

CI = A DEFINITION OF CAUSAL EFFECT

* Notation

• $A = \begin{cases} 1 & \text{if Treated} \\ 0 & \text{o/w} \end{cases} = \text{dichotomous treatment variable}$

(assuming here that there is a single version of the treatment; which is not always the same)

• $Y = \begin{cases} 1 & \text{if Survival} \\ 0 & \text{if death} \end{cases} = \text{dichotomous outcome variable}$

• $Y^{a=1} = Y^1 = \text{outcome variable that would have been observed under the treatment value } a=1$

• $Y^{a=0} = Y^0 = \text{--- " --- } a=0$

Y^0 & Y^1 are referred to as POTENTIAL OUTCOMES.

"either one of these may be potentially observed".

↳ one of them is actually factual i.e. observed.

If $A = a$, then $Y = Y^a$

"CONSISTENCY"

$$\Leftrightarrow Y = A Y^1 + (1 - A) Y^0$$

(Y_i^0, Y_i^1, A_i)

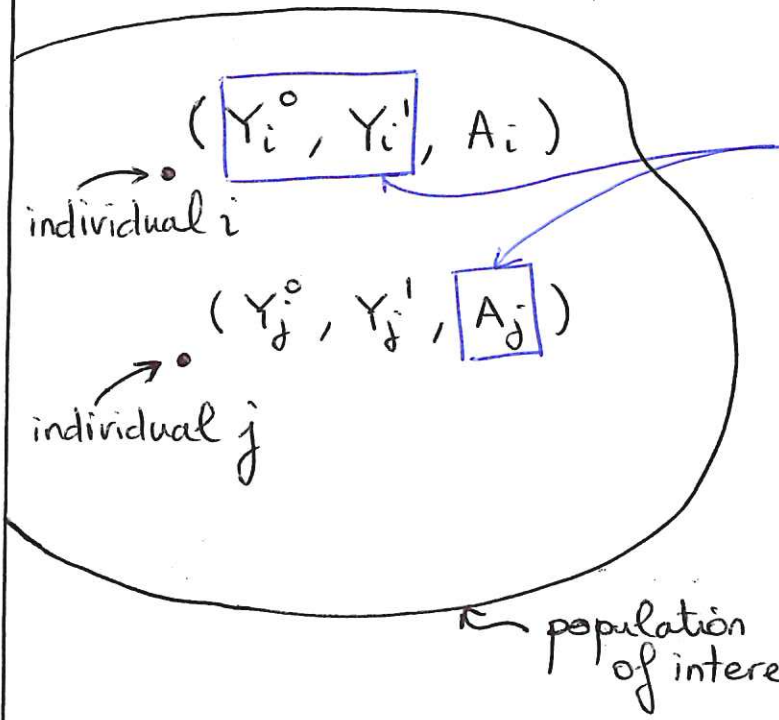
↑
individual i

population

The causal effect for individual i is
 $Y_i^1 - Y_i^0$

Individual effects cannot be identified since only one potential outcome is factual.

↓ instead, we look into aggregated effects i.e. average causal effects in a population of individuals.



The causal effect of individual i is well defined provided it does not depend on the treatment value of individual j : i.e. when there is no interference.

SUTVA assumption
Rubin (1980)

Under SUTVA, we can compute / estimate $P(Y^a = 1)$.

finite population,
all observed.

infinite / large population,
only a sample is observed

Def: An average causal effect of treatment A on outcome Y is present if $P(Y^1 = 1) \neq P(Y^0 = 1)$.

↖ An absence of average causal effect does not imply an absence of individual effects.

effect measures

$P(Y^1 = 1) - P(Y^0 = 1)$ [risk difference] (~ absolute numbers)

$P(Y^1 = 1) / P(Y^0 = 1)$ [risk ratio] (~ how much more likely)

* Remark = There are two sources of random error =
SAMPLING VARIABILITY + NONDETERMINISTIC P.O.

Causation vs Association.

For individual i , $Y_i = A_i Y_i^1 + (1 - A_i) Y_i^0$.

We observe the treatment A_i and the outcome Y_i .

We can compute / estimate $P(Y=1 | A=a)$

A measure of association between A and Y .

If $P(Y=1 | A=1) = P(Y=1 | A=0)$,

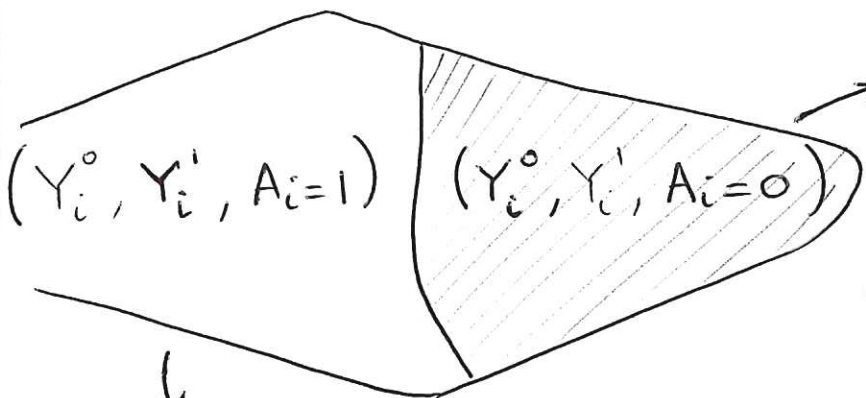
then $Y \perp A$ (independence)

Association Measures

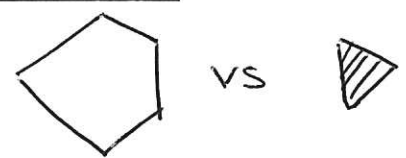
$P(Y=1 | A=1) - P(Y=1 | A=0)$

$\frac{P(Y=1 | A=1)}{P(Y=1 | A=0)}$

[Not unusual to find an association when there is no effect]



Association



$P(Y=1 | A=1)$ vs $P(Y=1 | A=0)$

These sub populations may not be representative of the whole population

Causation



$P(Y^1=1)$ vs $P(Y^0=1)$

CI = RANDOMISED EXPERIMENTS

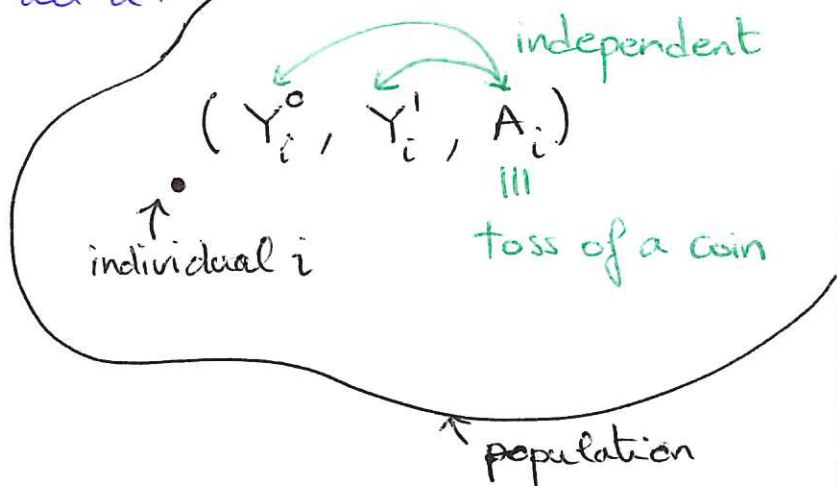
One of Y^0 or Y^1 is factual

↳ In a random experiment, which one of Y^0 or Y^1 is observed occurs just by chance.

⇒ Treatment allocation is independent of Y^a for all a .

This setup is referred to as EXCHANGEABILITY

$$Y^a \perp A \text{ for all } a$$



Consequence =

$$P(Y^0=1) = P(Y^0=1 | A=1) = P(Y^0=1 | A=0)$$

$$P(Y^1=1) = P(Y^1=1 | A=1) = P(Y^1=1 | A=0)$$

These are precisely the conditions of a Randomized Experiment.

• Remark: $Y^a \perp A$ vs $Y \perp A$

In cases where an average causal effect of treatment A on outcome Y exist, $P(Y^1=1) \neq P(Y^0=1)$

consistency

$$P(Y^1=1 | A=1) \stackrel{||}{=} P(Y=1 | A=1)$$

$$P(Y^0=1 | A=0) \stackrel{||}{=} P(Y=1 | A=0)$$

We conclude that $P(Y=1 | A=1) \neq P(Y=1 | A=0)$,
i.e. $Y \not\perp A$.

\Rightarrow Exchangeability condition is on the potential outcomes:
($Y^a \perp A$ for all a), not on the outcome variable Y .

Conditional Randomization.

Consider an additional variable $L \in \{0, 1\}$

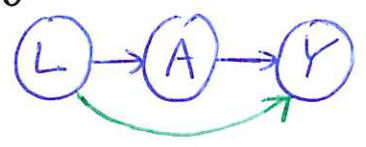
\rightarrow marginally randomized experiment
(toss one coin and allocate individuals)

\downarrow conditionally randomized experiment
(toss one coin for individuals with $L=0$, and
another coin for those with $L=1$)

\uparrow i.e. conditionally on L .

\downarrow equivalent to two separate marginally randomized experiments.

Thus $Y^a \perp A \mid L=0 \quad \forall a$
 $Y^a \perp A \mid L=1 \quad \forall a$



holds on \neq subsets / strata of the population.

In other words, $Y^a \perp A \mid L$

The average causal effect can be computed/estimated within each stratum, or across the whole population.

We discuss next two (equivalent) techniques to compute/estimate $P(Y^1=1) - P(Y^0=1)$ in a conditionally randomized experiment = STANDARDIZATION & INVERSE PROBABILITY WEIGHTING.

• STANDARDIZATION

Use the law of Total Probability:

$$\begin{aligned}
P(Y^a = 1) &= \sum_l P(Y^a = 1, L=l) \\
&= \sum_l \underbrace{P(Y^a = 1 \mid L=l)}_{\text{conditionally on } L, (Y^a \perp A \mid L)} P(L=l) \\
&= P(Y^a = 1 \mid L=l, A=a) \\
&= \sum_l P(Y^a = 1 \mid L=l, A=a) P(L=l) \\
&= \sum_l P(Y = 1 \mid L=l, A=a) P(L=l)
\end{aligned}$$

consistency ↗

(*) ↗

These can be computed / estimated from the data

• INVERSE PROBABILITY WEIGHTING

Note that $P(Y=1 \mid L=l, A=a) = \frac{P(L=l, A=a, Y=1)}{P(L=l, A=a)}$

⇒ From (*) we get

$$P(Y^a = 1) = \sum_l \left[\frac{1}{P(A=a \mid L=l)} \right] P(L=l, A=a, Y=1)$$

$\equiv w(a, l)$

$$P(Y^a = 1) = \sum_l w(a, l) P(L=l) P(A=a \mid L=l) \times P(Y=1 \mid L=l, A=a)$$

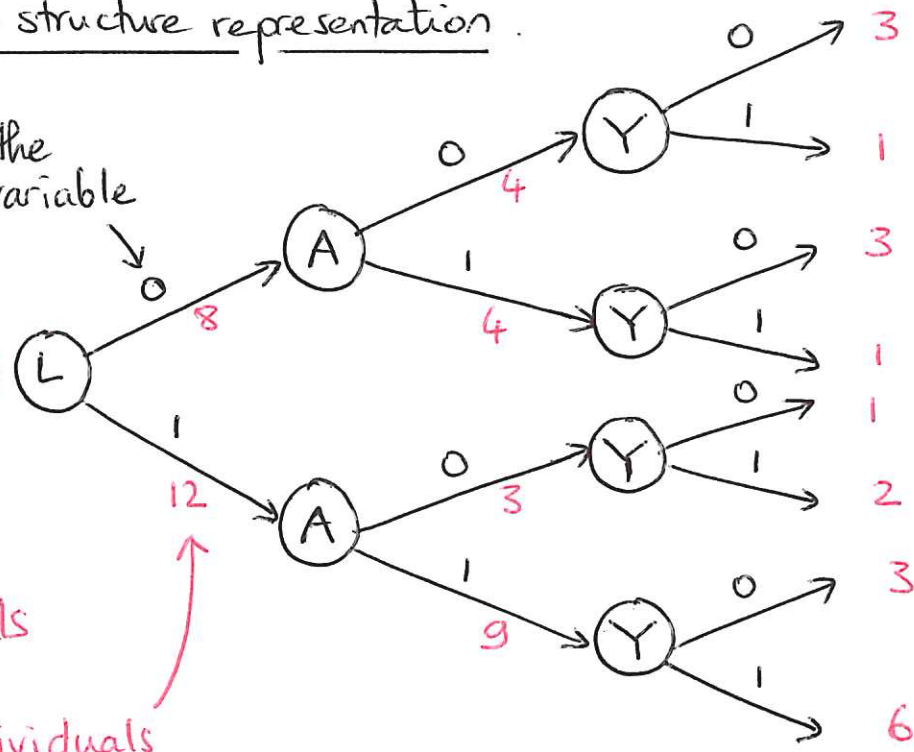
↖ We make sense of this expression on the next page, and in particular of the term $w(a, l)$.

Tree-structure representation.

value of the random variable

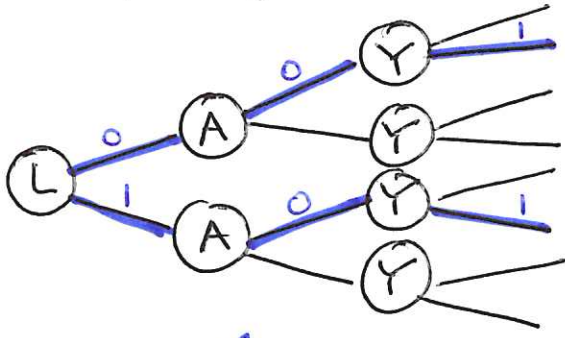
total # of individuals

of individuals with $L=1$



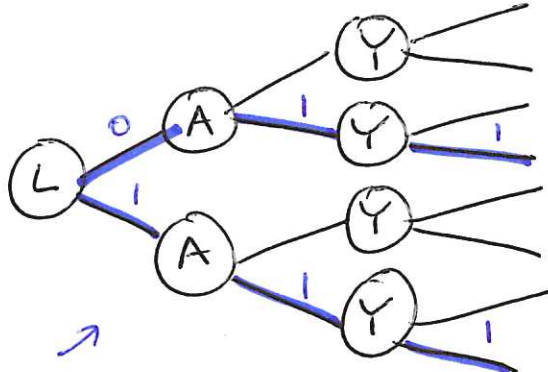
6 individuals have $(L=1, A=1, Y=1)$

To compute $P(Y^0=1)$, we need to look up branches corresponding to $A=0$ and $Y=1$



A situation where everyone (since inter-factual term) in the population remained untreated

Similarly, to compute $P(Y^1=1)$, we need to look up branches corresponding to $A=1$ and $Y=1$

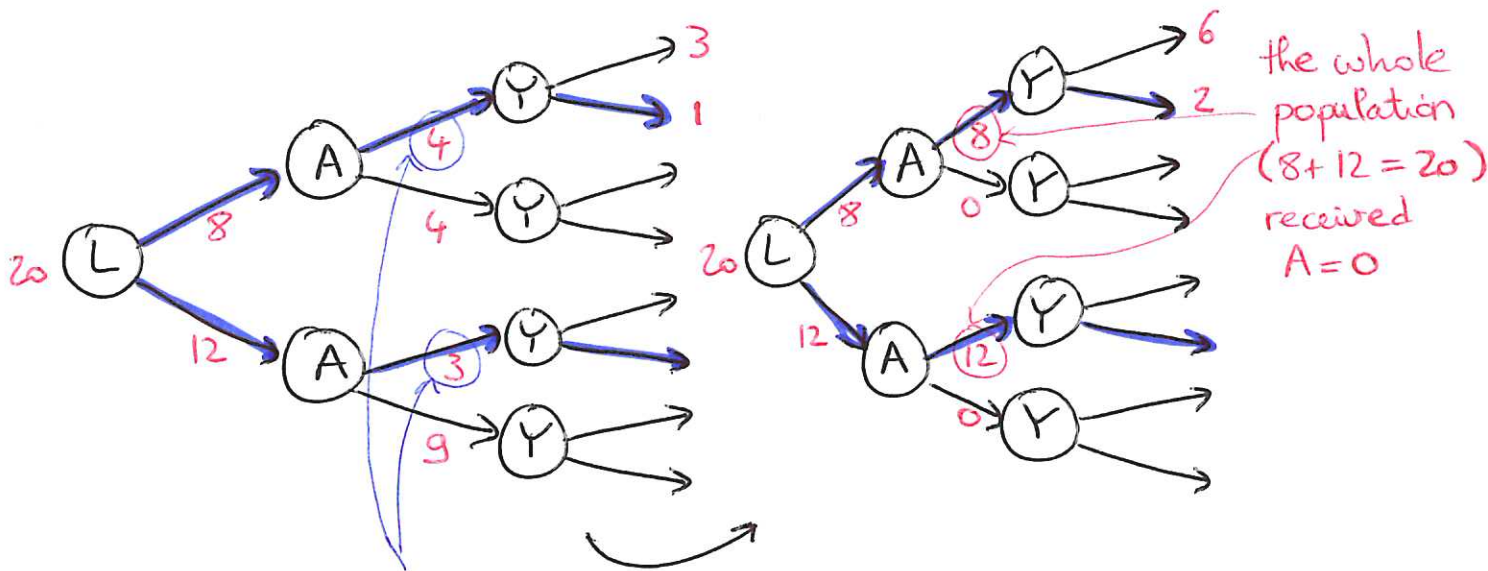


a situation where everyone in the population had been treated

Let's focus on $P(Y^0=1)$.

We need to

- transfer the 4 individuals who received $A=1 | L=0$ to non treatment
- transfer the 9 individuals who received $A=1 | L=1$ to non treatment



The number of individuals is scaled up

by $1/P(A=0 | L=l) = w(0, l)$

⇒ Individuals within each branch $L=0$ and $L=1$ are swapped or exchanged, as if they were interchangeable. But this is precisely the conditional exchangeability condition $Y^a \perp A | L \forall a$; that allow us to draw causal conclusions from observed values (A, L, Y) .

The scaling $w(0, l)$ translates all the way to the right hand side of the tree, where the number of individuals in each subbranch $(L=l, A=0, Y=1)$ is scaled up by $w(0, l)$.

$$P(Y^0=1) = \sum_l w(0, l) P(L=l, A=0, Y=1)$$

weight associated with obs $(L=l, A=0, Y=1)$

x Remark = Inverse Probability Estimator.

Consider a population of n individuals

$$\mathcal{L}_n = \{(L_1, A_1, Y_1), \dots, (L_n, A_n, Y_n)\}$$

$$\begin{aligned}
\bullet P(Y^1=1) &= \sum_l w(l, l) P(L=l, A=1, Y=1) \\
&= \sum_{l, a, y} w(a, l) a y P(L=l, A=a, Y=y) \\
&= E\{w(A, L) A Y\} \\
&\approx \frac{1}{n} \sum_{i=1}^n w(A_i, L_i) A_i Y_i
\end{aligned}$$

each observation (L_i, A_i, Y_i) is weighted by $w(A_i, L_i)$.

$$\begin{aligned}
\bullet P(Y^0=1) &= \sum_l w(0, l) P(L=l, A=0, Y=1) \\
&= \sum_{l, a, y} w(a, l) (1-a) y P(L=l, A=a, Y=y) \\
&= E\{w(A, L) (1-A) Y\} \\
&\approx \frac{1}{n} \sum_{i=1}^n w(A_i, L_i) (1-A_i) Y_i
\end{aligned}$$

$\Rightarrow P(Y^1=1) - P(Y^0=1)$ is estimated by

$$\frac{1}{n} \sum_{i=1}^n \left\{ \frac{A_i Y_i}{e(L_i)} - \frac{(1-A_i) Y_i}{1-e(L_i)} \right\}, \text{ where } e(L_i) = P(A_i=1 | L_i)$$

The Horvitz-Thompson estimator

Note that $E\{w(A, L)A\} = E\{w(A, L)(1-A)\} = 1$ (9a)

$$\begin{aligned}\Rightarrow P(Y^1=1) &= E\{w(A, L)AY\} \\ &= \frac{E\{w(A, L)AY\}}{E\{w(A, L)A\}} \\ &\approx \frac{\sum_{i=1}^n w(A_i, L_i) A_i Y_i}{\sum_{i=1}^n w(A_i, L_i) A_i} \\ &= \sum_{i|A_i=1} \left\{ \frac{w(A_i, L_i)}{\sum_{j|A_j=1} w(A_j, L_j)} \right\} Y_i \\ &= w'(A_i, L_i) \\ &\quad \text{(sum to 1)} \\ &\quad \text{(normalized weight)}\end{aligned}$$

& similarly for $P(Y^0=1)$

\Rightarrow leads to the modified HT estimator of the ATE
(Robins '98)

$$\widehat{ATE} = \sum_{i|A_i=1} w'(A_i, L_i) Y_i - \sum_{i|A_i=0} w'(A_i, L_i) Y_i$$

This estimator can be derived as the solution of a weighted LS problem.

Consider the (saturated) model $Y = \beta_0 + \beta_1 A + \varepsilon$, $A \in \{0, 1\}$
where ε is zero mean.

When weighting obs (A_i, Y_i, L_i) by $w(A_i, L_i)$,
we remove confounding by L , and $\beta_1 = ATE$.

$$(\hat{\beta}_0, \hat{\beta}_1) = \underset{\beta_0, \beta_1}{\text{argmin}} \sum_{i=1}^n \omega(A_i, L_i) (Y_i - \beta_0 - \beta_1 A_i)^2 \quad (9b)$$

$$\hat{\beta} = \begin{pmatrix} \hat{\beta}_0 \\ \hat{\beta}_1 \end{pmatrix} = (X^t W X)^{-1} X^t W Y,$$

where $W = \text{diag} \{ \omega(A_i, L_i) \}$
 $(n \times n)$

$$X = \begin{pmatrix} 1 & A_1 \\ \vdots & \vdots \\ 1 & A_n \end{pmatrix} \quad Y = \begin{pmatrix} Y_1 \\ \vdots \\ Y_n \end{pmatrix}$$

$(n \times 2)$ $(n \times 1)$

After calculations, $X^t W X = \begin{pmatrix} \sum_i \omega_i & \sum \omega_i A_i \\ \sum_i \omega_i A_i & \sum \omega_i A_i^2 \end{pmatrix}$

$\omega_i = \omega(A_i, L_i)$

$$X^t W Y = \begin{pmatrix} \sum_i \omega_i Y_i \\ \sum_i \omega_i A_i Y_i \end{pmatrix}$$

$$\Rightarrow \begin{pmatrix} \hat{\beta}_0 \\ \hat{\beta}_1 \end{pmatrix} = \begin{pmatrix} \frac{\sum \omega_i Y_i (1 - A_i)}{\sum \omega_i (1 - A_i)} \\ \frac{(\sum \omega_i)(\sum \omega_i A_i Y_i) - (\sum \omega_i A_i)(\sum \omega_i Y_i)}{(\sum \omega_i A_i)(\sum \omega_i (1 - A_i))} \end{pmatrix}$$

It follows that $\hat{\beta}_0 = \frac{\sum_{\substack{i | A_i=0 \\ j | A_j=0}} \omega(A_i, L_i) Y_i}{\sum_{j | A_j=0} \omega(A_j, L_j)}$

$$\hat{\beta}_0 + \hat{\beta}_1 = \frac{\sum_{i | A_i=1} \omega(A_i, L_i) Y_i}{\sum_{j | A_j=1} \omega(A_j, L_j)}$$

$\Rightarrow \hat{\beta}_1$ is the modified HT estimator (variance usually smaller than HT)

Normalized weights are often preferred over original weights, especially when $w(a, l)^{-1} = P(A=a | L=l)$ is close to 0 or 1, which typically increase the variance of the HT estimator. (9c)

→ Alternatively, truncate / trim weights to their 5% and 95% percentiles.

→ Another approach is to stabilize weights $w(a, l)$ by multiplying them by a quantity of the same order of magnitude as the denominator $P(A=a | L=l)$. A commonly used value is $P(A=a)$, leading to

the STABILIZED WEIGHTS $\tilde{w}(a, l) := \frac{P(A=a)}{P(A=a | L=l)}$

Note however that for a binary treatment $A \in \{0, 1\}$ and saturated model $Y = \beta_0 + \beta_1 A + \epsilon$, the weighted LS problem using weights \tilde{w} instead of w leads to exactly the same solution $(\hat{\beta}_0, \hat{\beta}_1)$ since $(\hat{\beta}_0, \hat{\beta}_1)$

$$\begin{aligned} \tilde{\beta}_0 &= \frac{\sum_{i|A_i=0} \tilde{w}_i(A_i, L_i) Y_i}{\sum_{j|A_j=0} \tilde{w}_j(A_j, L_j)} \\ &= \frac{\sum w_i(A_i, L_i) Y_i}{\sum w_i} \\ &= \hat{\beta}_0 \end{aligned}$$

Since the extra term $P(A=a)$ is common to all observations & cancels out

Stabilized weights must be used instead in more complex settings such as continuous or time-varying treatments.

x Remark: Instead of a saturated model

$Y = \beta_0 + \beta_1 A + \varepsilon$ fit using weights $w(A, L)$, we may consider equivalently the same linear model for counterfactual statements $Y^a = \beta_0 + \beta_1 a + \varepsilon$. Models involving Y^a quantities are commonly referred to as Marginal Structural Models.

x Remark: Back to standardization, Recall that

$$\mathbb{E} Y^a = \sum_l \mathbb{E}(Y \mid A=a, L=l) \mathbb{P}(L=l)$$

fit a regression/classification model to estimate the (mean of the) distribution of $Y \mid A=a, L=l$.

Once an estimate $\hat{\mathbb{E}}(Y \mid A=a, L=l)$ is available, we directly get an estimator of $\mathbb{E} Y^a$:

$$\mathbb{E} Y^a \approx \frac{1}{n} \sum_{i=1}^n \hat{\mathbb{E}}(Y \mid A=a, L=L_i)$$

Summary: • IP weighting requires the modeling of $A \mid L$
• standardization requires — " — $Y \mid A, L$.

↑ Parametric IP weighting & standardization are expected to differ. They can only agree asymptotically either (i) when non-parametric or (ii) when the parametric modeling for IP weighting & standardization are both correct. Model misspecification will introduce some bias.

* Revisiting conditional exchangeability with g-estimation.

(9)

Under the assumption $Y^a \perp A \mid L$,

$$\mathbb{P}(A=1 \mid L=1) = \mathbb{P}(A=1 \mid L, Y^0) \quad \begin{array}{l} \text{logistic regression} \\ \text{model} \end{array}$$
$$= \sigma^{-1}(\alpha_0 + \alpha_1 L + \alpha_2 Y^0) \quad (\diamond)$$

α_2 must be zero.

↓
Cannot fit this model using data since Y^0 is not observed. However, we can make the following additional assumption $\mathbb{E}(Y^a - Y^0) = \beta_1 a$ (*)

$$\Leftrightarrow \mathbb{E} Y^0 = \mathbb{E} Y^a - \beta_1 a \quad \begin{array}{l} \text{consistency} \\ \text{model} \end{array}$$
$$\Leftrightarrow Y^0 = Y - \beta_1 A + \varepsilon \quad ; \quad \mathbb{E} \varepsilon = 0.$$

The expression $Y^0 \approx Y - \beta_1 A$ can be used to reconstruct the counterfactual Y^0 under assumption (*).

→ Perform a grid search over β such that the reconstructed Y^0 yield a zero α_2 coefficient in (\diamond) .

* Remark: Condition (*) can be relaxed / weakened / generalized to include effect modifiers, aka Structural Nested Model (SNM)

Ex: $\mathbb{E}(Y^a - Y^0 \mid L) = \beta_1 a + \beta_2 a L$, which can in turn be used together with (\diamond) to recover the causal effect within each stratum $\{L=l\}$.

CI = OBSERVATIONAL STUDIES

In conditional randomized experiments, we showed (p.6) that

$$P(Y^a=1) = \sum_l w(a,l) P(L=l) P(A=a | L=l) \times P(Y=1 | L=l, A=a).$$

To derive this expression, we made use of the following.

a) CONSISTENCY If $A=a$, then $Y=Y^a$

$$P(Y^a=1 | L=l, A=a) = P(Y=1 | L=l, A=a)$$

b) CONDITIONAL EXCHANGEABILITY $Y^a \perp A | L$

$$P(Y^a=1 | L=l) = P(Y^a=1 | L=l, A=a)$$

c) POSITIVITY $P(A=a | L=l) > 0 \quad \forall l \mid P(L=l) > 0$

$$\text{since } w(a,l) = 1 / P(A=a | L=l)$$

b) and c) only are referred to as IGNORABILITY (weak & strong)

Need full exchangeability
 $(Y^0, Y^1) \perp A | L$

Rosenbaum & Rubin (1983)

Referred to as IDENTIFIABILITY conditions.
They hold by design.

Since "Causal Effects are identifiable".

⇒ To analyse observational data as if treatment had been randomly assigned conditionally on measured covariates L ,

We need these three properties to hold.

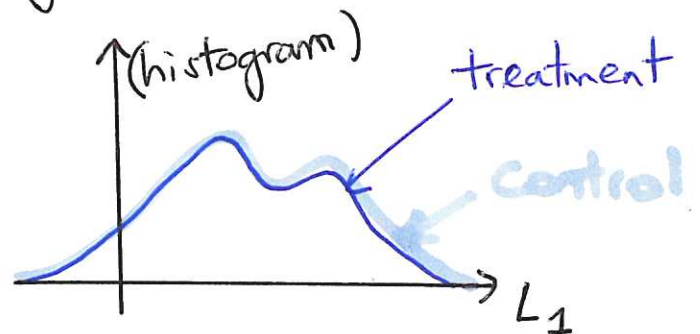
11

↳ these become assumptions and it is our responsibility to assess to which extent they hold / break down.

↳ conditional exchangeability.

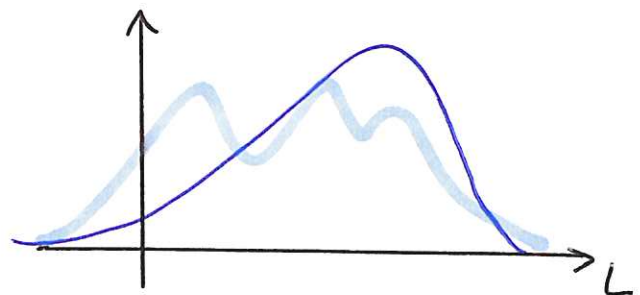
In a randomized experiment, $Y^a \perp A$, and all (measurable) features you can think of are balanced in the treatment and control group. This follows from the fact that the two subpopulations are representative samples of the whole population.

Practically, denoting L_1, L_2, \dots these covariates, a histogram of L_i conditionally on $A=0$ and $A=1$ should "look alike".



The same ideas apply in a conditional randomized experiment.

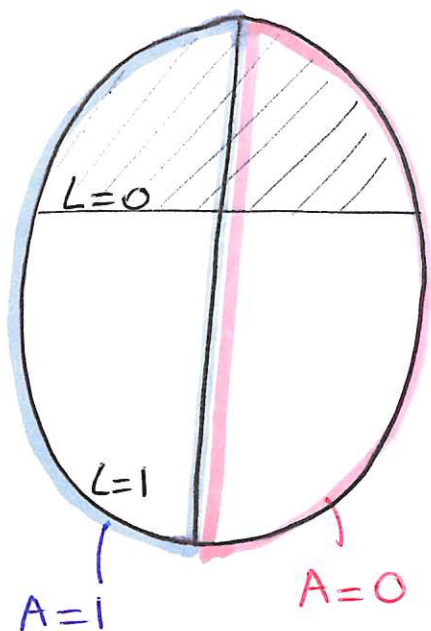
1) It is expected by design that L is unbalanced between the control and the treated groups, since $Y^a \perp A \mid L$.



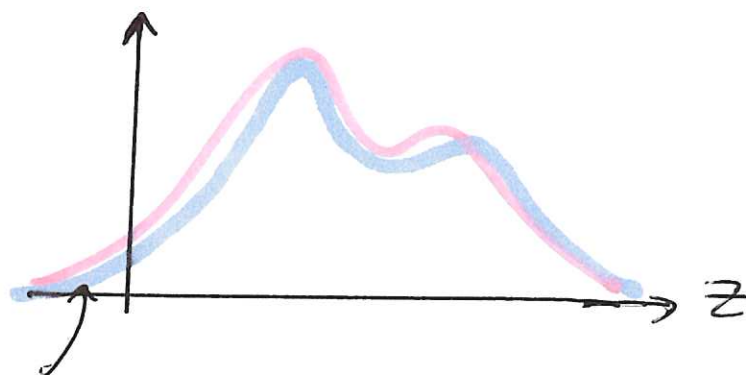
2) Within each stratum $L=l$, all other (measurable)

features Z must be balanced.

12



In each subpopulation $L=l$, plot the distribution (histogram) of another measured feature Z , conditionally on $A=a$



Conditionally on $L=l$, all other features Z must be balanced for $Y^a \perp A \mid L$ to hold.

If the histograms do not superpose, this is an indication that $Y^a \not\perp A \mid L$.

However...

- ↘ you may not be able to verify this in practice since you may not have collected the required data
- ↘ the set of features is unlimited & cannot be exhaustive in practice.
- ↘ some features are not measurable.

Conclusions:

- 1) Expert knowledge is needed.
- 2) We can easily gather information that conditional exchangeability breaks down, but we can never be fully sure it holds \Rightarrow risky task.

↳ positivity

In a randomized experiment, to compute the average causal effect $P(Y^1=1) - P(Y^0=1)$ using $P(Y=1 | A=1) - P(Y=1 | A=0)$, we need to collect data in both the treatment & the control groups.

The same idea applies in conditional randomized experiments. In order to estimate the average causal effect, the learning sample must be rich enough / contain enough information about the two groups. The mathematical way of translating this is that

$$\forall l \text{ s.t. } P(L=l) > 0 \Rightarrow P(A=a | L=l) > 0$$

↳ the positivity assumption / condition ↷

Let $\mathcal{A} = \{0, 1\}$ (binary treatment)

$\mathcal{L} = \{0, 1\}$ (binary feature)

$$\mathcal{L}(a) := \left\{ l \in \mathcal{L} \mid P(L=l) > 0 \right. \\ \left. \& P(A=a | L=l) > 0 \right\}$$

When positivity holds, $\mathcal{L}(0) = \mathcal{L}(1) = \mathcal{L}$

When positivity does not hold, $\mathcal{L}(0) \neq \mathcal{L}(1)$.

Recall the **Horvitz-Thompson estimator** (p. 9)

$$\frac{1}{n} \sum_{i=1}^n \left\{ \frac{A_i Y_i}{e(L_i)} - \frac{(1-A_i) Y_i}{1-e(L_i)} \right\}$$



$$e(l) = P(A=1 | L=l)$$

$$\approx \mathbb{E} \{ w(A, L) A Y \} - \mathbb{E} \{ w(A, L) (1-A) Y \}$$



$$w(a, l) = 1 / P(A=a | L=l)$$

We compute next $\mathbb{E}\{w(A, L)AY\}$ and $\mathbb{E}\{w(A, L)(1-A)Y\}$ (14)
 when positivity holds & does not hold.

$$\mathbb{E}\{w(A, L)AY\} = \sum_{y \in \{0, 1\}} \sum_{a \in A} \sum_{l \in \mathcal{L}(a)} w(a, l) ay \times P(L=l, A=a, Y=y)$$

the support of (L, A, Y)

terms vanish for $y=0$ and $a=0$

same as p. 6

$$= \sum_{l \in \mathcal{L}(1)} w(1, l) P(L=l, A=1, Y=1)$$

$$= \sum_{l \in \mathcal{L}(1)} P(Y=1 | L=l, A=1) P(L=l)$$

$$= \sum_{l \in \mathcal{L}(1)} P(Y'=1, L=l)$$

$$= \sum_{l \in \mathcal{L}} P(Y'=1, L=l, L \in \mathcal{L}(1))$$

$$= \sum_{l \in \mathcal{L}} P(Y'=1, L=l | L \in \mathcal{L}(1)) P(L \in \mathcal{L}(1))$$

LTP

$$= P(L \in \mathcal{L}(1)) P(Y'=1 | L \in \mathcal{L}(1)),$$

and similarly for $\mathbb{E}\{w(A, L)(1-A)Y\}$.

⇒ In general, the Horvitz-Thompson estimator is an appr. of

$$P(L \in \mathcal{L}(1)) P(Y'=1 | L \in \mathcal{L}(1)) - P(L \in \mathcal{L}(0)) P(Y^0=1 | L \in \mathcal{L}(0))$$

≠ $P(Y'=1) - P(Y^0=1)$ unless $\mathcal{L}(0) = \mathcal{L}(1) = \mathcal{L}$,
 i.e. unless the positivity assumption holds.

↳ consistency

If $A = a$, then $Y_i = Y_i^a$.



First, note that the treatment $A = a$ must be well defined (e.g. in experimental studies, the experimenter has control of the treatment, which must be administered the same way to all participants)

↳ usually ok in medical studies

↳ less ok in biological & social sciences.

↳ In obs studies, there may be different versions of the way a treatment is administered.

Let $r =$ version of a treatment.

We collect/analyze observations associated with all patients that received \neq versions of treatment a .

Then the consistency property holds provided the following assumption holds

If $A = a$, then $Y_i = Y_i^{a,r} \quad \forall r$



Potential Outcome associated with version r of treatment a

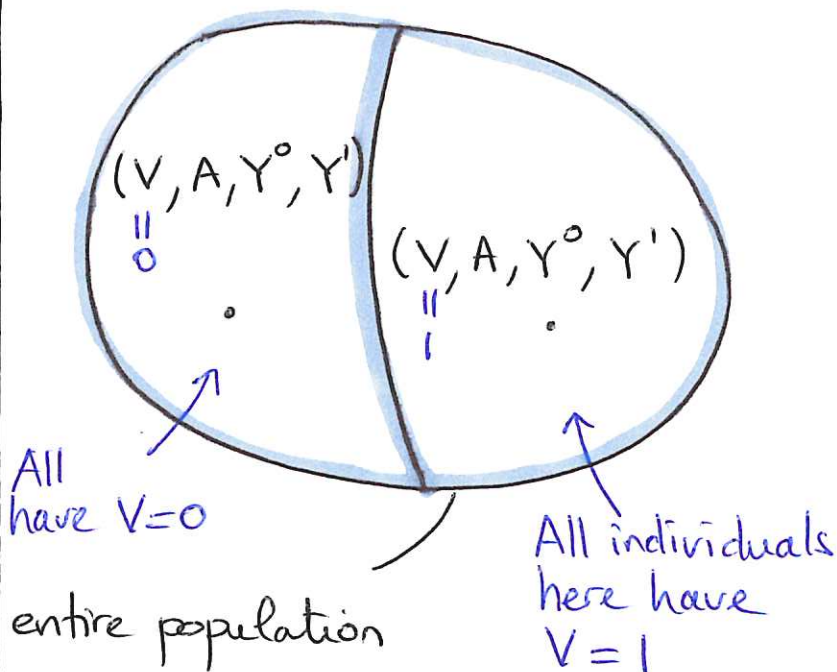
In other words, $Y_i^{a,r} = Y_i^{a,r'} \quad \forall r, r'$.

i.e. it is assumed that all versions of the treatment a produce the same causal effect.

CI = EFFECT MODIFICATION

We are interested here in the average causal effect about subsets of a population.

Let $V \in \{0, 1\}$ = binary indicator splitting an entire population into two subpopulations (e.g. Male vs Female)



We may be interested in comparing

$$P(Y^1=1) - P(Y^0=1)$$

with

$$P(Y^1=1 \mid V=1)$$

$$- P(Y^0=1 \mid V=1)$$

and

$$P(Y^1=1 \mid V=0)$$

$$- P(Y^0=1 \mid V=0)$$

Def: V is a MODIFIER of the effect of A on Y when the average causal effect of A on Y varies across levels of V .

↑ There can be a multiplicative but not additive effect modification by V .

To identify effect modification, a stratified analysis is a natural way to go:

• In an RCT:

- 1) Partition the population into strata $V = v$
- 2) For each strata, estimate the average causal effect $P(Y^1=1)_{V=v} - P(Y^0=1)_{V=v}$ using $P(Y=1 | A=1)_{V=v} - P(Y=1 | A=0)_{V=v}$.

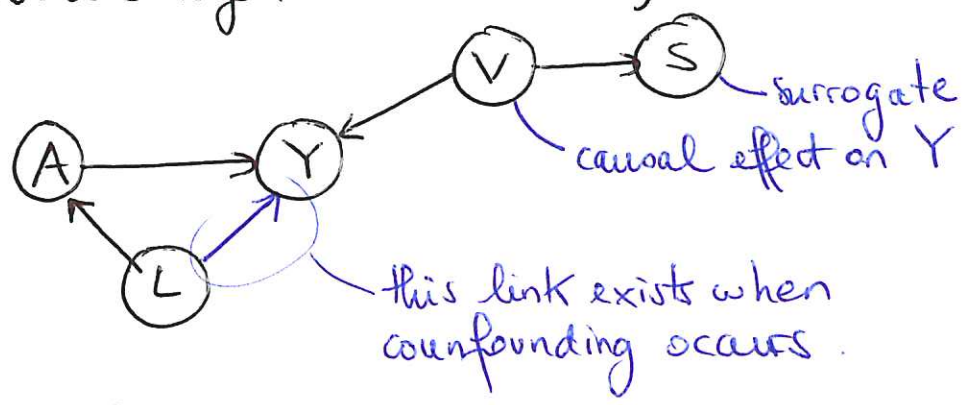
• In a conditional randomized experiment,

- 1) — " —
- 2) Use standardization / IP weighting under the exchangeability assumption $Y^a \perp A | L$

Remarks: • No claims are being made about the causal mechanisms. The variable V may be a SURROGATE effect modifier, while the CAUSAL effect modifier is unknown.

Eg: $V =$ carrying a lighter
 $Y =$ lung cancer

• graphical representation (difference between L , V , and a surrogate variable S)



• If V modifies the effect of treatment A on the outcome Y , then the average causal effect will differ btw populations with \neq prevalence of V .
 (\rightarrow question of TRANSPORTABILITY of causal

• Different forms of adjustments.

- ↳ stratification
- ↳ standardization / IP weighting
- ↳ matching

We review/introduce /compare these three techniques next.

Suppose we observe $d_n = \{(L_i, A_i, Y_i)\}_{i=1}^n$, where each $(L_i, A_i, Y_i) \stackrel{d}{\sim} (L, A, Y)$; $L \in \{0, 1\}$
 $A \in \{0, 1\}$
 $Y \in \{0, 1\}$
 ↑ i -th unit

a) STRATIFICATION splits the data into two subpopulations, according to the strata $L=0$ and $L=1$.

Within each stratum $L=l$, compute

$$P(Y=1 | A=a, L=l)$$

consistency $\hookrightarrow = P(Y^a=1 | \mathbf{A}=a, L=l)$

under the assumption of conditional exchangeability $\hookrightarrow = P(Y^a=1 | L=l)$
 $Y^a \perp A | L=l$

Under these two conditions,

$$P(Y^1=1 | L=l) - P(Y^0=1 | L=l)$$

$$= P(Y=1 | A=1, L=l) - P(Y=1 | A=0, L=l)$$

↑ One average causal estimate per stratum is computed.
 (generalizes well to non-binary L)

b) STANDARDIZATION / IP WEIGHTING

Adjust for L using the techniques detailed on page 6.

Under identifiability conditions, these provide an estimate of the ATE $\mathbb{P}(Y^1=1) - \mathbb{P}(Y^0=1)$

Only one estimate; for the entire population (\neq stratification; one per stratum).

x Remarks = (i) The risk difference and risk ratio for the entire population lies in the convex hull of the causal estimates computed within each stratum.

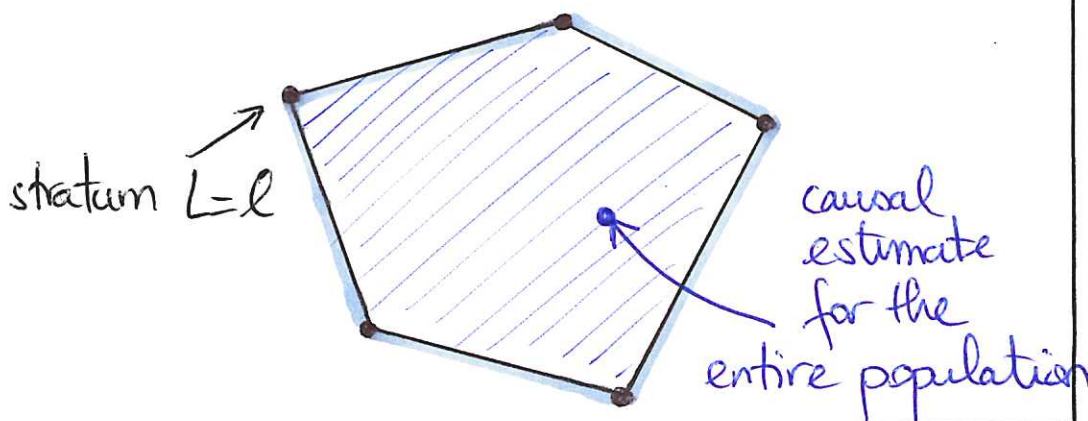
$$\text{proof} = \mathbb{P}(Y^1=1) - \mathbb{P}(Y^0=1) = \sum_l \{ \mathbb{P}(Y^1=1 | L=l) - \mathbb{P}(Y^0=1 | L=l) \} w_l$$

where $w_l = \mathbb{P}(L=l) \geq 0$; $\sum_l w_l = 1$.

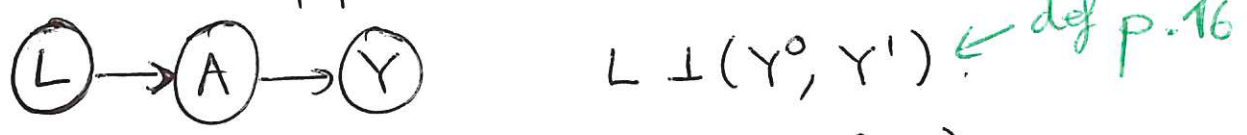
$$\mathbb{P}(Y^1=1) / \mathbb{P}(Y^0=1) = \sum_l \left\{ \frac{\mathbb{P}(Y^1=1 | L=l)}{\mathbb{P}(Y^0=1 | L=l)} \right\} w_l$$

where

$$w_l = \frac{\mathbb{P}(Y^0=1 | L=l) \mathbb{P}(L=l)}{\mathbb{P}(Y^0=1)} \geq 0 ; \sum_l w_l = 1$$

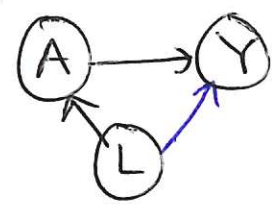


(ii). If L is not an effect modifier, then all average causal effects (within each stratum + at the entire population level) are equal.



Then $P(Y^a = 1 | L = l) = P(Y^a = 1)$; and the result follows from the expressions on page 19.

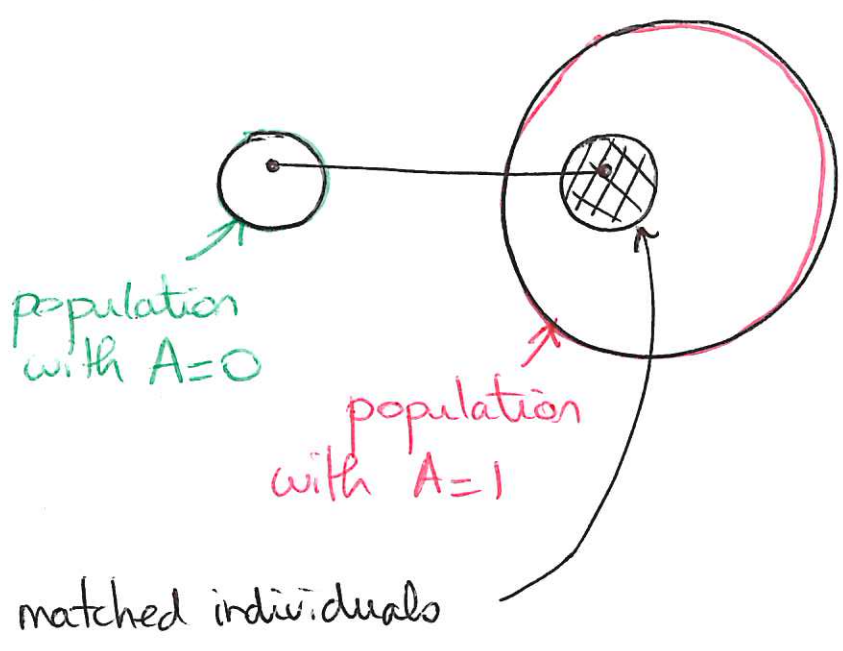
• If L is an effect modifier, then the average causal effects will be different





c) MATCHING.

The idea is to rebalance the distribution of the variable L within the two groups (to make them comparable)

• Suppose that we take the untreated ($A=0$) group as the reference group. Then, for each individual in $A=0$ with $L=l$, pick one / match it with an individual in $A=1$ with the same $L=l$.



By doing so, the two groups  and  are comparable (same distribution of L)

Provided conditional exchangeability holds in the original population $Y^a \perp A | L = l$, then

unconditional exchangeability holds for \bigcirc & \bigotimes . (21)

Since the reference population is $A=0$, a direct comparison of the two populations \bigcirc & \bigotimes yields an estimate of

$$\begin{aligned} & \mathbb{P}(Y^1=1 | A=0) - \mathbb{P}(Y^0=1 | A=0) \\ &= \mathbb{E}(Y^1 - Y^0 | A=0) \\ &= \text{ATNT} \end{aligned}$$

Average Treatment on the ^{Non-}Treated

The matching can be made by considering the group $A=1$ as the reference group; in which case we estimate the $\text{ATT} = \mathbb{E}(Y^1 - Y^0 | A=1)$.

x Remark: ATT and ATNT require the computation of quantities like $\mathbb{P}(Y^a=1 | A=a')$ $a \neq a'$. It can be done under the usual consistency & conditional exch. assumptions.

$$\mathbb{P}(Y^a=1 | A=a') = \sum_{\ell} \mathbb{P}(Y^a=1 | A=a', L=\ell) \mathbb{P}(L=\ell | A=a')$$

consistency $\hookrightarrow = \sum_{\ell} \mathbb{P}(Y^a=1 | A=a, L=\ell) \mathbb{P}(L=\ell | A=a')$

$Y^a \perp A | L$ $\hookrightarrow = \sum_{\ell} \mathbb{P}(Y=1 | A=a, L=\ell) \mathbb{P}(L=\ell | A=a')$

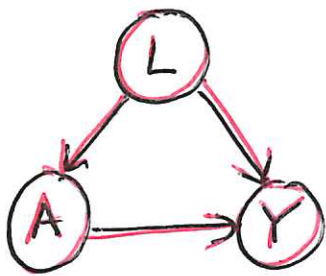
can be computed from data. \uparrow



CI = GRAPHICAL REPRESENTATION OF CAUSAL EFFECTS

A Structural Causal Model (SCM) $\mathcal{C} = (X, \mathbb{P}_{\mathcal{E}})$ consists of a collection of d assignments (X_1, \dots, X_d) $X_j = f_j(Pa_j, \mathcal{E}_j)$, where $Pa_j \subseteq \{X_1, \dots, X_d\}$ are the parents of X_j & a joint distribution $\mathbb{P}_{\mathcal{E}}$ over the noise variables $\mathcal{E} = (\mathcal{E}_1, \dots, \mathcal{E}_d)$. (unobserved)

Ex: Conditional RCT



$$\begin{aligned}
 X_1 = L & ; Pa_1 = \emptyset \\
 X_2 = A & ; Pa_2 = \{L\} \\
 X_3 = Y & ; Pa_3 = \{A, L\}
 \end{aligned}$$

↑ A Directed Acyclic Graph (DAG) with 3 nodes & 3 edges

[Together with structural equations $\begin{cases} X_1 = f_1(Pa_1, \mathcal{E}_1) \\ X_2 = f_2(Pa_2, \mathcal{E}_2) \\ X_3 = f_3(Pa_3, \mathcal{E}_3) \end{cases}$,

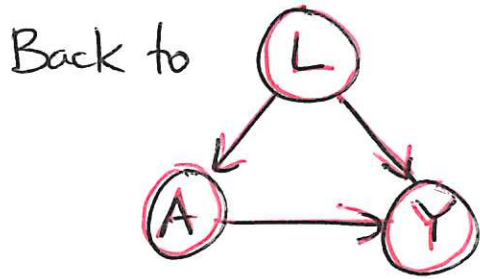
it constitutes a causal DAG]

↑
These allow us to link the DAG to obs & attach a probabilistic structure to it.

The structural equations show that time flows from \rightarrow to i.e. presence of a direct causal effect.

↘ No cycle = a variable cannot cause itself; either (23) directly or through another variable

↘ Conditionally on its direct causes (Pa), a variable is independent of any other variable for which it is not the cause.



$$\begin{aligned}
 L &= X_1 = f_1(\varepsilon_1) &= \varepsilon_1 \\
 A &= X_2 = f_2(L, \varepsilon_2) &= \alpha A + \varepsilon_2 \\
 Y &= X_3 = f_3(A, L, \varepsilon_3) &= \beta A + \gamma L + \varepsilon_3
 \end{aligned}$$

Assuming linear structural equations.

The joint distribution P_{ε} over $\varepsilon = (\varepsilon_1, \varepsilon_2, \varepsilon_3)$ induces a joint distribution $P = P_X$ over the vector $X = (X_1, X_2, X_3)$.

DAGs & Potential Outcomes.

This DAG representation does not naturally include the underlying counterfactual variables on the graph: the link between Y and (Y^0, Y') remain hidden.

Richardson & Robins (2013) found a graphical representation that explicitly include counterfactuals in the graph (Single World Intervention Graph SWIG); see next chapter.

However, the DAG representation can include potential outcomes

provided we intervene : any variable X_j on the DAG can be intervened on & set to a particular value x_j^a .
(assumption)

$$\begin{cases} L = X_1 = f_1(\epsilon_1) \\ A = X_2 = f_2(L, \epsilon_2) \\ Y = X_3 = f_3(A, L, \epsilon_3) \end{cases}$$

intervention \rightarrow
$$\begin{cases} L = X_1 = f_1(\epsilon_1) \\ A = X_2 = a \\ Y = X_3 = f_3(a, L, \epsilon_3) \end{cases}$$

↑ ↑
superscript indicating we intervened & set A to the value a; irrespectively of the structural graph.

The intervention modifies the joint distribution over (X_1, X_2, X_3)

Pearl's do calculus
"do(A=a)"

Intervention distribution

$$\mathbb{P}_{\text{do}(A=a)}$$

(randomized experiments)

↓
Joint distribution $P = P_X$
(observational studies)

x Remark = some invariants

$$\begin{aligned} \mathbb{P}(L=l) &= \mathbb{P}_{\text{do}(A=a)}(L=l) \\ \mathbb{P}(Y=1 | A=a, L=l) &= \mathbb{P}_{\text{do}(A=a)}(Y=1 | A=a, L=l) \end{aligned}$$

These invariants can be used to recover average causal effects:

$$\mathbb{P}_{\text{do}(A=a)}(Y^a=1) = \sum_l \mathbb{P}_{\text{do}(A=a)}(Y^a=1, A=a, L^a=l)$$

(what we want)

$$\begin{aligned}
 &= \sum_l P_{\text{do}(A=a)}(Y^a=1 \mid A=a, L=l) \\
 &\quad \parallel \\
 &P(Y=1 \mid A=a, L=l) \quad \times P_{\text{do}(A=a)}(A \neq a, L=l) \\
 &\quad \parallel \quad \text{(invariants)} \\
 &\quad \parallel \\
 &P(L=l) \\
 &= \sum_l P(Y=1 \mid A=a, L=l) P(L=l) \\
 &\quad \uparrow \\
 &\text{the } \underline{\text{STANDARDIZATION}} \text{ expression.}
 \end{aligned}$$

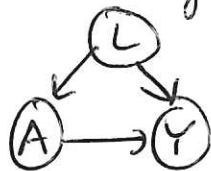
⇒ The techniques seen before have a graphical ~~equivalence~~; ~~equivalence~~ ^{cal} equivalence. Assumptions of positivity, consistency & conditional exchangeability must also be interpreted using DAGs.

↳ positivity = arrow from L to A is not deterministic

↳ consistency = arrow from A to Y corresponds to a possibly hypothetical but unambiguous intervention.

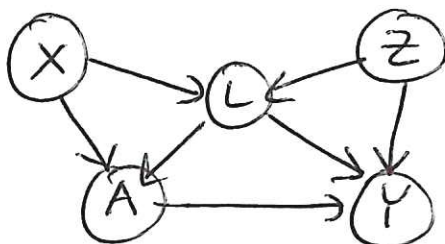
↳ cond. exch. = see next chapter. ■

• We recovered the average causal effect in a simple example where



(equivalent to standardization, see above)

For example, how should we generalize the expression above when the SCM looks something like

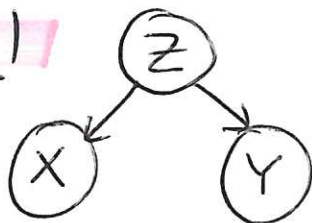


To answer this question, we must study indep & conditional indep statements

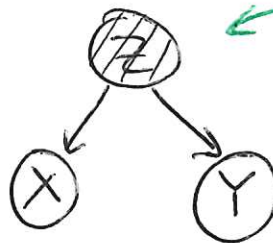
in DAGs for some simple structures. We identify 3 of them.

Conditional Independence & Graphs.

Example 1



vs



← an observed node is darkened

$$p(x, y, z) = p(x|z)p(y|z)p(z)$$

— Are X and Y independent? —

$$p(x, y) = \sum_z p(x|z)p(y|z)p(z) \neq p(x)p(y) \text{ in general}$$

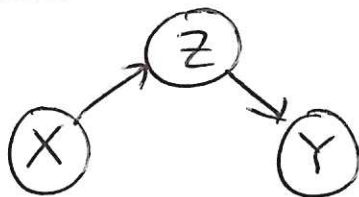
$$p(x, y|z) = p(x|z)p(y|z)$$

$$X \not\perp Y | \emptyset$$

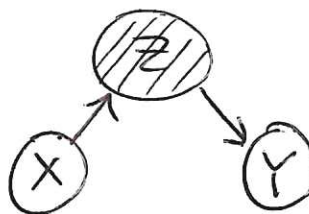
$$X \perp Y | z$$

We say that z blocks the path between X and Y.

Example 2



vs



$$p(x, y, z) = p(y|z)p(z|x)p(x)$$

— Are X and Y independent? —

$$p(x, y) = \sum_z p(y|z)p(z|x)p(x) = p(x)p(y|x) \neq p(x)p(y) \text{ in general}$$

$$p(x, y|z) = \frac{p(x, y, z)}{p(z)} = p(x|z)p(y|z)$$

$$X \perp Y | \emptyset$$

$$X \perp Y | z$$

z blocks the path between X & Y.

x Example 3



$$p(x, y, z) = p(x)p(y)p(z|x, y)$$

— Are X and Y independent? —

$$p(x, y) = \sum_z p(x)p(y)p(z|x, y) \quad p(x, y|z) = \frac{p(x)p(y)p(z|x, y)}{p(z)}$$

$$= p(x)p(y) \quad \checkmark \quad \neq p(x|z)p(y|z)$$

$X \perp Y | \emptyset$

$X \not\perp Y | z$

We say that z is a collider

z blocks the path between X and Y
 \Leftrightarrow z is tail-to-tail or head-to-tail wrt that path

This leads us to the following definition

• Def d-separation

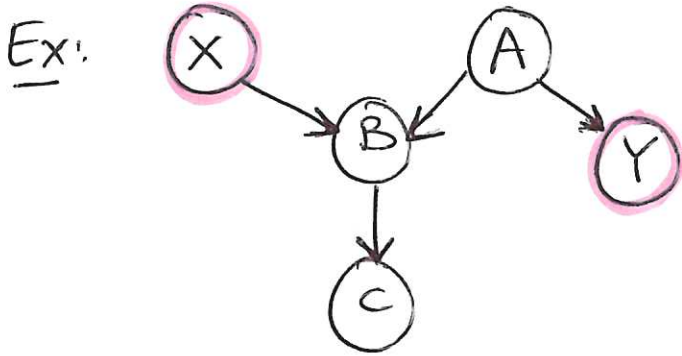
A path is blocked by a set of nodes Z iff

(i) The path contains a chain of nodes $A \rightarrow B \rightarrow C$ or a fork $A \leftarrow B \rightarrow C$ such that the middle node B is in Z (B is conditioned on)

(ii) The path contains a collider $A \rightarrow B \leftarrow C$ such that the collision node B is not in Z; and no descendent of B is in Z.

↓ If Z blocks every path between two nodes X & Y,

then X & Y are d-separated ~~by Z~~ by Z , (28)
 and thus are independent conditionally on Z .



• Is $X \perp Y | A$?

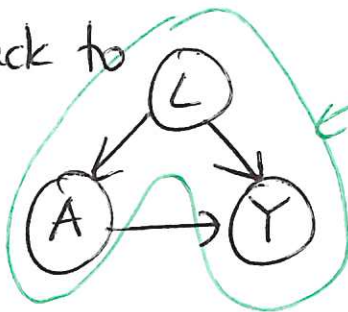
Conditioning on A blocks the path $X \rightarrow B \leftarrow A \rightarrow Y$.
 In addition, B (a collider) & its descendent C remain unobserved $\Rightarrow X$ and Y are d-separated by $\{A\}$.

• Is $X \perp Y | C$?

B is a collider and C is a descendent of $B \Rightarrow X$ and Y are **not** d-separated by $\{C\}$. ▣

Conditional independence for these three building blocks will allow us to identify which set of variables in a DAG we need to condition on (\equiv adjust for) in order to compute / identify the average causal effect of A on Y .

Back to



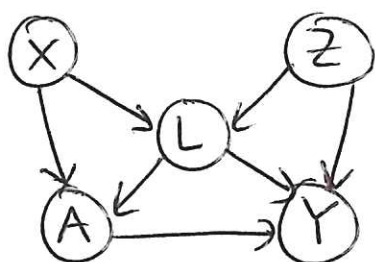
We are interested in the path

$$A \leftarrow L \rightarrow Y$$

We know that in this path, $A \not\perp Y$;
 while $A \perp Y | L$

\Rightarrow Provided L is observed, A and Y are independent along the path $A \leftarrow L \rightarrow Y$ (& this is the reason why in the standardization formula, we sum (& observe) over $L=l$)

We now have the tools to answer the question at the bottom of page 25 :



There are 4 paths to consider:

$$A \leftarrow X \rightarrow L \leftarrow Z \rightarrow Y \quad \text{--- (i)}$$

$$A \leftarrow L \rightarrow Y \quad \text{--- (ii)}$$

$$A \leftarrow L \leftarrow Z \rightarrow Y \quad \text{--- (iii)}$$

$$A \leftarrow X \rightarrow L \rightarrow Y \quad \text{--- (iv)}$$

We see from (i) that L is a collider so the set {L} does not block those paths.

However, {L, X} blocks all the 4 paths

Also {L, Z} and {L, Z, X}

We get back to this more extensively in subsequent chapters, where we discuss various types of bias preventing us to compute the average causal effect.

We say that there is a systematic bias whenever

$$P(Y^1=1) - P(Y^0=1) \neq P(Y=1 | A=1) - P(Y=1 | A=0);$$

which happens under a lack of exchangeability.

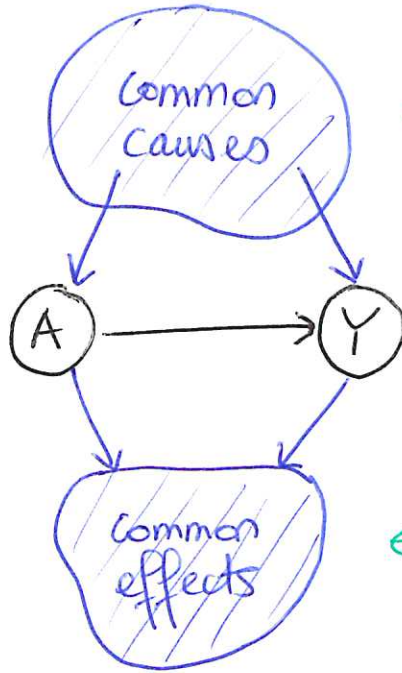
We will consider

- ↳ confounding (due to common causes)
- ↳ selection bias (cond. on common effects)
- ↳ measurement error.

CI = CONFOUNDING

• Task = estimate the causal effect from A to Y

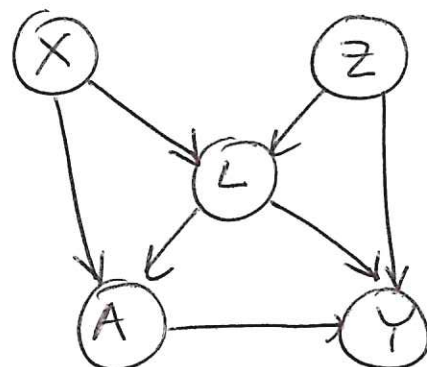
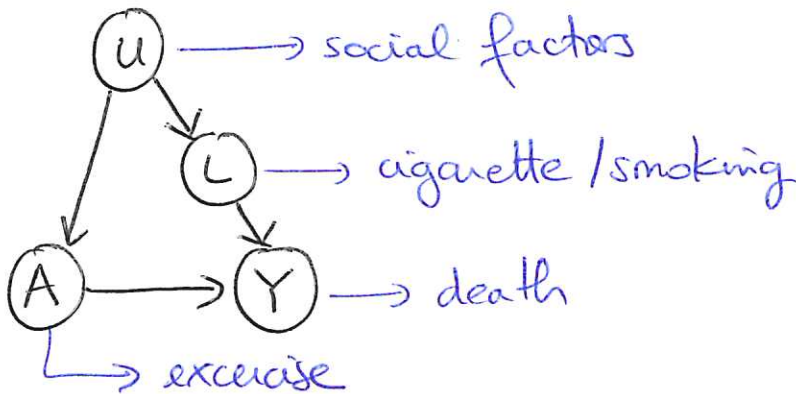
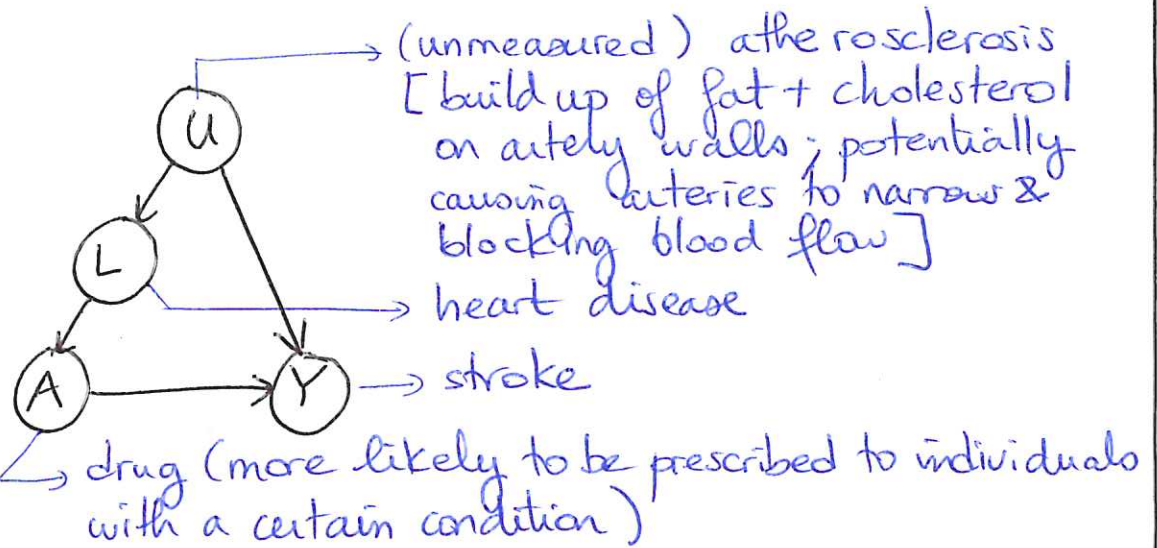
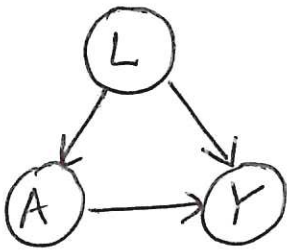
DA G =
+
Markov IP
on top of it



← In this chapter, we are interested in the structure of confounding i.e. bias due to common causes of treatment & outcome

← Assume here there are no common effects.

• Examples



Theorem (backdoor criterion)

A set of nodes L satisfy the backdoor criterion relative to $A \rightarrow Y$ if

- (i) no node in L is a descendent of A
(no common causes)
- (ii) L blocks every path between A and Y that contain an arrow into A .

If a set of nodes L satisfy the backdoor criteria for $A \rightarrow Y$, the causal effect from A to Y is $P(Y^1=1) - P(Y^0=1)$, where

$$P(Y^a=1) = \sum_l P(Y=y | L=l, A=a) P(L=l) \quad (*)$$

\Leftrightarrow all backdoor paths between A and Y are d-separated by the set of nodes L

x Sketch of proof / intuition

Expression (*) requires conditional exchangeability $Y^a \perp A | L$; see page 6. It remains to prove that

L satisfies the backdoor criterion for $A \rightarrow Y$	\Rightarrow	conditional exch. holds $Y^a \perp A L$
--	---------------	---

\Downarrow follows from the following lemma =

Lemma = The following statements are equivalent:

(a) Markov factorization property

$$p(x) = \prod_{j=1}^d p(x_j | pa_j)$$

(b) global Markov property

X_i & X_j are d-separated by $L \Rightarrow X_i \perp X_j \mid L$

(c) local Markov property

Each variable is independent of its non-descendants given its parents.

- \Rightarrow Intervening on the graph & removing all incoming arrows into A does not change the distribution of (A, Y) given L , whether $A=a$ or do($A=a$).
- \Rightarrow Causal effects can be identified & computed from observational conditional probabilities via the backdoor formula (*)

• Remark The equivalence

L satisfies the backdoor criterion for $A \rightarrow Y \iff$ cond. exchangeability holds $Y^a \perp A \mid L$

holds under the additional assumption of faithfulness

Faithfulness

X_i and X_j are d-sep by $L \iff X_i \perp X_j \mid L$

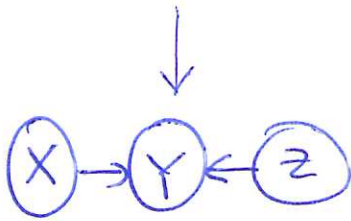
\leftarrow Reverse direction of the global Markov property (top of page)

A nice assumption since under Markov + Faithfulness, we have a 1-1 relation between statements about graphs & probability distributions (\hookrightarrow enables causal discovery)

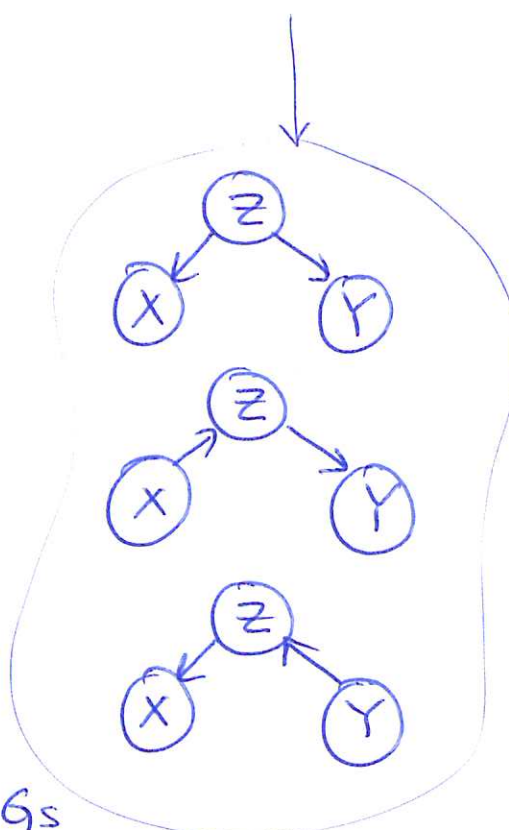
EX = Assume that \mathbb{P} is Markov & Faithful w.r.t. its DAG.

Recover all DAGs compatible with the following complete list of conditional independence statements.

(i) $X \perp Z$
(X, Y, Z)

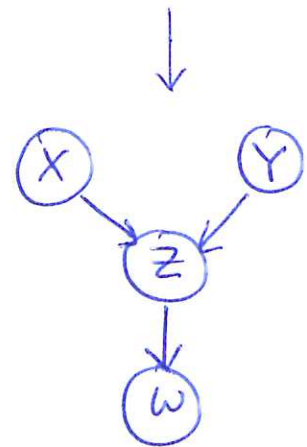


(ii) $X \perp Y | Z$



(iii) $X \perp Y$

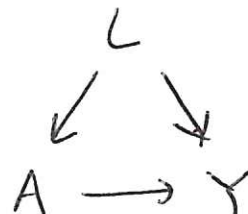
$X \perp W | Z$
 $Y \perp W | Z$
 $X \perp W | Z, Y$
 $Y \perp W | Z, X$



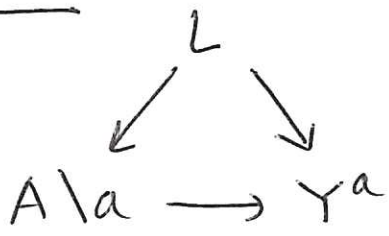
Enclose the same set of d-separated DAGs are called Markov Equivalent.

Remark = The backdoor criterion has a nice visual translation using Single-World Intervention Graphs (SWIG) that include counterfactuals as nodes
 ↳ The graph representing counterfactual worlds is created by a single intervention.

EX: original DAG
 (no counterfactual notation)



SWIG



— counterfactual notation since the outcome variable Y is a descendent of A

we intervene on A and do $(A=a)$
 two separate nodes

- a inherits all nodes coming out of A in the original DAG.
- A inherits all nodes incoming into A in the original DAG.



In this representation, the backdoor path $A \leftarrow L \rightarrow Y^a$ is d-separated by $L \Rightarrow Y^a \perp A \mid L$. ▣

To summarize, there are two classes of methods to consider in the presence of confounders L :

- (a) Methods requiring conditional exchangeability
 - g-methods (IP weighting / stand / backdoor)
 - stratification based methods / matching
- (b) Methods not requiring cond. exchangeability
 - front-door criterion
 - instrumental variables
 - difference-in-difference (DiD)



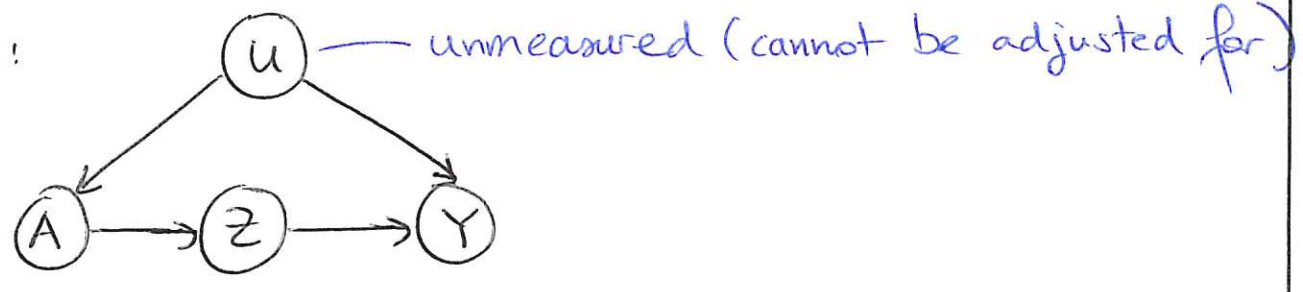
We introduce these next.

• Front-door-criterion.

↳ an alternative to the back-door criterion.

↳ typically when an additional variable Z fully mediates the effect of A on Y & that share no unmeasured causes with either A or Y.

Ex:



(•) Since $A \leftarrow U \rightarrow Y \leftarrow Z$ is blocked by Y,
 collider

A and Z are independent

$$\Rightarrow P_{do(A=a)}(Z=z) = P(Z=z | A=a)$$

(••) Since $Z \leftarrow A \leftarrow U \rightarrow Y$ is blocked cond. on A, the backdoor criterion relative to $Z \rightarrow Y$ applies:

$$P_{do(Z=z)}(Y=y) = \sum_a P(Y=y | Z=z, A=a) P(A=a)$$

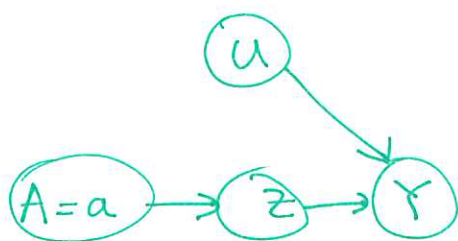
The law of total probability

$$P_{do(A=a)}(Y=y) = \sum_z P_{do(A=a)}(Y=y, Z=z)$$

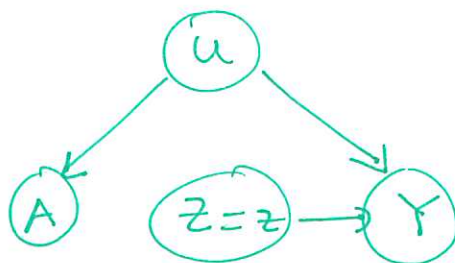
$$= \sum_z P_{do(A=a)}(Y=y | Z=z) P_{do(A=a)}(Z=z)$$

$$= P_{do(Z=z)}(Y=y)$$

since $Y = f(Z, U, N_Y)$ and ... ↘



vs



$$\begin{cases} u = f(N_u) \\ A = a \\ Z = f(A, N_Z) \\ Y = f(u, Z, N_Y) \end{cases}$$

int. prob $\mathbb{P}_{\text{do}(A=a)}(\cdot)$

$$\begin{cases} u = f(N_u) \\ A = f(u, N_A) \\ Z = z \\ Y = f(u, Z, N_Y) \end{cases}$$

int. prob $\mathbb{P}_{\text{do}(Z=z)}(\cdot)$

Expression of u and N_Y unchanged
+

Whether we do $(Z=z)$ or observe

$Z = f(a, N_Z) = z$, this has no influence on Y .

$$\Rightarrow \mathbb{P}_{\text{do}(A=a)}(Y=y | Z=z) = \mathbb{P}_{\text{do}(Z=z)}(Y=y)$$

⇓

$$\mathbb{P}_{\text{do}(A=a)}(Y=y) = \sum_z \underbrace{\mathbb{P}_{\text{do}(Z=z)}(Y=y)}_{(\cdot\cdot)} \underbrace{\mathbb{P}_{\text{do}(A=a)}(Z=z)}_{(\cdot)}$$

Front-door formula:

$$\mathbb{P}(Y^a = y) = \mathbb{P}_{\text{do}(A=a)}(Y=y)$$

$$= \sum_z \mathbb{P}(Z=z | A=a) \sum_{a'} \mathbb{P}(Y=y | Z=z, A=a') \times \mathbb{P}(A=a')$$

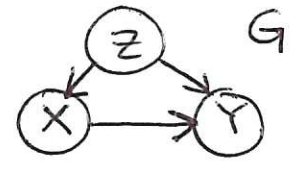
The backdoor criterion is not necessary to identify causal effects: if the frontdoor criterion is satisfied, we also have identifiability. However, the frontdoor criterion is itself a sufficient criterion. This naturally leads us to the following question: ~~can~~ can we identify causal estimates when the DAG satisfies neither the backdoor nor the frontdoor criterion?

↳ In other words, is there a general procedure to express intervention probabilities $P_{do}(X=x)(Y=y)$ using only observational proba $P(Y=y | X=x, \dots)$?

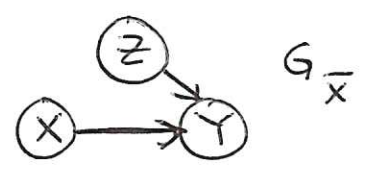
Pearl's do-calculus gives the answer to this question. The do-calculus provides tools to identify causal effects using only the structure of the causal graph.

Notation:

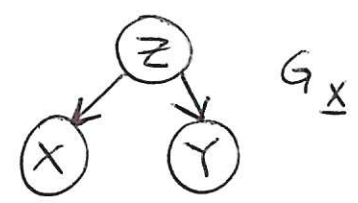
G = causal graph



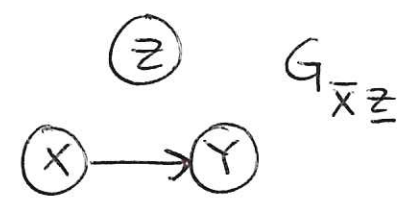
$G_{\bar{X}}$ = graph that we get if we take G and remove all incoming arrows into X



$G_{\underline{X}}$ = graph that we get if we take G and remove all outgoing edges from X



$G_{\bar{X}\underline{Z}}$ = graph that we get if we take G , and remove all incoming edges into X and remove all outgoing edges from Z



Thm: Rules of do-calculus

Simplified Rule

General Rule (Pearl)

• Rule 1 = adding / deleting an observation

$$p(y|z, w) = p(y|w)$$

if $Y \perp_G Z | W$

$$P_{do(x)}(y|z, w) = P_{do(x)}(y|w)$$

if $Y \perp_{G_x} Z | X, W$

this rule is generalized to intervention probabilities.

We need to consider G_x instead of G since we are interested in $P_{do(x)}(\cdot) \Rightarrow$ edges incoming to X must be removed. Also, $\{X=x\}$ is observed & thus conditioned on.

• Rule 2 = exchanging observation \leftrightarrow action

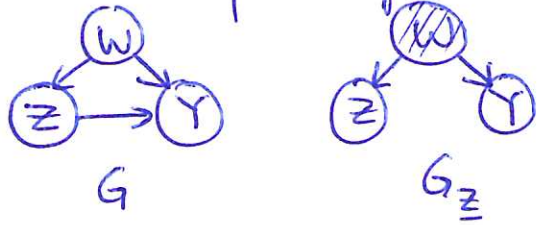
$$P_{do(z)}(y|w) = p(y|z, w)$$

if $Y \perp_{G_z} Z | W$

$$P_{do(x)}(y|w) = P_{do(z)}(y|z, w)$$

if $Y \perp_{G_{xz}} Z | X, W$

In other words, W blocks all backdoor paths from Z to Y



generalize by adding an extra "do" operator.

• Rule 3 = adding / deleting an action

$$P_{do(z)}(y) = p(y)$$

if $Y \perp_{G_z} Z$

(i) $P_{do(z)}(y|w) = p(y|w)$

if $Y \perp_{G_{z(w)}} Z | W$

where $Z(W)$ denotes the set of nodes of Z that are not ancestors of any node of W

$$(ii) P_{\text{do}(x)}^{\text{do}(z)}(y|w) = P_{\text{do}(z)}(y|w)$$

if $Y \perp_G Z \mid X, W$
 $\bar{X} \bar{Z}(w)$

Thm Shpitser & Pearl (2006)

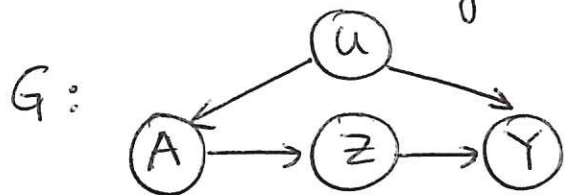
The do-calculus is complete.

↖ In other words, any identifiable quantity can be deduced from these 3 rules of do-calculus:
 The 3 rules are sufficient to derive all identifiable causal effects.
 (& their proof is constructive → identification can be obtained in polynomial time)

* Remark: All the above is about non-parametric identification
 The do-calculus tells us if we can identify causal effects using only the structure of the causal graph.
 Under additional parametric assumptions, one can recover more causal effects.

↘ cf instrumental variables p.65

Ex: Derivation of the frontdoor adjustment using the 3 rules of the do-calculus



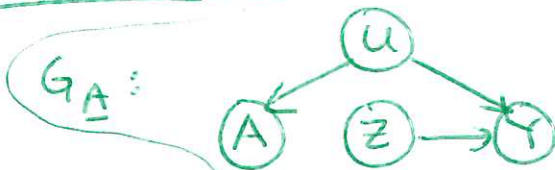
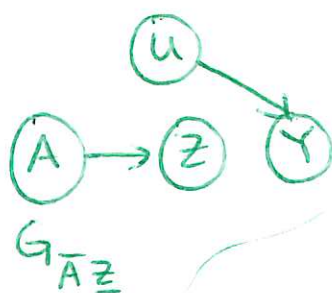
$$P(y | do(a)) = \sum_z P_{do(a)}(y, z) \quad (\text{Law of Total Proba.})$$

$$= \sum_z P_{do(a)}(y|z) P_{do(a)}(z)$$

$$= \sum_z P_{do(a)}(y|z) p(z|a)$$

Rule 2
 $Z \perp A$
 $G_{\bar{A}}$

Rule 2
 $Y \perp Z | A$
 $G_{\bar{A}Z}$



$$= \sum_z P_{do(a)}(y) p(z|a)$$

$$= \sum_z P_{do(z)}(y) p(z|a)$$

Rule 3
 $Y \perp A | Z$
 $G_{\bar{Z}A}$



$$P_{do(z)}(y) = \sum_a P_{do(z)}(y|a) P_{do(z)}(a)$$

$$= \sum_a p(y|a, z) p(a)$$

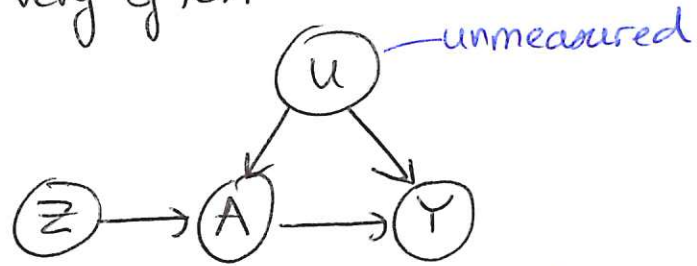
Rule 3
 $A \perp Z$
 $G_{\bar{Z}}$

Rule 2
 $Y \perp Z | A$
 $G_{\bar{Z}}$



• Instrumental variables

↳ used in economics very often
↳ typical situation

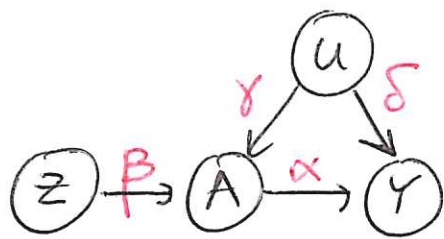


— Z is known as an instrumental variable —

To identify the causal effect of A on Y, we typically make further assumption on the structural equations.

Ex: Linear structural equations (continuous case)

$$\begin{cases} A = \beta Z + \gamma U + N_A \\ Y = \alpha A + \delta U + N_Y \end{cases}$$



↓ Since

$$Y = \alpha \left(\beta Z \right) + \left(\alpha \gamma + \delta \right) U + \alpha N_A + N_Y$$

new variable

noise

linear model

we can make use of the following procedure

- (a) Regress A on Z & get an unbiased estimate $\hat{\beta}$ of β .
- (b) Compute the fitted values $\hat{\beta} Z$
- (c) Regress Y on the fitted values $\hat{\beta} Z$ & get an unbiased estimate $\hat{\alpha}$ of α ; the causal effect of A on Y.

• difference-in-difference (DiD).

Recall: under the unconfoundedness assumption $(Y^0, Y^1) \perp A$, we can identify the $ATE = P(Y^1=1) - P(Y^0=1)$

$$\begin{aligned} (Y^0, Y^1) \perp A &\stackrel{\text{consistency}}{\hookrightarrow} = P(Y^1=1 | A=1) - P(Y^0=1 | A=0) \\ &\hookrightarrow = P(Y=1 | A=1) - P(Y=0 | A=0) \end{aligned}$$

To recover the ATT, we only need to assume that $Y^0 \perp A$ since

$$\begin{aligned} ATT &= P(Y^1=1 | A=1) - P(Y^0=1 | A=1) \\ &= P(Y^1=1 | A=1) - P(Y^0=1 | A=0) \\ &= P(Y=1 | A=1) - P(Y=1 | A=0) \end{aligned}$$

$ATE \Leftrightarrow (Y^0, Y^1) \perp A$
$ATT \Leftrightarrow Y^0 \perp A$

↑
When treatment and outcome are confounded by L , we need extra steps to recover the ATE & ATT. In addition, when the confounder is unmeasured, we cannot use the backdoor criterion. "Adjustment" can be made if a front-door criterion is applicable, or in the presence of an instrumental variable. Alternatively, we may recover the ATT (& potentially the ATE) provided pre-treatment measurements are available, and under similar / related

independence assumptions on the treatment variable and the potential outcomes.

Let τ = time indicator $\in \{ \overset{=0}{\text{pre-test}}, \overset{=1}{\text{test}} \}$

$Y^a(\tau)$ = potential outcome under treatment a at time τ .

We are interested in estimating the ATT

$$\text{ATT} = \mathbb{E}(Y^1(1) - Y^0(1) \mid A=1)$$

We will see that DiD is commonly used to identify the ATT; but additional assumptions the ATE can be identified as well.

We make the following assumptions.

(i) Consistency: $\forall \tau, A=a \Rightarrow Y(\tau) = Y^a(\tau)$.

(ii) Parallel trend: (the "essence" of DiD)

$$Y^0(1) - Y^0(0) \perp A$$

Compare with the traditional $Y^0 \perp A$ assumpt
* unconfoundedness is about the difference *

$$\text{csq: } \mathbb{E}(Y^0(1) - Y^0(0) \mid A=1) = \mathbb{E}(Y^0(1) - Y^0(0) \mid A=0)$$

(iii) no-pre-treatment effect.

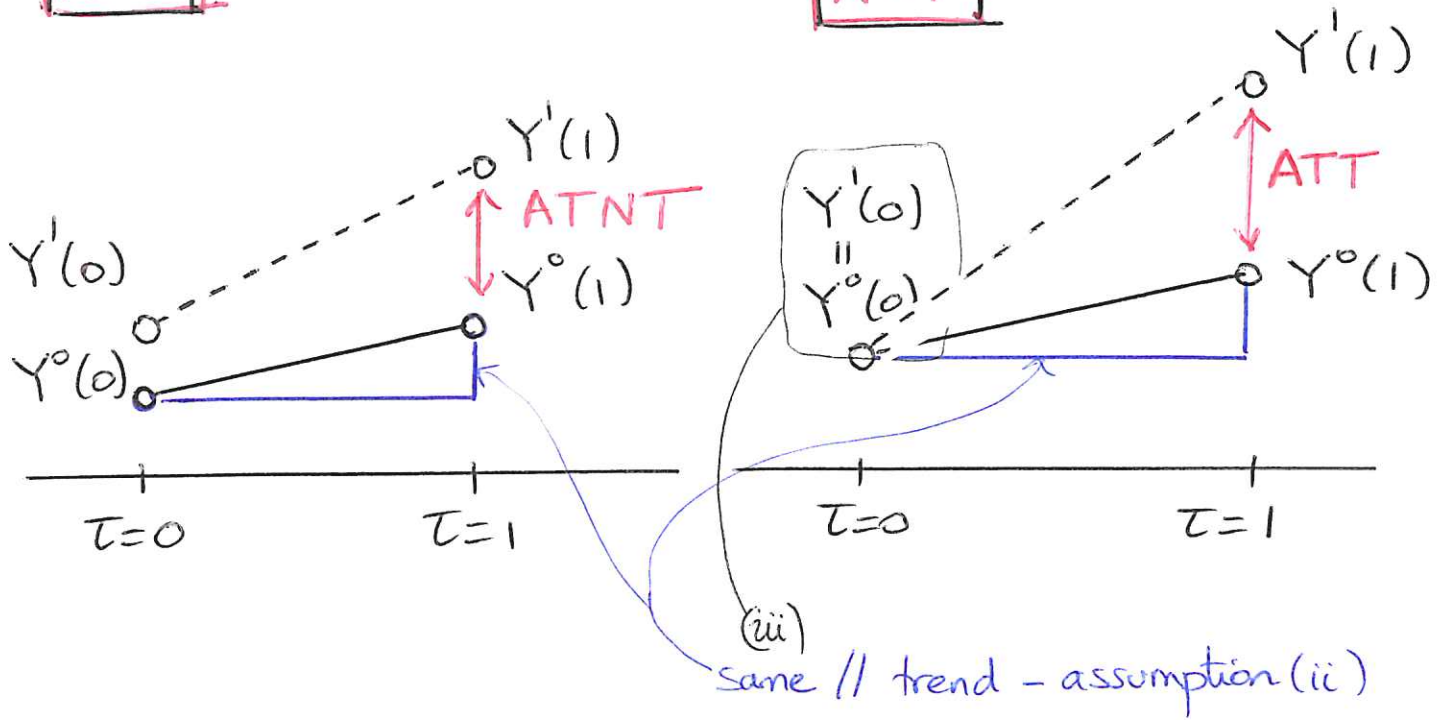
$$\mathbb{E}(Y^1(0) - Y^0(0) \mid A=1) = 0$$

usually ok, unless patients anticipate the treatment.

• Visually

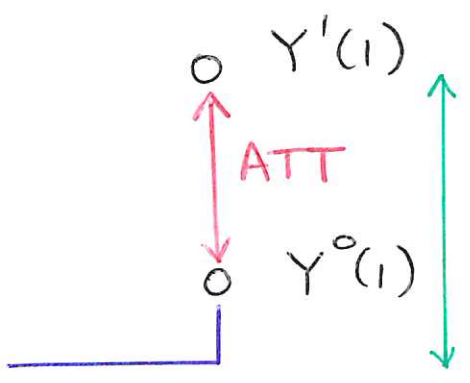
A=0

A=1



Not all quantities are observed: we only observe
 ↳ $Y^0(0)$ and $Y^0(1)$ in the group $A=0$
 ↳ $Y^1(0)$ and $Y^1(1)$ in the group $A=1$.
 (solid lines)

⇒ We immediately visualize that the parallel trend assumption allow us to identify the ATT. In addition,



$$ATT = E(Y(1) - Y(0) | A=1) - E(Y(1) - Y(0) | A=0)$$

• Formal proof = $ATT = E(Y^1(1) - Y^0(1) | A=1)$

$$= \mathbb{E}(Y'(1) | A=1) - \mathbb{E}(Y^{\circ}(1) | A=1)$$

(41)

$$(ii) \quad = \mathbb{E}(Y^{\circ}(0) | A=1) + \mathbb{E}(Y^{\circ}(1) | A=0) - \mathbb{E}(Y^{\circ}(0) | A=0)$$

$$= \left\{ \mathbb{E}(Y'(1) | A=1) - \mathbb{E}(Y^{\circ}(0) | A=1) \right\} - \left\{ \mathbb{E}(Y^{\circ}(1) | A=0) - \mathbb{E}(Y^{\circ}(0) | A=0) \right\}$$

$$\downarrow = \mathbb{E}(Y'(0) | A=1) \quad (iii)$$

$$= \mathbb{E}(Y'(1) - Y'(0) | A=1) - \mathbb{E}(Y^{\circ}(1) - Y^{\circ}(0) | A=0)$$

$$= \mathbb{E}(Y(1) - Y(0) | A=1) - \mathbb{E}(Y(1) - Y(0) | A=0)$$

x Remark : The ATNT = $\mathbb{E}(Y'(1) - Y^{\circ}(1) | A=0)$ can be identified replacing assumption (ii) and (iii) with

(ii)' $\underbrace{Y'(1) - Y'(0)}_{\text{on } Y' \text{ instead of } Y^{\circ}} \perp A$ (parallel trend)

(iii)' $\mathbb{E}(Y'(0) - Y^{\circ}(0) | A=0) = 0$ (no pre-tr. effect)

↓ Under (i), (ii)', (iii)', we get that

$$ATNT = \mathbb{E}(Y(1) - Y(0) | A=1) - \mathbb{E}(Y(1) - Y(0) | A=0)$$

same expression as before

Finally, under (i), (ii), (iii), (ii)', (iii)',

(42)

$$ATE = E(Y(1) - Y(0) | A=1)$$

$$- E(Y(1) - Y(0) | A=0)$$

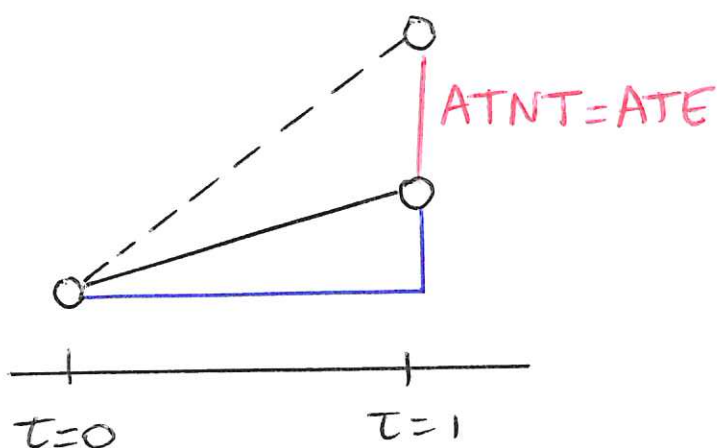
since

$$ATE = ATT \times P(A=1) + ATNT \times P(A=0)$$

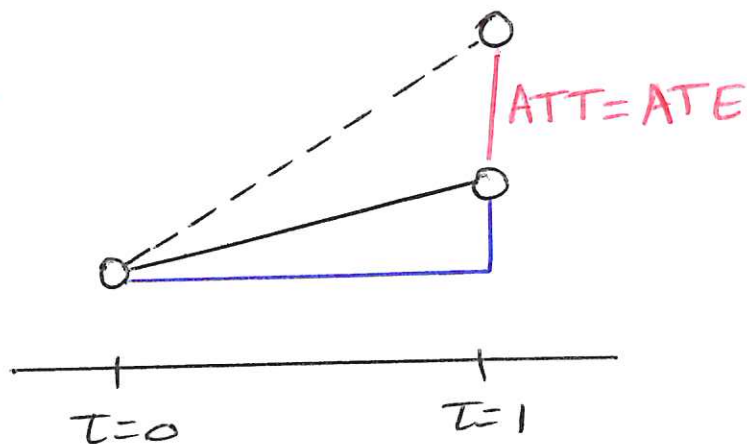
↑
same expression

- Visually, under (i), (ii), (iii), (ii)' and (iii)', all quantities are symmetrical

$A=0$

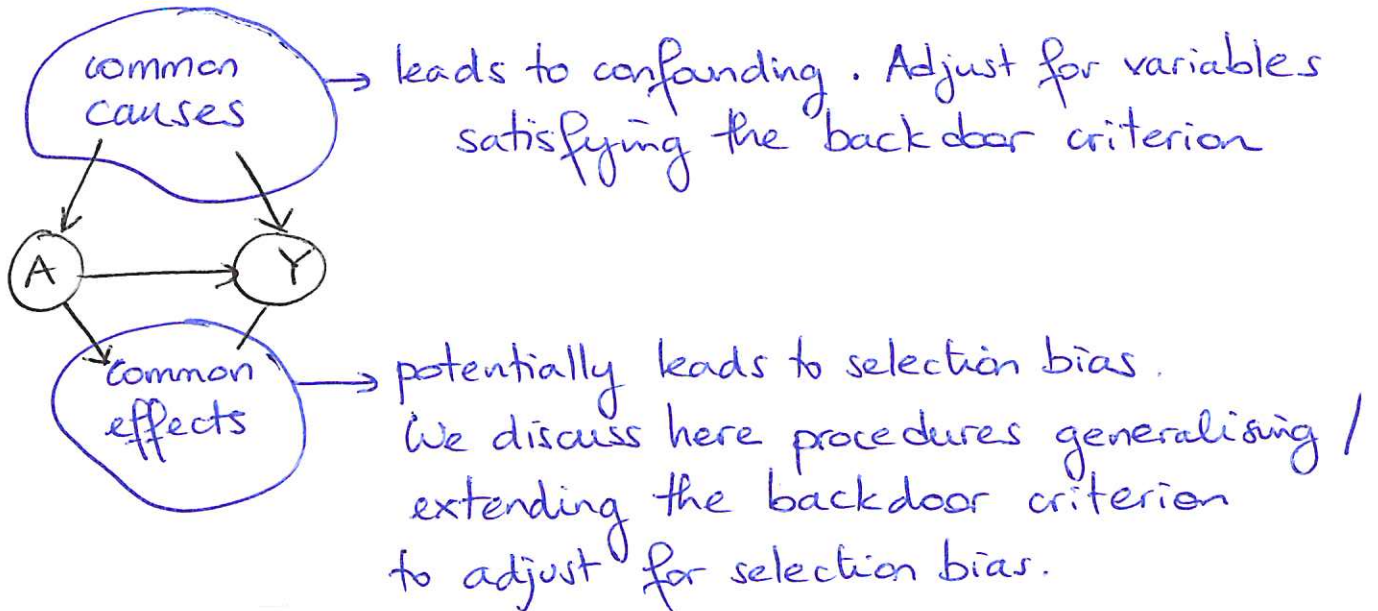


$A=1$

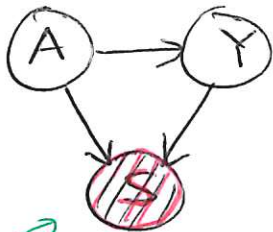


CI = SELECTION BIAS

When measuring the causal effect of A on Y, selection bias arises when conditioning on a common effect of both A and Y




Ex :



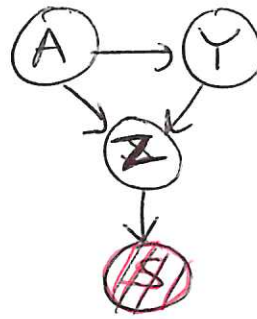
indicates that S is conditioned on

A = folic acid supplement intakes
Y = cardiac malformation during the first 2 months of pregnancy
S = death before birth.

↳ The study selects individuals who survived until birth.

Note that without explicitly including the sampling process; the study, assuming carried under an RCT, would look like $A \rightarrow Y$ \Rightarrow this DAG is oblivious to selection bias \Rightarrow highlights the importance to include a node  encoding the sampling procedure.

Alternatively, we may consider



A & Y as before

Z = death before birth

S = parental grief.

Selection bias remains when conditioning on a descendent of a common cause Z of both A and Y.

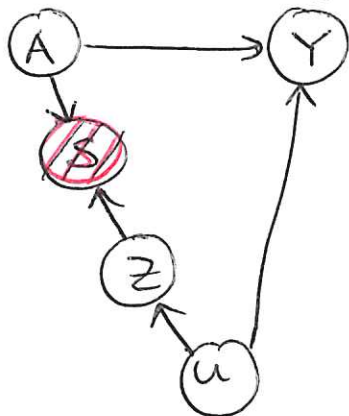
If the sampling process was completely random, S would be independent of all other variables in the analysis (no arrow coming in or out of S). When samples are collected preferentially, the causal effect must be recovered from

$P(\underline{X} = \underline{x} \mid S = 1)$ instead of $P(\underline{X} = \underline{x})$

↳ encodes selected individuals

↳ $\underline{X} = (X_1, \dots, X_d)$ = variables appearing in the DAG; including A and Y.

Ex: HIV study



A = antiretroviral drug intake

Y = 3-year death

U = 1 (high level of immunosuppression)

Z = presence of symptoms (fever, weight loss)
CD4 counts

viral load in plasma

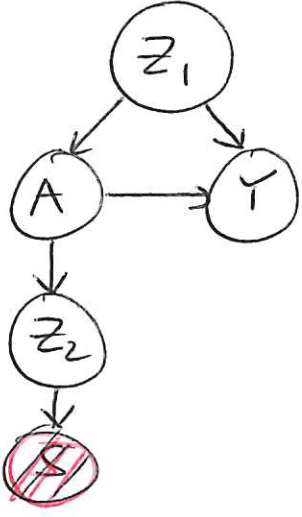
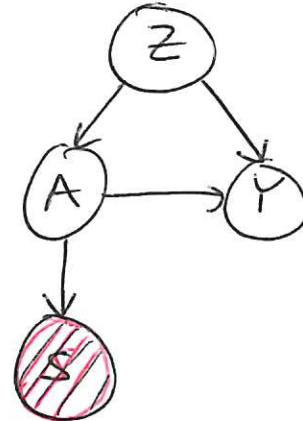
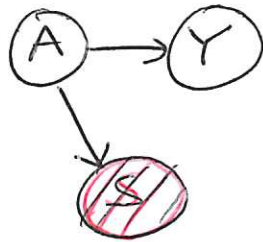
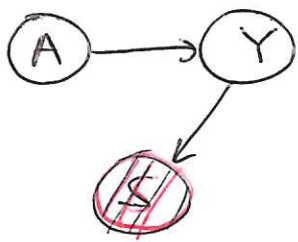
S = 1 (individual remain in the study)

Individuals with $U=1$ are more likely to be censored ($S=0$) due to the severity of the disease. Individuals

receiving the treatment are more likely to drop out due to the presence of heavy side effects.

The causal effect of A on Y suffers from the sampling process, restricted to individuals not dropping out of the study.

• Other examples: do we have selection bias in the following cases?



and if yes, is there a procedure to correct / adjust for it?

x Remark = Selection bias occurs in observational studies as well as in RCT.

Randomization protects against confounding, but not against selection bias.

no selection bias = sampling process at random
 no confounding = treatment allocation at random

Consequences of selection bias are worse than with confounding = we are not even sure to recover conditional probabilities $P(Y=y | A=a)$ from the dataset collected under selection bias $P(\underline{X}=\underline{x} | S=1)$; let alone intervention probabilities $P(Y^a=y) = P(Y=y | do(A=a))$

↑ This is well understood intuitively: conditional probabilities are related to the frequency of events in a population. When some individuals are preferred, their characteristics may not be representative of the whole population.

Following [Barainboim, Tian & Pearl \(2014\)](#), we consider 3 problems

(a) Selection without external data.

dataset $\sim P(\underline{X}=\underline{x} | S=1)$

Q: Under which conditions is $P(Y=y | A=a)$ recoverable?

(b) Selection with external data

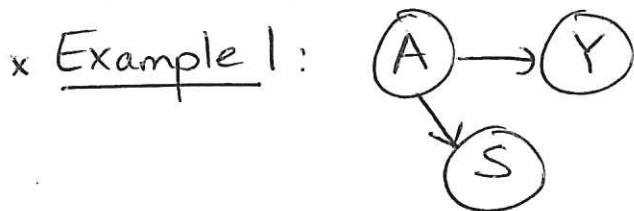
dataset $\sim P(\underline{X}=\underline{x} | S=1) +$ unbiased samples $P(\underline{T}=\underline{t})$
for $\underline{T} \subset \underline{X}$

Q: Under which conditions is $P(Y=y | A=a)$ recoverable?

(c) Selection in causal inference

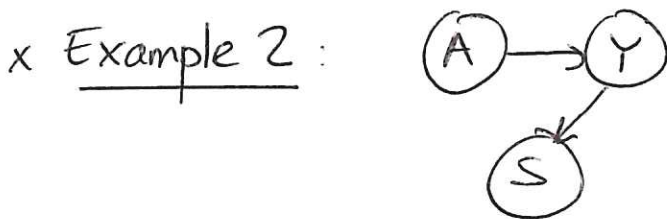
dataset $\sim P(\underline{X}=\underline{x} | S=1) +$ unbiased samples $P(\underline{T}=\underline{t})$
for $\underline{T} \subset \underline{X}$

Q: Under which conditions is $P(Y^a=y) = P(Y=y | do(A=a))$ recoverable?



Since $Y \perp S \mid A$, $P(Y=y \mid A=a, S=1)$
 $\stackrel{||}{=} P(Y=y \mid A=a)$ (*)
 $\stackrel{||}{=} P(Y^a=y)$ no confounding

\Rightarrow No need to resort to additional data external to the biased study to recover conditional effects (& causal effects).

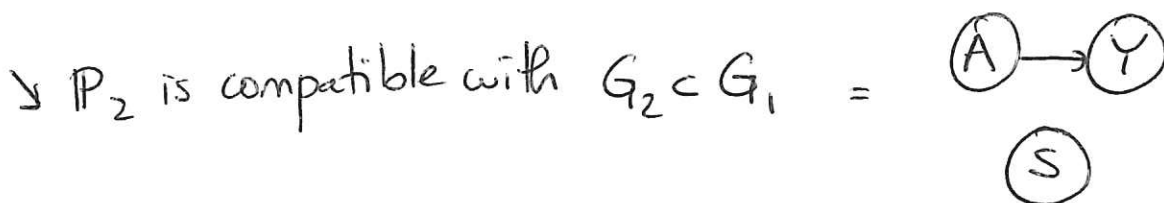
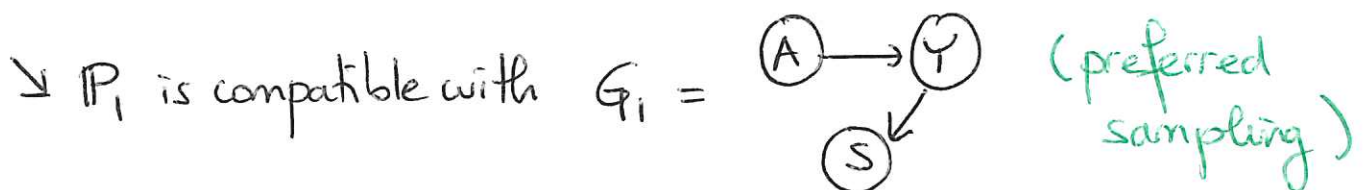


$Y \not\perp S \mid A$ so (*) is not true in general.

In fact, we show next

that no matter how many samples $(A, Y, S=1)$ are drawn $\sim P(A=a, Y=y \mid S=1)$, or how sophisticated the estimation technique is, the estimator of $P(Y=y \mid A=a)$ will never converge to its true value.

Consider P_1 and P_2 such that



(sampling at random)

Task = Construct P_1 and P_2 such that

$$P_1(A=a, Y=y | S=1) = P_2(A=a, Y=y | S=1)$$

they agree on the biased sample

but for which

$$P_1(Y=y | A=a) \neq P_2(Y=y | A=a)$$

↑ The assumptions embedded in the causal model are not sufficient to express conditional probabilities in terms of the distribution under selected bias.

Put $p_i(a, y, s) = P_i(A=a, Y=y, S=s)$, $i=1, 2$.

- P_1 compatible with $G_1 \Rightarrow p_1(a, y, s) = p_1(s|y)p_1(y|a)p_1(a)$
- P_2 " " " $G_2 \Rightarrow p_2(a, y, s) = p_2(s)p_2(y|a)p_2(a)$

we need to define explicitly all these terms so that the task above is fulfilled.

$$p_1(a, y | S=1) = p_2(a, y | S=1) = p_2(a, y)$$

↑ enforced by our task $A, Y \perp S$ in G_2

Take $p_i(a) = p_2(a) = \text{something arbitrary}$, $\forall a$

dividing all members by $p_1 = p_2$ yields

$$p_1(y | a, S=1) = p_2(y | a, S=1) = p_2(y | a)$$

Thus

$$\begin{aligned}
 P_2(y|a) &= p_1(y|a, S=1) = \frac{p_1(y, a, S=1)}{p_1(a, S=1)} \\
 &= \frac{p_1(S=1|y) p_1(y|a) p_1(a)}{p_1(S=1|a) p_1(a)} \\
 &= \frac{p_1(S=1|y) p_1(y|a)}{p_1(S=1|y, a) p(y|a) + p_1(S=1|1-y, a) \times p_1(1-y|a)}
 \end{aligned}$$

Put $\alpha = p_1(S=1|y)$

$\beta = p_1(S=1|1-y)$; $\alpha \neq \beta$

$p_1(y|a) = p_1(1-y|a) = \frac{1}{2}$

(all terms p_1
are now
defined)

⇓

$$P_2(y|a) = \frac{\alpha}{\alpha + \beta} \neq \frac{1}{2} p_1(y|a), \text{ as required } \blacksquare$$

Definition: S-RECOVERABILITY.

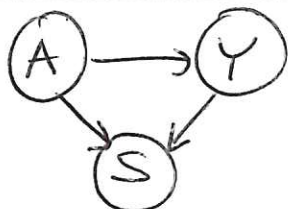
The distribution $P(Y=y | A=a)$ is said to be S -recoverable from the selection biased data if \forall two distributions P_1 and P_2 compatible with the DA G_1 ,

$$P_1(\underline{X}=\underline{x} | S=1) = P_2(\underline{X}=\underline{x} | S=1) > 0$$

⇓

$$P_1(Y=y | A=a) = P_2(Y=y | A=a)$$

x Example 3



(50)

$P(Y=y | A=a)$ is not s -recoverable (corollary from example 2).

Theorem Bareinboim, Tian & Pearl (2014)

$P(Y=y | A=a)$ is s -recoverable from its DAG
 \iff
 $S \perp Y | A$.

\Uparrow Sufficiency is immediate

\Downarrow Necessity is tedious: need to show that for all DAGs in which $S \not\perp Y | A$ we can construct a counterexample showing agreement on $P(\underline{X}=\underline{x} | S=1)$ & disagreement on $P(Y=y | A=a)$.

We turn our attention to the case where external additional measurements are available at the population level (e.g. from census data).

\hookrightarrow The definition of s -recoverability needs to be redefined to account for external data.

Suppose we collect

\rightarrow data under selection bias $P_i(\underline{M}=\underline{m} | S=1)$

\rightarrow external data $P_i(\underline{T}=\underline{t})$

where $\underline{M} \subset \{X_1, \dots, X_d\}$; $\underline{T} \subset \{X_1, \dots, X_d\}$

The distribution $P(Y=y | A=a)$ is said to be s-recoverable from selection bias if for any two proba distrib P_1 and P_2 compatible with the DAG,

$$P_1(M=m | S=1) = P_2(M=m | S=1) > 0$$

$$\& P_1(T=t) = P_2(T=t) > 0$$

$$\Rightarrow P_1(Y=y | A=a) = P_2(Y=y | A=a)$$

We identify conditions under which s-recoverability holds:

$$P(Y=y | A=a) = \sum_z P(Y=y, Z=z | A=a)$$

For some variables Z that are measured both in the biased study with A and Y , and at the population level together with A .

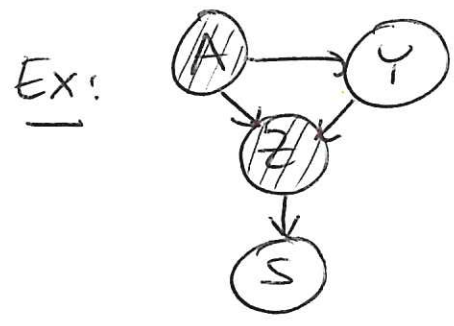
$$= \sum_z P(Y=y | Z=z, A=a) P(Z=z | A=a)$$

if $S \perp Y | Z, A$

$$= \sum_z \underbrace{P(Y=y | Z=z, A=a, S=1)}_{\text{Estimated from selection biased data}} \times \underbrace{P(Z=z | A=a)}_{\text{Estimated from external data (population level)}}$$

This term can be estimated from selection biased data

Estimated from external data (population level)



Ex: $S \perp Y | A, Z$ so we can recover the cond. proba $P(Y=y | A=a)$ using the expression above.

Summary = If there is a set Z that is measured in the biased study together with $\{A, Y\}$ and at the population level with A such that $S \perp Y | Z, A$, then $P(Y=y | A=a)$ is S -recoverable and

$$P(Y=y | A=a) = \sum_z P(Y=y | Z=z, A=a, S=1) \times P(Z=z | A=a)$$

• Recovering causal effects $P(Y^a=y)$

Notation: $Z^* =$ ~~set~~ set of variables = $Z^+ \cup Z^-$
 Z^+ contains all non-descendants of A
 Z^- — " ——— descendants of A

x Case I = No confounding ($Z^+ = \emptyset$)

$$P(Y^a=y) = P(Y=y | do(A=a)) = P(Y=y | A=a) \quad \text{no confounding}$$

$$= \sum_{z^-} \underbrace{P(Y=y | A=a)}_{\parallel} P(Z^-=z^-)$$

if $Z^- \perp Y | A$

$$\underbrace{P(Y=y | A=a, Z^-=z^-)}_{\parallel}$$

if $S \perp Y | A, Z^-$

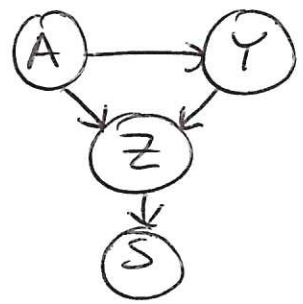
$$P(Y=y | A=a, Z^-=z^-, S=1)$$

$$P(Y^a=y) = \sum_{z^-} P(Y=y | A=a, Z^-=z^-, S=1) P(Z^-=z^-)$$

↑ Expression obtained under $S \perp Y | A, Z^-$ ($\equiv S$ -recov.) + additional assumption $Z^- \perp Y | A$; the price

to pay for no external measurement of both A and Z at the population level (compare the weighting by $P(Z=z)$ instead of $P(Z=z | A=a)$)

• Ex (continued)

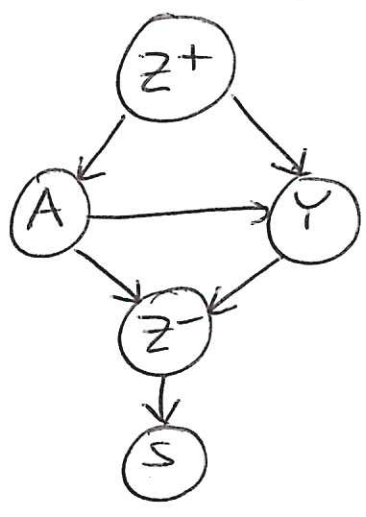


$Z \not\perp Y | A$
 $S \perp Y | A, Z$

$$P(Y^a = y) = \sum_z P(Y=y | A=a, Z=z, S=1) P(Z=z | A=a)$$

≠ $\sum_z \text{---} \text{---} \text{---} \times P(Z=z)$ □

• Case II = general case



$$P(Y=y | do(A=a))$$

$$= \sum_{z^+} P(Y=y | A=a, z^+=z^+) P(z^+=z^+)$$

(backdoor criterion adjusting for z^+)

$$= \sum_{z^+, z^-} P(Y=y | A=a, z^+=z^+) \times P(z^+=z^+, z^-=z^-)$$

$$= P(Y=y | A=a, z^+=z^+, z^-=z^-)$$

provided $z^- \perp Y | A, z^+$

$$= P(Y=y | A=a, z=z)$$

$$= P(Y=y | A=a, z=z, S=1)$$

provided $S \perp Y | A, z$

$$P(Y=y | do(A=a)) = \sum_z P(Y=y | A=a, Z=z, S=1) P(Z=z).$$

Under (i) Z^+ blocks all backdoor paths from A to Y.

(i) $Z^- \perp Y | A, Z^+$

(ii) $S \perp Y | A, Z$

(iii) $Z \cup \{A, Y\} \underset{Z}{\sim}$ selection biased
 $\underset{Z}{\sim}$ population level.



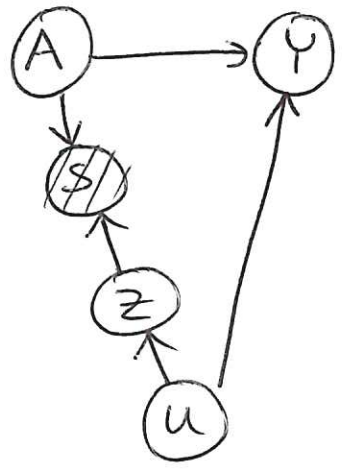
These 4 conditions are called selection-backdoor criterion.
 & the expression above the selection-backdoor adjustment.

* Remark = The selection-backdoor adjustment coincide with the backdoor formula derived in the presence of confounding. Note however that we are adjusting for $Z = Z^+ \cup Z^-$, where Z^- contains descendants of A, which was forbidden in the backdoor criterion.

There are no contradictions. The reason is that the backdoor criterion is a sufficient condition. There is extra room to adjust for descendants under additional assumptions, and this is exactly what the selection-backdoor criterion is.

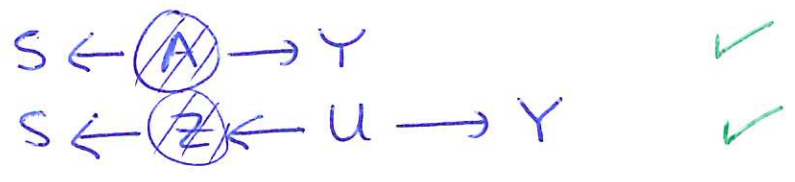
Note that the selection-backdoor criterion is a sufficient condition to recover the causal effect, as illustrated in the next example.

• Ex = HIV study

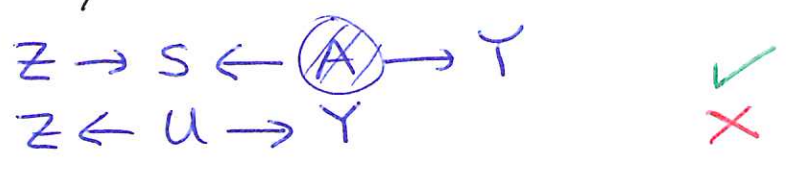


No confounding $Z_+ = \emptyset$
 Take $Z_- = Z$.

(a) $S \perp Y \mid A, Z$ since



(b) $Z \not\perp Y \mid A$ since



⇒ Condition (ii) of the selection-backdoor criterion fails to hold. However, we can derive the same adjustment formula from scratch, making use of alternative independence structure emerging from the graph.

$$\begin{aligned}
 P(Y=y \mid \text{do}(A=a)) &= P(Y=y \mid A=a) \\
 &= \sum_z P(Y=y \mid Z=z, A=a) \underbrace{P(Z=z \mid A=a)}_{P(Z=z)}
 \end{aligned}$$

since $Z \perp A$:



since $S \perp Y \mid A, Z$

$$= \sum_z P(Y=y \mid Z=z, A=a, S=1) P(Z=z)$$



x Remark : Selection-backdoor adjustment and exchangeability

(56)

Assuming that the selection-backdoor criterion holds,

$$\mathbb{P}(Y^a = y) = \sum_z \mathbb{P}(Y = y \mid A = a, Z = z, S = 1) \mathbb{P}(Z = z)$$

$$= \sum_z \frac{\mathbb{P}(Y = y, A = a, Z = z, S = 1)}{\mathbb{P}(A = a, Z = z, S = 1)} \mathbb{P}(Z = z)$$

$$= \sum_z \frac{1}{\mathbb{P}(A = a, S = 1 \mid Z = z)} \mathbb{P}(Y = y, A = a, Z = z, S = 1)$$

$$= \sum_z \frac{1}{\mathbb{P}(A = a \mid Z = z)} \frac{1}{\mathbb{P}(S = 1 \mid Z = z, A = a)} \mathbb{P}(\text{-joint-})$$

adjusts for confounding adjusts for selection bias

$$= \sum_z \underbrace{\omega(a, z) \omega^s(a, z)}_{=: \omega^c(a, z)} \mathbb{P}(\text{-joint-})$$

⇒ Each observation_i receives a weight $\omega^c(A_i, Z_i)$
 $(A_i, Z_i, Y_i, S_i = 1)$

[We may consider as well stabilized weights

$$\tilde{\omega}^c(a, z) := \tilde{\omega}(a, z) \tilde{\omega}^s(a, z), \text{ where}$$

$$\tilde{\omega}(a, z) := \mathbb{P}(A = a) / \mathbb{P}(A = a \mid Z = z) \quad (\text{p. 9c})$$

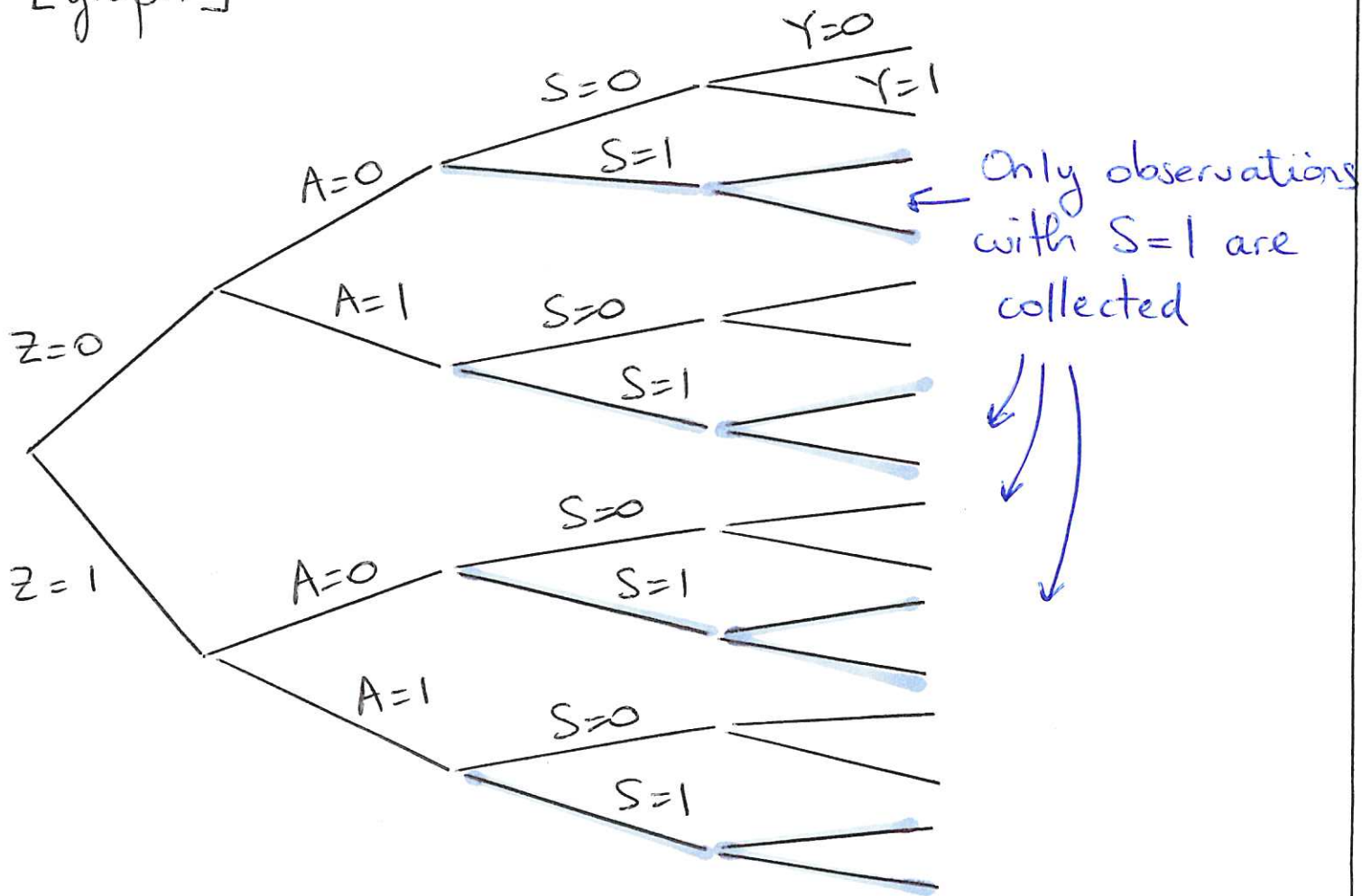
$$\tilde{\omega}^s(a, z) := \mathbb{P}(S = 1 \mid A = a) / \mathbb{P}(S = 1 \mid A = a, Z = z)]$$

↑
 The numerator must be independent of Z

The modified HT estimator of the ATE is obtained by fitting the (saturated) linear model

$Y | A, S=1 = \beta_0 + \beta_1 A + \varepsilon$ using a weighted LS approach, where each observation $(A_i, Z_i, Y_i, S_i=1)$ receives weight $w^c(A_i, Z_i)$.

[graph]

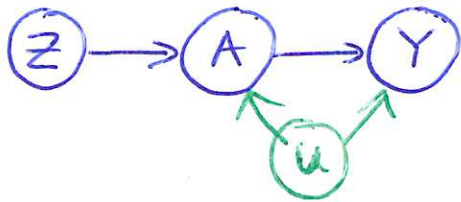


When interested in estimating $P(Y^a=1)$, we must reallocate all users with $A=0$ to the branch $A=1$ [that's the job of $w^a(a, z)$] & reallocate users with $S=0$ to the branch $S=1$ [the job of $w^s(a, z)$]

The appropriate scalings are $1/P(A=a|z=z)$ and $1/P(S=1|A=a, z=z)$, respectively. (see explanation p.7/8 with no censoring)

CI = MEASUREMENT BIAS

- The main contribution of this chapter is to introduce the Intent-To-Treat (ITT) effect under partial compliance to treatment assignment.



Z = treatment allocation $\in \{0, 1\}$

A = treatment status $\in \{0, 1\}$

Y = outcome $\in \{0, 1\}$

U = unmeasured confounders.

↑

In a RCT with partial compliance, the treatment assignment is independent of any confounder U , so that

$$\mathbb{P}(Y^{z=1} = 1) = \mathbb{P}(Y=1 \mid Z=1)$$

$$\& \mathbb{P}(Y^{z=0} = 1) = \mathbb{P}(Y=1 \mid Z=0).$$

↑

The difference between these two quantities is precisely the Intent-To-Treat (ITT):

$$\text{ITT} = \mathbb{E}(Y^{z=1} - Y^{z=0}).$$

- The notation $Y^{z=z}$ stands for Y under the intervention $do(Z=z)$. Similarly, we denote $A^{z=1}$ and $A^{z=0}$ the potential outcomes of A under the interventions $do(Z=1)$ and $do(Z=0)$, so that $A = Z A^{z=1} + (1-Z) A^{z=0}$.
- In addition, we introduce $Y^{a=1}$ and $Y^{a=0}$ to denote Y under the intervention $do(A=1)$ and $do(A=0)$.

Then, $Y = AY^{a=1} + (1-A)Y^{a=0}$.

individual i

indpt of Z

$(Z, A^{z=0}, A^{z=1}, Y^{a=0}, Y^{a=1})$

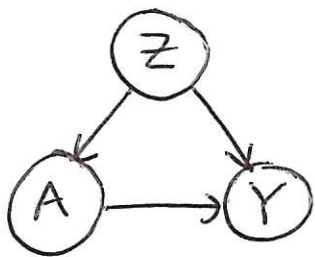
$A = ZA^{z=1} + (1-Z)A^{z=0}$

$Y = AY^{a=1} + (1-A)Y^{a=0}$

↳ observe (Z_i, A_i, Y_i) for each individual i .

x Remark = The DAG $Z \rightarrow A \rightarrow Y$ implicitly assumes the exclusion principle: no direct arrow from Z to Y .

Consider the more general scenario:



Ex: $Z \rightarrow Y$ encodes ^{that} the awareness of the assigned treatment might lead to changes in the behavior of study participants; having an effect on the outcome.

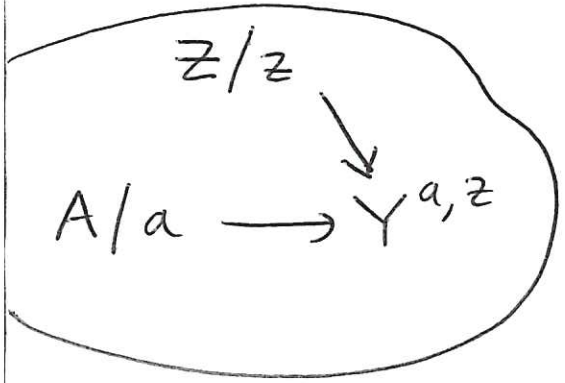
$Z \rightarrow A$: degree of adherence to the assigned treatment

$A \rightarrow Y$: effect of the treatment

$Z \rightarrow Y$: concurrent behavioral change

↳ to break this, perform a double-blind RCT.

In the general case, we need to consider potential outcomes $Y^{a,z}$, $a \in \{0,1\}$ $z \in \{0,1\}$ under the double intervention $do(Z=z)$ & $do(A=a)$.



Observe

$$Y = AZY^{1,1} + (1-A)Y^{0,1} + A(1-Z)Y^{1,0} + (1-A)(1-Z)Y^{0,0}$$

When the exclusion principle holds (no arrow from Z to Y),

$$Y^{a=0, z=0} = Y^{a=0, z=1} = Y^{a=0}$$

$$Y^{a=1, z=0} = Y^{a=1, z=1} = Y^{a=1}$$

and Y simplifies to $Y = AY^{a=1} + (1-A)Y^{a=0}$, as required.

Assumption 1: $A^{z=0} \sim B(\pi_0)$
 $A^{z=1} \sim B(\pi_1)$

If allocated to treatment, the proba to receive the treatment is π_0 .

If allocated to control, the proba to receive the treatment is π_1 .

The total population can be partitioned into 4 groups, based on their values of $A^{z=0}$ and $A^{z=1}$.

		$A^{z=1} =$	
		0	1
$A^{z=0} =$	0	Never Takers	Compliers
	1	Defiers	Always Takers

Assumption 2 $= \pi_0 = 0$

Subjects in the control group will not take the treatment
 $\Leftrightarrow A^{z=0} = 0$ a.s.

Under assumption 2, the population is split between never takers & compliers:

$$\text{Never Takers} = \{ \omega \in \Omega \mid A^{z=1}(\omega) = 0 \}$$

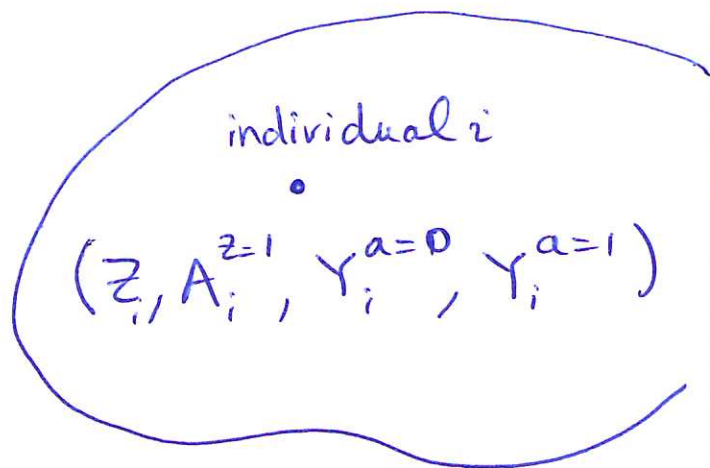
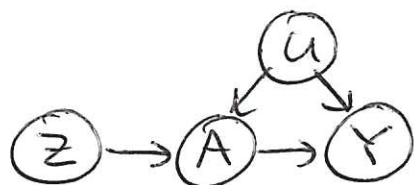
$$\text{Compliers} = \{ \omega \in \Omega \mid A^{z=1}(\omega) = 1 \}$$

• Summary: Under assumptions 1 and 2, and under the exclusion principle,

$$A^{z=1} \sim B(n_i)$$

$$A = z A^{z=1}$$

$$Y = A Y^{a=1} + (1-A) Y^{a=0}$$



Note that under the intervention $do(Z=0)$ and $do(Z=1)$,

62

• When $do(Z=0)$, then $A = A^{Z=0} \equiv 0$ (assumption 2),
we
and $Y = Y^{a=0}$.

• When we $do(Z=1)$, then $A = A^{Z=1}$, and
 $Y = A^{Z=1} Y^{a=1} + (1 - A^{Z=1}) Y^{a=0}$.

$$\begin{aligned} ITT &= E(Y^{Z=1} - Y^{Z=0}) \\ &= E(A^{Z=1} Y^{a=1} + (1 - A^{Z=1}) Y^{a=0} - Y^{a=0}) \\ &= E(\text{---} \text{---} \mid A^{Z=1} = 1) P(A^{Z=1} = 1) \\ &\quad + \underbrace{E(\text{---} \text{---} \mid A^{Z=1} = 0)}_{=0} P(A^{Z=1} = 0) \end{aligned}$$

$$ITT = E(Y^{a=1} - Y^{a=0} \mid A^{Z=1} = 1) P(A^{Z=1} = 1)$$

treatment effect of the
compliers
aka the

$= \pi_1 =$ proportion
of compliers.

Local Average Treatment Effect (LATE)
or the Complier Average Causal Effect (CACE)

* Remark: When $\pi_1 = 1$, then $ITT = E(Y^{a=1} - Y^{a=0}) = ATE$.

x Remark = Identification of the LATE under monotonicity (Imbens & Angrist '94)

$$\begin{aligned} ITT &= E(Y|Z=1) - E(Y|Z=0) \\ &= E(A^{z=1} Y^{a=1} + (1-A^{z=1}) Y^{a=0} | Z=1) \\ &\quad - E(A^{z=1} Y^{a=1} + (1-A^{z=1}) Y^{a=0} | Z=0) \end{aligned}$$

Assuming $Z \perp (A^{z=0}, A^{z=1}, Y^{a=0}, Y^{a=1})$ (p.59)

$$\begin{aligned} &= E[(A^{z=1} - A^{z=0})(Y^{a=1} - Y^{a=0})] \\ &= E(Y^{a=1} - Y^{a=0} | A^{z=1} - A^{z=0} = 1) P(A^{z=1} - A^{z=0} = 1) \\ &\quad + E(\text{---} | \text{---} = -1) P(\text{---} = -1) \end{aligned}$$

Without further assumptions, we may construct proba on $(A^{z=0}, A^{z=1}, Y^{a=0}, Y^{a=1})$ such that the ITT is zero, while the causal effect of A on Y is > 0 .

However, assuming $A^{z=1} \geq A^{z=0}$ (monotonicity), the second term in the RHS vanishes and

$$ITT = E(Y^{a=1} - Y^{a=0} | A^{z=1} \neq A^{z=0}) P(A^{z=1} \neq A^{z=0})$$

Population made of ~~of~~ compliers. Under the previous/earlier assumption that $A^{z=0} \equiv 0$, we recover conditioning on compliers as well.

"linear latent model"

Ex: $Z \in \{0, 1\}$ $A = \mathbb{1}(\gamma_0 + \gamma_1 Z + \varepsilon > 0)$ $Y = \beta_0 + \beta_1 A + \eta$

Re-arranging the terms,

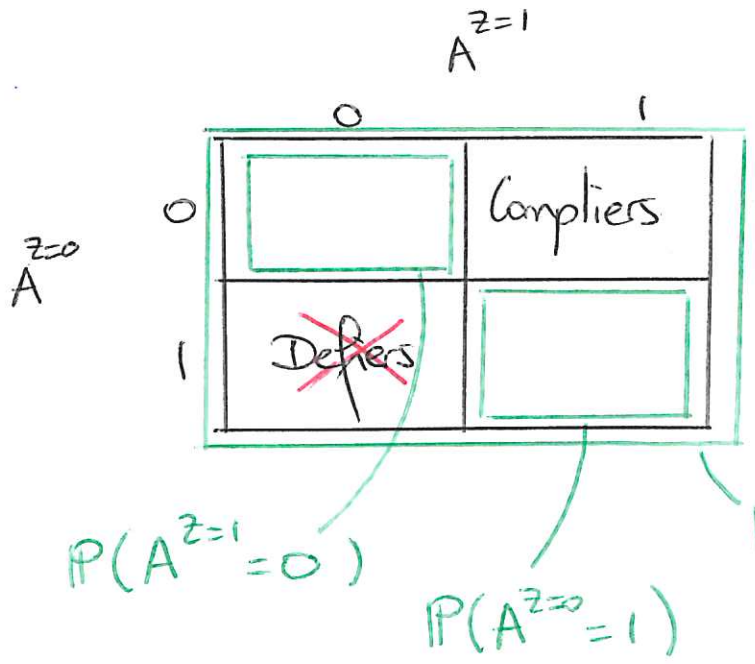
$$E(Y^{a=1} - Y^{a=0} | A^{z=1} \neq A^{z=0}) = \frac{E(Y|z=1) - E(Y|z=0)}{P(A^{z=1} \neq A^{z=0})}$$

This is the CATE i.e. the ATE within compliers. Under monotonicity, we have indeed that $\{A^{z=1} \neq A^{z=0}\} = \{A^{z=0} = 0, A^{z=1} = 1\} = \{\text{compliers}\}$

denominateur = $P(A^{z=0} = 0, A^{z=1} = 1)$, which can be re-expressed as $1 - P(A^{z=1} = 0) - P(A^{z=0} = 1)$ since there are no defiers.

$$= 1 - P(A=1 | z=0) - P(A=0 | z=1)$$

$$= P(A=1 | z=1) - P(A=1 | z=0)$$

$$= E(A|z=1) - E(A|z=0)$$


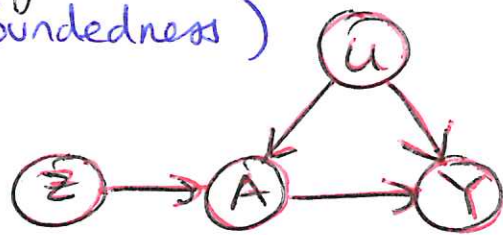
Summary = _____

$$CATE = E(Y^{a=1} - Y^{a=0} | A^{z=1} \neq A^{z=0}) = \frac{E(Y|z=1) - E(Y|z=0)}{E(A|z=1) - E(A|z=0)}$$

↑ Non-parametric identification of the CATE under monotonicity assumptions.

When there is partial compliance to treatment assignment, (65)
the variable Z is called an INSTRUMENT. It
satisfies 3 key properties:

- (i) Z has a causal effect on A (relevance)
- (ii) Z has a causal effect on Y that is
fully mediated by A (exclusion restriction)
- (iii) There are no backdoor paths from
 Z to Y (instrumental unconfoundedness)



Previously, we established a non-parametric expression
of the CATE. It turns out that there are no
non-parametric expressions for the ATE, since the
effect of A on Y is confounded by U .

⇒ With instrument variables, we must make assumptions
about the parametric form to identify causal effects.

"do" probabilities can be expressed as
functions of observation probabilities
(eg. backdoor adjustment, frontdoor
adjustment)

Ex 1: Binary linear

$$A, Z \in \{0, 1\}$$

$$Y = \delta A + \alpha U$$

$$\delta = \text{ATE}$$

$$E(Y|Z=1) - E(Y|Z=0)$$

(66)

$$= E(\delta A + \alpha U | Z=1) - E(\delta A + \alpha U | Z=0)$$

$$= \delta [E(A|Z=1) - E(A|Z=0)]$$

$$+ \alpha [E(U|Z=1) - E(U|Z=0)]$$

= 0 since $U \perp Z$

$$\Rightarrow \delta = ATE = \frac{E(Y|Z=1) - E(Y|Z=0)}{E(A|Z=1) - E(A|Z=0)}$$

= same expression as the CATE under monotonicity. ▣

EX 2: Continuous linear $A, Z \in \mathbb{R}$

Consider instead $\text{Cov}(Y, Z) = E(YZ) - E(Y)E(Z)$

$$= E(\delta A + \alpha U)Z$$

$$- E(\delta A + \alpha U)E(Z)$$

$$= \delta \{E(AZ) - (EA)(EZ)\}$$

$$+ \alpha \{E(UZ) - (EU)(EZ)\}$$

$$= \delta \text{Cov}(A, Z) + \alpha \underbrace{\text{Cov}(U, Z)}_{=0}$$

$$\Rightarrow \delta = ATE = \frac{\text{Cov}(Y, Z)}{\text{Cov}(A, Z)}$$

▣

x Remark: When there is a weak association between A and Z , the denominators $E(A|Z=1) - E(A|Z=0)$ & $\text{Cov}(A, Z)$ are close to 0, inflating the numerator = ITT [potential for explosive bias].

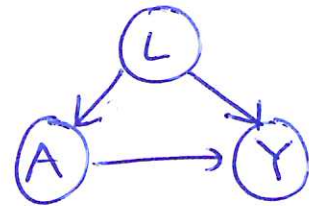
CI = PROPENSITY SCORES

When L satisfies the backdoor criterion (\equiv conditional exchangeability), so that $(Y^0, Y^1) \perp A \mid L$, we may wonder if it is necessary to condition on the whole vector L , especially when L is high dimensional.

Rosenbaum & Rubin (1983) showed that it is enough to condition on $P(A=1 \mid L)$. The quantity $e(L) := P(A=1 \mid L)$ is commonly referred to as the PROPENSITY SCORE (PS) since it represents the propensity (probability) of receiving the treatment given L .

Thm $(Y^0, Y^1) \perp A \mid L \Rightarrow (Y^0, Y^1) \perp A \mid e(L)$

proof: In the simple setting where



the arrow from L to A symbolizes the probabilistic relationship between L and A i.e. $P(A=1 \mid L)$ for a binary treatment variable, which is precisely the PS.

More formally, we show that $(Y^0, Y^1) \perp A \mid e(L)$ by showing that $P(A=1 \mid Y^a, e(L))$ does not depend on Y^a , where $a=0$ or $a=1$.

$$\mathbb{E}(A=1 | Y^a, e(L)) = \mathbb{E}(A | Y^a, e(L)) \quad (65)$$

$$= \mathbb{E} \left[\underbrace{\mathbb{E}\{A | Y^a, \cancel{e(L)}, L\}}_{\text{average on something finer.}} | Y^a, e(L) \right]$$

$$= \mathbb{E} \left[\underbrace{\mathbb{E}\{A | Y^a, L\}}_{=} | Y^a, e(L) \right]$$

$$\mathbb{E}\{A | L\} \text{ since } Y^a \perp A | L$$

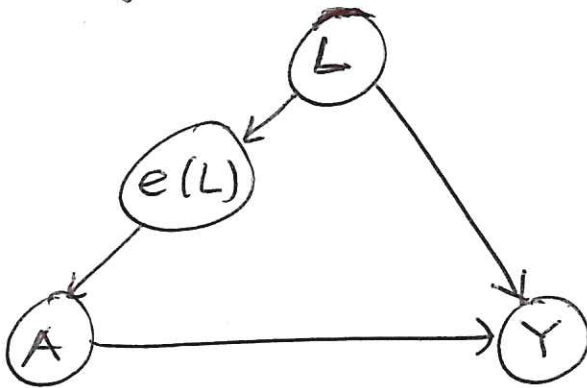
$$\mathbb{P}(A=1 | L) = e(L)$$

$$= \mathbb{E}[e(L) | Y^a, e(L)]$$

$$= e(L)$$

← independent of Y^a .

Graphically, this means that $e(L)$ is a full mediator of the effect of L on A



And we see that $A \perp L | e(L)$
 i.e. the PS balances the covariates between the treated & the untreated.

vs Randomization which balances measured & unmeasured covariates [& thus preferred over PS]

Csq: We can adjust for $e(L)$ in place of L .

Remark: balancing scores

Rosenbaum & Rubin (1983) introduced the concept of a balancing score $b(L)$.

Def: $b(L)$ is a balancing score

$$A \perp L \mid b(L)$$

↓ the PS is the simplest example of a balancing score.

They showed the following result

Thm: $A \perp L \mid b(L)$

$$e(L) = f(b(L)) \text{ for some function } f$$

Ex: $b(L) = L$ is the finest balancing score
 $b(L) = e(L)$ is the coarsest.

proof \uparrow Suppose $\exists f$ s.t. $e(L) = f(b(L))$

We want to show that $A \perp L \mid b(L)$.

It is sufficient to show that

$$\mathbb{E}(A \mid b(L)) = e(L)$$

$$\text{since } e(L) = \mathbb{E}(A \mid L)$$

$$= \mathbb{E}(A \mid L, b(L)).$$

This would then imply that

$$\mathbb{E}(A \mid b(L)) = \mathbb{E}(A \mid L, b(L))$$

i.e. $A \perp L \mid b(L)$.

$$\begin{aligned} \mathbb{E}(e(L) | b(L)) &= \mathbb{E}\left(\mathbb{E}(A | L) | b(L)\right) \quad (67) \\ &= \mathbb{E}\left(\underbrace{\mathbb{E}(A | L, b(L))}_{\text{finer averaging}} | b(L)\right) \\ &= \mathbb{E}(A | b(L)) \end{aligned}$$

On the other hand,

$$\begin{aligned} \mathbb{E}(e(L) | b(L)) &= \mathbb{E}(f(b(L)) | b(L)) \\ &= f(b(L)) = e(L), \end{aligned}$$

so that $\mathbb{E}(A | b(L)) = e(L)$ indeed.

□ Suppose that $L \perp A | b(L)$.

By contradiction, suppose that $\exists l_1, l_2$ s.t.

$$e(l_1) \neq e(l_2)$$

$$\text{but } b(l_1) = b(l_2)$$

By definition of e ,

$$\mathbb{P}(A=1 | L=l_1) \neq \mathbb{P}(A=1 | L=l_2)$$

$$\mathbb{P}(A=1 | L=l_1, b(L)=b(l_1)) \quad \parallel$$

$$\mathbb{P}(A=1 | L=l_2, b(L)=b(l_2)) \quad \parallel$$

$\Rightarrow L \not\perp A | b(L)$; a contradiction. □

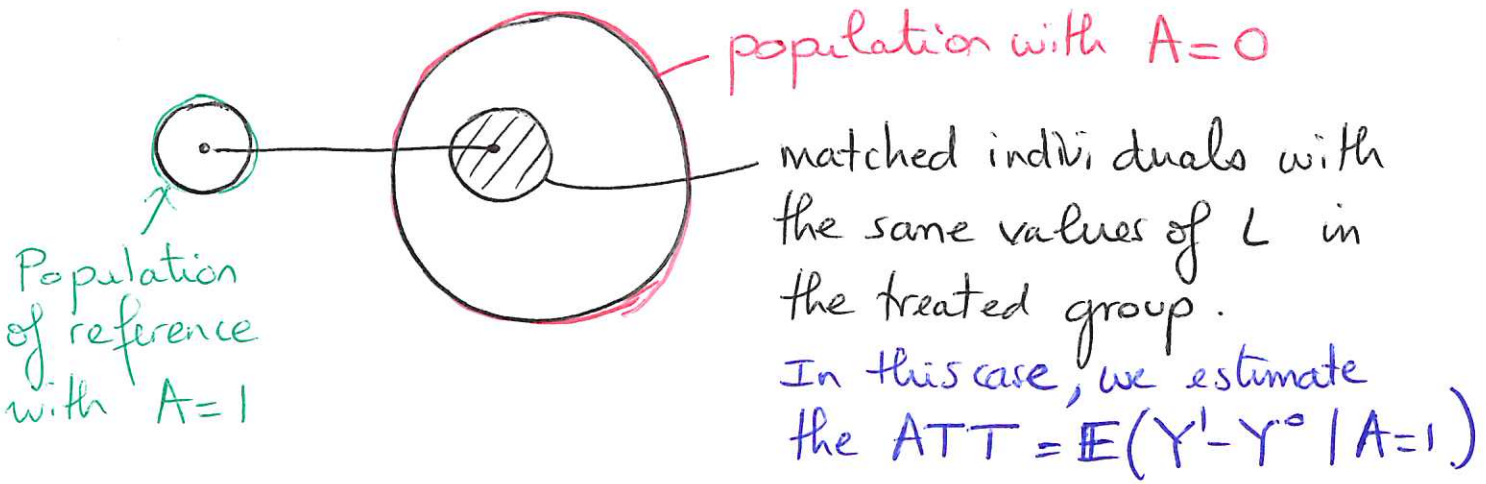
• Propensity Standardization = Adjust for $e(L)$ instead of L .

$$\mathbb{E}Y^a = \sum_e \mathbb{P}(Y=y | A=a, e(L)=e) \mathbb{P}(e(L)=e)$$

(if discrete)

• Propensity Matching

Recall that matching rebalances the distribution of the variable L within the two groups. Depending on how the matching is performed, the quantity of interest is the ATT or ATNT (see p.20)



Instead of matching based on L , we may match based on $e(L)$

$$ATT = E(Y^1 | A=1) - E(Y^0 | A=1) \quad \text{consistency}$$

$$= \underbrace{E(Y | A=1)}_{\text{easily estimated}} - \underbrace{E(Y^0 | A=1)}_{\text{counterfactual term}}$$

$$\downarrow$$

$$\frac{1}{n} \sum_{i=1}^n Y_i \mathbb{1}(A_i=1)$$

$i = i$ -th individual

$$\downarrow$$

$$\frac{1}{n} \sum_{i=1}^n \hat{Y}_i^0 \mathbb{1}(A_i=1)$$

where \hat{Y}_i^0 is an estimate of the potential outcome for individual i .

$$\hat{Y}_i^0 = \frac{1}{K} \sum_{k \in K_i} Y_k,$$

where K_i contains the K closest individuals in the control group based on their $e(L)$ value.

• Stratification.

↳ After computing the PS $e(L)$, the sample is divided into strata / blocks such that within each block, the PS is approximately constant.

Consider J intervals $(b_{j-1}, b_j]$. We define

$$s(L) = j \quad \text{if} \quad \log \left\{ \frac{e(L)}{1-e(L)} \right\} \in (b_{j-1}, b_j]$$

Rubin (2001)

$$\begin{aligned} \text{ATT} &= \mathbb{E}(Y^1 - Y^0 \mid A=1) \\ &= \mathbb{E}(\mathbb{E}(Y^1 - Y^0 \mid A=1, s(L)) \mid A=1) \\ &\quad \parallel \\ &\quad \text{ATT}(s(L)) \\ &= \mathbb{E}(\text{ATT}(s(L)) \mid A=1) \\ &= \sum_s \underbrace{\text{ATT}(s)}_{\substack{\uparrow \\ \text{estimated using}}} \underbrace{\mathbb{P}(S(L) = s \mid A=1)}_{\substack{\downarrow \\ \approx \text{number of individuals} \\ \text{in the treated group} \\ \text{with } s(L) = s \text{ divided} \\ \text{by the number of} \\ \text{individuals in the treated} \\ \text{group.}}} \end{aligned}$$

estimated using

$$\begin{aligned} &\frac{1}{n_{1s}} \sum_{i=1}^n Y_i \mathbb{1}(A_i=1) \mathbb{1}(s(L_i)=s) \\ &- \frac{1}{n_{0s}} \sum_{i=1}^n Y_i \mathbb{1}(A_i=0) \mathbb{1}(s(L_i)=s) \end{aligned}$$

↑ direct comparison of the two groups since the PS are almost the same, rebalancing the confounders L .

$$n_{as} = \sum_{i=1}^n \mathbb{1}(A_i=a) \mathbb{1}(s(L_i)=s) = \# \text{ observations with treatment allocation } a \text{ within stratum } s.$$

x Summary of part I

Interested in identifying the causal effect of A on Y

$$ATE = \mathbb{P}(Y^{a=1} = 1) - \mathbb{P}(Y^{a=0} = 1)$$

We may also be interested in the average causal effect amongst the treated individuals

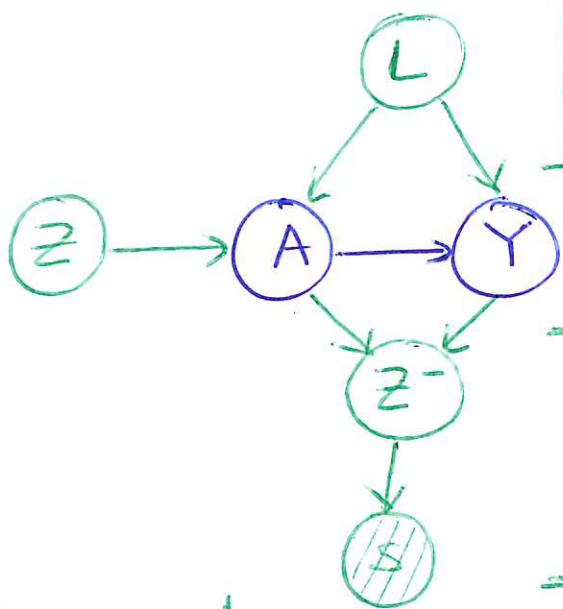
$$ATT = \mathbb{P}(Y^{a=1} = 1 | A=1) - \mathbb{P}(Y^{a=0} = 1 | A=1)$$

In presence of confounders, adjust via the backdoor formula

$$\mathbb{P}(Y^a = y | A=1) \quad (\text{p. 6/21})$$

$$= \sum_l \mathbb{P}(Y=y | A=a, L=l) \mathbb{P}(L=l | A=1)$$

if L blocks all backdoor paths between A and Y. (p. 31)



Under partial complier, consider instead the ITT (p. 58) or the

LATE/CACE (p. 62)

In presence of selection bias, causal effects may be recovered with or without external data. When external data are available, and when the selection ~~is~~ backdoor criterion is satisfied (p. 54),

$$\mathbb{P}(Y^a = y | A=1)$$

$$= \sum_{\text{shaded}} \mathbb{P}(Y=y | A=a, \text{shaded}, S=1) \times \mathbb{P}(\text{shaded} | A=1)$$