced with positivity violations, one should either re-define the estimand so that there is no positivity violation, choose an estimator that is less affected by positivity problems, or both (Petersen et al., 2012).²¹ Alternative estimands include the effect of treatment on the treated or untreated, various types of stochastic effects [including incremental propensity score effects (Kennedy, 2019), which do not require that positivity hold (Naimi et al., 2021)], or “blip” effects that are encoded in structural nested models, and can be estimated with g estimation. One can also use collaborative targeted minimum loss-based estimation,²² and the parametric g formula, which tend to be less sensitive to positivity violations (Cole et al., 2013; Porter et al., 2011; Ju et al., 2017).

There are a number of different procedures one can use to evaluate whether positivity is a problem. Among these include propensity score overlap plots. Consider again our data from the last section. To get the propensity score for

²¹ Keep in mind: one cannot simply “avoid” positivity. In extreme setting, nonpositivity means that those who were unexposed in the sample are very unlikely to be exposed (and vice versa). In such a situation, it may not make sense to estimate the average treatment effect, because there is a subset of the population who may never realistically be exposed (or unexposed). In this case, g estimation, cTMLE, and the parametric g formula can actually estimate parameters that differ slightly or profoundly from the ATE.

²² there is mounting evidence that standard (not collaborative) TMLE is very sensitive to positivity violations.

a binary exposure, we can fit a logistic model to the exposure data, conditional on confounders. Here, we use the Lalonde dataset, which is well known in econometric circles. This dataset was originally obtained from a study used to evaluate the effect of a training program (treat) on income:

```
library(MatchIt)
data("lalonde")

head(lalonde)
```

```
##      treat age educ  race married nodegree re74 re75      re78
## NSW1     1  37  11 black        1         1  0  0 9930.0460
## NSW2     1  22   9 hispan       0         1  0  0 3595.8940
## NSW3     1  30  12 black        0         0  0  0 24909.4500
## NSW4     1  27  11 black        0         1  0  0  7506.1460
## NSW5     1  33   8 black        0         1  0  0   289.7899
## NSW6     1  22   9 black        0         1  0  0 4056.4940
```

```
propensity_score <- glm(treat ~ age + educ +
  re75 + re78, data = lalonde, family = binomial(link = "logit"))$fitted.values

## by appending a '$fitted.values' to
## the end of this glm function, we are
## keeping the predicted values from
## the model under the observed data
## settings.
```

We can now plot the density of this propensity score for each exposure group to see how they overlap:

```
exposure <- lalonde$treat

plot_data <- data.frame(propensity_score,
  Exposure = as.factor(lalonde$treat))
```

```
p1 <- ggplot(data = plot_data) + scale_y_continuous(expand = c(0,
  0)) + scale_x_continuous(expand = c(0,
  0)) + ylab("Density") + xlab("Propensity Score") +
  scale_color_manual(values = c("#000000",
    "#D55E00")) + geom_density(aes(x = propensity_score,
    group = Exposure, color = Exposure)) +
  xlim(0, 1)

ggsave("../figures/2022_01_10-ps_overlap.pdf",
  plot = p1)
```

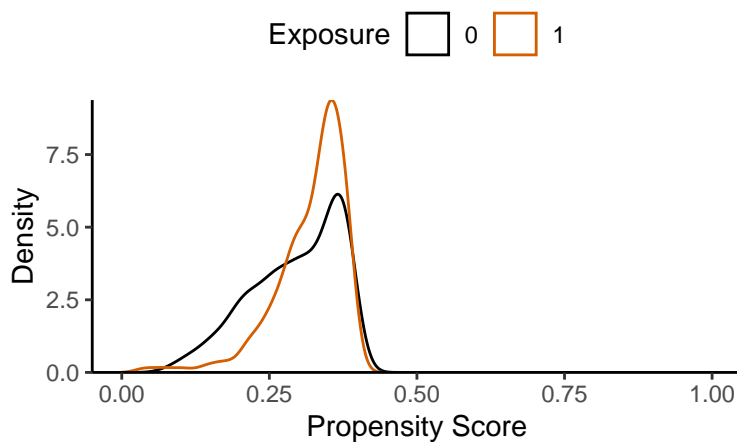


Figure 1: Propensity score overlap plot for the training intervention in 614 individuals in the Lalonde dataset.

Since the mass of the density for the exposed occurs in the same place as the density mass for the unexposed, positivity does not seem to be much of an issue here. Another way to check positivity is to create stabilized inverse probability weights²³ and look at their descriptive statistics.

```
sw <- (mean(exposure)/propensity_score) *
  exposure + ((1 - mean(exposure))/(1 -
    propensity_score)) * (1 - exposure)

summary(sw)
```

```
##      Min. 1st Qu.  Median    Mean 3rd Qu.    Max.
```

²³ We won't get too deep into the theory for / definition of weights here. But here is some code for creating stabilized weights and evaluating positivity.

```
## 0.7575 0.8811 0.9823 1.0080 1.0864 7.7241
```

The mean of the stabilized weights is 1, and the max weight is not large at all, suggesting very well-behaved weights. Thus, in this particular case, we are not concerned with violations of the positivity assumption.


Technical Note:

In a large body of methods literature, particularly econometrics, you are likely to encounter these causal assumptions articulated in different ways. Most commonly, researchers will often invoke **ignorability** as a core assumption in causal inference. There are at least two versions of ignorability: strong and weak. **Strong ignorability** is defined as the combination of the conditional independence assumption, and the positivity assumption. Technically, strong ignorability holds if, for individual i with a binary exposure $X \in \{0, 1\}$:

$$(Y_i^{x=0}, Y_i^{x=1}) \perp\!\!\!\perp X_i \mid \mathbf{C}_i, \text{ and} \\ 0 < P(X_i = 1 \mid \mathbf{C}_i) < 1,$$

where the $\perp\!\!\!\perp$ symbol denotes independence (in this case, conditional independence since we include $\mid \mathbf{C}_i$). In this case, the ignorability is “strong” because the independence is assumed to exist *jointly* between both potential outcomes $(Y_i^{x=0}, Y_i^{x=1})$ for individual i , and the exposure X_i , conditional on \mathbf{C}_i . Sometimes, a weaker version of the assumption is made:

$$(Y_i^x) \perp\!\!\!\perp X_i \mid \mathbf{C}_i, \text{ and} \\ 0 < P(X_i = 1 \mid \mathbf{C}_i) < 1,$$

This version of the assumption is weaker in that we need not worry about whether the potential outcomes are jointly independent of the observed exposure. Rather, we only need each potential outcome (Y_i^x) to be independent of the observed exposure.

Even still, these assumptions are stronger than what we need to identify the causal risk difference, risk ratio, odds ratio, or other typical summary contrasts we often quantify in epidemiology. In the above proof, we demonstrated that identifiability is obtained under **mean exchangeability** $E(Y^x \mid X, \mathbf{C}) = E(Y^x \mid \mathbf{C})$. This assumption is even weaker than those articulated in the formalization of strong and weak ignorability.

The distinctions between strong ignorability, weak ignorability, and mean exchangeability are of little practical consequence. While it is important to recognize that ignorability typically consists of the combination of exchangeability and positivity. Additional details on and intuition behind the different versions of exchangeability can be found in Technical Point 2.1 of ?.

7 Identifiability: Instrumental Variables and Local Average Treatment Effects

Average treatment effects (ATEs) are extremely common estimands in the empirical sciences. Researchers will often conduct randomized or observational studies, quantify treatment effects, and then interpret these treatment effects as ATEs. As we just encountered, the validity of this interpretation depends heavily on the above stated assumptions: counterfactual consistency, no interference, (conditional) exchangeability, correct model specification, and positivity. Each of these assumptions provides an opportunity for our causal interpretation to break down.

This has prompted a large body of research to develop on both what happens when these assumptions do not hold, and what we can do about it. For example, when counterfactual consistency is violated, we can consider re-evaluating our exposure or research question, and examine the possible “versions of treatment” ([VanderWeele, 2009](#), [Hernán and VanderWeele \(2011\)](#)). When interference is violated, we can look to evaluate direct, indirect, and spillover effects, which are defined specifically for treatments that result in interference ([Hudgens and Halloran, 2008](#)). Violations of the correct model specification assumption has lead to work on using machine learning for semi-parametric or nonparametric causal inference ([Naimi et al., 2022](#), [Kennedy \(2016\)](#)). When positivity is violated, there are several options one can pursue to either mitigate its impact ([Petersen et al., 2012](#)) or avoid the need for assuming positivity altogether ([Naimi et al., 2021](#), [Kennedy \(2019\)](#)).

Exchangeability can be violated when there is unadjusted confounding, selection bias, or information bias. By far, violations of exchangeability that result from unadjusted (unmeasured) confounding have taken up much of the literature’s focus. Fortunately, work on measurement error is increasing (?), and the impact of selection bias is arguably not as large as confounding and information bias ([Greenland, 2003](#)).

A natural question that one may ask is whether there is a way we can conduct an analysis with observational data without having to worry at all about the “no measured or unmeasured confounding” assumption. It turns out, we can use instrumental variables to do this. Because of this, instrumental variable estimators might be considered an “epidemiologist’s dream” ([Hernán and](#)

Robins, 2006). However, as with (in my experience, at least) pretty much everything, this strong benefit of the IV approach comes with important tradeoffs to consider we make when we seek to analyse data using IVs (Hernán and Robins, 2006).

7.1 Randomization and ITT

To begin our explanation of IV methods, let's revisit the randomized trial example from the section on identifying the ATE. Consider again the situation where we want to estimate the effect of aspirin on headache in a cohort of individuals aged 18-40 years. Again, we randomize 50% of the individuals in our sample to aspirin, and the other 50% to placebo. From this study, we can collect information on a number of things, beyond just the outcome of interest, and whether an individual was assigned to the treatment/placebo group.

In addition, we can measure whether or not the individual took the pill they were actually assigned too (aspirin or placebo). For instance, some may have decided to throw the pill away, others may have gotten distracted and forgot to take it, and others still may have had some GI discomfort, and opted simply not to take it. We can also measure variables that might predict whether they took the pill or not, and that predict whether they will experience a headache or not.

For example, we can again consider sleep quantity in the night previous. An individual who slept very little the night before may not feel up to taking their assigned treatment, or may have simply forgotten. Similarly, this same individual may be at a higher risk of experiencing headache. For this reason, we would collect information on sleep quantity in the night previous to randomization.

The example given in this scenario allows us to draw a DAG summarizing the relations between the treatment assignment indicator, taking the treatment/placebo, the sleep quantity confounder, and the outcome of interest. This DAG is shown in Figure 2:

This Figure shows some important elements that result from the design of a randomized trial. First, and most important, is the fact that X is an *exogenous*²⁴ variable. Under such exogeneity, we can perform a simple regression analysis to quantify the effect of X on Y . This is the primary reason why randomized trials and the intention-to-treat principle are often considered the "gold standard" for estimating causal effects.

²⁴ "Exogenous" is meant to connote a variable that is "outside of the system." In other words, the variable's influence does not depend on the distribution of other variables in the system. Exogenous variables can be used as a sort of "lever" against the complexity of a particular system. Randomization, in particular, leads to the presence of an exogenous variable because (in expectation) random allocation "breaks" any potential pre-existing relation between the randomized variable and other variables.

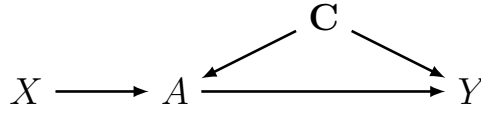


Figure 2: Directed Acyclic Graph for the relation between being assigned to treatment (X), adhering to treatment (A), a confounder (C) of the adherence-outcome relation, and the outcome (Y).

Unfortunately, as can be seen in Figure 2 the benefits of this exogeneity only directly extend to the *treatment assignment* variable X , and not to the indicator of whether they actually took the drug (in this example, aspirin A).²⁵ That is, we can obtain unconfounded estimates of the effect of X on Y without issue. We cannot do the same for A . Instrumental variable estimation is a technique developed to leverage the exogeneity of X to answer questions about the effect of A on Y .

²⁵ Recall, under intention-to-treat, we quantify the effect of being *assigned* to treatment, not the effect of treatment itself.

7.2 The Instrumental Variable Estimator

The IV estimator is based on a realization that information on the effect of A on Y defined as $E(Y^{a=1} - Y^{a=0})$ is contained in the following two effects:

$$\begin{aligned} E(Y^{x=1} - Y^{x=0}) \\ E(A^{x=1} - A^{x=0}) \end{aligned}$$

One can visualize this referring to Figure 2 by noting that the path $X \rightarrow Y$ is composed of paths $X \rightarrow A$ and $A \rightarrow Y$. Instrumental variables use information in $X \rightarrow Y$ and $X \rightarrow A$, both of which are identified without resorting to “no unmeasured confounding” assumptions, to estimate the effect of $A \rightarrow Y$.

To do this, we use the IV estimator, defined as:

$$\frac{E(Y \mid X = 1) - E(Y \mid X = 0)}{E(A \mid X = 1) - E(A \mid X = 0)}$$

Conceptually, this equation estimates the effect of X on A , and then removes this component from the estimated effect of X on Y . What remains is the effect of A on Y .

7.3 Local Average Treatment Effects

Instrumental variable estimators were introduced early in the 20th century, and were used extensively in economics. In 1996, Angrist, Imbens, and Rubin (Angrist et al., 1996) used potential outcomes notation to clearly define what we estimate when we implement IVs. They showed that, under a set of assumptions that differ from those for the ATE, the IV estimator quantifies (Angrist et al., 1996, Angrist and Pischke (2009)):

$$E(Y^{a=1} - Y^{a=0} \mid A^{x=1} > A^{x=0})$$

This effect is referred to as the *local average treatment effect*. It is “local” because it estimates the average treatment effect of A on Y in a subset of the population defined by $A^{x=1} > A^{x=0}$. This conditioning statement can be best understood by referring to the Table in Figure 3.

		$A^{x=0}$	
		0	1
$A^{x=1}$	0	$Y^{[x=1, A^{x=1}=0]} - Y^{[x=0, A^{x=0}=0]} = 0$	$Y^{[x=1, A^{x=1}=0]} - Y^{[x=0, A^{x=0}=1]} = -[Y^{x=1} - Y^{x=0}]$
		Never-Taker	Defier
$A^{x=1}$	1	$Y^{[x=1, A^{x=1}=1]} - Y^{[x=0, A^{x=0}=0]} = Y^{x=1} - Y^{x=0}$	$Y^{[x=1, A^{x=1}=1]} - Y^{[x=0, A^{x=0}=1]} = 0$
		Complier	Always-Taker

Referring back to our RCT example of taking aspirin and headache risk, the subset of the population defined by $A^{x=1} > A^{x=0}$ refers to individuals who would have taken aspirin if they were assigned to it, and who would have not taken aspirin if not assigned. These individuals are referred to as “compliers”. To see why, note that for compliers, $A^{x=1} = 1$ and $A^{x=0} = 0$, which, for a binary treatment A , is the only way in which the condition $A^{x=1} > A^{x=0}$ can be satisfied.²⁶ Thus, the LATE quantifies the average treatment effect among a subset of the sample who would have taken aspirin had they been assigned to it, AND who would have taken placebo had they been assigned to it.

Here we encounter one of the first problems with IV methods: the local average treatment effect. The problem here is that the conditioning statement is defined as a function of potential outcomes. Thus, the question that always arises is: who are these individuals? There is no way to empirically identify which individuals in the sample are “compliers.” As a result, it is often hard to generalize the LATE, because we don’t know who it applies to beyond the vague

Figure 3: Table of response types showing never-takers, always-takers, compliers, and defiers, as well as the effects (defined via potential outcomes) that result for each scenario. Table reproduced from Angrist et al 1996 JASA.

²⁶ We can also define the conditioning statement as: $A^{x=1} = 1, A^{x=0} = 0$. Note also that, though we’ve not said it yet, we really to require binary treatments X and A for this to work. While IV estimation has been generalized to the continuous setting (see Kennedy et al 2019 JRSS-B), these required slightly different assumptions and result in different interpretations.

abstraction that we refer to as “compliers.”

The Table in Figure 3 shows the compliers individuals in the bottom left cell. That is, under the most common formulation of the IV estimator,²⁷ the local average treatment effect is identical with the *complier average causal effect*.

²⁷ We will see precisely what I mean by this “formulation” next.

7.4 IV Identifiability Assumptions

Before we discuss how the IV estimators quantifies the LATE, and precisely why the LATE is the effect of treatment among compliers (CACE). We need to discuss precisely what constitutes an instrumental variable, and what assumptions are needed to identify the LATE using IVs. By definition, a variable X is referred to as an “instrumental variable” if:

Assumption 1) X has a causal effect on A .²⁸ Formally, this assumption can be represented as: $A^{x=1} - A^{x=0} \neq 0$

²⁸ Note that Hernán and Robins (2006) generalized this condition to be that X is associated with A .

Assumption 2) X affects Y only through A . Formally, this assumption can be represented as: $Y^{a,x=1} = Y^{a,x=0} = Y^a$, and is often referred to as the *exclusion restriction* assumption.

Assumption 3) X does not share common causes with the outcome (i.e., no confounding of the effect of X on Y). Formally, this is often represented as $Y^{A^{x=1},x=1}, Y^{A^{x=0},x=0}, A^{x=1}, A^{x=0} \perp\!\!\!\perp X$, which implies mean exchangeability [i.e., $E(Y^{A^{x=1},x=1} | X) = E(Y^{A^{x=1},x=1})$ and $E(A^{x=1} | X) = E(A^{x=1})$]²⁹

²⁹ See technical note on strong and weak ignorability.

With these assumptions, we can begin the process of re-writing the IV estimator above so that we obtain the LATE (CACE). In effect, we want to prove that³⁰:

³⁰ The material in this section comes from Angrist et al. (1996), Hernán and Robins (2006), and Angrist and Pischke (2009)

$$\frac{E(Y | X = 1) - E(Y | X = 0)}{E(A | X = 1) - E(A | X = 0)} = E(Y^{a=1} - Y^{a=0} | A^{x=1} > A^{x=0})$$

If we assume counterfactual consistency and no interference, we can re-write the terms in the numerator of the IV estimator as:

$$E(Y | X = 1) = E(Y^{a=A,x=1} | X = 1)$$

Recall that Assumption 2, the exclusion restriction, states that $Y^{a,x=1} = Y^{a,x=0} = Y^a$. Thus, under exclusion restriction, the right hand side of the above equation can be simplified as:

$$E(Y^{a=A, x=1} \mid X = 1) = E(Y^{a=A} \mid X = 1)$$

If we make Assumption 3, we can assume mean exchangeability, which allows us to write:

$$E(Y^{a=A} \mid X = 1) = E(Y^{a=A})$$

At this point, it's important to understand what $Y^{a=A}$ refers to (as opposed to Y^a). This is the potential outcome that would be observed if the exposure was set to their observed value. But in the actual dataset, this consists of two possible potential outcomes: $Y^{a=1}$ for those with $A = 1$, and $Y^{a=0}$ for those with $A = 0$. With this in mind, we can now slightly re-formulate the right hand side of the above equation as:

$$E(Y^{a=A}) = E[Y^{a=0} + (Y^{a=1} - Y^{a=0})A^{x=1}]$$

Note that if $A^{x=1} = 1$, then the expectation on the right hand side resolves to $Y^{a=1}$, whereas if $A^{x=1} = 0$ it resolves to $Y^{a=0}$. Thus, the equality holds.

What we've done so far is to show that, under counterfactual consistency and no interference, exclusion restriction (Assumption 2), and exchangeability (Assumption 3), we can re-write the first term in the numerator of the IV estimator as:

$$E(Y \mid X = 1) = E[Y^{a=0} + (Y^{a=1} - Y^{a=0})A^{x=1}]$$

The numerator of the IV estimator can thus be written as:

$$E[Y^{a=0} + (Y^{a=1} - Y^{a=0})A^{x=1}] - E[Y^{a=0} + (Y^{a=1} - Y^{a=0})A^{x=0}]$$

From here, we can do some basic re-arrangement of the terms so that the numerator of the IV estimator is:

$$E[(Y^{a=1} - Y^{a=0})(A^{x=1} - A^{x=0})]$$

which holds only under the assumptions listed (SUTVA + Assumptions 2 and 3). This latter equation states that the causal effect of X on Y is the product of the effect of A on Y and X on A .

Consider that the difference in potential outcomes ($A^{x=1} - A^{x=0}$) can only

yield three possible values: $\{1,0,-1\}$. Under Assumption 1, we have that $A^{x=1} - A^{x=0} \neq 0$, which means that we can make the following decomposition:

$$\begin{aligned} E[(Y^{a=1} - Y^{a=0})(A^{x=1} - A^{x=0})] \\ = E(Y^{a=1} - Y^{a=0} \mid A^{x=1} - A^{x=0} = 1)P(A^{x=1} - A^{x=0} = 1) \\ + E(Y^{a=1} - Y^{a=0} \mid A^{x=1} - A^{x=0} = 0)P(A^{x=1} - A^{x=0} = 0) \\ + E(Y^{a=1} - Y^{a=0} \mid A^{x=1} - A^{x=0} = -1)P(A^{x=1} - A^{x=0} = -1) \end{aligned} \quad (11)$$

In effect, this states that product of the effect of A on Y and X on A (the numerator of the IV estimator) is a weighted average of the ATE among the **compliers** and the **defiers**.

To complete the proof and show that the IV estimator can be connected to the CACE, we make one additional assumption: monotonicity: $A^{x=1} \geq A^{x=0}$. This assumption essentially states that there are no defiers in the population. Adding this assumption allows us to complete the proof. In effect, we have:

$$E[(Y^{a=1} - Y^{a=0})(A^{x=1} - A^{x=0})] = E(Y^{a=1} - Y^{a=0} \mid A^{x=1} - A^{x=0} = 1)P(A^{x=1} - A^{x=0} = 1)$$

Noting that $P(A^{x=1} - A^{x=0} = 1) = E(A^{x=1} - A^{x=0})$, we get:

$$\frac{E(Y \mid X = 1) - E(Y \mid X = 0)}{E(A \mid X = 1) - E(A \mid X = 0)} = \frac{E[(Y^{a=1} - Y^{a=0})(A^{x=1} - A^{x=0})]}{E(A^{x=1} - A^{x=0})} = E(Y^{a=1} - Y^{a=0} \mid A^{x=1} \geq A^{x=0})$$

7.5 Weak Instruments

7.6 Homogeneity Makes the LATE \rightarrow ATE

8 Non-Identifiability and Partial Identifiability: Bounding Effects

What happens when the effect we want to estimate is not identifiable? Suppose, for example, exchangeability is violated because we could not randomize our exposure and were aware of the absence of key (unmeasured) confounders? Or perhaps there was some loss to follow-up that could not be accounted for with absolute certainty? More likely there is both unmeasured confounding and loss to follow-up. When this happens, we get a point estimate for the causal effect of interest, but it could either be smaller or larger

in magnitude due to the influence of the unmeasured confounder and loss to follow-up.

In order to get a precise measure of **all the values the point estimate can possibly take** as a result of unmeasured confounding and loss to follow-up, we can estimate bounds for the point estimate of interest.³¹ Confidence intervals are bounds on the point estimate of interest that capture the uncertainty that results from random variation (Wasserman, 2004). In contrast, identification bounds capture the uncertainty that results from potential violations in some of the identification conditions (Manski, 2003).

³¹ Another way of phrasing this is: what range of point estimate values is compatible with the data?

Consider a study by Cole et al. (2019) in which they sought to quantify the effect of injection drug use on time to AIDS or death in a cohort of 1164 adult HIV-positive, AIDS-free women. These women were followed for AIDS or death up to 10 years from 12/6/95 in the Women's Interagency HIV Study (Barkan et al., 1998). Overall, 127 of 1164 women (11%) were lost to follow up. Adjusted risk differences were obtained via inverse probability weighting. Adjustment was made for age, race and nadir CD4 cell count.

Figure 4 shows the results from the analysis (obtained via personal communication with Stephen R. Cole; only a subset of these were presented in the manuscript). The top left panel shows the unadjusted risk difference over follow-up. The top right panel shows the corresponding risk difference after adjusting for loss to follow-up and measured confounders. The bottom left panel shows the identification bounds that result from loss to follow-up. And the bottom right panel shows the identification bounds that result from both loss to follow-up and unmeasured confounding. Specifically, the black area shows all possible risk differences that could arise given the data.

The bottom right panel in Figure 4 tells us something critically important that we often fail to consider when conducting an empirical study. Without assumptions, data alone rarely provide much information about a causal effect of interest. Rather, when we interpret that a point estimate from a statistical model as a causal effect estimate, we are invoking a whole set of assumptions (knowingly or unknowingly) that allow us to get a single number out of our data, rather than a range of possible values. One of these sets of assumptions we discussed here (counterfactual consistency, no interference, positivity, exchangeability, correct model specification). Nonparametric bounds such as those depicted in the study by Cole et al. (2019) help us understand exactly

Figure 4: Difference in risk of AIDS or death by injection drugs use, as a function of time on study, Women's Interagency HIV Study, 1995 to 2006. Courtesy of Stephen R. Cole.

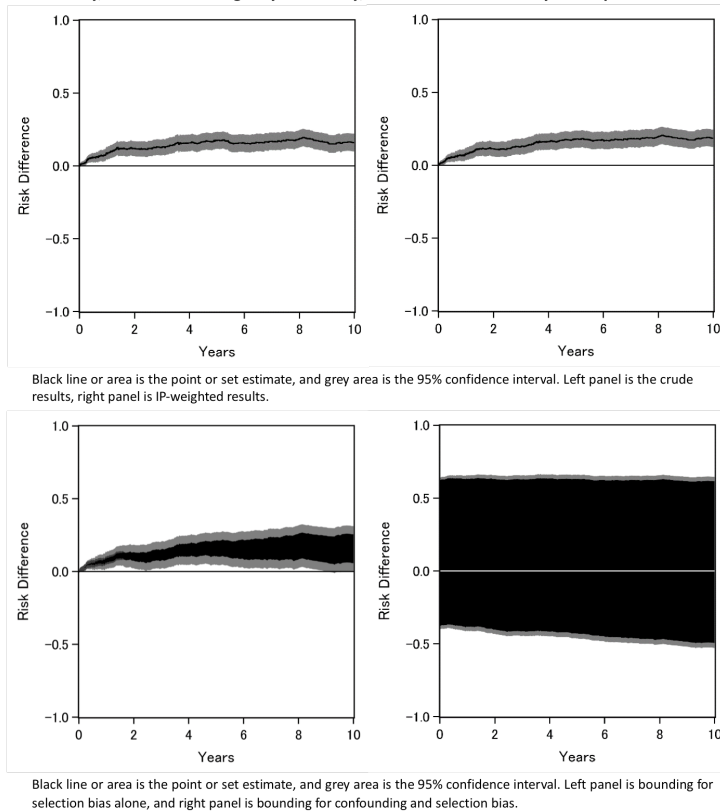


Figure 4: Bounds figure

how much support our data provide for an effect of interest, and how much of our results rely on unverifiable assumptions.

9 Takeaway

References

- Joshua D. Angrist and Jörn-Steffen Pischke. *Mostly Harmless Econometrics: An Empiricist's Companion*. Princeton University Press, January 2009. ISBN 978-0-691-12035-5. Google-Books-ID: YSAzEAAAQBAJ.
- Joshua D. Angrist, Guido W. Imbens, and Donald B. Rubin. Identification of causal effects using instrumental variables. *J Am Stat Assoc*, 91(434): 444–455, 1996.
- S E Barkan, S L Melnick, S Preston-Martin, K Weber, L A Kalish, P Miotti, M Young, R Greenblatt, H Sacks, and J Feldman. The women's interagency hiv study. wihs collaborative study group. *Epidemiology*, 9(2):117–125, Mar 1998.
- George Casella and Roger L. Berger. *Statistical Inference*. Duxbury Press, Pacific Grove, CA, 2nd edition, 2002.
- Stephen R. Cole, David B. Richardson, Haitao Chu, and Ashley I. Naimi. Analysis of occupational asbestos exposure and lung cancer mortality using the g formula. *Am J Epidemiol*, 177(9):989–996, 2013.
- Stephen R. Cole, Michael G Hudgens, Jessie K Edwards, M Alan Brookhart, David B Richardson, Daniel Westreich, and Adaora A Adimora. Nonparametric bounds for the risk function. *American Journal of Epidemiology*, 188(4): 632–636, 2019.
- Scott Cunningham. *Causal Inference: The Mixtape*. Yale University Press, New Haven, CT, 2021.
- Constantine E. Frangakis and Donald B. Rubin. Principal stratification in causal inference. *Biometrics*, 58(1):21–29, 2002.
- Sander Greenland. Quantifying biases in causal models: Classical confounding vs collider-stratification bias. *Epidemiol*, 14(3):300–306, 2003.

- Sander Greenland and James Robins. Identifiability, exchangeability and confounding revisited. *Epidemiol Perspect Innov*, 6(1):4, 2009.
- Sander Greenland and JM Robins. Identifiability, exchangeability, and epidemiological confounding. *Int J Epidemiol*, 15(3):413–419, 1986.
- Sander Greenland, James M. Robins, and Judea Pearl. Confounding and collapsibility in causal inference. *Stat Sci*, 14(1):29–46, 1999.
- M Elizabeth Halloran and Michael G Hudgens. Dependent happenings: A recent methodological review. *Curr Epidemiol Rep*, 3(4):297–305, Dec 2016.
- M. A. Hernán and JM Robins. *Causal Inference*. Chapman/Hall, Boca Raton, FL, Forthcoming.
- M A Hernan and S L Taubman. Does obesity shorten life? the importance of well-defined interventions to answer causal questions. *Int J Obes*, 32(S3):S8–S14, 2008.
- M. A. Hernán, S. R. Cole, J. Margolick, M. Cohen, and J. M. Robins. Structural accelerated failure time models for survival analysis in studies with time-varying treatments. *Pharmacoepidemiol Drug Saf*, 14(7):477–91, 2005.
- Miguel A. Hernán. Invited commentary: Hypothetical interventions to define causal effects—afterthought or prerequisite? *Am J Epidemiol*, 162(7):618–620, 2005.
- Miguel A Hernán and James M Robins. Estimating causal effects from epidemiological data. *J Epidemiol Community Health*, 60(7):578–586, 2006.
- Miguel A Hernán and Tyler J VanderWeele. Compound treatments and transportability of causal inference. *Epidemiol*, 22(3):368–377, May 2011. DOI: 10.1097/EDE.0b013e3182109296.
- M. G. Hudgens and M. E. Halloran. Toward causal inference with interference. *J Am Stat Assoc*, 103(482):832–842, 2008.
- Guido W Imbens and Donald B Rubin. *Causal Inference for Statistics, Social, and Biomedical Sciences: An Introduction*. Cambridge University Press, New York, NY, 2015.

- Cheng Ju, Susan Gruber, Samuel D Lendle, Antoine Chambaz, Jessica M Franklin, Richard Wyss, Sebastian Schneeweiss, and Mark J van der Laan. Scalable collaborative targeted learning for high-dimensional data. *Statistical Methods in Medical Research*, 28(2):532–554, 2017.
- Edward H Kennedy. Semiparametric theory and empirical processes in causal inference. In Hua He, Pan Wu, and Ding-Geng (Din) Chen, editors, *Statistical Causal Inferences and Their Applications in Public Health Research*. Springer International, Switzerland, 2016.
- Edward H. Kennedy. Nonparametric causal effects based on incremental propensity score interventions. *Journal of the American Statistical Association*, 114(526):645–656, 2019.
- Charles F Manski. *Partial identification of probability distributions*. Springer Science & Business Media, New York, NY, 2003.
- Kathleen M Mortimer, Romain Neugebauer, Mark van der Laan, and Ira B Tager. An application of model-fitting procedures for marginal structural models. *Am J Epidemiol*, 162(4):382–388, Aug 2005. doi: 10.1093/aje/kwi208.
- Ivan Diaz Munoz and Mark van der Laan. Population intervention causal effects based on stochastic interventions. *Biometrics*, 68(2):541–549, 2012.
- A I Naimi. The Counterfactual Implications of Fundamental Cause Theory. *Curr Epidemiol Reports*, In Press, 2016.
- AI Naimi, A Mishler, and Edward H. Kennedy. Challenges in obtaining valid causal effect estimates with machine learning algorithms. *Am J Epidemiol*, kwab201, 2022.
- Ashley I. Naimi and Jay S. Kaufman. Counterfactual theory in social epidemiology: Reconciling analysis and action for the social determinants of health. *Curr Epidemiol Reports*, 2(1):52–60, 2015.
- Ashley I. Naimi and Daniel J. Westreich. Big data: A revolution that will transform how we live, work, and think. *American Journal of Epidemiology*, 179(9): 1143–1144, 2014.
- Ashley I Naimi and Brian W Whitcomb. Estimating risk ratios and risk differences using regression. *American Journal of Epidemiology*, 189(6):508–510, 2020.

- Ashley I. Naimi, E Rudolph, H Kennedy, A Cartus, SI Kirkpatrick, DM Haas, H Simhan, and LM Bodnar. Incremental propensity score effects for time-fixed exposures. *Epidemiology*, 32(2):202–208, 2021.
- Judea Pearl, Madelyn R Glymour, and Nicholas Jewell. *Causal Inference in Statistics: A Primer*. Wiley, United Kingdom, 2016.
- Karl Pearson. *The Grammar of Science*. London, J.M. Dent & sons Ltd, 3rd edition, 1911.
- Maya L Petersen, Kristin E Porter, Susan Gruber, Yue Wang, and Mark J van der Laan. Diagnosing and responding to violations in the positivity assumption. *Stat Methods in Med Res*, 21(1):31–54, 2012.
- Kristin E Porter, Susan Gruber, Mark J van der Laan, and Jasjeet S Sekhon. The relative performance of targeted maximum likelihood estimators. *Int J Biostat*, 7(1), 2011.
- James M. Robins and Sander Greenland. Adjusting for differential rates of prophylaxis therapy for pcp in high-versus low-dose azt treatment arms in an aids randomized trial. *J Am Stat Assoc*, 89(427):737–749, 1994.
- Donald B. Rubin. Randomization analysis of experimental data: The fisher randomization test comment. *J Am Stat Assoc*, 75(371):591–593, 1980.
- Donald B Rubin. Causal inference using potential outcomes. *J Am Stat Assoc*, 100(469):322–331, 2005.
- Ian Shrier, Russell J Steele, Evert Verhagen, Rob Herbert, Corinne A Riddell, and Jay S Kaufman. Beyond intention to treat: what is the right question? *Clin Trials*, 11(1):28–37, Feb 2014. DOI: 10.1177/1740774513504151.
- Eric J. Tchetgen Tchetgen. Identification and estimation of survivor average causal effects. *Stat Med*, 33(21):3601–3628, 2014.
- Eric J Tchetgen Tchetgen and Tyler J VanderWeele. On causal inference in the presence of interference. *Stat Methods in Med Res*, 21(1):55–75, 2012.
- T. J. VanderWeele. Concerning the consistency assumption in causal inference. *Epidemiol*, 20(6):880–883, 2009.
- Tyler J VanderWeele and Miguel Ángel Hernán. Causal inference under multiple versions of treatment. *Journal of Causal Inference*, 1(1):1–20, 2013.

Larry Wasserman. *All of statistics: a concise course in statistical inference*. Springer, New York, 2004.

S. Weisberg and A. H. Welsh. Adapting for the missing link. *The Annals of Statistics*, 22(4):1674–1700, 1994.

Daniel Westreich and Stephen R. Cole. Invited commentary: Positivity in practice. *Am J Epidemiol*, 171(6):674–677, 2010.