Part 6 Extensions

The development above has mainly focussed on

- binary treatment,
- linear models,
- average treatment effects,
- single time-point studies,

but extensions can be developed to handle each case.

For continuous treatments, if necessary we may carry out adjustment via the conditional *expectation*

$$b(X) = \mathbb{E}_{Z|X}^{\mathcal{O}}[Z|X]$$

rather than (for example) the propensity score which is based on the conditional *probability* model

 $f^{\mathcal{O}}_{Z|X}(Z|X).$

Example: G-estimation

The model

$$Y = Z\psi + \mu_0(X;\beta) + \varepsilon$$

which forms the basis of the G-estimation procedure can be utilized if Z is *continuous*.

This relies on the construction of a model for $\mathbb{E}^{\mathcal{O}}_{Z|X}[Z|X]$

- justified by a focus on *orthogonality* under the *covariance* inner product;
- focus on conditional *uncorrelatedness* rather than conditional *independence*.

For IPW estimation, the construction

$$\widetilde{\mu}_{\scriptscriptstyle \mathrm{IPW}}(\mathsf{Z}) = rac{1}{n}\sum_{i=1}^n rac{\mathbbm{1}_{\{\mathsf{Z}\}}(Z_i)Y_i}{f^{\mathcal{O}}_{Z|X}(Z_i|X_i)}.$$

also works in the continuous case as

$$\mathbb{E}^{\mathcal{O}}_{X,Y,Z}\left[\frac{\mathbbm{1}_{\{\mathbf{Z}\}}(Z)Y}{f^{\mathcal{O}}_{Z|X}(Z|X)}\right] = \mu(\mathbf{Z})$$

but in practice this estimator can be variable in finite sample.

The focus on $\mathbb{E}^{\mathcal{O}}_{Z|X}[Z|X]$ emphasizes the role of *regression* approaches to constructing the treatment model: we denote

$$\xi(X) \equiv \xi(X; \alpha) = \mathbb{E}_{Z|X}^{\mathcal{O}}[Z|X].$$

In most cases, some form of *generalized linear model* will suffice, with estimating equation such as

$$\sum_{i=1}^{n} \mathbf{x}_{i}^{\top}(z_{i} - \xi(\mathbf{x}_{i}\alpha)) = \mathbf{0}$$

for parameter $\boldsymbol{\alpha}$ so that the fitted values

$$\xi(\mathbf{x}_i\widehat{\alpha}) \quad i=1,\ldots,n$$

can be computed.

In any subsequent procedure, these fitted values are used in estimating the causal effects: for example, in *G*-estimation, we consider the *plug-in* estimator

$$\widehat{\delta}_{\scriptscriptstyle G} = \frac{\displaystyle\sum_{i=1}^n Y_i(Z_i - \xi(\mathbf{X}_i \widehat{\alpha}))}{\displaystyle\sum_{i=1}^n Z_i(Z_i - \xi(\mathbf{X}_i \widehat{\alpha}))}$$

In IPW estimation, for binary treatment, we use

$$\widehat{\mu}_{\scriptscriptstyle \mathrm{IPW}}(\mathbf{1}) = rac{1}{n}\sum_{i=1}^n rac{Z_i Y_i}{e(X_i;\widehat{lpha})}.$$

The application of many of the causal adjustment methods rely on regression modelling for

the outcome mean model

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

the propensity or balancing score

$$\mathbf{e}(\mathbf{x}) = \Pr[Z = 1 | X = \mathbf{x}] \qquad \mathbf{b}(\mathbf{x}) = \mathbb{E}_{Z|X}[Z|X = Z]$$

In either case we can utilize

- Inear/generalized linear models
- flexible models (eg splines)
- prediction-based approaches (eg machine learning methods, regression trees, neural networks)
- ensemble methods (eg model averaging, boosting)

to construct fitted versions of each model.

The advanced methods can be effective, but there are potential pitfalls:

- 1. The quantification of *uncertainty*;
 - no ready analytic answers, typically relies on bootstrap;
 - large computational burden
- 2. Positivity violations.
 - prediction (of treatment mechanism) is not the fundamental goal;
 - can be overcome using methods that target balance/overlap.

If the outcome variable Y is discrete, or it is contended that the outcome mean model

$$\mathbb{E}^{\mathcal{O}}_{Y|X,Z}[Y|X=x,Z=z] = \mu(x,z;\beta,\psi)$$

is not linear (in the treatment or parameters), it is necessary to extend some of the previous concepts.

In the following,

- X is the random vector of confounders, and x its observed values,
- ► x is the vector of functions of x that form the linear predictor,
- X is the random variable version of x.

Suppose that a *log-linear* model is deemed appropriate:

$$\log \mu(\mathbf{x}, \mathbf{z}; \beta, \psi) = \mathbf{x}_{\beta}\beta + \mathbf{z}\mathbf{x}_{\psi}\psi$$

Regarding this as the structural model, we then have that

$$\mathbb{E}[Y(\mathsf{z})] \equiv \mathbb{E}_{Y|Z}^{\varepsilon}[Y|Z = \mathsf{z}] = \mathbb{E}_{X}^{\varepsilon}[\exp\{\mathbf{X}_{\beta}\beta + \mathsf{z}\mathbf{X}_{\psi}\psi\}]$$

which would then invoke the estimator

$$\widetilde{\mu}(\mathbf{z}) = \frac{1}{n} \sum_{i=1}^{n} \exp\{\mathbf{x}_{\beta i} \widehat{\beta} + \mathbf{z} \mathbf{x}_{\psi i} \widehat{\psi}\}.$$

where $(\widehat{\beta},\widehat{\psi})$ are estimated from a correctly specified model.

Then in the binary treatment case

$$\frac{\mathbb{E}[Y(1)]}{\mathbb{E}[Y(0)]} = \frac{\mathbb{E}_X^{\varepsilon}[\exp\{\mathbf{X}_{\beta}\beta + \mathbf{X}_{\psi}\psi\}]}{\mathbb{E}_X^{\varepsilon}[\exp\{\mathbf{X}_{\beta}\beta\}]}$$

but notice that unlike in the linear case

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

depend on both β and ψ . Also

$$\frac{\mathbb{E}[Y(1)]}{\mathbb{E}[Y(0)]} \neq \mathbb{E}\left[\frac{Y(1)}{Y(0)}\right]$$

in general.

The above model assumes that

$$\mu(\mathbf{x}, \mathbf{1}; \beta, \psi) = \mu(\mathbf{x}, \mathbf{0}; \beta, \psi) \exp\{\mathbf{x}_{\psi}\psi\}.$$

with the effect of Z represented *conditional* on X; *marginally*, ψ alone does not capture the effect of treatment.

We could formulate a model for the potential outcomes where

 $Y(\mathbf{1}) = Y(\mathbf{0}) \exp\{\mathbf{x}_{\psi}\psi\}$

that is

$$\log Y(\mathbf{1}) - \log Y(\mathbf{0}) = \mathbf{x}_{\psi} \psi.$$

If we can estimate β and ψ consistently, then we can recover APO and ATE estimators from them. By standard regression arguments, we know that *correct specification* of the outcome model is needed.

For the log-linear model,

$$\log \mu(\mathbf{x}, \mathbf{z}; \beta, \psi) = \mathbf{x}_{\beta}\beta + \mathbf{z}\mathbf{x}_{\psi}\psi$$

the standard estimating equation takes the form

$$\sum_{i=1}^{n} \begin{pmatrix} \mathbf{x}_{\beta i}^{\top} \\ z_{i} \mathbf{x}_{\psi i}^{\top} \end{pmatrix} (y_{i} - \exp\{\mathbf{x}_{\beta i}\beta + z_{i} \mathbf{x}_{\psi i}\psi\}) = \mathbf{0}.$$

We might try to make inference robust to mis-specification using a G-estimation-like strategy, and modify the estimating equation to be

$$\sum_{i=1}^{n} \begin{pmatrix} \mathbf{x}_{\beta i}^{\top} \\ (\mathbf{z}_{i} - \mathbf{e}(\mathbf{x}_{i}))\mathbf{x}_{\psi_{i}}^{\top} \end{pmatrix} (\mathbf{y}_{i} - \exp\{\mathbf{x}_{\beta i}\beta + \mathbf{z}_{i}\mathbf{x}_{\psi i}\psi\}) = \mathbf{0}$$

where $e(x_i)$ is the propensity score; if necessary, this can be estimated using a further parametric model

$$\sum_{i=1}^{n} \mathbf{x}_{i}^{\top}(z_{i} - \boldsymbol{e}(\boldsymbol{x}_{i}; \alpha)) = \mathbf{0}.$$

In the *linear* case, this modification led to consistent estimation of ψ as the resulting estimating equation is *unbiased*, that is

$$\mathsf{E}_{X,Y,Z}[(Z-\mathbf{e}(X))\mathbf{X}_{\psi}^{\top}(Y-\mathbf{X}_{\beta}\beta-Z\mathbf{X}_{\psi}\psi)]=\mathbf{0}$$

provided $\mathbf{x}_{\psi}\psi$ correctly captures the effect of treatment, and either

- propensity score e(X), or
- treatment-free mean component $\mathbf{x}_{\beta}\beta$

is correctly specified.

In the *log-linear* case, the equivalent requirement would be

$$\mathbb{E}_{X,Y,Z}[(Z - \boldsymbol{e}(X))\mathbf{X}_{\psi}^{\top}(Y - \exp\{\mathbf{X}_{\beta}\beta + Z\mathbf{x}_{\psi}\psi)\}] = \mathbf{0}.$$

However, this requirement is *not* met if the treatment-free mean component is mis-specified, even if the propensity score model is correctly specified.

Suppose in reality

$$\mathsf{E}_{Y|X,Z}[Y|X=x,Z=z] = \mu(x,z) = \mu_0(x) \exp\{z\mathbf{x}_{\psi}\psi\}.$$

Then

$$\mathbb{E}_{Y|X,Z}[Y - \exp\{\mathbf{x}_{\beta}\beta + Z\mathbf{x}_{\psi}\psi)\}|X = \mathbf{x}, Z = \mathbf{z}]$$

= $\mu_0(\mathbf{x}) \exp\{\mathbf{z}\mathbf{x}_{\psi}\psi\} - \exp\{\mathbf{x}_{\beta}\beta + \mathbf{z}\mathbf{x}_{\psi}\psi)\}$
= $\exp\{\mathbf{z}\mathbf{x}_{\psi}\psi\}(\mu_0(\mathbf{x}) - \exp\{\mathbf{x}_{\beta}\beta\})$

which cannot be made independent of z unless

$$\mu_0(\mathbf{x}) \equiv \exp\{\mathbf{x}_\beta\beta\}.$$

We must modify the estimating function to be

$$\varphi(X, Y, Z) = (Z - e(X)) \mathbf{X}_{\psi}^{\top} \exp\{-Z \mathbf{X}_{\psi} \psi\} (Y - \exp\{\mathbf{X}_{\beta}\beta + Z \mathbf{X}_{\psi} \psi\})$$
so that

$$\begin{split} \mathbb{E}_{Y|X,Z}[\varphi(X,Y,Z)|X = x, Z = z] \\ &= (z - e(x))\mathbf{x}_{\psi}^{\top} \exp\{-z\mathbf{x}_{\psi}\psi\} \exp\{z\mathbf{x}_{\psi}\psi\}(\mu_0(x) - \exp\{\mathbf{x}_{\beta}\beta\}) \\ &= (z - e(x))\mathbf{x}_{\psi}^{\top}(\mu_0(x) - \exp\{\mathbf{x}_{\beta}\beta\}) \end{split}$$

It then follows that

$$\begin{split} \mathsf{E}_{Z|X}[\mathbb{E}_{Y|X,Z}[\varphi(X,Y,Z)|X=x,Z=z] \mid X=x] \\ &= \mathbf{x}_{\psi}^{\top}(\mu_0(\mathbf{x}) - \exp\{\mathbf{x}_{\beta}\beta\})\mathbb{E}_{Z|X}[(Z-e(X))|X=x] = \mathbf{0} \end{split}$$

if the propensity model is correctly specified. Thus for consistent estimation, we need to solve

$$\sum_{i=1}^{n} \begin{pmatrix} \mathbf{x}_{\beta i}^{\top} \\ (z_i - \mathbf{e}(\mathbf{x}_i)) \mathbf{x}_{\psi_i}^{\top} \exp\{-z_i \mathbf{x}_{\psi_i} \psi\} \end{pmatrix} (y_i - \exp\{\mathbf{x}_{\beta i} \beta + z_i \mathbf{x}_{\psi_i} \psi\}) = \mathbf{0}$$

If outcome Y is *binary*, we might attempt to use a *logistic* model and estimate using

$$\sum_{i=1}^{n} \begin{pmatrix} \mathbf{x}_{\beta i}^{\top} \\ z_{i} \mathbf{x}_{\psi i}^{\top} \end{pmatrix} (y_{i} - \operatorname{expit}\{\mathbf{x}_{\beta i}\beta + z_{i} \mathbf{x}_{\psi i}\psi\}) = \mathbf{0}$$

where

$$\operatorname{expit}(\mathbf{x}) = \frac{e^{\mathbf{x}}}{1 + e^{\mathbf{x}}}.$$

This method is not robust to mis-specification of the mean model, and it cannot be rescued using the previous trick.

A solution is found by noting that

$$\begin{split} \exp\{\mathbf{x}_{\psi}\psi\} &= \frac{\Pr[Y=1|Z=1,X=x]/\Pr[Y=0|Z=1,X=x]}{\Pr[Y=1|Z=0,X=x]/\Pr[Y=0|Z=0,X=x]}\\ &\equiv \frac{\Pr[Z=1|Y=1,X=x]/\Pr[Z=0|Y=1,X=x]}{\Pr[Z=1|Y=0,X=x]/\Pr[Z=0|Y=0,X=x]} \end{split}$$

that is, ψ parameterizes the conditional *log-odds ratio*.

Therefore we can construct a doubly robust estimator assuming that **either**

- the model for *Y* given *X* and *Z*, **or**
- the model for Z given X and Y

is correctly specified, given that the log-odds model is correctly specified.

The R package drgee uses the estimating procedure based on first estimating parameters in the two systems

$$\sum_{i=1}^{n} \begin{pmatrix} \mathbf{x}_{\gamma i}^{\top} \\ y_{i} \mathbf{x}_{\psi i}^{\top} \end{pmatrix} (z_{i} - \operatorname{expit}(\mathbf{x}_{\gamma i} \gamma + y_{i} \mathbf{x}_{\psi i} \psi^{\dagger})) = \mathbf{0}$$

and

$$\sum_{i=1}^{n} \begin{pmatrix} \mathbf{x}_{\beta i}^{\top} \\ z_{i} \mathbf{x}_{\psi i}^{\top} \end{pmatrix} (y_{i} - \operatorname{expit}(\mathbf{x}_{\beta i}\beta + z_{i} \mathbf{x}_{\psi i}\psi^{\ddagger})) = \mathbf{0}$$

utilizing two additional nuisance parameters ψ^{\dagger} and ψ^{\ddagger} . This yields estimates of β and γ .

Then ψ is estimated using

$$\sum_{i=1}^{n} (z_i - e^*(\mathbf{x}_i; \psi, \hat{\beta}, \hat{\gamma})) \mathbf{x}_{\psi i}^{\top} \left(y_i - \operatorname{expit}(\mathbf{x}_{\beta i} \hat{\beta} + z_i \mathbf{x}_{\psi i} \psi) \right) = \mathbf{0}$$

where

$$\mathbf{e}^*(\mathbf{x};\psi,\widehat{\beta},\widehat{\gamma}) = \left[1 + \frac{1 - \operatorname{expit}(\mathbf{x}_{\gamma}\widehat{\gamma})}{\operatorname{expit}(\mathbf{x}_{\gamma}\widehat{\gamma})} \frac{\operatorname{expit}(\mathbf{x}_{\beta}\widehat{\beta})}{\operatorname{expit}(\mathbf{x}_{\beta}\widehat{\beta} + \mathbf{x}_{\psi}\psi)}\right]^{-1}$$

It can be shown that this is an unbiased estimating equation.

.

In the binary case, we have mainly focussed on the estimation of the average treatment effect (ATE)

$$\mathbb{E}[Y(\mathbf{1}) - Y(\mathbf{0})] = \mathbb{E}_{Y|Z}^{\varepsilon}[Y \mid Z = \mathbf{1}] - \mathbb{E}_{Y|Z}^{\varepsilon}[Y \mid Z = \mathbf{0}]$$

The average treatment effect on the treated (ATT) is

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

using the potential outcome notation.

That is, the ATT aims to identify the causal effect on intervening to change Z = 0 to Z = 1 but *only* in the subpopulation of individuals who are *observed* to receive treatment

Note that

$$\mathbb{E}[Y(1) - Y(0)] = \mathbb{E}[Y(1) - Y(0) \mid Z = 0] \Pr[Z = 0]$$
$$+ \mathbb{E}[Y(1) - Y(0) \mid Z = 1] \Pr[Z = 1]$$

In this calculation, we imagine

- ▶ the observational distribution $f_{X,Y,Z}^{\mathcal{O}}$ generating the observed data $\{(x_i, y_i, z_i), i = 1, \dots, n\}$
- in the subgroup observed to have Z = 1, we then consider a second (hypothetical) experimental intervention to change Z to z which over-rides the original Z if z = 0,
- we then consider comparison of hypothetical outcomes between the two hypothetical subgroups indexed by z.



 $f_X(\mathbf{x}) f_{Z|X}(z|\mathbf{x}) f_{Y(\mathbf{0}),Y(\mathbf{1})|X}(y_0,y_1|\mathbf{x}) f_{Y|Z,Y(\mathbf{0}),Y(\mathbf{1})}(y|z,y_0,y_1)$

defines the observational distribution, where

$$f_{Y|Z,Y(0),Y(1)}(y|z,y_0,y_1) = \left\{egin{array}{cc} 1 & y = (1-z)y_0 + zy_1 \ 0 & ext{otherwise} \end{array}
ight.$$

Then

$$\begin{split} f_{Y|X,Z}(y|x,z) &= \int f_{Y|X,Z,Y(0),Y(1)}(y|x,z,y_0,y_1) f_{Y(0),Y(1)|X,Z}(y_0,y_1|x,z) \ dy_0 \ dy_1 \\ &= \int f_{Y|Z,Y(0),Y(1)}(y|z,y_0,y_1) f_{Y(0),Y(1)|X}(y_0,y_1|x) \ dy_0 \ dy_1 \\ &= f_{Y(z)|X}(y|x) \end{split}$$

as the integral reduces to a evaluation at a single point where (y_0,y_1) satisfy

$$y=(1-z)y_0+zy_1.$$

For the ATT, we can represent the quantity of interest using a modified DAG that proposes a second hypothetical binary treatment, A.

We allow Z to cause A, and then allow Z to act as a selection mechanism, but ensure that

 $X \perp\!\!\!\perp A \mid Z.$



where the implied joint model is

 $f_X(x)f_{Z|X}(z|x)f_{Y(\textbf{0}),Y(\textbf{1})|X}(y_0,y_1|x)f_{A|Z}(a|z)f_{Y|A,Y(\textbf{0}),Y(\textbf{1})}(y|a,y_0,y_1)$

This is the new 'experimental' distribution \mathcal{E} .

Again

$$f^arepsilon_{Y|A,Y(m{0}),Y(m{1})}(y|a,y_0,y_1) = \left\{egin{array}{cc} 1 & y=(1-a)y_0+ay_1\ 0 & ext{otherwise} \end{array}
ight.$$

As A is binary, the model $f_{A|Z}(a|z)$ must take the form

$$f^{arepsilon}_{A|Z}(a|z) = p^{\mathsf{a}}_{z}(1-p_{z})^{1-a} \qquad a,z\in\{0,1\}$$

for $0 \leq p_z \leq 1$ for z = 0, 1.

We can then express the ATT via the new DAG as

$$\mathbb{E}^{\mathcal{E}}_{Y|A,Z}[Y|A=\mathbf{1},Z=1] - \mathbb{E}^{\mathcal{E}}_{Y|A,Z}[Y|A=\mathbf{0},Z=1]$$

We have that from the DAG that

$$f^{\mathcal{E}}_{Y|A,X,Z}(y|a,x,z) \equiv f^{\mathcal{E}}_{Y|A,X}(y|a,x)$$

and hence as before

$$\begin{split} & f_{Y|A,X}^{\varepsilon}(y|a,x) \\ & = \int f_{Y|A,Y(0),Y(1)}^{\varepsilon}(y|a,y_0,y_1) f_{Y(0),Y(1)|A,X}^{\varepsilon}(y_0,y_1|a,x) \ dy_0 \ dy_1 \\ & = f_{Y(a)|X}^{\varepsilon}(y|x). \end{split}$$

Also from the DAG, $X \perp \!\!\!\perp A \mid Z$, so for all a, x, z

$$f_{X|A,Z}^{\mathcal{E}}(\mathbf{x}|\mathbf{a},z) \equiv f_{X|Z}^{\mathcal{E}}(\mathbf{x}|z).$$

For a = 0, 1, we therefore have

$$\begin{split} \mathsf{E}^{\varepsilon}_{Y\mid A,Z}[Y\mid A = \mathsf{a}, Z = z] \\ = \iint y \, f^{\varepsilon}_{Y\mid A,X}(y\mid \mathsf{a}, x) f^{\varepsilon}_{X\mid Z}(x\mid z) \, dy \, dx. \end{split}$$

Note that there is a potential *incompatibility* in the conditioning between

$$f^{\mathcal{E}}_{Y \mid A, X}(y \mid \mathsf{a}, \mathbf{x}) \qquad \text{and} \qquad f^{\mathcal{E}}_{X \mid Z}(\mathbf{x} \mid z)$$

when we try to write the integral in terms of the data generating mechanism.
As before, choosing the form of

$$f^{arepsilon}_{Y\mid A,X}(y\mid \mathsf{a},x)$$

is to be avoided if possible.

We seek to resolve the incompatibility using the importance sampling trick, and write the expectation with respect to the observational model

$$f^{\mathcal{O}}_{Y|X,Z}(y|x,z)f^{\mathcal{O}}_{Z|X}(z|x)f^{\mathcal{O}}_X(x).$$

First note that

$$f^arepsilon_{X\mid Z}(x\mid z) = rac{f^arepsilon_{Z\mid X}(z\mid x) f^arepsilon_X(x)}{f^arepsilon_Z(z)}$$

so the integral can be rewritten

$$\frac{1}{f^{\varepsilon}_{Z}(z)} \iint y \, f^{\varepsilon}_{Y\mid A, X}(y \mid \mathbf{a}, x) f^{\varepsilon}_{Z\mid X}(z \mid x) f^{\varepsilon}_{X}(x) \, dy \, dx.$$

For t = 0, 1, we can re-write the integrand using the importance sampling trick as

$$y\, f^{\mathcal{E}}_{Y\mid A,X}(y\mid \mathbf{a},x) \frac{f^{\mathcal{E}}_{Z\mid X}(z\mid x)}{f^{\mathcal{E}}_{Z\mid X}(\mathbf{a}\mid x)}\, f^{\mathcal{E}}_{Z\mid X}(\mathbf{a}\mid x) f^{\mathcal{E}}_{X}(\mathbf{x})$$

which can be rearranged to

$$\left\{ y \; \frac{f^{\varepsilon}_{Z\mid X}(z\mid \mathbf{x})}{f^{\varepsilon}_{Z\mid X}(\mathbf{a}\mid \mathbf{x})} \right\} \; f^{\varepsilon}_{Y\mid A, X}(y\mid \mathbf{a}, \mathbf{x}) \; f^{\varepsilon}_{Z\mid X}(\mathbf{a}\mid \mathbf{x}) f^{\varepsilon}_{X}(\mathbf{x}),$$

Comparing the observational and experimental DAGs, we see that for all \boldsymbol{x} and \boldsymbol{z}

$$f_{Z|X}^{\mathcal{E}}(z \mid x) \equiv f_{Z|X}^{\mathcal{O}}(z \mid x) \qquad f_{X}^{\mathcal{E}}(x) \equiv f_{X}^{\mathcal{O}}(x) \qquad f_{Z}^{\mathcal{E}}(z) \equiv f_{Z}^{\mathcal{O}}(z).$$

Also, we have for any t and y that

$$f^{\mathcal{E}}_{Y\mid A,X}(y\mid t,x) \equiv f^{\mathcal{O}}_{Y\mid X,Z}(y\mid x,t).$$

Therefore we have

$$f^{\mathcal{E}}_{Y\mid A,X}(y\mid \mathbf{a},x) \ f^{\mathcal{E}}_{Z\mid X}(\mathbf{a}\mid x) f^{\mathcal{E}}_{X}(x) \equiv f^{\mathcal{O}}_{Y\mid X,Z}(y\mid x,\mathbf{a}) \ f^{\mathcal{O}}_{Z\mid X}(\mathbf{a}\mid x) f^{\mathcal{O}}_{X}(x).$$

Average Treatment Effect on the Treated

Thus

$$\begin{split} \mathbb{E}_{Y|A,Z}^{\mathcal{E}}[Y \mid A = \mathbf{a}, Z = z] \\ &= \frac{1}{f_Z^{\mathcal{O}}(z)} \iint \left\{ y \frac{f_{Z|X}^{\mathcal{O}}(z \mid x)}{f_{Z|X}^{\mathcal{O}}(\mathbf{a} \mid x)} \right\} \ f_{X,Y,Z}^{\mathcal{O}}(x, y, \mathbf{a}) \ dy \ dx \\ &= \frac{1}{f_Z^{\mathcal{O}}(z)} \iiint \left\{ \mathbbm{1}_{\{\mathbf{a}\}}(t) y \frac{f_{Z|X}^{\mathcal{O}}(z \mid x)}{f_{Z|X}^{\mathcal{O}}(t \mid x)} \right\} \ f_{X,Y,Z}^{\mathcal{O}}(x, y, t) \ dy \ dx \ dt \\ &= \frac{1}{f_Z^{\mathcal{O}}(z)} \mathbb{E}_{X,Y,Z}^{\mathcal{O}}\left[\mathbbm{1}_{\{\mathbf{a}\}}(Z) Y \ \frac{f_{Z|X}^{\mathcal{O}}(z \mid X)}{f_{Z|X}^{\mathcal{O}}(Z \mid X)} \right]. \end{split}$$

For the ATT, we are interested only in z = 1. The momentbased estimator is therefore

$$\widehat{\mathbb{E}}_{Y|A,Z}^{\varepsilon}[Y \mid A = a, Z = 1] = \frac{\sum\limits_{i=1}^{n} \mathbb{1}_{\{a\}}(Z_i)w_1(X_i, Z_i)Y_i}{\sum\limits_{i=1}^{n} \mathbb{1}_{\{1\}}(Z_i)}$$

where

$$w_z(X_i, Z_i) = \frac{f_{Z|X}^{\mathcal{O}}(z \mid X_i)}{f_{Z|X}^{\mathcal{O}}(Z_i \mid X_i)}$$

When a = 1,

$$\mathbb{1}_{\{a\}}(Z_i)w_1(X_i, Z_i) = \mathbb{1}_{\{1\}}(Z_i) = Z_i$$
 w.p. 1

as the weight is identically 1, so therefore

$$\widehat{\mathbb{E}}_{Y|A,Z}^{\varepsilon}[Y \mid A = 1, Z = 1] = \frac{\sum\limits_{i=1}^{n} Z_i Y_i}{\sum\limits_{i=1}^{n} Z_i}$$

that is, the mean in the treated group.

When a = 0,

$$\begin{split} \mathbb{1}_{\{\mathbf{a}\}}(Z_i) w_1(X_i, Z_i) &= \mathbb{1}_{\{\mathbf{0}\}}(Z_i) \frac{f_{Z|X}^{\mathcal{O}}(1 \mid X_i)}{f_{Z|X}^{\mathcal{O}}(Z_i \mid X_i)} \\ &= (1 - Z_i) \frac{f_{Z|X}^{\mathcal{O}}(1 \mid X_i)}{f_{Z|X}^{\mathcal{O}}(\mathbf{0} \mid X_i)}. \end{split}$$

Average Treatment Effect on the Treated

Therefore

$$\hat{\mathsf{E}}^{\varepsilon}_{Y|A,Z}[Y \mid A = \mathbf{0}, Z = 1] = \frac{\sum_{i=1}^{n} (1 - Z_i) w(X_i) Y_i}{\sum_{i=1}^{n} Z_i}$$

where

$$w(X_i) = rac{f^{\mathcal{O}}_{Z|X}(1 \mid X_i)}{f^{\mathcal{O}}_{Z|X}(0 \mid X_i)} = rac{e(X_i)}{1 - e(X_i)}$$

That is, this estimator is a weighted sum of contributions from the *untreated* individuals.

Thus the estimator for the ATT is

$$rac{\sum\limits_{i=1}^n (Z_i-(1-Z_i)w(X_i))Y_i}{\sum\limits_{i=1}^n Z_i}.$$

Under the standard assumptions, this estimator is consistent for the ATT and asymptotically normally distributed if e(x) is correctly specified; that is, it is *singly robust*. To achieve *double robustness*, we proceed in the usual fashion and augment the estimand for a = 0 as follows:

$$\mathbb{E}[Y(0) \mid Z = 1] = \mathbb{E}[Y(0) - \mu(X, 0) \mid Z = 1] + \mathbb{E}[\mu(X, 0) \mid Z = 1]$$

where

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

is the modelled conditional mean for Y.

Average Treatment Effect on the Treated

To estimate the first term, we use

$$\widehat{\mathbb{E}}[Y(\mathbf{0}) - \mu(X, \mathbf{0}) \mid Z = 1] = \frac{\sum_{i=1}^{n} (1 - Z_i) w(X_i) (Y_i - \mu(X_i, \mathbf{0}))}{\sum_{i=1}^{n} Z_i}$$

as in the singly robust case. For the second term, we have

$$\widehat{\mathbb{E}}[\mu(X,\mathbf{0}) \mid Z=1] = \frac{\sum\limits_{i=1}^{n} Z_i \mu(X_i,\mathbf{0})}{\sum\limits_{i=1}^{n} Z_i}$$

Therefore

$$\widehat{E}[Y(\mathbf{0}) \mid Z = 1] = \frac{\sum_{i=1}^{n} (1 - Z_i) w(X_i) (Y_i - \mu(X_i, \mathbf{0})) + Z_i \mu(X_i, \mathbf{0})}{\sum_{i=1}^{n} Z_i}$$

which yields the augmented ATT estimator

$$\frac{\sum_{i=1}^{n} (Z_i - (1 - Z_i) w(X_i))(Y_i - \mu(X_i, \mathbf{0}))}{\sum_{i=1}^{n} Z_i}.$$

Suppose we have the following system:



where $U \perp W$, and $Y \perp W \mid Z$.

We aim to find out the causal effect of Z on Y (both scalars).

- There is an open biasing path ZUY
- ➤ We cannot condition on U to block this path as it is unmeasured.

Suppose the structural model underlying this system is linear

$$Y_i = \beta_0 + Z_i \psi_0 + \mathbf{u}_i \zeta + \varepsilon_i$$

where

$$\mathbb{E}[\varepsilon_i|\mathbf{u}_i]=0$$

and for all *i*, $Z_i \perp \varepsilon_i$.

Clearly, regressing Y on Z alone will lead to inconsistent estimation of $\psi_0.$

Suppose now that W (also scalar) is observed and we in fact have the two-equation structural system

$$Y_i = \beta_0 + Z_i \psi_0 + \mathbf{u}_i \zeta + \varepsilon_i$$
$$Z_i = \gamma_0 + W_i \gamma_1 + \epsilon_i$$

with ϵ_i uncorrelated with W_i and ε_i . Then, by the usual trick, we have

$$W_i Y_i = W_i \beta_0 + W_i Z_i \psi_0 + W_i \mathbf{u}_i \zeta + W_i \varepsilon_i$$
$$W_i Z_i = W_i \gamma_0 + W_i^2 \gamma_1 + W_i \epsilon_i$$

Taking expectations and standardizing, using the usual arguments from linear regression, we have

 $Cov[W, Y] = \psi_0 Cov[W, Z]$ $Cov[W, Z] = \gamma_1 Var[W]$

as W_i and U_i are *independent* by assumption.

Hence on rearrangement we have that

$$\psi_0 = \frac{\operatorname{Cov}[W, Y]}{\operatorname{Cov}[W, Z]} = \frac{\operatorname{Cov}[W, Y] / \operatorname{Var}[W]}{\operatorname{Cov}[W, Z] / \operatorname{Var}[W]} = \frac{\operatorname{Cov}[W, Y] / \operatorname{Var}[W]}{\gamma_1}$$

This suggests a procedure to estimate ψ_0 :

- (i) regress Z on W to obtain estimate $\widehat{\gamma}_1$
- (ii) regress *Y* on *W* to obtain estimate $\widehat{\lambda}_1$
- (iii) form the estimate of ψ_0 as

$$\widehat{\psi}_0 = \frac{\widehat{\lambda}_1}{\widehat{\gamma}_1}$$

This is a an *instrumental variable ratio estimator* using the *instrument* W.



Note

- For this method to work, we cannot have an *interaction* between Z and U in the structural model
- If there is only *weak association* between W and Z, so that γ_1 is close to zero, this estimator can have undesirable properties.
- Statistical properties of $\hat{\psi}_0$ not trivial to derive.

A common situation in which the linear model is not quite appropriate is when W and Z are both *binary*. For example, we might have in a randomized study

▶ *W* is the randomized *allocated* treatment indicator; eg

 $W = 1 \implies$ "allocated to treatment arm"

► *Z* is the actual treatment *received* indicator; eg

 $Z = 1 \implies$ "took treatment"

- Could have W = 1, Z = 0 (*non-compliance*);
- In some situations, could have W = 0, Z = 1.

Then we have by iterated expectation, for w = 0, 1

$$\begin{split} \mathbb{E}[Y|W = w] = \mathbb{E}[Y|W = w, Z = 0] \Pr[Z = 0|W = w] \\ + \mathbb{E}[Y|W = w, Z = 1] \Pr[Z = 1|W = w] \\ = \mathbb{E}[Y|Z = 0] \Pr[Z = 0|W = w] \\ + \mathbb{E}[Y|Z = 1] \Pr[Z = 1|W = w] \end{split}$$

as $Y \perp \!\!\!\perp W \mid Z$.

Therefore

$$\begin{split} \mathbb{E}[Y|W = 1] &= \mathbb{E}[Y|Z = 0] \Pr[Z = 0|W = 1] \\ &+ \mathbb{E}[Y|Z = 1] \Pr[Z = 1|W = 1] \\ \mathbb{E}[Y|W = 0] &= \mathbb{E}[Y|Z = 0] \Pr[Z = 0|W = 0] \\ &+ \mathbb{E}[Y|Z = 1] \Pr[Z = 1|W = 0] \end{split}$$

Hence subtracting the second equation from the first

$$\begin{split} \mathbb{E}[Y|W = 1] - \mathbb{E}[Y|W = 0] \\ &= \mathbb{E}[Y|Z = 0](\Pr[Z = 0|W = 1] - \Pr[Z = 0|W = 0]) \\ &+ \mathbb{E}[Y|Z = 1](\Pr[Z = 1|W = 1] - \Pr[Z = 1|W = 0]) \end{split}$$

But

$$Pr[Z = 0|W = 1] - Pr[Z = 0|W = 0]$$
$$= Pr[Z = 1|W = 0] - Pr[Z = 1|W = 1]$$

Therefore

$$\begin{split} \mathbb{E}[Y|W = 1] - \mathbb{E}[Y|W = 0] \\ &= (\Pr[Z = 1|W = 1] - \Pr[Z = 1|W = 0]) \\ &\times (\mathbb{E}[Y|Z = 1] - \mathbb{E}[Y|Z = 0]) \end{split}$$

That is, provided $\Pr[Z = 1 | W = 1] \neq \Pr[Z = 1 | W = 0]$

$$\begin{split} \mathbb{E}[Y|Z = 1] - \mathbb{E}[Y|Z = 0] \\ &= \frac{\mathbb{E}[Y|W = 1] - \mathbb{E}[Y|W = 0]}{\Pr[Z = 1|W = 1] - \Pr[Z = 1|W = 0]} \end{split}$$

or equivalently

$$\mathbb{E}[Y|Z=1] - \mathbb{E}[Y|Z=0] = \frac{\mathbb{E}[Y|W=1] - \mathbb{E}[Y|W=0]}{\mathbb{E}[Z|W=1] - \mathbb{E}[Z|W=0]}.$$

As W is observed, we can deploy the moment-based estimator

$$\frac{\sum_{i=1}^{n} W_i Y_i}{\sum_{i=1}^{n} W_i} - \frac{\sum_{i=1}^{n} (1 - W_i) Y_i}{\sum_{i=1}^{n} (1 - W_i)}$$
$$\frac{\sum_{i=1}^{n} W_i Z_i}{\sum_{i=1}^{n} W_i} - \frac{\sum_{i=1}^{n} (1 - W_i) Z_i}{\sum_{i=1}^{n} (1 - W_i)}$$

This is the binary analogue of the earlier regression-based ratio estimator.

In terms of potential outcomes, we can consider potential outcomes for treatment received

 $\{Z({\color{black} 0}),Z({\color{black} 1})\}$

corresponding to two settings $\{0,1\}$ of W, and potential outcomes for outcome

$$\{Y(w, z), w = 0, 1, z = 0, 1\}$$

corresponding to the four possible combinations of W and Z.

We have the "consistency" requirement relating the potential outcomes to the observed data:

$$Z = (1 - W)Z(0) + WZ(1)$$

and

$$Y = \sum_{\mathsf{z}=0}^{1} \sum_{\mathsf{w}=0}^{1} \mathbb{1}_{\{\mathsf{w}\}}(W) \mathbb{1}_{\{\mathsf{z}\}}(Z) Y(\mathsf{w},\mathsf{z}).$$

We must have that for z = 0, 1

$$Y(\mathbf{1},\mathbf{Z})=Y(\mathbf{0},\mathbf{Z})$$

as $Y \perp W | Z$ by assumption. We have *ignorability*

$$\{Y(\mathbf{w},\mathbf{z}),Z(\mathbf{w}):\forall \mathbf{w},\mathbf{z}\} \perp Z.$$

Consequently

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

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

which simplifies the causal quantity of interest.

We are interested in the contrast

$$\mathbb{E}[Y(1) - Y(0)] \equiv \mathbb{E}[Y(W, 1) - Y(W, 0)].$$

We have that

$$Y = (1 - Z)Y(0) + ZY(1)$$
$$= Y(0) + (Y(1) - Y(0))Z$$

which suggests the linear model for the observed data

$$Y_i = \beta_0 + \psi_0 Z_i + \varepsilon_i$$

We also need that $W \not\perp Z$, that is

 $\mathbb{E}[Z(1)] \neq \mathbb{E}[Z(0)] \quad \text{or} \quad \Pr[Z(1) = 1] \neq \Pr[Z(0) = 1]$

- ▶ Manipulation of *W* does change *Z*;
- ▶ From the assumptions and the DAG, we can state

$$\mathbb{E}[Z(1) - Z(0)] \equiv \mathbb{E}[Z|W = 1] - \mathbb{E}[Z|W = 0].$$

Recall the earlier formula

$$\mathbb{E}[Y|Z=1] - \mathbb{E}[Y|Z=0] = \frac{\mathbb{E}[Y|W=1] - \mathbb{E}[Y|W=0]}{\mathbb{E}[Z|W=1] - \mathbb{E}[Z|W=0]}.$$

We now have that the denominator on the right hand side is non-zero.

For the numerator, we have

$$\mathbb{E}[Y|W = 1] = \mathbb{E}[Y(0) + (Y(1) - Y(0))Z|W = 1]$$
$$= \mathbb{E}[Y(0) + (Y(1) - Y(0))Z(1)]$$

 $\mathbb{E}[Y|W = 0] = \mathbb{E}[Y(0) + (Y(1) - Y(0))Z(0)]$

Thus, taking the difference we have

$$\mathbb{E}[Y|W = 1] - \mathbb{E}[Y|W = 0] = \mathbb{E}[(Y(1) - Y(0))(Z(1) - Z(0))]$$

so therefore

$$\mathbb{E}[Y|Z = \mathbf{1}] - \mathbb{E}[Y|Z = \mathbf{0}] = \frac{\mathbb{E}[(Y(\mathbf{1}) - Y(\mathbf{0}))(Z(\mathbf{1}) - Z(\mathbf{0}))]}{\mathbb{E}[Z(\mathbf{1}) - Z(\mathbf{0})]}.$$

We now try to understand the causal implications of this identity.

Stratum	<i>Z</i> (1)	Z(0)	Description
Compliers	1	0	<i>Follows</i> treatment protocol, takes assigned treatment w.
Defiers	0	1	Defies treatment protocol, takes other treatment $1 - w$.
Always Takers	1	1	<i>Ignores</i> treatment protocol, takes treatment 1.
Never Takers	0	0	<i>Ignores</i> treatment protocol, takes treatment 0.
In the *observed* data:

W	Ζ	Possible strata		
1	0	Defier, Never Taker		
0	1	Defier, Always Taker		
1	1	Complier, Always Taker		
0	0	Complier, Never Taker		

We can never conclude with certainty to which stratum an individual belongs from the observed data.

It is common (and reasonable in most cases) to assume that

$$\Pr[Z(1) \ge Z(0)] = \Pr[Z(1) - Z(0) \ge 0] = 1$$

that is, being assigned W = 1 cannot decrease the value of $Z(\mathbf{w})$ compared with being assigned W = 0.

i.e. being allocated treatment increases or leaves unchanged the received treatment status variable.

This *monotonicity assumption* rules out the possibility of the *Defier* stratum.

Assuming $\Pr[Z(1) - Z(0) \ge 0] = 1$

W Z Possible strata

- 1 0 Defier, Never Taker
- 0 1 Defier, Always Taker
- 1 1 Complier, Always Taker
- 0 0 Complier, Never Taker

To get a further simplification in the expression

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

we use iterated expectation using the partitioning event

Z(1) > Z(0)

which is an event occurs if and only if

$$Z(1) = 1$$
 and $Z(0) = 0$.

We have in the numerator

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

= $\mathbb{E}[(Y(1) - Y(0))(Z(1) - Z(0))|Z(1) > Z(0)]\Pr[Z(1) > Z(0)]$
+ $\mathbb{E}[(Y(1) - Y(0))(Z(1) - Z(0))|Z(1) \le Z(0)]\Pr[Z(1) \le Z(0)]$

and similarly in the denominator

 $\mathbb{E}[(Z(1) - Z(0))]$ = $\mathbb{E}[(Z(1) - Z(0))|Z(1) > Z(0)]\Pr[Z(1) > Z(0)]$ + $\mathbb{E}[(Z(1) - Z(0))|Z(1) \le Z(0)]\Pr[Z(1) \le Z(0)]$ _

Stratum	Z(1) Z(0)		Z(1) - Z(0)	
Compliers	1	0	1	
Defiers	0	1	-1	
Always Takers	1	1	0	
Never Takers	0	0	0	

We therefore have in the numerator:

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

= $\mathbb{E}[(Y(1) - Y(0))|Z(1) > Z(0)]\Pr[Z(1) > Z(0)]$

as only the contribution Z(1) = 1 and Z(0) = 0 remains. Similarly in the denominator:

$$\mathbb{E}[(Z(1) - Z(0))] = \Pr[Z(1) > Z(0)]$$

Therefore

$$\mathbb{E}[Y|Z = 1] - \mathbb{E}[Y|Z = 0] = \frac{\mathbb{E}[(Y(1) - Y(0))(Z(1) - Z(0))]}{\mathbb{E}[Z(1) - Z(0)]}$$
$$\equiv \mathbb{E}[(Y(1) - Y(0))|Z(1) > Z(0)].$$

It is therefore evident that the IV estimation procedure estimates the causal contrast

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

As

$$Z(1) > Z(0) \quad \Longleftrightarrow \quad Z(1) = 1, Z(0) = 0$$

the quantity

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

is termed the

- complier average treatment effect (CATE), or
- complier average causal effect (CACE), or
- local average treatment effect (LATE),

as Z(1) > Z(0) identifies the *Complier* stratum.

Note

- If we observe Z = 1, the subject could be a Complier or Always Taker; we do not know who the Compliers are.
- The Compliers may be different depending on the instrument *W* chosen.
- In general the three quantities

ATE : $\mathbb{E}[(Y(1) - Y(0))]$ ATT : $\mathbb{E}[(Y(1) - Y(0))|Z = 1]$ CATE : $\mathbb{E}[(Y(1) - Y(0))|Z(1) > Z(0)]$

are not equal.

Note

• It is often in practice realistic to make the assumption

$$\Pr[Z=1|W=0]=0$$

which asserts that the Always Takers stratum is empty, and that all subjects for whom Z = 1 are Compliers.

It follows that

$$\frac{\mathbb{E}[Y|Z=1] - \mathbb{E}[Y|Z=0]}{\Pr[Z=1|W=1]} = \mathbb{E}[(Y(1) - Y(0))|Z=1]$$

which equates to the ATT.

Extension: with measured confounders



We have $U \perp W$, and $Y \perp W \mid Z, X$.

We can extend the previous logic by appealing to *conditional* ignorability

$$\{Y(\mathsf{w},\mathsf{z}),Z(\mathsf{w}):\forall \mathsf{w},\mathsf{z}\} \perp Z | X \}$$

and develop similar methods.

Suppose we have the following confounded system:



No adjustment method is feasible.

Including the propensity score, we have



This does not directly help with unmeasured confounding.

However, consider a variable W, with $\mathbb{E}[W] = 0$, such that

$$Z = e(X) + W$$

so that

$$W = Z - e(X) = Z - \mathbb{E}[Z|X]$$

We have

$$\mathbb{E}[W \ e(X)] = \mathbb{E}[(Z - e(X))e(X)] = \mathbb{E}[(e(X) - e(X))e(X)] = 0$$

by iterated expectation. Therefore W and e(X) must be uncorrelated.

We also have that

$$\mathbb{E}[WZ] = \mathbb{E}[(Z - e(X))Z] = \mathbb{E}[Z - Ze(X)] = \mathbb{E}[e(X)(1 - e(X))]$$

as $Z^2 = Z$ (w.p. 1). Therefore as $\mathbb{E}[W] = 0$,
$$\operatorname{Cov}[W, Z] = \mathbb{E}[e(X)(1 - e(X))] > 0$$

under the usual positivity assumption 0 < e(x) < 1. Hence $W \not\perp Z$. If we consider W a manipulable variable, we could propose the following DAG:



It is evident that $Y \perp W \mid X, Z$, and also that $U \perp W$.

Thus

$$W = Z - e(X) = Z - \mathbb{E}[Z|X].$$

is an instrumental variable. We can compute the ratio estimator based on the earlier identity

$$\psi_0 = \frac{\operatorname{Cov}[W, Y]}{\operatorname{Cov}[W, Z]}$$

that is

$$\hat{\psi}_{0} = \frac{\sum_{i=1}^{n} W_{i}Y_{i}}{\sum_{i=1}^{n} W_{i}Z_{i}} = \frac{\sum_{i=1}^{n} (Z_{i} - e(X_{i}))Y_{i}}{\sum_{i=1}^{n} (Z_{i} - e(X_{i}))Z_{i}}$$

which is identical to the *singly robust G-estimator*.

Statistical methods for causal *mediation* analysis can also be formulated using potential outcomes:

$$Z \longrightarrow M \longrightarrow Y$$

where z and m will represent potential levels of treatment Z and mediator M respectively, with *potential outcome*

Y(z, m)

and *potential mediator* value in light of Z = z

 $M(\mathbf{z}).$

In general, both Z and M can take arbitrary values, but we most commonly consider the case when

- *Z* is binary {0, 1}
- ► *M* is
 - binary,
 - discrete, or
 - continuous.

Then following Imai et al. (*Statistical Science*, 2010)

▶ the *total causal effect* (TCE) of treatment is based on

 $Y(\mathbf{Z}, M(\mathbf{Z}))$

i.e. set Z = z and observe the consequence.

For example

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

▶ the *causal mediation effect* (CME) is based on

 $Y(\mathbf{Z}, M(\mathbf{Z^*}))$

i.e. set Z = z and observe the consequence of changing the mediator as if $Z = z^*$

For example

 $\delta(\mathsf{z}) = \mathbb{E}[Y(\mathsf{z}, M(1)) - Y(\mathsf{z}, M(0))] \qquad \mathsf{z} = \mathsf{0}, \mathsf{1}.$

▶ the *direct effect* (DE) of treatment is based on

$$Y(\mathbf{z}^{*}, M(\mathbf{z}))$$

i.e. set the mediator as if Z = z, and observe the consequence of changing treatment according to z^* .

For example

$$\zeta(\mathsf{z}) = \mathbb{E}[Y(\mathsf{1}, M(\mathsf{z})) - Y(\mathsf{0}, M(\mathsf{z}))] \qquad \mathsf{z} = \mathsf{0}, \mathsf{1}.$$

Terminology:

		Imai et al.	Pearl	Robins
Total effect	au	ATCE	Total	Total
Indirect effect	$\delta(\mathbf{z})$	ACME	Natural	Pure/Total
Direct effect	$\zeta(\mathbf{Z})$	ADE	Natural	Pure/Total
with				

Pure : z = 0 Total : z = 1.

That is,

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

= $\mathbb{E}[Y(1, M(1)) - Y(1, M(0))]$
+ $\mathbb{E}[Y(1, M(0)) - Y(0, M(0))]$

so that

$$\tau = \delta(\mathbf{1}) + \zeta(\mathbf{0}) = \delta(\mathbf{0}) + \zeta(\mathbf{1}).$$

We have in the observed data that

$$M = (1-Z)M(\mathbf{0}) + ZM(\mathbf{1})$$

and, if M is binary say,

$$Y = (1 - Z)(1 - M) Y(0, 0)$$
$$+ (1 - Z)M Y(0, 1)$$
$$+ Z(1 - M) Y(1, 0)$$
$$+ Z M Y(1, 1)$$

etc.



The DAG encompasses two ignorability assumptions

(i)
$$\{Y(z^*, m), M(z)\} \perp Z \mid X = x;$$

(ii) $Y(z^*, m) \perp M(z) \mid Z = z, X = x;$

for all z, z^*, m and x, and relies upon two positivity assumptions

$$0 < \Pr[Z = z | X = x] < 1$$

and

$$0 < \Pr[M(z) = m | X = x, Z = z] < 1.$$

Imai et al (2010) gives the following identities:

$$\begin{split} \delta(\mathbf{Z}) &= \iiint y f_{Y|M,X,Z}(y|m,\mathbf{x},\mathbf{Z}) f_{M|X,Z}(m|\mathbf{x},1) f_X(\mathbf{x}) \ dy \ dm \ dx \\ &- \iiint y f_{Y|M,X,Z}(y|m,\mathbf{x},\mathbf{Z}) f_{M|X,Z}(m|\mathbf{x},0) f_X(\mathbf{x}) \ dy \ dm \ dx \end{split}$$

$$\begin{split} \zeta(\mathbf{Z}) &= \iiint y f_{Y|M,X,Z}(y|m,x,1) f_{M|X,Z}(m|x,\mathbf{Z}) f_X(x) \ dy \ dm \ dx \\ &- \iiint y f_{Y|M,X,Z}(y|m,x,0) f_{M|X,Z}(m|x,\mathbf{Z}) f_X(x) \ dy \ dm \ dx \end{split}$$

As for the ATT case, the challenge in estimating these quantities using moment-based methods is the mis-match in the terms

$$f_{Y|M,X,Z}(y|m,x,\mathsf{z})f_{M|X,Z}(m|x,\mathsf{z}^*)$$

whenever $z \neq z^*$.

Model-based estimation via the mean-model

$$\mu(m,x,z) = \int y \; f_{Y|M,X,Z}(y|m,x,z) \; dy$$

and the mediator model

 $f_{M|X,Z}(m|x,z).$

is possible.

It is quite common to encounter longitudinal data with

- ▶ treatment
- confounders
- ▶ outcome

recorded over several time points for each individual.

Longitudinal Data

For example



In such cases, and for the analysis of direct effects (within an interval) the methods already established can be utilized sequentially.

For *cumulative* (across interval) effects of treatment on a single terminal outcome, adjustment methods need careful application.

Longitudinal Data

For a two time point setting:



In this formulation, the time ordering

$$X_1 \longrightarrow Z_1 \longrightarrow X_2 \longrightarrow Z_2 \longrightarrow Y$$

delimits the possible causal pathways.
We can consider the expected counterfactual outcomes associated with treatment *patterns*

 $\mathbb{E}[Y(\mathbf{Z_1},\mathbf{Z_2})]$

or equivalently

$$\mathbb{E}_{Y|Z_1,Z_2}^{\mathcal{E}}[Y|Z_1 = \mathbf{z_1}, Z_2 = \mathbf{z_2}]$$

where the experimental distribution $\ensuremath{\mathcal{E}}$ assumes randomized treatments.

To learn about the APOs for different treatment patterns from observational data is not straightforward.

- Z₁ has a *direct* effect on Y, but also has a *mediated* effect via X₂ and Z₂;
- ► Z₂ has a *direct* effect on Y, but it is *confounded* by X₂; to remove this confounding we need to condition on X₂;
- ► However, conditioning on X₂ blocks the directed path from Z₁ to Y and hence affects the causal effect.

Therefore, we cannot break the confounding by blocking paths by conditioning to get at the aggregate effect.

We may use *inverse weighting* to break the confounding as in the single interval case. For example, for APO

$$\mu(\mathbf{Z_1},\mathbf{Z_2}) = \mathbb{E}_{Y|Z_1,Z_2}^{\mathcal{E}}[Y|Z_1 = \mathbf{Z_1}, Z_2 = \mathbf{Z_2}]$$

we may use the estimator

$$\widetilde{\mu}(\mathsf{Z}_{1},\mathsf{Z}_{2}) = \frac{1}{n} \sum_{i=1}^{n} \frac{\mathbb{1}_{\{\mathsf{Z}_{1}\}}(Z_{1i})\mathbb{1}_{\{\mathsf{Z}_{2}\}}(Z_{2i})}{f^{\mathcal{O}}_{Z_{1},Z_{2}|X_{1},X_{2}}(Z_{1i},Z_{2i}|X_{1i},X_{2i})} Y_{i}$$

Each outcome data point is re-weighted by the IPW weight across the whole treatment sequence.

In the re-weighted data, *model-based* analysis can also be used: for example, we could propose a marginal model

$$\mathbb{E}_{Y|Z_{1},Z_{2}}^{\varepsilon}[Y|Z_{1}=\mathbf{Z}_{1},Z_{2}=\mathbf{Z}_{2}]=\beta_{0}+\psi_{1}\mathbf{Z}_{1}+\psi_{2}\mathbf{Z}_{2}$$

or, using the *total* treatment

$$\mathbb{E}_{Y|Z_1,Z_2}^{\mathcal{E}}[Y|Z_1 = \mathbf{Z}_1, Z_2 = \mathbf{Z}_2] = \beta_0 + \psi_0(\mathbf{Z}_1 + \mathbf{Z}_2)$$

and then perform a *weighted least squares analysis* (WLS) to estimate (ψ_1, ψ_2) or ψ_0 .

Such a model is termed a *marginal structural model* (MSM).

That is for example

$$(\hat{\beta}_0, \hat{\psi}_1, \hat{\psi}_2) = \arg \min_{(\beta_0, \psi_1, \psi_2)} \sum_{i=1}^n w_i (y_i - \beta_0 - \psi_1 z_{1i} - \psi_2 z_{2i})^2$$

where

$$w_i = \frac{1}{f_{Z_1, Z_2 | X_1, X_2}^{\mathcal{O}}(z_{1i}, z_{2i} | x_{1i}, x_{2i})}$$

where

$$f^{\mathcal{O}}_{Z_1,Z_2|X_1,X_2}(z_1,z_2|x_1,x_2)=f^{\mathcal{O}}_{Z_1|X_1}(z_1|x_1)f^{\mathcal{O}}_{Z_2|X_1,X_2,Z_1}(z_2|x_1,x_2,z_1).$$

An alternative weight

$$w_i = \frac{f_{Z_1,Z_2}^{\mathcal{O}}(z_{1i}, z_{2i})}{f_{Z_1,Z_2|X_1,X_2}^{\mathcal{O}}(z_{1i}, z_{2i}|\mathbf{x}_{1i}, \mathbf{x}_{2i})}$$

where

$$f^{\mathcal{O}}_{Z_1,Z_2}(z_{1i},z_{2i}) = f^{\mathcal{O}}_{Z_1}(z_1) f^{\mathcal{O}}_{Z_2|Z_1}(z_2|z_1)$$

is modelled could be used. This generalizes the earlier form of IPW estimator.

If a *non-parametric* model for $f_{Z_1,Z_2}^{\mathcal{O}}(z_{1i}, z_{2i})$ is adopted, then the new weight essentially reduces to the original weight.

Using a *parametric* model,

$$w_i = \frac{f^{\mathcal{O}}_{Z_1,Z_2}(z_{1i}, z_{2i}; \hat{\alpha})}{f^{\mathcal{O}}_{Z_1,Z_2|X_1,X_2}(z_{1i}, z_{2i}|x_{1i}, x_{2i}; \hat{\gamma})}$$

where

• $\alpha = (\alpha_1, \alpha_2)$ is estimated from the model

$$f^{\mathcal{O}}_{Z_{1},Z_{2}}(z_{1},z_{2};\alpha) = f^{\mathcal{O}}_{Z_{1}}(z_{1};\alpha_{1})f^{\mathcal{O}}_{Z_{2}|Z_{1}}(z_{2}|z_{1};\alpha_{2})$$

• $\gamma = (\gamma_1, \gamma_2)$ is estimated from the model

$$\begin{split} f^{\mathcal{O}}_{Z_1,Z_2|X_1,X_2}(z_1,z_2|x_1,x_2;\gamma) &= f^{\mathcal{O}}_{Z_1|X_1}(z_1|x_1;\gamma_1) \\ &\times f^{\mathcal{O}}_{Z_2|X_1,X_2,Z_1}(z_2|x_1,x_2,z_1;\gamma_2) \end{split}$$

It is also possible to carry out a *conditional* analysis given baseline predictors X_1 which are not subject to the influence of any treatment: for example

$$\mathbb{E}_{Y|Z_1,Z_2,X_1}^{\mathcal{E}}[Y|Z_1 = \mathbf{Z_1}, Z_2 = \mathbf{Z_2}, X_1 = x_1] = \beta_0 + x_1\beta_1 + \psi_1\mathbf{Z_1} + \psi_2\mathbf{Z_2}$$

for which the so-called *stabilized* weights

$$w_i = \frac{f_{Z_1, Z_2 | X_1}^{\mathcal{O}}(z_{1i}, z_{2i} | x_{1i})}{f_{Z_1, Z_2 | X_1, X_2}^{\mathcal{O}}(z_{1i}, z_{2i} | x_{1i}, x_{2i})}$$

should be used.

Separate models are needed for numerator and denominator

Denominator:

$$f^{\mathcal{O}}_{Z_1,Z_2|X_1,X_2}(z_1,z_2|x_1,x_2) = f^{\mathcal{O}}_{Z_1|X_1}(z_1|x_1)f^{\mathcal{O}}_{Z_2|X_1,X_2,Z_1}(z_2|x_1,x_2,z_1)$$

Numerator:

$$f^{\mathcal{O}}_{Z_1,Z_2|X_1,X_2}(z_1,z_2|x_1) = f^{\mathcal{O}}_{Z_1|X_1}(z_1|x_1) f^{\mathcal{O}}_{Z_2|X_1,Z_1}(z_2|x_1,z_1)$$

parameterized by α and γ respectively, say.

That is, after cancelling terms,

$$w_i = rac{f^{\mathcal{O}}_{Z_2|X_1,Z_1}(z_2|\mathbf{x}_1,z_1;\widehat{lpha})}{f^{\mathcal{O}}_{Z_2|X_1,X_2,Z_1}(z_2|\mathbf{x}_1,\mathbf{x}_2,z_1;\widehat{\gamma})}$$

This weight may be less extreme than the unstabilized counterpart.

Note that conditioning on X_1 in the outcome model is necessary to account for possible confounding.

Note

The term 'stabilized' is slightly misleading; the introduction of the numerator term *changes the estimand*, so the original and stabilized versions of the MSM estimate *different* quantities.

In many cases the stabilized weights will be more uniform, and this has the effect of reducing estimator variance, but the estimation target is changed. For treatment sequence (z_1, z_2) , we seek the *optimal* treatment sequence (z_1^{opt}, z_2^{opt}) to *maximize* expected response.

- The objective is to establish a sequence of decision rules to be applied at each treatment interval;
- For interval k = 1, 2, the decision rule for interval k must utilize the information available only up to and including (the start of) that interval;
- The decision rule should be *hyperopic*, that is, prioritize *long-term* outcomes over short-term (*myopic*) outcomes;
- The rules will be learned from observational data, where optimal treatment is not necessarily guaranteed.

To solve this problem, a recursive approach is adopted

- optimal treatment sequence will be *personalized*, that is, dependent on an individual's characteristics;
- the computation will use reverse (or Bellman) optimization;
- first optimize Stage 2 treatment, and then optimize Stage 1 treatment assuming optimal treatment at Stage 2.

Decompose the expected (counterfactual) response as

$$\begin{split} \mathbb{E}[Y(\mathsf{Z}_1,\mathsf{Z}_2)] &= \mathbb{E}[Y(z_1^{\text{opt}},z_2^{\text{opt}})] \\ &- \left\{ \mathbb{E}[Y(z_1^{\text{opt}},z_2^{\text{opt}}) - Y(\mathsf{Z}_1,z_2^{\text{opt}})] \right\} \\ &- \left\{ \mathbb{E}[Y(\mathsf{Z}_1,z_2^{\text{opt}}) - Y(\mathsf{Z}_1,\mathsf{Z}_2)] \right\} \end{split}$$

This is the basis of the *Structural Nested Mean Model* (SNMM).

Optimal Dynamic Treatment Regimes

Further

$$\begin{split} \mathbb{E}[Y(z_1^{\text{opt}}, z_2^{\text{opt}}) - Y(\mathbf{Z}_1, z_2^{\text{opt}})] &= \mathbb{E}[Y(z_1^{\text{opt}}, z_2^{\text{opt}}) - Y(0, z_2^{\text{opt}})] \\ &- \mathbb{E}[Y(\mathbf{Z}_1, z_2^{\text{opt}}) - Y(0, z_2^{\text{opt}})] \end{split}$$

$$\begin{split} \mathbb{E}[Y(\mathsf{z}_1, z_2^{\text{opt}}) - Y(\mathsf{z}_1, \mathsf{z}_2)] &= \mathbb{E}[Y(\mathsf{z}_1, z_2^{\text{opt}}) - Y(\mathsf{z}_1, 0)] \\ &- \mathbb{E}[Y(\mathsf{z}_1, \mathsf{z}_2) - Y(\mathsf{z}_1, 0)] \end{split}$$

We specify models

$$\begin{split} & \mathbb{E}[Y(\mathsf{Z}_1,\mathsf{Z}_2) - Y(\mathsf{Z}_1,0)] = (\mathbf{X}_2\psi_2)\mathsf{Z}_2 & \text{Stage 2} \\ & \mathbb{E}[Y(\mathsf{Z}_1,z_2^{\text{opt}}) - Y(0,z_2^{\text{opt}})] = (\mathbf{X}_1\psi_1)\mathsf{Z}_1 & \text{Stage 1} \end{split}$$

where

- ► X₂ can depend on all data including observed treatments – observed up to Stage 2.
- X_1 can depend on all data observed up to Stage 1.

Optimal Dynamic Treatment Regimes

1. Use *G*-estimation on *Y* conditioning on $(\mathbf{x}_1, z_1, \mathbf{x}_2, z_2)$ at second stage to estimate ψ_2 using a proposed mean model

$$Y = \mathbf{X}_{21}\beta_2 + z_2\mathbf{X}_{22}\psi_2 + \varepsilon$$

2. For each individual infer the optimal Stage 2 treatment

$$z_2^{ ext{opt}} = \mathbb{1}\{\mathbf{X}_{22}\widehat{\psi}_2 > 0\}$$

3. Form *pseudo-outcome* Y_1

$$Y_1 = Y - (\mathbf{X}_{22}\widehat{\psi}_2)(z_2^{\text{opt}} - z_2)$$

4. Use G-estimation on Y_1 conditioning on (\mathbf{x}_1, z_1) at first stage to estimate ψ_1 using a proposed mean model

$$Y_1 = \mathbf{X}_{11}\beta_1 + z_1\mathbf{X}_{12}\psi_1 + \varepsilon$$

5. For each individual infer the optimal Stage 1 treatment

$$z_1^{\text{opt}} = \mathbb{1}\{\mathbf{X}_{11}\widehat{\psi}_1 > 0\}$$

The method is robust to mis-specification of the *nuisance mean model*, provided the treatment model is correctly specified.

Can infer optimized potential outcome

$$Y + (\mathbf{X}_{12} \widehat{\psi}_1) (z_1^{ ext{opt}} - z_1) + (\mathbf{X}_{22} \widehat{\psi}_2) (z_2^{ ext{opt}} - z_2)$$

which takes the observed outcome Y and adds

the additional benefit of optimal treatment at stage 1

$$(\mathbf{X}_{12}\widehat{\psi}_1)(z_1^{\mathrm{opt}}-z_1)$$

▶ the additional benefit of optimal treatment at stage 2

$$(\mathbf{X}_{22}\widehat{\psi}_2)(z_2^{\text{opt}}-z_2).$$

Many of the adjustment methods rely upon the construction of estimators assuming knowledge of

- confounder structure
- time-ordering
- propensity score construction

all of which depend on knowledge of the data generating DAG.

In many applications, the DAG is *unknown*.

There are several approaches to *causal discovery*, that is, learning the DAG from the observed data

- PC algorithm: pcalg in R
- SGS algorithm

Both rely on recursive identification of conditional independencies, identification of colliders and edge orientation procedures based on statistical tests.

Can be effective, but typically need large sample sizes.