

Propensity Score Methods, Models and Adjustment

Dr David A. Stephens

Department of Mathematics & Statistics
McGill University
Montreal, QC, Canada.

d.stephens@math.mcgill.ca
www.math.mcgill.ca/dstephens/SISCR2016/

Part 1: Introduction

- 1.1 The central causal question
- 1.2 Notation
- 1.3 Causal estimands
- 1.4 Basics of estimation
- 1.5 The Monte Carlo paradigm
- 1.6 Collapsibility
- 1.7 The randomized study
- 1.8 Confounding
- 1.9 Statistical modelling

Part 2: The Propensity Score

- 2.1 Manufacturing balance
- 2.2 The propensity score for binary exposures
- 2.3 Matching via the propensity score
- 2.4 The Generalized Propensity Score
- 2.5 Propensity score regression
- 2.6 Adjustment by weighting
- 2.7 Augmentation and double robustness

Part 3: Implementation and Computation

3.1 Statistical tools

3.2 Key considerations

Part 4: Extensions

4.1 Longitudinal studies

4.2 The Marginal Structural Model (MSM)

Part 5: New Directions

5.1 New challenges

5.2 Bayesian approaches

Part 1

Introduction

The central causal question

In many research domains, the objective of an investigation is to quantify the effect on a measurable outcome of changing one of the conditions under which the outcome is measured.

- ▶ in a health research setting, we may wish to discover the benefits of a new therapy compared to standard care;
- ▶ in economics, we may wish to study the impact of a training programme on the wages of unskilled workers;
- ▶ in transportation, we may attempt to understand the effect of embarking upon road building schemes on traffic flow or density in a metropolitan area.

The central statistical challenge is that, unless the condition of interest is changed independently, the inferred effect may be subject to the influence of other variables.

The central causal question

Example: The effect of nutrition on health

In a large cohort, the relationship between diet and health status is to be investigated. Study participants are queried on the nutritional quality of their diets, and their health status in relation to key indicators is assessed via questionnaires.

For a specific outcome condition of interest, incidence of cardiovascular disease (CVD), the relation to a specific dietary component, vitamin E intake, is to be assessed.

In the study, both incidence of disease and vitamin E intake were dichotomized

- ▶ Exposure: Normal/Low intake of vitamin E.
- ▶ Outcome: No incidence/Incidence of CVD in five years from study initiation.

The central causal question

Example: The effect of nutrition on health

		Outcome	
		No CVD	CVD
Exposure	Normal	27	8020
	Low	86	1879

Question: does a diet lower in vitamin E lead to higher chance of developing CVD ? More specifically, is this a causal link ?

- ▶ that is, if we were to intervene to change an individual's exposure status, by how much would their risk of CVD change ?

The language of causal inference

We seek to quantify the effect on an outcome of changes in the value of an exposure or treatment.

- ▶ Outcome: could be
 - ▶ binary;
 - ▶ integer-valued;
 - ▶ continuous-valued.
- ▶ Exposure: could be
 - ▶ binary;
 - ▶ integer-valued;
 - ▶ continuous-valued.
- ▶ Study: could be
 - ▶ cross-sectional (single time point);
 - ▶ longitudinal (multiple time points), with single or multiple exposures.

We consider an intervention to change exposure status.

Notation

We adopt the following notation: let

- ▶ i index individuals included in the study;
- ▶ Y_i denote the outcome for individual i ;
- ▶ Z_i denote the exposure for individual i ;
- ▶ X_i denote the values of other predictors or covariates.

For a cross-sectional study, Y_i and Z_i will be scalar-valued; for the longitudinal case, Y_i and Z_i may be vector valued. X_i is typically vector-valued at each measurement time point.

We will treat these variables as random quantities, and regard them as samples from an infinite population, rather than a finite population.

Counterfactual or Potential Outcomes

In order to phrase causal questions of interest, it is useful to consider certain hypothetical outcome quantities that represent the possible outcomes under different exposure alternatives.

We denote by

$$Y_i(\mathbf{z})$$

the hypothetical outcome for individual i if we intervene to set exposure to \mathbf{z} .

$Y_i(\mathbf{z})$ is termed a counterfactual or potential outcome.

Counterfactual or Potential Outcomes

If exposure is binary, the pair of potential outcomes

$$\{Y_i(0), Y_i(1)\}$$

represent the outcomes that would result for individual i if that subject was not exposed, or exposed, respectively.

The observed outcome, Y_i , may be written in terms of the potential outcomes and the observed exposure, Z_i , as

$$Y_i = (1 - Z_i)Y_i(0) + Z_iY_i(1).$$

Counterfactual or Potential Outcomes

If exposure is multi-valued, the potential outcomes

$$\{Y_i(\mathbf{z}_1), Y_i(\mathbf{z}_2), \dots, Y_i(\mathbf{z}_d)\}$$

represent the outcomes that would result for individual i if that subject exposed to exposure level $\mathbf{z}_1, \mathbf{z}_2, \dots, \mathbf{z}_d$ respectively.

The observed outcome, Y_i , may then be written in terms of the potential outcomes and the observed exposure, Z_i , as

$$Y_i = \sum_{j=1}^d \mathbb{1}_{\mathbf{z}_j}(Z_i) Y_i(\mathbf{z}_j).$$

where $\mathbb{1}_{\mathcal{A}}(Z)$ is the indicator random variable for the set \mathcal{A} , with $\mathbb{1}_{\mathcal{A}}(Z) = 1$ if $Z \in \mathcal{A}$, and zero otherwise.

Counterfactual or Potential Outcomes

If exposure is continuous-valued, the potential outcomes

$$\{Y_i(\mathbf{z}), \mathbf{z} \in \mathcal{Z}\}$$

represent the outcomes that would result for individual i if that subject exposed to exposure level \mathbf{z} which varies in the set \mathcal{Z} .

Counterfactual or Potential Outcomes

Note 1

It is rare that we can ever observe more than one of the potential outcomes for a given subject in a given study, that is, for binary exposures it is rare that we will be able to observe both

$$Y_i(0) \quad \text{and} \quad Y_i(1)$$

in the same study.

In the previous example, we cannot observe the CVD outcome under both the assumption that the subject did and simultaneously did not have a low vitamin E diet.

This is the first fundamental challenge of causal inference.

Causal Estimands

The central question of causal inference relates to comparing the (expected) values of different potential outcomes.

We consider the causal effect of exposure to be defined by differences in potential outcomes corresponding to different exposure levels.

Note 2

This is a statistical, rather than necessarily mechanistic, definition of causality.

Binary Exposures

For a binary exposure, we define the causal effect of exposure by considering contrasts between $Y_i(0)$ and $Y_i(1)$; for example, we might consider

- ▶ Additive contrasts

$$Y_i(1) - Y_i(0)$$

- ▶ Multiplicative contrasts

$$Y_i(1)/Y_i(0)$$

Continuous Exposures

For a continuous exposure, we might consider the path tracing how $Y_i(z)$ changes as z changes across some relevant set of values.

This leads to a causal dose-response function.

Example: Occlusion Therapy for Amblyopia

We might seek to study the effect of occlusion therapy (patching) on vision improvement of amblyopic children. Patching ‘doses’ are measured in terms of time for which the fellow (normal functioning) eye is patched.

As time is measured continuously, we may consider how vision improvement changes for any relevant dose of occlusion.

Expected counterfactuals

In general, we are interested in population or subgroup, rather than individual level causal effects. The potential outcomes are random quantities. Therefore, we more typically consider expected potential outcomes

$$\mathbb{E}[Y_i(\mathbf{z})]$$

or contrasts of these quantities.

We might also consider subgroup conditional expected quantities

$$\mathbb{E}[Y_i(\mathbf{z})|i \in \mathcal{S}]$$

where \mathcal{S} is some stratum of interest in the general population.

We typically assume that subject i is randomly sampled from the population or stratum, so that these individual-level expectations are representative of the population.

Expected counterfactuals: binary exposure

For a binary exposure, we might consider the average effect of exposure (or average treatment effect, ATE) defined as

$$\mathbb{E}[Y_i(1) - Y_i(0)] = \mathbb{E}[Y_i(1)] - \mathbb{E}[Y_i(0)]$$

If the outcome is also binary, we note that

$$\mathbb{E}[Y_i(\mathbf{z})] \equiv \Pr[Y_i(\mathbf{z}) = 1]$$

so may also consider odds or odds ratios quantities

$$\frac{\Pr[Y_i(\mathbf{z}) = 1]}{\Pr[Y_i(\mathbf{z}) = 0]} \qquad \frac{\Pr[Y_i(1) = 1] / \Pr[Y_i(1) = 0]}{\Pr[Y_i(0) = 1] / \Pr[Y_i(0) = 0]}.$$

Expected counterfactuals: binary exposure

We may also consider quantities such as the

average treatment effect on the treated, ATT

defined as

$$\mathbb{E}[Y_i(1) - Y_i(0)|Z_i = 1]$$

although such quantities can be harder to interpret.

The utility of the potential outcomes formulation is evident in this definition.

Example: antidepressants and autism

Antidepressants are quite widely prescribed and for a variety of mental health concerns. However, patients can be reluctant to embark on a course of antidepressants during pregnancy. We might wish to investigate, in a population of users (and potential users) of antidepressants, the incidence of autism-spectrum disorder in early childhood and to assess the causal influence of antidepressant use on this incidence.

- ▶ Outcome: binary, recording the a diagnosis of autism-spectrum disorder in the child by age 5;
- ▶ Exposure: antidepressant use during 2nd or 3rd trimester of pregnancy.

Then we may wish to quantify

$$\mathbb{E}[Y_i(\text{antidepressant}) - Y_i(\text{no antidepressant}) | \text{Antidep. actually used}].$$

Estimation of average potential outcomes

We wish to obtain estimates of causal quantities of interest based on the available data, which typically constitute a random sample from the target population.

We may apply to standard statistical principles to achieve the estimation. Typically, we will use sample mean type quantities: for a random sample of size n , the sample mean

$$\frac{1}{n} \sum_{i=1}^n Y_i$$

is an estimator of the population mean and so on.

Estimation of average potential outcomes

In a typical causal setting, we wish to perform estimation of average potential outcome (APO) values.

Consider first the situation where all subjects in a random sample receive a given exposure \mathbf{z} ; we wish to estimate $\mathbb{E}[Y(\mathbf{z})]$.

In terms of a formal probability calculation, we write this as

$$\begin{aligned}\mathbb{E}[Y(\mathbf{z})] &= \int y f_{Y(\mathbf{z})}(y) \, dy \\ &= \int y f_{Y(\mathbf{z}),X}(y, x) \, dy \, dx\end{aligned}$$

where the second line recognizes that in the population, the values of the predictors X vary randomly according to some probability distribution

Estimation of average potential outcomes

- ▶ $f_{Y(\mathbf{z})}(y)$ is the marginal distribution of the potential outcome when we set the exposure to \mathbf{z} .
- ▶ $f_{Y(\mathbf{z}),X}(y, x)$ is the joint distribution of $(Y(\mathbf{z}), X)$ in the population where we set the exposure to \mathbf{z} .

Note that we may also write

$$\mathbb{E}[Y(\mathbf{z})] = \int y \mathbb{1}_{\mathbf{z}}(z) f_{Y(\mathbf{z}),X}(y, x) dy dz dx$$

that is, imagining an exposure distribution degenerate at $z = \mathbf{z}$. Our random sample is from the population with density

$$\mathbb{1}_{\mathbf{z}}(z) f_{Y(\mathbf{z}),X}(y, x) \equiv \mathbb{1}_{\mathbf{z}}(z) f_{Y|Z,X}(y|z, x) f_X(x).$$

Estimation of average potential outcomes

Then we may estimate the relevant APO $\mathbb{E}[Y(\mathbf{z})]$ by

$$\widehat{\mathbb{E}}[Y(\mathbf{z})] = \frac{1}{n} \sum_{i=1}^n Y_i = \bar{Y}.$$

Note 3

To estimate functions of the sample mean, we may use simple transformations of the estimator; for example, if the outcome is binary, we estimate the odds

$$\frac{\Pr[Y_i(\mathbf{z}) = 1]}{\Pr[Y_i(\mathbf{z}) = 0]} \quad \text{by} \quad \frac{\bar{Y}}{1 - \bar{Y}}.$$

Monte Carlo methods

Causal quantities are typically average measures across a given population, hence we often need to consider integrals with respect to probability distributions.

Recall a simplified version of the calculation above: for any function $g(\cdot)$, we have

$$\begin{aligned}\mathbb{E}[g(Y)] &= \int g(y) f_Y(y) \, dy \\ &= \int g(y) f_{Y,X}(y, x) \, dy \, dx\end{aligned}$$

Rather than performing this calculation analytically using integration, we may consider approximating it numerically using Monte Carlo.

Monte Carlo methods

Monte Carlo integration proceeds as follows:

- ▶ generate a sample of size n from the density

$$f_Y(y)$$

to yield y_1, \dots, y_n ; there are standard techniques to achieve this.

- ▶ approximate $\mathbb{E}[g(Y)]$ by

$$\widehat{\mathbb{E}}[g(Y)] = \frac{1}{n} \sum_{i=1}^n g(y_i).$$

- ▶ if n is large enough, $\widehat{\mathbb{E}}[g(Y)]$ provides a good approximation to $\mathbb{E}[g(Y)]$

Note 4

This calculation is at the heart of frequentist methods in statistics:

- ▶ we collect a random sample of data of size n , and then form estimates based on this sample which often correspond to sample averages.
- ▶ if our sample is large enough, we are confident in our results.

Importance sampling

A variation on standard Monte Carlo is importance sampling:
by noting that

$$\begin{aligned}\mathbb{E}[g(Y)] &= \int g(y) f_Y(y) \, dy \\ &= \int g(y) \frac{f_Y(y)}{f_Y^*(y)} f_Y^*(y) \, dy\end{aligned}$$

we have that

$$\mathbb{E}_{f_Y}[g(Y)] \equiv \mathbb{E}_{f_Y^*} \left[g(Y) \frac{f_Y(Y)}{f_Y^*(Y)} \right].$$

Here $f_Y^*(y)$ is some other density from we are able to sample.

Importance sampling

Thus importance sampling proceeds as follows:

- ▶ generate a sample of size n from the density

$$f_Y^*(y)$$

to yield y_1, \dots, y_n ;

- ▶ approximate $\mathbb{E}[g(Y)]$ by

$$\widehat{\mathbb{E}}[g(Y)] = \frac{1}{n} \sum_{i=1}^n g(y_i) \frac{f_Y(y_i)}{f_Y^*(y_i)}.$$

Importance sampling

This means that even if we do not have a sample from the distribution of interest, f_Y , we can still compute averages with respect to f_Y if we have access to a sample from a related distribution, f_Y^* .

Clearly, for the importance sampling computation to work, we need that

$$\frac{f_Y(y_i)}{f_Y^*(y_i)}$$

is finite for the required range of Y , which means that we must have

$$f_Y^*(y) > 0 \quad \text{whenever} \quad f_Y(y) > 0.$$

Marginal and conditional measures of effect

Many of the causal measures described above are marginal measures, that is, they involve averaging over the distribution of predictors: for example

$$\begin{aligned}\mathbb{E}[Y(\mathbf{z})] &= \int y f_{Y(\mathbf{z}),X}(y, x) \, dy \, dx \\ &= \int y f_{Y(\mathbf{z})|X}(y|x) f_X(x) \, dy \, dx \\ &= \int y f_{Y|\mathbf{Z},X}(y|\mathbf{z}, x) f_X(x) \, dy \, dx\end{aligned}$$

Marginal and conditional measures of effect

Marginal measures are not typically the same as the equivalent measure defined for the conditional model

$$f_{Y|Z,X}(y|z, x)$$

Marginal measures that do not have the same interpretation in the conditional model are termed non-collapsible.

Example: logistic regression

Consider the binary response, binary exposure regression model, where

$$\Pr[Y = 1|Z = z, X = x] = \frac{\exp\{\beta_0 + \beta_1 z + \beta_2 x\}}{1 + \exp\{\beta_0 + \beta_1 z + \beta_2 x\}} = \mu(z, x; \beta)$$

say. We then have that in this conditional (on x) model, the parameter

$$\beta_1 = \log \left(\frac{\Pr[Y = 1|Z = 1, X = x]/\Pr[Y = 0|Z = 1, X = x]}{\Pr[Y = 1|Z = 0, X = x]/\Pr[Y = 0|Z = 0, X = x]} \right)$$

is the log odds ratio comparing outcome probabilities with for $Z = 1$ and $Z = 0$ respectively.

Example: logistic regression

In the marginal model, we wish to consider

$$\Pr[Y = 1|Z = z]$$

directly, and from the specified conditional model we have

$$\Pr[Y = 1|Z = z] = \int \Pr[Y = 1|Z = z, X = x]f_X(x) \, dx$$

assuming, for the moment, that Z and X are independent.
Explicitly,

$$\Pr[Y = 1|Z = z] = \int \mu(z, x; \beta)f_X(x) \, dx$$

Example: logistic regression

Typically, the integral that defines $\Pr[Y = 1|Z = z]$ in this way is not tractable. However, as, Y is binary, we may still consider a logistic regression model in the marginal distribution, say parameterized as

$$\Pr[Y = 1|Z = z] = \frac{\exp\{\theta_0 + \theta_1 z\}}{1 + \exp\{\theta_0 + \theta_1 z\}}$$

where θ_1 is the marginal log odds ratio.

In general, $\beta_1 \neq \theta_1$.

The randomized study

The approach that intervenes to set exposure equal to z for all subjects, however, does not facilitate comparison of APOs for different values of z .

Therefore consider a study design based on randomization; consider from simplicity the binary exposure case. Suppose that a random sample of size $2n$ is obtained, and split into two equal parts.

- ▶ the first group of n are assigned the exposure and form the ‘treated’ sample,
- ▶ the second half are left ‘untreated’.

The randomized study

For both the treated and untreated groups we may use the previous logic, and estimate the ATE

$$\mathbb{E}[Y_i(1) - Y_i(0)] = \mathbb{E}[Y_i(1)] - \mathbb{E}[Y_i(0)]$$

by the difference in means in the two groups, that is

$$\frac{1}{n} \sum_{i=1}^n Y_i - \frac{1}{n} \sum_{i=n+1}^{2n} Y_i.$$

The key idea here is that the two halves of the original sample are exchangeable with respect to their properties; the only systematic difference between them is due to exposure assignment.

The randomized study

In a slightly modified design, suppose that we obtain a random sample of size n from the study population, but then assign exposure randomly to subjects in the sample: subject i receives treatment with probability p .

- ▶ if $p = 1/2$, then there is an equal chance of receiving treatment or not;
- ▶ we may choose any value of $0 < p < 1$.

In the final sample, the number treated, n_1 , is a realization of a random variable N_1 where

$$N_1 \sim \text{Binomial}(n, p).$$

The randomized study

This suggests the estimators

$$\widehat{\mathbb{E}}[Y(\mathbf{z})] = \frac{\sum_{i=1}^n \mathbf{1}_{\mathbf{z}}(Z_i) Y_i}{\sum_{i=1}^n \mathbf{1}_{\mathbf{z}}(Z_i)} \quad \mathbf{z} = 0, 1 \quad (1)$$

where the indicators $\mathbf{1}_{\mathbf{z}}(Z_i)$ identify those individuals that received treatment \mathbf{z} .

The randomized study

Note that for the denominator,

$$\sum_{i=1}^n \mathbb{1}_1(Z_i) \sim \text{Binomial}(n, p)$$

so we may consider replacing the denominators by their expected values

$$np \quad \text{and} \quad n(1-p)$$

respectively for $z = 0, 1$. This yields the estimators

$$\widehat{\mathbb{E}}[Y(1)] = \frac{1}{np} \sum_{i=1}^n \mathbb{1}_1(Z_i) Y_i \quad \widehat{\mathbb{E}}[Y(0)] = \frac{1}{n(1-p)} \sum_{i=1}^n \mathbb{1}_0(Z_i) Y_i. \quad (2)$$

The randomized study

Note 5

The estimators in (1) are more efficient than the estimators in (2), that is, they have lower variances.

It is more efficient to use an estimated value of p

$$\hat{p} = \frac{N_1}{n}$$

than p itself.

The randomized study

We have that

$$\mathbb{E}[Y(\mathbf{z})] = \frac{\int y \mathbf{1}_{\mathbf{z}}(z) f_{Y|Z,X}(y|z, x) f_X(x) f_Z(z) \, dy \, dz \, dx}{\int \mathbf{1}_{\mathbf{z}}(z) f_Z(z) \, dz}$$

and we have taken the random sample from the joint density

$$f_{Y|Z,X}(y|z, x) f_X(x) f_Z(z)$$

which demonstrates that the estimators in (1) are akin to Monte Carlo estimators.

The challenge of confounding

The second main challenge of causal inference is that for non-randomized (or observational, or non-experimental) studies, exposure is not necessarily assigned according to a mechanism independent of other variables.

For example, it may be that exposure is assigned dependent on one or more of the measured predictors. If these predictors also predict outcome, then there is the possibility of *confounding* of the causal effect of exposure by those other variables.

The challenge of confounding

Specifically, in terms of densities, if predictor(s) X

- ▶ predicts outcome Y in the presence of Z :

$$f_{Y|Z,X}(y|z, x) \neq f_{Y|Z}(y|z)$$

and

- ▶ predicts exposure Z :

$$f_{Z|X}(z|x) \neq f_Z(z)$$

then X is a confounder.

Confounding: example

Example: The effect of nutrition on health: revisited

The relationship between low vitamin E diet and CVD incidence may be confounded by socio-economic status (SES); poorer individuals may have worse diets, and also may have higher risk of cardiovascular incidents via mechanisms other than those determined by diet:

- ▶ smoking;
- ▶ pollution;
- ▶ access to preventive measures/health advice.

Confounding

Confounding is a central challenge as it renders the observed sample unsuitable for causal comparisons unless adjustments are made:

- ▶ in the binary case, if confounding is present, the treated and untreated groups are not directly comparable;
- ▶ the effect of confounder X on outcome is potentially different in the treated and untreated groups.
- ▶ direct comparison of sample means does not yield valid insight into average treatment effects;

Causal inference is fundamentally about comparing exposure subgroups on an equal footing, where there is no residual influence of the other predictors. This is possible in the randomized study as randomization breaks the association between Z and X .

Note 6

Confounding is not the same as non-collapsibility.

- ▶ Non-collapsibility concerns the measures of effect being reported, and the parameters being estimated; parameters in a marginal model do not in general correspond to parameters in a conditional model.

Non-collapsibility is a property of the model, not the study design. It may be present even for a randomized study.

- ▶ Confounding concerns the inter-relationship between outcome, exposure and confounder. It is not model-dependent, and does depend on the study design.

Simple confounding example

Suppose that Y, Z and X are all binary variables. Suppose that the true (structural) relationship between Y and (Z, X) is given by

$$\mathbb{E}[Y|Z = z, X = x] = \Pr[Y = 1|Z = z, X = x] = 0.2 + 0.2z - 0.1x$$

with $\Pr[X = 1] = q$. Then, by iterated expectation

$$\mathbb{E}[Y(z)] = 0.2 + 0.2z - 0.1q$$

and

$$\mathbb{E}[Y(1) - Y(0)] = 0.2.$$

Simple confounding example

Suppose also that in the population from which the data are drawn

$$\Pr[Z = 1|X = x] = \begin{cases} p_0 & x = 0 \\ p_1 & x = 1 \end{cases} = (1 - x)p_0 + xp_1.$$

in which case

$$\Pr[Z = 1] = (1 - q)p_0 + qp_1.$$

Simple confounding example

If we consider the estimators in (2)

$$\widehat{\mathbb{E}}[Y(1)] = \frac{1}{np} \sum_{i=1}^n \mathbf{1}_1(Z_i)Y_i \quad \widehat{\mathbb{E}}[Y(0)] = \frac{1}{n(1-p)} \sum_{i=1}^n \mathbf{1}_0(Z_i)Y_i$$

and set $p = (1-q)p_0 + qp_1$, we see that for the first term

$$\begin{aligned}\mathbb{E}_{Y,Z,X}[\mathbf{1}_1(Z)Y] &= \mathbb{E}_{Z,X}[\mathbf{1}_1(Z)\mathbb{E}_{Y|Z,X}[Y|Z,X]] \\ &= \mathbb{E}_{Z,X}[\mathbf{1}_1(Z)(0.2 + 0.2Z - 0.1X)] \\ &= 0.2\mathbb{E}_X[\mathbb{E}_{Z|X}[\mathbf{1}_1(Z)|X]] \\ &\quad + 0.2\mathbb{E}_X[\mathbb{E}_{Z|X}[\mathbf{1}_1(Z)Z|X]] \\ &\quad - 0.1\mathbb{E}_X[X\mathbb{E}_{Z|X}[\mathbf{1}_1(Z)|X]]\end{aligned}$$

Simple confounding example

Now

$$\begin{aligned}\mathbb{E}_{Z|X}[\mathbf{1}_1(Z)|X] &= \mathbb{E}_{Z|X}[\mathbf{1}_1(Z)Z|X] \\ &\equiv \Pr[Z = 1|X] = (1 - X)p_0 + Xp_1\end{aligned}$$

and

$$\begin{aligned}\mathbb{E}_X[\mathbb{E}_{Z|X}[\mathbf{1}_1(Z)|X]] &= (1 - q)p_0 + qp_1 = p \\ \mathbb{E}_X[\mathbb{E}_{Z|X}[\mathbf{1}_1(Z)Z|X]] &= (1 - q)p_0 + qp_1 = p \\ \mathbb{E}_X[X\mathbb{E}_{Z|X}[\mathbf{1}_1(Z)|X]] &= qp_1\end{aligned}$$

and therefore

$$\mathbb{E}_{Y,Z,X}[\mathbf{1}_1(Z)Y] = 0.4p - 0.1qp_1.$$

Simple confounding example

$$\therefore \mathbb{E} \left[\frac{1}{np} \sum_{i=1}^n \mathbb{1}_1(Z_i) Y_i \right] = \frac{0.4p - 0.1p_1}{p}$$

By a similar calculation, as $\mathbb{1}_0(Z) = 1 - \mathbb{1}_1(Z)$,

$$\mathbb{E}_X[\mathbb{E}_{Z|X}[\mathbb{1}_0(Z)|X]] = 1 - p$$

$$\mathbb{E}_X[\mathbb{E}_{Z|X}[\mathbb{1}_0(Z)Z|X]] = 0$$

$$\mathbb{E}_X[X\mathbb{E}_{Z|X}[\mathbb{1}_0(Z)|X]] = q(1 - p_1)$$

so

$$\mathbb{E} \left[\frac{1}{n(1-p)} \sum_{i=1}^n \mathbb{1}_0(Z_i) Y_i \right] = \frac{0.2(1-p) - 0.1q(1-p_1)}{1-p}$$

Simple confounding example

Finally, therefore ATE estimator

$$\widehat{\mathbb{E}}[Y(1)] - \widehat{\mathbb{E}}[Y(0)]$$

has expectation

$$\frac{0.4p - 0.1qp_1}{p} - \frac{0.2(1-p) - 0.1q(1-p_1)}{1-p}$$

which equals

$$0.2 - 0.1q \left\{ \frac{p_1}{p} - \frac{1-p_1}{1-p} \right\}$$

and therefore the unadjusted estimator based on (2) is biased.

Simple confounding example

The bias is caused by the fact that the two subsamples with

$$Z = 0 \quad \text{and} \quad Z = 1$$

are not directly comparable - they have a different profile in terms of X ; by Bayes theorem

$$\Pr[X = 1|Z = 1] = \frac{p_1q}{p} \quad \Pr[X = 1|Z = 0] = \frac{(1 - p_1)q}{1 - p}$$

so, here, conditioning on $Z = 1$ and $Z = 0$ in turn in the computation of (2), leads to a different composition of X values in the two subsamples.

As X structurally influences Y , this renders the resulting Y values not directly comparable.

Instruments

If predictor \tilde{Z} predicts Z , but does not predict Y in the presence of Z , then \tilde{Z} is termed an instrument.

Example: Non-compliance

In a randomized study of a binary treatment, if Z_i records the treatment actually received by individual i , suppose that there is non-compliance with respect to the treatment; that is, if \tilde{Z}_i records the treatment assigned by the experimenter, then possibly

$$\tilde{z}_i \neq z_i.$$

Then \tilde{Z}_i predicts Z_i , but is not associated with outcome Y_i given Z_i .

Instruments

Instruments are not confounders as they do not predict outcome once the influence of the exposure has been accounted for.

Suppose in the previous confounding example, we had

$$\mathbb{E}[Y|Z = z, X = 0] = \Pr[Y = 1|Z = z, X = 1] = 0.2 + 0.2z$$

for the structural model, but

$$\Pr[Z = 1|X] = (1 - X)p_0 + Xp_1.$$

Then X influences Z , and there is still an imbalance in the two subgroups indexed by Z with respect to the X values, but as X does not influence Y , there is no bias if the ATE estimator based on (2) is used.

Critical Assumption

An important assumption that is commonly made is that of

No unmeasured confounding

that is, the measured predictors X include (possibly as a subset) all variables that confound the effect of Z on Y .

We must assume that all variables that simultaneously influence exposure and outcome have been measured in the study.

This is a strong (and possibly unrealistic) assumption in practical applications. It may be relaxed and the influence of unmeasured confounders studied in sensitivity analyses.

Model-based analysis

So far, estimation based on the data via (1) and (2) has proceeded in a nonparametric or model-free fashion.

- ▶ models such as

$$f_{Y(z),X}(y, x)$$

have been considered, but not modelled parametrically.

We now consider semiparametric specifications, specifically models where parametric models for example for

$$\mathbb{E}[Y(z)|X]$$

are considered.

Correct model specification

Suppose we posit an outcome conditional mean model

$$\mathbb{E}[Y|Z, X] = \mu(Z, X)$$

that may be parametric in nature, say

$$\mathbb{E}[Y|Z, X; \beta] = \mu(Z, X; \beta)$$

which perhaps might be linear in β , or a monotone transform of a linear function of β .

The importance of ‘no unmeasured confounders’

An important consequence of the no unmeasured confounders assumption is that we have the equivalence of the conditional mean structural and observed-data outcome models, that is

$$\mathbb{E}[Y(\mathbf{z})|X] \quad \text{and} \quad \mathbb{E}[Y|X, Z = \mathbf{z}]$$

when this model is correctly specified.

Inference under correct specification

We might (optimistically) assume that the model $\mathbb{E}[Y|Z, X]$ is correctly specified, and captures the true relationship.

If this is, in fact, the case, then

No special (causal) techniques are needed to estimate the causal effect.

That is, we may simply use regression of Y on (Z, X) using mean model $\mathbb{E}[Y|Z, X]$.

To estimate the APO, we simply set

$$\hat{\mathbb{E}}[Y(\mathbf{z})] = \frac{1}{n} \sum_{i=1}^n \mu(\mathbf{z}, X_i) \quad (3)$$

and derive other estimates from this: if $\mu(\mathbf{z}, x)$ correctly captures the relationship of the outcome to the exposure and confounders, then the estimator in (3) is consistent.

Inference under correct specification

The third challenge of causal inference is that

correct specification cannot be guaranteed.

- ▶ we may not capture the relationship between Y and (Z, X) correctly.

Part 2

The Propensity Score

Constructing a balanced sample

Recall the randomized trial setting in the case of a binary exposure.

- ▶ we obtain a random sample of size n of individuals from the target population, and measure their X values;
- ▶ according to some random assignment procedure, we intervene to assign treatment Z to individuals, and measure their outcome Y ;
- ▶ the link between X and Z is broken by the random allocation.

Constructing a balanced sample

Recall that this procedure led to the valid use of the estimators of the ATE based on (1) and (2).

The important feature of the randomized study is that we have, for confounders X (indeed all predictors)

$$f_{X|Z}(x|1) \equiv f_{X|Z}(x|0) \quad \text{for all } x,$$

or equivalently, in the case of a binary confounder,

$$\Pr[X = 1|Z = 1] = \Pr[X = 1|Z = 0].$$

Constructing a balanced sample

The distribution of X is balanced across the two exposure groups; this renders direct comparison of the outcomes possible. Probabilistically, X and Z are independent.

In a non-randomized study, there is a possibility that the two exposure groups are not balanced

$$f_{X|Z}(x|1) \neq f_{X|Z}(x|0) \quad \text{for some } x,$$

or in the binary case

$$\Pr[X = 1|Z = 1] \neq \Pr[X = 1|Z = 0].$$

If X influences Y also, then this imbalance renders direct comparison of outcomes in the two groups impossible.

Constructing a balanced sample

Whilst global balance may not be present, it may be that ‘local’ balance, within certain strata within the sample, may be present.

That is, for $x \in \mathcal{S}$ say, we might have balance; within \mathcal{S} , X is independent of Z .

$$f_{X|Z:\mathcal{S}}(x|1 : x \in \mathcal{S}) = f_{X|Z}(x|0 : x \in \mathcal{S})$$

Then, for individuals who have X values in \mathcal{S} , there is the possibility of direct comparison of the treated and untreated groups.

We might then restrict attention to causal statements relating to stratum \mathcal{S} .

Constructing a balanced sample

For discrete confounders, we might consider defining strata where the X values are precisely matched, and then comparing treated and untreated within those strata.

Consider matching strata $\mathcal{S}_1, \dots, \mathcal{S}_K$. We would then be able to compute the ATE by noting that

$$\mathbb{E}[Y(1) - Y(0)] = \sum_{k=1}^K \mathbb{E}[Y(1) - Y(0) | X \in \mathcal{S}_k] \Pr[X \in \mathcal{S}_k]$$

- ▶ $\mathbb{E}[Y(1) - Y(0) | X \in \mathcal{S}_k]$ may be estimated nonparametrically from the data by using (1) or (2) for data restricted to have $x \in \mathcal{S}_k$.
- ▶ $\Pr[X \in \mathcal{S}_k]$ may be estimated using the empirical proportion of x that lie in \mathcal{S}_k .

Constructing a balanced sample

For continuous confounders, we might consider the same strategy: consider matching strata $\mathcal{S}_1, \dots, \mathcal{S}_K$. Then the formula

$$\mathbb{E}[Y(1) - Y(0)] = \sum_{k=1}^K \mathbb{E}[Y(1) - Y(0) | X \in \mathcal{S}_k] \Pr[X \in \mathcal{S}_k]$$

still holds. However

- ▶ we must assume a model for how $\mathbb{E}[Y(1) - Y(0) | X \in \mathcal{S}_k]$ varies with x for $x \in \mathcal{S}_k$.

In both cases, inference is restricted to the set of X space contained in

$$\bigcup_{k=1}^K \mathcal{S}_k.$$

Constructing a balanced sample

In the continuous case, the above calculations depend on the assumption that the treatment effect is similar for x values that lie ‘close together’ in predictor (confounder) space. However

- I. Unless we can achieve exact matching, then the term ‘close together’ needs careful consideration.
- II. If X is moderate or high-dimensional, there may be insufficient data to achieve adequate matching to facilitate the estimation of the terms

$$\mathbb{E}[Y(1) - Y(0)|X \in \mathcal{S}_k];$$

recall that we need a large enough sample of treated and untreated subjects in stratum \mathcal{S}_k .

Nevertheless, matching in this fashion is an important tool in causal comparison.

Balance via the propensity score

We now introduce the important concept of the propensity score that facilitates causal comparison via a balancing approach.

Recall that our goal is to mimic the construction of the randomized study that facilitates direct comparison between treated and untreated groups. We may not be able to achieve this globally, but possibly can achieve it locally in strata of X space.

The question is how to define these strata.

Balance via the propensity score

Recall that in the binary exposure case, balance corresponds to being able to state that within \mathcal{S} , X is independent of Z :

$$f_{X|Z:\mathcal{S}}(x|1 : x \in \mathcal{S}) = f_{X|Z}(x|0 : x \in \mathcal{S})$$

This can be achieved if \mathcal{S} is defined in terms of a statistic, $e(X)$ say. That is, we consider the conditional distribution

$$f_{X|Z,e(X)}(x|z, e)$$

and attempt to ensure that, given $e(X) = e$, Z is independent of X , so that within strata of $e(X)$, the treated and untreated groups are directly comparable.

Balance via the propensity score

By Bayes theorem, for $z = 0, 1$, we have that

$$f_{X|Z,e(X)}(x|z, e) = \frac{f_{Z|X,e(X)}(z|x, e)f_{X|e(X)}(x|e)}{f_{Z|e(X)}(z|e)} \quad (4)$$

Now, as Z is binary, we must be able to write the density in the denominator as

$$f_{Z|e(X)}(z|e) = p(e)^z(1 - p(e))^{1-z} \quad z \in 0, 1$$

where $p(e)$ is a probability, a function of the fixed value e , and where $p(e) > 0$.

Balance via the propensity score

Therefore, in order to make the density $f_{X|Z,e(X)}(x|z, e)$ functionally independent of z , and achieve the necessary independence, we must have that

$$f_{Z|X,e(X)}(z|x, e) = p(e)^z(1 - p(e))^{1-z} \quad z \in 0, 1$$

also. But $e(X)$ is a function of X , so automatically we have that

$$f_{Z|X,e(X)}(z|x, e) \equiv f_{Z|X}(z|x).$$

Therefore, we require that

$$f_{Z|X}(z|x) = f_{Z|X}(z|x, e) = p(e)^z(1 - p(e))^{1-z} \equiv f_{Z|e(X)}(z|e)$$

for all relevant z, x .

Balance via the propensity score

This can be achieved by choosing the statistic

$$e(x) = f_{Z|X}(1|x) = \Pr_{Z|X}[Z = 1|x]$$

and setting $p(\cdot)$ to be the identity function, so that

$$f_{Z|X}(z|x) = e^z(1 - e)^{1-z} \quad z = 0, 1.$$

More generally, choosing $e(x)$ to be some monotone transform of $f_{Z|X}(1|x)$ would also achieve the same balance.

The corresponding random variable $e(X)$ defines the strata via which the causal calculation can be considered.

Balance via the propensity score

The function $e(x)$ defined in this way is the propensity score¹. It has the following important properties

- (i) as seen above, it is a balancing score; conditional on $e(X)$, X and Z are independent.
- (ii) it is a scalar quantity, irrespective of the dimension of X .
- (iii) in noting that for balance we require that

$$f_{Z|X}(z|x) \equiv f_{Z|e(X)}(z|e),$$

the above construction demonstrates that if $\tilde{e}(X)$ is another balancing score, then $e(X)$ is a function of $\tilde{e}(X)$;

- ▶ that is, $e(X)$ is the ‘coarsest’ balancing score.

¹ see Rosenbaum & Rubin (1983), *Biometrika*

Evaluating the propensity score

To achieve balance we must have

$$e(X) = \Pr[Z = 1|X]$$

correctly specified; that is, for confounders X , we must precisely specify the model $\Pr[Z = 1|X]$.

- ▶ If X comprises entirely discrete components, then we may be able to estimate $\Pr[Z = 1|X]$ entirely nonparametrically, and satisfactorily if the sample size is large enough.
- ▶ If X has continuous components, it is common to use parametric modelling, with

$$e(X; \alpha) = \Pr[Z = 1|X; \alpha].$$

Balance then depends on correct specification of this model.

Unconfoundedness given the propensity score

The assumption of ‘no unmeasured confounders’ amounts to assuming that the potential outcomes are jointly independent of exposure assignment given the confounders, that is

$$\{Y(0), Y(1)\} \perp Z \mid X$$

that is, in terms of densities

$$\begin{aligned} f_{Y(z), Z|X}(y, z|x) &= f_{Y(z)|X}(y|x) f_{Z|X}(z|x) \\ &= f_{Y|Z, X}(y|z, x) f_{Z|X}(z|x). \end{aligned}$$

Unconfoundedness given the propensity score

Now consider conditioning on propensity score $e(X)$ instead of X : we have by factorization that

$$f_{Y(z), Z|e(X)}(y, z|e) = \frac{1}{f_{e(X)}(e)} \int_{\mathcal{S}_e} f_{Y(z), Z, X}(y, z, x) \, dx$$

where \mathcal{S}_e is the set of x values

$$\mathcal{S}_e \equiv \{x : e(x) = e\}$$

that yield a propensity score value equal to the value e .

Unconfoundedness given the propensity score

Now we have by unconfoundedness given X that

$$f_{Y(z),Z,X}(y, z, x) = f_{Y(z)|X}(y|x)f_{Z|X}(z|x)f_X(x)$$

and on the set \mathcal{S}_e , we have

$$f_{Z|X}(z|x) = e^z(1 - e)^{1-z} \equiv f_{Z|e(X)}(z|e).$$

Unconfoundedness given the propensity score

Therefore, recalling the \mathcal{S}_e is defined via the fixed constant e ,

$$\begin{aligned}\int_{\mathcal{S}_e} f_{Y(z),Z,X}(y,z,x) \, dx &= \int_{\mathcal{S}_e} f_{Y(z)|X}(y|x) e^z (1-e)^{1-z} f_X(x) \, dx \\ &= e^z (1-e)^{1-z} \int_{\mathcal{S}_e} f_{Y(z)|X}(y|x) f_X(x) \, dx \\ &= f_{Z|e(X)}(z|e) f_{Y(z)|e(X)}(y|e).\end{aligned}$$

Hence

$$f_{Y(z),Z|e(X)}(y,z|e) = \frac{1}{f_{e(X)}(e)} f_{Z|e(X)}(z|e) f_{Y(z)|e(X)}(y|e)$$

and so

$$Y(z) \perp Z \mid e(X) \quad \text{for all } z.$$

Estimation using the propensity score

We now consider the same stratified estimation strategy as before, but using $e(X)$ instead X to stratify.

Consider strata $\mathcal{S}_1, \dots, \mathcal{S}_K$ defined via $e(X)$. In this case, recall that

$$0 < e(X) < 1$$

so we might consider an equal quantile partition, say using quintiles.

Then we have

$$\mathbb{E}[Y(1) - Y(0)] = \sum_{k=1}^K \mathbb{E}[Y(1) - Y(0) | e(X) \in \mathcal{S}_k] \Pr[e(X) \in \mathcal{S}_k]$$

still holds approximately if the \mathcal{S}_k are small enough.

Estimation using the propensity score

This still requires us to be able to estimate

$$\mathbb{E}[Y(1) - Y(0) | e(X) \in \mathcal{S}_k]$$

which requires us to have a sufficient number of treated and untreated individuals with $e(X) \in \mathcal{S}_k$ to facilitate the ‘direct comparison’ within this stratum.

If the expected responses are constant across the stratum, the formulae (1) and (2) may be used.

Matching

The derivation of the propensity score indicates that it may be used to construct matched individuals or groups that can be compared directly.

- ▶ if two individuals have precisely the same value of $e(x)$, then they are exactly matched;
- ▶ if one of the pair is treated and the other untreated, then their outcomes can be compared directly, as any imbalance between their measured confounder values has been removed by the fact that they are matched on $e(x)$;
- ▶ this is conceptually identical to the standard procedure of matching in two-group comparison.

Matching

For an exactly matched pair (i_1, i_0) , treated and untreated respectively, the quantity

$$y_{i_1} - y_{i_0}$$

is an unbiased estimate of the ATE

$$\mathbb{E}[Y(1) - Y(0)];$$

more typically we might choose m such matched pairs, usually with different $e(x)$ values across pairs, and use the estimate

$$\frac{1}{m} \sum_{i=1}^m (y_{i_1} - y_{i_0})$$

Matching

Exact matching is difficult to achieve, therefore we more commonly attempt to achieve approximate matching

- ▶ May match one treated to M untreated ($1 : M$ matching)
- ▶ caliper matching;
- ▶ nearest neighbour/kernel matching;
- ▶ matching with replacement.

Most standard software packages have functions that provide automatic matching using a variety of methods.

Beyond binary exposures

The theory developed above extends beyond the case of binary exposures.

Recall that we require balance to proceed with causal comparisons; essentially, with strata defined using X or $e(X)$, the distribution of X should not depend on Z .

We seek a scalar statistic such that, conditional on the value of that statistic, X and Z are independent. In the case of general exposures, we must consider balancing scores that are functions of both Z and X .

Beyond binary exposures

For a balancing score $b(Z, X)$, we require that

$$X \perp Z \mid b(Z, X).$$

We denote $B = b(Z, X)$ for convenience.

Consider the conditional distribution $f_{Z|X,B}(z|x, b)$: we wish to demonstrate that

$$f_{Z|X,B}(z|x, b) = f_{Z|B}(z|b) \quad \text{for all } z, x, b.$$

That is, we require that B completely characterizes the conditional distribution of Z given X .

Beyond binary exposures

This can be achieved by choosing the statistic

$$b(z, x) = f_{Z|X}(z|x)$$

in line with the choice in the binary case.

The balancing score defined in this way is termed the

Generalized Propensity Score

which is a balancing score for general exposures.

Beyond binary exposures

Note, however, that this choice that mimics the binary exposure case is not the only one that we might make. The requirement

$$f_{Z|X,B}(z|x,b) = f_{Z|B}(z|b)$$

for all relevant z, x is met if we define $b(Z, X)$ to be any sufficient statistic that characterizes the conditional distribution of Z given X .

It may be possible, for example, to choose functions purely of X .

Beyond binary exposures

Example: Normally distributed exposures

Suppose that continuous valued exposure Z is distributed as

$$Z|X = x \sim \text{Normal}(x\alpha, \sigma^2)$$

for row-vector confounder X . We have that

$$f_{Z|X}(z|x) = \frac{1}{\sqrt{2\pi\sigma^2}} \exp \left\{ -\frac{1}{2\sigma^2} (z - x\alpha)^2 \right\}$$

Beyond binary exposures

Example: Normally distributed exposures

We might therefore choose

$$b(Z, X) = \frac{1}{\sqrt{2\pi\sigma^2}} \exp \left\{ -\frac{1}{2\sigma^2} (Z - X\alpha)^2 \right\}.$$

However, the linear predictor

$$b(X; \alpha) = X\alpha$$

also characterizes the conditional distribution of Z given X ; if we know that $x\alpha = b$, then

$$Z|X = x \equiv Z|B = b \sim \text{Normal}(b, \sigma^2).$$

In both cases, parameters α are to be estimated.

Beyond binary exposures

The generalized propensity score inherits all the properties of the standard propensity score;

- ▶ it induces balance;
- ▶ if the potential outcomes and exposure are independent given X under the unconfoundedness assumption, they are also independent given $b(Z, X)$.

However, how exactly to use the generalized propensity score in causal adjustment for continuous exposures is not clear.

Propensity Score Regression

Up to this point we have considered using the propensity score for stratification, that is, to produce directly comparable groups of treated and untreated individuals.

Causal comparison can also be carried out using regression techniques: that is, we consider building an estimator of the APO by regressing the outcome on a function of the exposure and the propensity score.

Regressing on the propensity score is a means of controlling the confounding.

Propensity Score Regression

If we construct a model

$$\mathbb{E}[Y|Z = z, b(Z, X) = b] = \mu(z, b)$$

then because potential outcomes $Y(\mathbf{z})$ and Z are independent given $b(Z, X)$, we have

$$\mathbb{E}[Y(\mathbf{z})|b(Z, X) = b] = \mathbb{E}[Y|Z = \mathbf{z}, b(\mathbf{z}, X) = b] = \mu(\mathbf{z}, b)$$

and therefore

$$\mathbb{E}[Y(\mathbf{z})] = \mathbb{E}_{b(\mathbf{z}, X)}[\mathbb{E}[Y|Z = \mathbf{z}, b(\mathbf{z}, X)]] = \mathbb{E}_{b(\mathbf{z}, X)}[\mu(\mathbf{z}, b(\mathbf{z}, X))].$$

Propensity Score Regression

That is, to estimate the APO, we might

- ▶ fit the propensity score model $b(Z, X)$ to the observed exposure and confounder data by regressing Z on X ;
- ▶ fit the conditional outcome model $\mu(z, b)$ using the fitted $b(Z, X)$ values, $\hat{b}(z_i, x_i)$;
- ▶ for each z of interest, estimate the APO by

$$\frac{1}{n} \sum_{i=1}^n \hat{\mu}(z, \hat{b}(z, x_i)).$$

Propensity Score Regression

If the propensity function $b(Z, X) \equiv b(X)$, we proceed similarly, and construct a model

$$\mathbb{E}[Y|Z = z, b(X) = b] = \mu(z, b)$$

then

$$\mathbb{E}[Y(\mathbf{z})|b(X) = b] = \mathbb{E}[Y|Z = \mathbf{z}, b(X) = b] = \mu(\mathbf{z}, b)$$

and therefore

$$\mathbb{E}[Y(\mathbf{z})] = \mathbb{E}_{b(X)}[\mathbb{E}[Y|Z = \mathbf{z}, b(X)]] = \mathbb{E}_{b(X)}[\mu(\mathbf{z}, b(X))].$$

Propensity Score Regression

To estimate the APO:

- ▶ fit the propensity score model $b(X)$ to the observed exposure and confounder data by regressing Z on X ;
- ▶ fit the conditional outcome model $\mu(z, b)$ using the fitted $b(X)$ values, $\hat{b}(x_i)$;
- ▶ for each z of interest, estimate the APO by

$$\frac{1}{n} \sum_{i=1}^n \hat{\mu}(z, \hat{b}(x_i)).$$

Example: Binary Exposure

We specify

- ▶ $e(X; \alpha) = \Pr[Z = 1|X, \alpha]$ then regress Z on X to obtain $\hat{\alpha}$ and fitted values $\hat{e}(X) \equiv e(X; \hat{\alpha})$.
- ▶ $\mathbb{E}[Y|Z = z, e(X) = e; \beta] = \mu(z, e; \beta)$ and estimate this model by regressing y_i on z_i and $\hat{e}(x_i)$. For example, we might have that

$$\mathbb{E}[Y|Z = z_i, e(X_i) = e_i; \beta] = \beta_0 + \beta_1 z_i + \beta_2 e_i.$$

This returns $\hat{\beta}$.

We finally compute the predictions under this model, and average them to obtain the APO estimate

$$\hat{\mathbb{E}}[Y(\mathbf{z})] = \frac{1}{n} \sum_{i=1}^n \mu(\mathbf{z}, \hat{e}(x_i); \hat{\beta}).$$

Example: Continuous Exposure

In the case of a continuous exposure, we have a parametric probability density for the exposure

$$b(Z, X; \alpha) = f_{Z|X}(Z|X; \alpha)$$

for which we estimate α by regressing Z on X to obtain $\hat{\alpha}$ and fitted values $\hat{b}(Z, X) \equiv b(Z, X; \hat{\alpha})$.

Then we specify outcome model

$$\mathbb{E}[Y|Z = z, b(X) = b; \beta] = \mu(z, b; \beta)$$

and estimate this model by regressing y_i on z_i and $\hat{b}(z_i, x_i)$. Again, we might have that

$$\mathbb{E}[Y|Z = z_i, b(Z_i, X_i) = b_i; \beta] = \beta_0 + \beta_1 z_i + \beta_2 b_i.$$

This returns $\hat{\beta}$.

Example: Binary Exposure

We then compute the predictions under this model, and average them to obtain the APO estimate

$$\widehat{\mathbb{E}}[Y(\mathbf{z})] = \frac{1}{n} \sum_{i=1}^n \mu(\mathbf{z}, \widehat{b}(\mathbf{z}, x_i); \widehat{\beta}).$$

Note that here the propensity terms that enter into μ are computed at the target \mathbf{z} values

not the observed exposure values.

Propensity Score Regression

These procedures require us to make two modelling choices:

- ▶ the propensity model, $b(Z, X)$ or $b(X)$;
- ▶ the outcome mean model $\mu(z, b)$.

Unfortunately, both models must be correctly specified for consistent inference.

Misspecification of the outcome mean model will lead to bias; this model needs to capture the outcome to exposure and propensity function relationship correctly.

Weighting approaches

For a causal quantity of interest, we focus on the APO

$$\mathbb{E}[Y(\mathbf{z})] = \int y f_{Y(\mathbf{z}),X}(y, x) dy dx$$

that is, the average outcome, over the distribution of the confounders and predictors, if we hypothesize that the intervention sets the exposure to \mathbf{z} .

We now study methods that utilize the components already described, including the propensity score, but in a different fashion;

- ▶ instead of accounting for confounding by balancing through matching, we aim to achieve balance via weighting

Average potential outcome

If we could intervene at the population level to set $Z = \mathbf{z}$ for all individuals independently of their X value, we might rewrite this as

$$\mathbb{E}[Y(\mathbf{z})] = \int y \mathbb{1}_{\mathbf{z}}(z) f_{Y(\mathbf{z}),X}(y, x) \, dy \, dz \, dx$$

and take a random sample from the population with density

$$\mathbb{1}_{\mathbf{z}}(z) f_{Y(\mathbf{z}),X}(y, x) \equiv \mathbb{1}_{\mathbf{z}}(z) f_{Y|Z,X}(y|z, x) f_X(x).$$

We could then construct the ‘Monte Carlo’ estimator

$$\widehat{\mathbb{E}}[Y(\mathbf{z})] = \frac{1}{n} \sum_{i=1}^n Y_i$$

as $Z_i = \mathbf{z}$ for all i .

Average potential outcome: Experimental study

In a randomized (experimental) study, suppose that exposure $Z = z$ is assigned with probability determined by $f_Z(z)$. Then

$$\begin{aligned}\mathbb{E}[Y(z)] &= \frac{\int y \mathbf{1}_z(z) f_{Y(z),X}(y, x) f_Z(z) \, dy \, dz \, dx}{\int \mathbf{1}_z(z) f_{Y(z),X}(y, x) f_Z(z) \, dy \, dz \, dx} \\ &= \frac{\int y \mathbf{1}_z(z) f_{Y|Z,X}(y|z, x) f_X(x) f_Z(z) \, dy \, dz \, dx}{\int \mathbf{1}_z(z) f_Z(z) \, dz}\end{aligned}$$

Average potential outcome: Experimental study

This suggests the Monte Carlo estimators

$$\widehat{\mathbb{E}}[Y(\mathbf{z})] = \frac{\sum_{i=1}^n \mathbf{1}_{\mathbf{z}}(Z_i) Y_i}{\sum_{i=1}^n \mathbf{1}_{\mathbf{z}}(Z_i)} \quad \text{or} \quad \widehat{\mathbb{E}}[Y(\mathbf{z})] = \frac{1}{n f_Z(\mathbf{z})} \sum_{i=1}^n \mathbf{1}_{\mathbf{z}}(Z_i) Y_i.$$

Average potential outcome: Observational study

Denote by $P_{\mathcal{E}}$ the probability measure for samples drawn under the experimental measure corresponding to the density

$$f_{Y|Z,X}^{\mathcal{E}}(y|z, x) f_X^{\mathcal{E}}(x) f_Z^{\mathcal{E}}(z).$$

Now consider the case where the data arise from the observational (non-experimental) measure $P_{\mathcal{O}}(dy, dz, dx)$.

We have

$$\begin{aligned} \mathbb{E}[Y(\mathbf{z})] &= \frac{1}{f_Z^{\mathcal{E}}(\mathbf{z})} \int y \mathbf{1}_{\mathbf{z}}(z) P_{\mathcal{E}}(dy, dz, dx) \\ &= \frac{1}{f_Z^{\mathcal{E}}(\mathbf{z})} \int y \mathbf{1}_{\mathbf{z}}(z) \underbrace{\frac{P_{\mathcal{E}}(dy, dz, dx)}{P_{\mathcal{O}}(dy, dz, dx)}}_{\textcircled{1}} P_{\mathcal{O}}(dy, dz, dx). \end{aligned}$$

Average potential outcome: Observational study

In terms of densities ① becomes

$$\begin{aligned} & \frac{f_{Y|Z,X}^{\mathcal{E}}(y|z,x) f_Z^{\mathcal{E}}(z) f_X^{\mathcal{E}}(x)}{f_{Y|Z,X}^{\mathcal{O}}(y|z,x) f_{Z|X}^{\mathcal{O}}(z|x) f_X^{\mathcal{O}}(x)} \\ &= \frac{f_{Y|Z,X}^{\mathcal{E}}(y|z,x)}{f_{Y|Z,X}^{\mathcal{O}}(y|z,x)} \times \frac{f_Z^{\mathcal{E}}(z)}{f_{Z|X}^{\mathcal{O}}(z|x)} \times \frac{f_X^{\mathcal{E}}(x)}{f_X^{\mathcal{O}}(x)} \end{aligned}$$

- ▶ for the first term, we assume that

$$\frac{f_{Y|Z,X}^{\mathcal{E}}(y|z,x)}{f_{Y|Z,X}^{\mathcal{O}}(y|z,x)} = 1 \quad \text{for all } y, z, x;$$

this is essentially a no unmeasured confounders assumption.

- ▶ the third term equals 1 by assumption.

Experimental vs observational sampling

The second term

$$\frac{f_Z^{\mathcal{E}}(z)}{f_{Z|X}^{\mathcal{O}}(z|x)}$$

constitutes a weight that appears in the integral that yields the desired APO; the term

$$\frac{1}{f_{Z|X}^{\mathcal{O}}(z|x)}$$

accounts for the imbalance that influences the confounding and measures the difference between the observed sample and a hypothetical idealized randomized sample.

This suggests the (nonparametric) estimators

$$\widehat{\mathbb{E}}[Y(\mathbf{z})] = \frac{1}{n} \sum_{i=1}^n \frac{\mathbf{1}_{\mathbf{z}}(Z_i) Y_i}{f_{Z|X}^{\mathcal{O}}(Z_i|X_i)} \quad (\text{IPW0})$$

which is unbiased, or

$$\widehat{\mathbb{E}}[Y(\mathbf{z})] = \frac{\sum_{i=1}^n \frac{\mathbf{1}_{\mathbf{z}}(Z_i) Y_i}{f_{Z|X}^{\mathcal{O}}(Z_i|X_i)}}{\sum_{i=1}^n \frac{\mathbf{1}_{\mathbf{z}}(Z_i)}{f_{Z|X}^{\mathcal{O}}(Z_i|X_i)}} \quad (\text{IPW})$$

which is consistent, each provided $f_{Z|X}^{\mathcal{O}}(\cdot|\cdot)$ correctly specifies the conditional density of Z given X for all (z, x) .

Note 7

Inverse weighting constructs a pseudo-population in which there are no imbalances on confounders between the exposure groups. The pseudo-population is balanced, as required for direct comparison of treated and untreated groups.

Note 8

The term in the denominator of the components of the sum is $f_{Z|X}^O(Z_i|X_i)$, that is, the probability model that captures the conditional model for Z_i given X_i . If Z_i is binary, this essentially reduces to

$$e(X_i)^{Z_i}(1 - e(X_i))^{1-Z_i}$$

where $e(\cdot)$ is the propensity score as defined previously.

Note 9

It is evident that we must have

$$f_{Z|X}^O(Z_i|X_i) > 0$$

with probability 1 for this calculation to be valid.

This is commonly assumed, and is termed the positivity or experimental treatment assignment assumption.

Note 10

The inverse weighting procedure can also be justified from a weighted likelihood perspective.

Estimation via Augmentation

We may write

$$\mathbb{E}[Y(\mathbf{z})] = \mathbb{E}[Y(\mathbf{z}) - \mu(\mathbf{z}, X)] + \mathbb{E}[\mu(\mathbf{z}, X)]$$

where $\mu(\mathbf{z}, x) = \mathbb{E}[Y|Z = \mathbf{z}, X = x]$.

We have the alternate estimator

$$\widehat{\mathbb{E}}[Y(\mathbf{z})] = \frac{1}{n} \sum_{i=1}^n \frac{\mathbf{1}_{\mathbf{z}}(Z_i)(Y_i - \mu(Z_i, X_i))}{f_{Z|X}^{\mathcal{O}}(Z_i|X_i)} + \frac{1}{n} \sum_{i=1}^n \mu(\mathbf{z}, X_i) \quad (\text{AIPW})$$

and

$$\text{Var}_{\text{AIPW}} \leq \text{Var}_{\text{IPW}}.$$

Furthermore, (AIPW) is doubly robust (i.e. consistent even if one of $f_{Z|X}^{\mathcal{O}}(z|x)$ and $\mu(\mathbf{z}, x)$ is misspecified).

Properties under misspecification

Implementing (AIPW) relies on specification of the components

$$f_{Z|X}^{\mathcal{O}}(z|x) \quad \mu(z, x).$$

Suppose that, in reality, the correct specifications are

$$\tilde{f}_{Z|X}(z|x) \quad \tilde{\mu}(z, x).$$

Then the bias of (AIPW) is

$$\mathbb{E} \left[\frac{(f_{Z|X}^{\mathcal{O}}(z|X) - \tilde{f}_{Z|X}(z|X))(\mu(z, X) - \tilde{\mu}(z, X))}{f_{Z|X}^{\mathcal{O}}(z|X)} \right] \quad (5)$$

which is zero if

$$f_{Z|X}^{\mathcal{O}} \equiv \tilde{f}_{Z|X} \quad \text{or} \quad \mu(z, x) \equiv \tilde{\mu}(z, x)$$

that is, (AIPW) is doubly robust.

Properties under misspecification

Asymptotically, for estimators that are sample averages, the variance of the estimator converges to zero under standard conditions.

Therefore in large samples it is the magnitude of the bias as given by (5) that determines the quality of the estimator.

- ▶ equation (5) demonstrates that misspecification in the functions $\mu(z, x)$ and $f_{Z|X}^{\mathcal{O}}$ play equal roles in the bias.

Parametric modelling: two-stage approach

In the formulation, the nonparametric models

$$f_{Z|X}^{\mathcal{O}}(z|x) \quad \mu(z, x)$$

are commonly replaced by parametric models

$$f_{Z|X}^{\mathcal{O}}(z|x; \alpha) \quad \mu(z, x; \beta) = \int y f_{Y|Z,X}^{\mathcal{O}}(y|z, x; \beta) dy.$$

Parameters (α, β) are estimated from the observed data by regressing

- ▶ Stage I: Z on X using $(z_i, x_i), i = 1, \dots, n$,
 - ▶ Stage II: Y on (Z, X) using $(y_i, z_i, x_i), i = 1, \dots, n$
- and using plug-in version of (IPW) and (AIPW).

The estimated propensity score

Note 11

It is possible to conceive of situations where the propensity-type model

$$f_{Z|X}^{\mathcal{O}}(z|x) \quad \text{or} \quad f_{Z|X}^{\mathcal{O}}(z|x; \alpha)$$

is known precisely and does not need to be estimated.

This is akin to the randomized study where the allocation probabilities are fixed by the experimenter. It can be shown that using estimated quantities

$$\widehat{f}_{Z|X}^{\mathcal{O}}(z|x) \quad \text{or} \quad f_{Z|X}^{\mathcal{O}}(z|x; \widehat{\alpha})$$

yields lower variances for the resulting estimators than if the known quantities are used.

Alternative view of augmentation

Scharfstein et al. (1999), Bang & Robins (2005) write the estimating equation yielding (AIPW) as

$$\sum_{i=1}^n \frac{\mathbb{1}_{\mathbf{z}}(Z_i)}{f_{Z|X}^{\mathcal{O}}(Z_i|X_i)} (Y_i - \mu(Z_i, X_i)) + \sum_{i=1}^n \{\mu(\mathbf{z}, X_i) - \mu(\mathbf{z})\} = 0$$

Alternative view of augmentation

The first summation is a component of the score obtained when performing OLS regression for Y with mean function

$$\mu(z, x) = \mu_0(z, x) + \epsilon \frac{\mathbf{1}_z(z)}{f_{Z|X}^{\mathcal{O}}(z|x)}.$$

and $\mu_0(z, x)$ is a conditional mean model, and ϵ is a regression coefficient associated with the derived predictor

$$\frac{\mathbf{1}_z(z)}{f_{Z|X}^{\mathcal{O}}(z|x)}.$$

Alternative view of augmentation

Therefore, estimator (AIPW) can be obtained by regressing Y on (X, Z) for fixed z using the mean specification $\mu(z, x)$, and forming the estimator

$$\frac{1}{n} \sum_{i=1}^n \left\{ \mu_0(Z_i, X_i) + \hat{\epsilon} \frac{\mathbf{1}_z(Z_i)}{f_{Z|X}^{\mathcal{O}}(Z_i|X_i)} \right\}.$$

In a parametric model setting, this becomes

$$\frac{1}{n} \sum_{i=1}^n \left\{ \mu_0(Z_i, X_i; \hat{\beta}) + \hat{\epsilon} \frac{\mathbf{1}_z(Z_i)}{f_{Z|X}^{\mathcal{O}}(Z_i|X_i; \hat{\alpha})} \right\}$$

where α is estimated from Stage (I), and β is estimated along with ϵ in the secondary regression.

Augmentation and contrasts

The equivalent to (AIPW) for estimating the ATE for binary treatment

$$\mathbb{E}[Y(1)] - \mathbb{E}[Y(0)]$$

is merely $\widehat{\mathbb{E}}[Y(1)] - \widehat{\mathbb{E}}[Y(0)]$ or

$$\begin{aligned} \frac{1}{n} \sum_{i=1}^n \left[\frac{\mathbf{1}_1(Z_i)}{f_{Z|X}^{\mathcal{O}}(1|X_i)} - \frac{\mathbf{1}_0(Z_i)}{f_{Z|X}^{\mathcal{O}}(0|X_i)} \right] (Y_i - \mu(Z_i, X_i)) \\ + \frac{1}{n} \sum_{i=1}^n \{\mu(1, X_i) - \mu(0, X_i)\}. \end{aligned}$$

Augmentation and contrasts

Therefore we can repeat the above argument and base the contrast estimator on the regression of Y on (X, Z) using the mean specification

$$\mu(z, x) = \mu_0(z, x) + \epsilon \left[\frac{\mathbb{1}_1(z)}{f_{Z|X}^{\mathcal{O}}(1|x)} - \frac{\mathbb{1}_0(z)}{f_{Z|X}^{\mathcal{O}}(0|x)} \right]$$

Part 3

Implementation and Computation

Causal inference typically relies on reasonably standard statistical tools:

1. **Standard distributions:**

- ▶ Normal;
- ▶ Binomial;
- ▶ Time-to-event distributions (Exponential, Weibull etc.)

2. **Regression tools:**

- ▶ linear model/ordinary least squares;
- ▶ generalized linear model, typically linear regression;
- ▶ survival models.

Pooled logistic regression

For a survival outcome, pooled logistic regression is often used.

The usual continuous survival time outcome is replaced by a discrete, binary outcome;

- ▶ this is achieved by partitioning the outcome space into short intervals,

$$(0, t_1], (t_1, t_2], \dots$$

and assuming that the failure density is approximately constant in each interval.

- ▶ using a hazard parameterization, we have that

$$\Pr[\text{Failure in } (t_{j-1}, t_j] | \text{No failure before } t_{j-1}] = q_j$$

which converts each single failure time outcome into a series of binary responses, with 0 recording ‘no failure’ and 1 recording ‘failure’.

Semiparametric estimation

Semiparametric models based on estimating equations are typically used:

- ▶ such models make no parametric assumptions about the distributions of the various quantities, but instead make moment restrictions;
- ▶ resulting estimators inherit good asymptotic properties;
- ▶ variance of estimators typically estimated in a ‘robust’ fashion using the sandwich estimator of the asymptotic variance.

Key considerations

In light of the previous discussions, in order to facilitate causal comparisons, there are several key considerations that practitioners must take into account.

1. **The importance of no unmeasured confounding.**

When considering the study design, it is essential for valid conclusions to have measured and recorded all confounders.

2. Model construction for the outcome regression.

- ▶ ideally, the model for the expected value of Y given Z and X , $\mu(z, x)$, should be correctly specified, that is, correctly capture the relationship between outcome and the other variables.
- ▶ if this can be done, then no causal adjustments are necessary.
- ▶ conventional model building techniques (variable selection) can be used; this will prioritize predictors of outcome and therefore will select all confounders;
- ▶ however, in finite sample, this method may omit weak confounders that may lead to bias.

3. Model construction for the propensity score.

Ideally, the model for the (generalized) propensity score, $e(x)$ or $b(z, x)$, should be correctly specified, that is, correctly capture the relationship between the exposure and the confounders. We focus on

- 3.1 identifying the confounders;
- 3.2 ignoring the instruments: instruments do not predict the outcome, therefore cannot be a source of bias (unless there is unmeasured confounding) - however they can increase the variability of the resulting propensity score estimators.
- 3.3 the need for the specified propensity model to induce balance;
- 3.4 ensuring positivity, so that strata constructed from the propensity score have sufficient data within them to facilitate comparison;
- 3.5 effective model selection.

Key considerations

Note 12

Conventional model selection techniques (stepwise selection, selection via information criteria, sparse selection) should not be used when constructing the propensity score.

This is because such techniques prioritize the accurate prediction of exposure conditional on the other predictors; however, this is not the goal of the analysis.

These techniques may merely select strong instruments and omit strong predictors of outcome that are only weakly associated with exposure.

Note 13

An apparently conservative approach is to build rich (highly parameterized) models for both $\mu(z, x)$ and $e(x)$.

This approach prioritizes bias elimination at the cost of variance inflation.

4. **The required measure of effect.**

Is the causal measure required

- ▶ a risk difference ?
- ▶ a risk ratio ?
- ▶ an odds ratio ?
- ▶ an ATT, ATE or APO ?

Part 4

Extensions

Longitudinal studies

It is common for studies to involve multiple longitudinal measurements of exposure, confounders and outcomes.

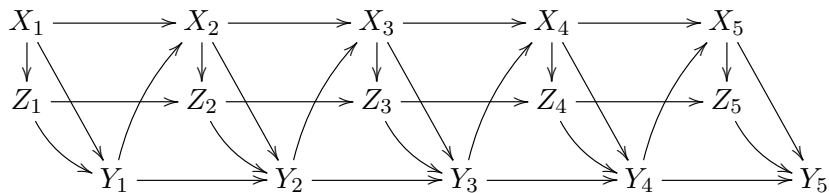
In this case, the possible effect of confounding of the exposure effect by the confounders is more complicated.

Furthermore, we may be interested in different types of effect:

- ▶ the direct effect: the effect of exposure in any given interval on the outcome in that interval, or the final observed outcome;
- ▶ the total effect: the effect of exposure aggregated across intervals final observed outcome;

Illustration

Possible structure across five intervals:



Mediation and time-varying confounding

- ▶ The effect of exposure on later outcomes may be mediated through variables measured at intermediate time points
 - ▶ for example, the effect of exposure Z_1 may have a direct effect on Y_1 that is confounded by X_1 ; however, the effect of Z_1 on Y_2 may also be non-negligible. This effect is mediated via X_2 .
- ▶ There may be time-varying confounding;

Multivariate versions of the propensity score

The propensity score may be generalized to the multivariate setting. We consider longitudinal versions of the measured variables: for $j = 1, \dots, m$, consider

- ▶ exposure: $\tilde{Z}_{ij} = (Z_{i1}, \dots, Z_{ij})$;
- ▶ outcome: $\tilde{Y}_{ij} = (Y_{i1}, \dots, Y_{ij})$;
- ▶ confounders: $\tilde{X}_{ij} = (X_{i1}, \dots, X_{ij})$.

Sometimes the notation

$$Z_{1:m} = (Z_1, \dots, Z_m)$$

will be useful.

Multivariate versions of the propensity score

We consider vectors of potential outcomes corresponding to these observed quantities, that is, we consider a potential sequence of interventions up to time j

$$\tilde{\mathbf{z}}_{ij} = (\mathbf{z}_{i1}, \dots, \mathbf{z}_{ij})$$

and then the corresponding sequence of potential outcomes

$$\tilde{Y}(\tilde{\mathbf{z}}_{ij}) = (Y(\mathbf{z}_{i1}), \dots, Y(\mathbf{z}_{ij})).$$

Multivariate versions of the propensity score

We define the multivariate (generalized) propensity score by

$$b_j(z, x) = f_{Z_j|X_j, \tilde{Z}_{j-1}, \tilde{X}_{j-1}}(z|x, \tilde{z}_{j-1}, \tilde{x}_{j-1})$$

that is, using the conditional distribution of exposure at interval j , given the confounder at interval j , and the historical values of exposures and confounders.

Under the sequential generalizations of the ‘no unmeasured confounders’ and positivity assumptions, this multivariate extension of the propensity score provides the required balance, and provides a means of estimating the direct effect of exposure.

The use of mixed models

The multivariate generalization above essentially builds a joint model for the sequence of exposures, and embeds this in a full joint distribution for all measured variables.

An alternative approach uses mixed (or random effect) models to capture the joint structure.

- ▶ such an approach is common in longitudinal data analysis;
- ▶ here we consider building a model for the longitudinal exposure data that encompasses a random effect.

The use of mixed models

Suppose first we have a continuous exposure: we consider the mixed effect model where for time point j

$$Z_{ij} = \tilde{X}_{ij}\alpha + \tilde{Z}_{i,j-1}\vartheta + \xi_i + \epsilon_{ij}$$

where

- ▶ $\tilde{X}_{ij}\alpha$ captures the fixed effect contribution of past and current confounders;
- ▶ $\tilde{Z}_{i,j-1}\vartheta$ captures the fixed effect contribution of past exposures;
- ▶ ξ_i is a subject specific random effect;
- ▶ ϵ_{ij} is a residual error.

The use of mixed models

The random effect ξ_i helps to capture unmeasured time-invariant confounding.

The distributional assumption made about ϵ_{ij} determine the precise form of a generalized propensity score that can again be used to estimate the direct effect of exposure.

The use of mixed models

For binary or other discrete exposures, the random effect model is built on the linear predictor scale, with say

$$\eta_{ij} = \tilde{X}_{ij}\alpha + \tilde{Z}_{i,j-1}\vartheta + \xi_i$$

determining the required conditional mean for the exposure at interval j .

Full-likelihood based inference may be used, but also generalized estimating approaches may be developed.

Estimation of Total Effects

The estimation of the total effect of exposure in longitudinal studies is more complicated as the need to acknowledge mediation and time-varying confounding renders standard likelihood-based approaches inappropriate.

The Marginal Structural Model is a semiparametric inverse weighting methodology designed to estimate total effects of functions of aggregate exposures that generalizes conventional inverse weighting.

The Marginal Structural Model

We observe for each individual i a sequence of exposures

$$Z_{i1}, Z_{i2}, \dots, Z_{im}$$

and confounders

$$X_{i1}, X_{i2}, \dots, X_{im}$$

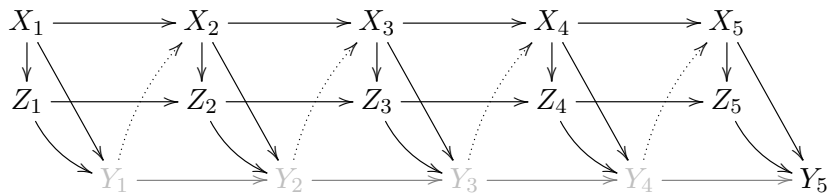
along with outcome $Y_i \equiv Y_{im}$ measured at the end of the study.

Intermediate outcomes $Y_{i1}, Y_{i2}, \dots, Y_{i,m-1}$ also possibly available.

We might also consider individual level frailty variables $\{v_i\}$, which are determinants of both the outcome and the intermediate variables, but can be assumed conditionally independent of the exposure assignments.

The Marginal Structural Model

For example, with $m = 5$:



Common example: pooled logistic regression

- ▶ discrete time survival outcome
- ▶ outcome is binary, intermediate outcomes monotonic
- ▶ length of follow-up is random, or event time is censored.

The Marginal Structural Model

We seek to quantify the causal effect of exposure pattern

$$\tilde{\mathbf{z}} = (\mathbf{z}_1, \mathbf{z}_2, \dots, \mathbf{z}_m)$$

on the outcome. If the outcome is binary, we might consider²

$$\log \left(\frac{f(Y_{im} = 1 | \tilde{\mathbf{z}}, \theta)}{f(Y_{im} = 0 | \tilde{\mathbf{z}}, \theta)} \right) = \theta_0 + \theta_1 \sum_{j=1}^m z_j$$

as the true, structural model. Note that this is a marginal model.

To avoid complicated notation in what follows, all probability distributions will be generically denoted $p(.|.)$.

² We might also consider structural models in which the influence of co-variates/confounders is recognized.

The Marginal Structural Model

However, this model is expressed for data presumed to be collected under an experimental design, \mathcal{E} .

In reality, it is necessary to adjust for the influence of

- ▶ time-varying confounding due to the observational nature of exposure assignment
- ▶ mediation as past exposures may influence future values of the confounders, exposures and outcome.

The adjustment can be achieved using inverse weighting via a marginal structural model.

The Marginal Structural Model

Causal parameter θ may be estimated via the weighted pseudo-likelihood

$$q(\theta; \tilde{x}, y, \tilde{z}, \gamma_0, \alpha_0) \equiv \prod_{i=1}^n f(y_i | \tilde{z}_i, \theta)^{w_{i0}},$$

where

$$w_{i0} = \frac{\prod_{j=1}^m f(z_{ij} | \tilde{z}_{i(j-1)}, \alpha_{0j})}{\prod_{j=1}^m f(z_{ij} | \tilde{z}_{i(j-1)}, \tilde{x}_{ij}, \gamma_{0j})}$$

defines ‘stabilized’ case weights in which the true parameter values (γ_0, α_0) are (for now) taken to be known.

The Marginal Structural Model: The logic

- ▶ Inference is required under target population \mathcal{E} but a sample from the population of interest is not directly available
- ▶ Samples from observational design \mathcal{O} relevant for learning about this target population are available.
- ▶ In population \mathcal{E} , the conditional independence $z_{ij} \perp \tilde{x}_{ij} \mid \tilde{z}_{i(j-1)}$ holds true.
- ▶ The weights w_{i0} convey information on how much \mathcal{O} resembles \mathcal{E} : this information is contained in the parameters γ .
- ▶ \mathcal{E} has the same marginal exposure assignment distribution as \mathcal{O} , characterized by α .

The Marginal Structural Model: Implementation

- ▶ Parameters α typically must (and should) be estimated from the exposure and confounder data;
- ▶ The usual logic of the propensity score applies here: in constructing the terms that enter the conditional models that enter into the stabilized weights, we use confounders but omit instruments.
- ▶ Inference using the weighted likelihood typically proceeds using robust (sandwich) variance estimation, or the bootstrap.

Antiretroviral therapy (ART) has reduced morbidity and mortality due to nearly all HIV-related illnesses, apart from mortality due to end-stage liver disease, which has increased since ART treatment became widespread.

In part, this increase may be due to improved overall survival combined with Hepatitis C virus (HCV) associated hepatic liver fibrosis, the progress of which is accelerated by immune dysfunction related to HIV-infection.

The Canadian Co-infection Cohort Study is one of the largest projects set up to study the role of ART on the development of end-stage liver disease in HIV-HCV co-infected individuals.

Given the importance of ART in improving HIV-related immunosuppression, it is hypothesized that liver fibrosis progression in co-infected individuals may be partly related to adverse consequences of ART interruptions.

Study comprised $N = 474$ individuals with at least one follow-up visit (scheduled at every six months) after the baseline visit, and 2066 follow-up visits in total (1592 excluding the baseline visits). The number of follow-up visits m_i ranged from 2 to 16 (median 4).

We adopt a pooled logistic regression approach:

- ▶ a single binary outcome (death at study termination)
- ▶ longitudinal binary exposure (adherence to ART)
- ▶ possible confounders
 - ▶ baseline covariates: female gender, hepatitis B surface antigen (HBsAg) test and baseline APRI, as well as
 - ▶ time-varying covariates: age, current intravenous drug use (binary), current alcohol use (binary), duration of HCV infection, HIV viral load, CD4 cell count, as well as ART interruption status at the previous visit.
- ▶ need also a model for informative censoring.

Real Data Example : ART interruption in HIV/HCV co-infected individuals

- ▶ We included co-infected adults who were not on HCV treatment and did not have liver fibrosis at baseline.
- ▶ The outcome event was defined as aminotransferase-to-platelet ratio index (APRI), a surrogate marker for liver fibrosis, being at least 1.5 in any subsequent visit
- ▶ We included visits where the individuals were either on ART or had interrupted therapy ($Z_{ij} = 1$), based on self-reported medication information, during the 6 months before each follow-up visit.

Real Data Example : ART interruption in HIV/HCV co-infected individuals

- ▶ Individuals suspected of having spontaneously cleared their HCV infection (based on two consecutive negative HCV viral load measurements) were excluded as they are not considered at risk for fibrosis progression.
- ▶ To ensure correct temporal order in the analyses, in the treatment assignment model all time-varying covariates (x_{ij}), including the laboratory measurements (HIV viral load and CD4 cell count), were lagged one visit.
- ▶ Follow-up was terminated at the outcome event ($Y_{ij} = 1$), while individuals starting HCV medication during the follow-up were censored.

We considered the structural model

$$\log \left(\frac{f(Y_{ij} = 1 | \tilde{z}_{ij}, \theta)}{f(Y_{ij} = 0 | \tilde{z}_{ij}, \theta)} \right) = \theta_0 + \theta_1 z_j$$

so that θ_1 measures the total effect of exposure in an interval, allowing for mediation.

Real Data Example : ART interruption in HIV/HCV co-infected individuals

Results:

Estimator	$\hat{\theta}_1$	SE	z
Naive	4.616	0.333	13.853
MSM	0.354	0.377	0.937
Bootstrap	0.308	0.395	0.780

After adjustment for confounding, the effect of exposure is non-significant.

Part 5

New Challenges and Approaches

New challenges

The main challenge for causal adjustments using the propensity score is the nature of the observational data being recorded.

The data sets and databases being collected are increasingly complex and typically originate from different sources. The benefits of ‘Big Data’ come with the costs of more involved computation and modelling.

There is always an important trade off between the sample size n and the dimension of the confounder (and predictor) set.

Examples

- ▶ pharmacoepidemiology;
- ▶ electronic health records and primary care decision making;
- ▶ real-time health monitoring;

Data synthesis and combination

For observational databases, the choice of inclusion/exclusion criteria for analysis can have profound influence on the ultimate results:

- ▶ different databases can lead to different conclusions for the same effect of interest purely because of the methodology used to construct the raw data, irrespective of modelling choices.
- ▶ the key task of the statistician is to report uncertainty in a coherent fashion, ensuring that all sources of uncertainty are reflected. This should include uncertainty introduced due to lack of compatibility of data sources.

Classic challenges

Alongside the challenges of modern quantitative health research are more conventional challenges:

- ▶ missing data: many causal adjustment procedures are adapted forms of procedures developed for handling informative missingness (especially inverse weighting);
- ▶ length-bias and left truncation in prevalent case studies: selection of prevalent cases is also a form of ‘selection bias’ that causes bias in estimation if unadjusted;
- ▶ non-compliance: in randomized and observational studies there is the possibility of non- or partial compliance which is again a potential source of selection bias.

The Bayesian version

The Bayesian paradigm provides a natural framework within which decision-making under uncertainty can be undertaken.

Much of the reasoning on causal inference, and many of the modelling choices we must make for causal comparison and adjustment, are identical under Bayesian and classical (frequentist, semiparametric) reasoning.

The advantages of Bayesian thinking

With increasingly complex data sets in high dimensions, Bayesian methods can be useful as they

- ▶ provide a means of informed and coherent decision making in the presence of uncertainty;
- ▶ yield interpretable variability estimates in finite sample at the cost of interpretable modelling assumptions;
- ▶ allow the statistician to impose structure onto the inference problem that is helpful when information is sparse;
- ▶ naturally handle prediction, hierarchical modelling, data synthesis, and missing data problems.

Typically, these advantages come at the cost of more involved computation.

Bayesian causal inference: recent history

- ▶ D.B. Rubin formulated the modern foundations for causal inference from a largely Bayesian (missing data) perspective:
 - ▶ revived potential outcome concept to define causal estimand
 - ▶ inference through Bayesian (model-based) predictive formulation
 - ▶ focus on matching
- ▶ Semiparametric frequentist formulation pre-dominant from mid 80s
- ▶ Recent Bayesian approaches largely mimic semiparametric approach, but with explicit probability models.

Bayesian inference for two-stage models

- ▶ Full Bayes: full likelihood in two parametric models
 - ▶ needs correct specification;
 - ▶ two component models are treated independently.
- ▶ Quasi-Bayes: use semiparametric estimating equation approach for Stage II, with Stage I parameters treated in a fully Bayesian fashion.
 - ▶ possibly good frequentist performance;
 - ▶ difficult to understand frequentist properties.
- ▶ Pseudo-Bayes: use amended likelihood to avoid feedback between Stage I and Stage II
 - ▶ not fully Bayesian, no proper probability model

Five Considerations

1. The causal contrast
2. Do we really need potential outcomes ?
3. ‘Observables’ implies ‘Prediction’
4. The Fundamental Theory of Bayesian Inference.
5. The Bayesian Causal Specification

Part 6

Conclusions

Conclusions

- ▶ Causal inference methods provide answers to important questions concerning the impact of hypothetical exposures;
- ▶ Standard statistical methods are used;
- ▶ Balance is the key to accounting for confounding;
- ▶ The propensity score is a tool for achieving balance;
- ▶ The propensity score can be used for
 - ▶ matching,
 - ▶ weighting, and
 - ▶ as part of regression modelling.
- ▶ Bayesian methods are not widely used, but are generally applicable.

Key remaining challenges

- ▶ Model selection;
- ▶ Scale and complexity of observational data;

Collaborators & Acknowledgments

- ▶ McGill: Erica Moodie, Michael Wallace, Marina Klein
- ▶ Toronto: Olli Saarela
- ▶ Imperial College London: Dan Graham, Emma McCoy



Part 7

Appendix: Bayesian methods

Five Considerations

1. The causal contrast
2. Do we really need potential outcomes ?
3. 'Observables' implies 'Prediction'
4. The Fundamental Theory of Bayesian Inference.
5. The Bayesian Causal Specification

What is the causal contrast to work with ?

The causal effect of changing exposure from \mathbf{z}_1 to \mathbf{z}_2 is

$$\delta(\mathbf{z}_1, \mathbf{z}_2) = \mathbb{E}[Y_i(\mathbf{z}_2) - Y_i(\mathbf{z}_1)]$$

that is, an expected difference between potential outcomes for the same individual, which is presumed to be the same for all i .

There is no meaningful inferential difference³ – in most settings – if we work instead with

$$\delta(\mathbf{z}_1, \mathbf{z}_2) = \mathbb{E}[Y(\mathbf{z}_2)] - \mathbb{E}[Y(\mathbf{z}_1)]$$

imagining the exposures acting on the whole population.

³ In a randomized trial, we do not assign the same individual multiple exposures, we compare different individuals who have been randomly allocated to different, yet comparable, exposure groups.

Do we really need potential outcomes ?

The introduction of the potential outcome random variables is not strictly necessary for inference.

The notation and conceptualization is useful, but not necessary.

A joint distribution for potential outcomes $\{Y_i(\mathbf{z}) : \mathbf{z} \in \mathcal{Z}\}$ is imagined, but never utilized.

‘Observables’ implies ‘Prediction’

The quantity

$$\delta(\mathbf{z}_1, \mathbf{z}_2) = \mathbb{E}[Y(\mathbf{z}_2)] - \mathbb{E}[Y(\mathbf{z}_1)]$$

concerns the expected value of observable⁴ quantities.

In an inferential setting, we wish to make statements about μ in light of observed data $(y_i, z_i, x_i), i = 1, \dots, n$.

Hence, in a Bayesian setting we should be examining predictive quantities⁵.

⁴ albeit hypothetical

⁵ this is Rubin’s argument

Exchangeability and de Finetti's Representation

Suppose that U_1, \dots, U_n, \dots are an infinite sequence of observable random variables such that, for all $n \geq 1$,

$$P(U_{1:n} \in \mathcal{B}_{1:n}) = P(U_{\rho(1:n)} \in \mathcal{B}_{1:n})$$

where

- ▶ $U_{1:n} = (U_1, \dots, U_n)$,
- ▶ $\mathcal{B}_{1:n} = \mathcal{B}_1 \times \dots \times \mathcal{B}_n$,
- ▶ $\rho(\cdot)$ is a permutation operator.

Then the $\{U_n\}$ are exchangeable.

Exchangeability and de Finetti's Representation

If $\{U_n\}$ are exchangeable, then there exists a measure π_0 on the space of probability measures such that

$$P(U_{1:n} \in \mathcal{B}_{1:n}) = \int \left\{ \prod_{i=1}^n \mathbb{Q}(U_i \in \mathcal{B}_i) \right\} \pi_0(d\mathbb{Q}) \quad (\blacklozenge)$$

- ▶ $\pi_0(\cdot)$ is a distribution function for the limiting empirical probability measure \mathbb{P} , that is, $\mathbb{P} \sim \pi_0$ and

$$\mathbb{P}(\mathcal{B}) = \lim_{n \rightarrow \infty} \sum_{i=1}^n \mathbb{P}_n(\mathcal{B}) = \lim_{n \rightarrow \infty} \frac{1}{n} \sum_{i=1}^n \mathbb{1}_{\mathcal{B}}(U_i)$$

- ▶ The term in brackets defines a ‘likelihood’

$$P(U_{1:n} \in \mathcal{B}_{1:n} | \mathbb{P} = \mathbb{Q}) = \prod_{i=1}^n \mathbb{Q}(U_i \in \mathcal{B}_i).$$

Exchangeability and de Finetti's Representation

That is, we have a Bayesian model with ‘parameter’ \mathbb{P} : to generate realizations from the model, we

- ▶ specify π_0 as a ‘prior’ distribution for \mathbb{P} ,
- ▶ draw \mathbb{Q} from π_0 ,
- ▶ draw $U_1, U_2, \dots, U_n, \dots$ independently from \mathbb{Q} .

Prediction:

$$P(U_{(n+1):(n+m)}|U_{1:n}) = \frac{P(U_{1:(n+m)})}{P(U_{1:n})} \quad (\blacksquare)$$

and use (\blacklozenge) in numerator and denominator.

Inference:

$$P(\mathbb{P}|U_{1:n}) = \lim_{m \rightarrow \infty} \mathbb{P}_{m|n}(U_{(n+1):(n+m)}|U_{1:n}) \quad (\star)$$

The Bayesian Causal Specification

We treat $U_i = (Y_i, Z_i, X_i), i = 1, \dots, n$ as a realization of an exchangeable sequence, and obtain de Finetti representations for experimental and observational measures $P_{\mathcal{E}}(U_{1:n})$ and $P_{\mathcal{O}}(U_{1:n})$ which facilitate inference.

We might consider parametric versions, and denote by

$$\vartheta_{\mathcal{E}} = (\theta_{\mathcal{E}}, \alpha, \psi_{\mathcal{E}}) \quad \vartheta_{\mathcal{O}} = (\theta_{\mathcal{O}}, \gamma, \psi_{\mathcal{O}})$$

the different parameters that appear in the ‘likelihood’ part of
(♦)

$$f_{\mathcal{E}}(y_i, x_i, z_i; \vartheta_{\mathcal{E}}) = f_{\mathcal{E}}(y_i|x_i, z_i; \theta_{\mathcal{E}})f_{\mathcal{E}}(z_i; \alpha)f_{\mathcal{E}}(x_i; \psi_{\mathcal{E}})$$

$$f_{\mathcal{O}}(y_i, x_i, z_i; \vartheta_{\mathcal{O}}) = f_{\mathcal{O}}(y_i|x_i, z_i; \theta_{\mathcal{O}})f_{\mathcal{O}}(z_i|x_i; \gamma)f_{\mathcal{O}}(x_i; \psi_{\mathcal{O}})$$

under the two assumed design settings.

The Bayesian Causal Solution

We can therefore adopt the fully Bayesian causal inference strategy:

1. Propose models for observed data, presumed exchangeable, under the experimental and observational assumptions;
2. Formulate the causal effect estimation problem as a prediction problem, and resolve to examine posterior predictive expectations under the experimental setting;
3. As data collected under the experimental setting are not available, use the observational posterior and change of measure/importance sampling ideas to re-express the posterior predictive expectation of interest in terms of the observational posterior measure;
4. Perform computations using Monte Carlo and MCMC.

Point binary treatment/binary response

Suppose $Y, Z, z \in \{0, 1\}$. A natural estimand is

$$\Pr[Y(z) = 1] = \mathbb{E}[Y(z)] \quad z = 0, 1.$$

We construct de Finetti representations for data

$$U_{1:n} = \{Y_{1:n}, Z_{1:n}, X_{1:n}\}$$

under the experimental and observation assumptions

$$P_{\mathcal{E}}(U_{1:n}) \quad P_{\mathcal{O}}(U_{1:n})$$

and define the estimand as

$$\lim_{n \rightarrow \infty} \mathbb{E}_n^{\mathcal{E}}[Y^* | Z^* = z, u_{1:n}]$$

where $\mathbb{E}_n^{\mathcal{E}}[\cdot]$ denotes expectation with respect to the posterior predictive distribution computed under \mathcal{E} , $p_n^{\mathcal{E}}(\cdot)$.

Point binary treatment/binary response

We justify the estimator

$$\mathbb{E}_n^{\mathcal{E}}[Y^* | Z^* = \mathbf{z}, u_{1:n}] \quad (\text{PPE})$$

via a maximum expected utility argument; we define utility $\mathcal{U}(\cdot, \cdot)$ for predictions as

$$\mathcal{U}(y^*, g(u_{1:n})) = -(y^* - g(u_{1:n}))^2.$$

for functions $g(\cdot)$ of the data, and solve

$$\widehat{\mathbb{E}}_n^{\mathcal{E}}[Y(\mathbf{z})] = \arg \max_{g(\cdot)} \int \mathcal{U}(y^*, g(y)) p_n^{\mathcal{E}}(y|\mathbf{z}) \, dy \quad (\text{E})$$

This yields the posterior predictive expectation for Y^* under the setting $Z^* = \mathbf{z}$, marginally over X^* , as in (PPE)

Point binary treatment/binary response

In principle, the integral in (E) could be approximated using Monte Carlo using a sample from conditional posterior predictive $p_n^{\mathcal{E}}(y|z)$, computed from

$$p_n^{\mathcal{E}}(y, z, x).$$

However, we do not have access to data collected under \mathcal{E} ; we must rewrite the computation in terms of the posterior predictive computed under \mathcal{O}

$$p_n^{\mathcal{O}}(y, z, x).$$

To do this we re-purpose the earlier frequentist calculation.

$$\begin{aligned}\mathbb{E}_n^{\mathcal{E}}[Y(\mathbf{z})] &= \int y \mathbf{1}_{\mathbf{z}}(z) p_n^{\mathcal{E}}(dy, dz, dx) \\ &= \int y \mathbf{1}_{\mathbf{z}}(z) \frac{p_n^{\mathcal{E}}(dy, dz, dx)}{p_n^{\mathcal{O}}(dy, dz, dx)} p_n^{\mathcal{O}}(dy, dz, dx) \\ &= \dots \\ &= \int \frac{y \mathbf{1}_{\mathbf{z}}(z)}{p_n^{\mathcal{O}}(z|x)} p_n^{\mathcal{O}}(dy, dz, dx)\end{aligned}$$

Point binary treatment/binary response

We estimate the integral by Monte Carlo using a sample from the posterior predictive

$$p_n^{\mathcal{O}}(y, z, x)$$

obtained by

- ▶ bootstrap resampling from the original data: this corresponds to a nonparametric posterior predictive procedure; or
- ▶ MCMC sampling using a parametric model.

Bayesian nonparametric procedure

The Bayesian nonparametric approximation is

$$\hat{p}_n^{\mathcal{O}}(y, z, x) = \sum_{i=1}^n \omega_i \delta_{x_i, y_i, z_i}(x, y, z)$$

where

$$(\omega_1, \dots, \omega_n) \sim \text{Dirichlet}(1, \dots, 1)$$

The non-parametric posterior predictive distribution is thus a discrete random measure with support equal to the original data.

Bayesian nonparametric procedure

To produce a sample from this distribution, we first sample the ω from the Dirichlet, and then (x^*, y^*, z^*) from the Multinomial distribution on the original data with probabilities ω .

We then repeat this N times to get a sample of size N from the posterior predictive.

Bayesian nonparametric procedure

Now the expression

$$\widehat{p}_n^{\mathcal{O}}(y, z, x) = \sum_{i=1}^n \omega_i \delta_{x_i, y_i, z_i}(x, y, z)$$

implies an expression for the conditional $\widehat{p}_n^{\mathcal{O}}(z|x)$ of

$$\widehat{p}_n^{\mathcal{O}}(z|x) = \frac{\sum_{i=1}^n \omega_i \delta_{x_i, z_i}(x, z)}{\sum_{i=1}^n \omega_i \delta_{x_i}(x)}$$

where this expression may be evaluated for $x = x_i$ $i = 1, \dots, n$, for each z .

Bayesian nonparametric procedure

This can be rewritten

$$\widehat{p}_n^{\mathcal{O}}(z|x) = \sum_{i=1}^n \left(\frac{\omega_i \delta_{x_i}(x)}{\sum_{j=1}^n \omega_j \delta_{x_j}(x)} \right) \delta_{z_i}(z)$$

This is also a random measure, and it will vary with each draw of ω .

Parametric version

In the parametric version

$$\begin{aligned} p_n^{\mathcal{O}}(y, z, x) &= \int f_{\mathcal{O}}(y, x, z | \beta, \alpha) \pi_n^{\mathcal{O}}(\beta, \alpha) \, d\beta \, d\alpha \\ &= \left\{ \int f_{\mathcal{O}}(y | x, z, \beta) \pi_n^{\mathcal{O}}(\beta) \, d\beta \right\} \left\{ \int f_{\mathcal{O}}(z | x, \alpha) \pi_n^{\mathcal{O}}(\alpha) \, d\alpha \right\} p_n^{\mathcal{O}}(x) \\ &= p_n^{\mathcal{O}}(y | z, x) p_n^{\mathcal{O}}(z | x) p_n^{\mathcal{O}}(x) \end{aligned}$$

In the calculation, the posterior predictive density

$$p_n^{\mathcal{O}}(z | x) = \int f_{Z|X}^{\mathcal{O}}(z | x, \alpha) \pi_n^{\mathcal{O}}(\alpha) \, d\alpha$$

is typically computed using MCMC.

Real Data Example : ART interruption in HIV/HCV co-infected individuals

The following table shows the results for alternative estimators for the interruption effect in a marginal model

$$f_{\mathcal{E}}(Y_{ij} = 1 \mid \mathbf{z}, \theta) = \text{expit}\{\theta_0 + \theta_1 \mathbf{z}\},$$

and the corresponding standard errors.

- ▶ The weights in the multiple imputation type estimator, as well as in the two Bayesian estimators were calculated from samples of size 2500 from the posterior distributions
- ▶ Flat improper priors were used for all parameters.
- ▶ Multinomial, Dirichlet and bootstrap estimates were calculated from 2500 replications.

Real Data Example : ART interruption in HIV/HCV co-infected individuals

Estimator	$\hat{\theta}_1$	SE	z
Naive	4.616	0.333	13.853
MSM	0.354	0.377	0.937
Quasi-Bayes MI	0.316	0.529	0.597
Bootstrap	0.308	0.395	0.780
Dirichlet	0.366	0.375	0.976
Multinomial	0.361	0.400	0.902

Real Data Example : ART interruption in HIV/HCV co-infected individuals

The five alternative estimates are similar, with the exception of the MI-type estimator, which appears to inflate the standard error. This indicates that sampling from a quasi-posterior in the two-step approach introduces excess variability to the estimation.

In contrast, the Multinomial and Dirichlet sampling standard errors are close to the bootstrap standard error, without involving re-estimation of the treatment assignment and censoring models in each replication.