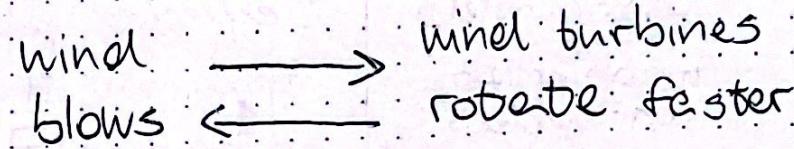


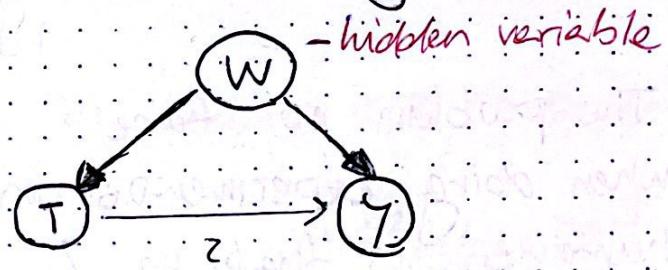
# Methods for Causal Inference

## Reverse causation:

if we don't see the bigger picture - knowing what causes what - then we may conclude the reverse causation!



## Confounding factor



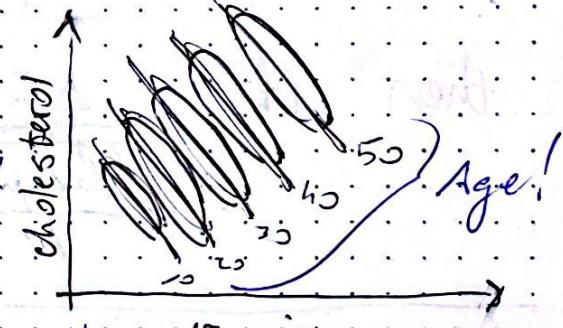
## Simpson's Paradox

- concluding causality from purely associational measure (correlation) can be

very wrong → would be better not to draw any conclusions at all!!

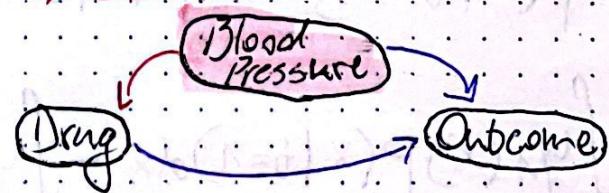


we can see that there is a factor we haven't accounted for



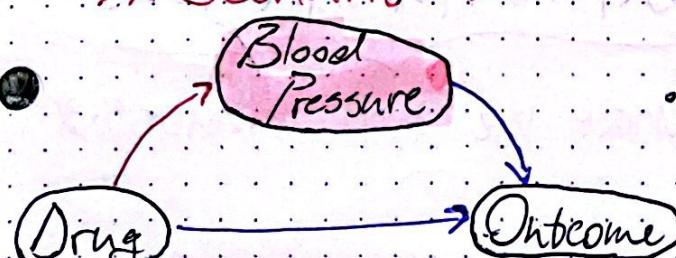
Consider all variables affecting the system of interest ~~and~~ and the role they play

### Scenario 1



• blood pressure is a confounder

here



• blood pressure is a mediator

T - treatment

Y - outcome

X - confounders

U - unobservable confounders

$$\mathbb{E}[Y_{t=1}(x) - Y_{t=0}(x)] =$$

Expected value of  $Y$  when we apply the treatment vs. when we don't given some confounders (= features)

$$= \int (Y_1(x) - Y_0(x)) P(x) dx$$

↑ probability distribution of subjects given confounder  $X$  (like age or sex or smth)

The problem we face

when doing experiments on

humans is that no 2 humans

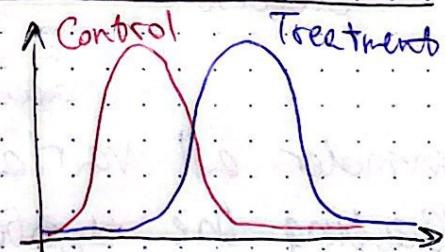
are genetically the same  $\Rightarrow$  let's take average across the whole population

$$\Delta\bar{\mu} = \bar{\mathbb{E}}[Y_{t=1}(x) - Y_{t=0}(x)] = \frac{1}{N} \sum_{i=1}^N (y_1^{(i)}(x) - y_0^{(i)}(x))$$

then if  $\frac{\Delta\bar{\mu}}{\sqrt{\frac{(G\Delta\mu)^2}{N}}} > 6^*$  is big enough then we gucci gucci

~~What's wrong?~~ What's wrong?

$$P(x|t=1) \neq P(x|t=0) \Rightarrow$$



$$\int Y_1(x) P(x|t=1) dx - \int Y_0(x) P(x|t=0) dx$$

if we naively measure it, then we will get false result!

$$\neq \int (Y_1(x) - Y_0(x)) P(x) dx$$

what we are interested in

possible solution is **STRATIFICATION**

↳ measure outcome within each of the groups separately

$$\begin{aligned} \mathbb{E}(\text{healed} | t=1) &= \mathbb{E}(\text{healed} | t=1, \text{young}) p(\text{young}) \\ &\quad + \mathbb{E}(\text{healed} | t=1, \text{old}) p(\text{old}) \end{aligned}$$

## Disadvantages:

- ① all possible outcomes need to be observed
- ② samples can become non representative

2 main frameworks:

↳ causal estimation!

\*1 Potential Outcomes  $\Rightarrow$  Rubin

\*2 Structural Causal Models  $\Rightarrow$  Pearl

↳ causal discovery

Probability axioms for two mutually exclusive events:

$$P(A \text{ or } B) = P(A) + P(B)$$

$$P(A \text{ and } B) = P(A) \cdot P(B)$$

$$P(A) = P(A, \text{not } B) + P(A, B)$$

### Conditional Probability

$$P(X, Y) = P(X|Y) P(Y)$$

### Bayes rule

$$P(X_i|Y) = \frac{P(Y|X_i) P(X_i)}{P(Y)}$$

$$= \frac{P(X_i) P(Y|X_i)}{P(X_1) P(Y|X_1) + \dots + P(X_n) P(Y|X_n)}$$

Incorporating knowledge about the process that generated the data  $\Rightarrow$  **FIRST STEP TOWARDS CAUSAL INFERENCE**

Two events are said to be mutually exclusive when their occurrence is not simultaneous

Two events are said to be independent when occurrence of one cannot control occurrence of other

$$P(X, Y) = P(X) P(Y)$$

$$P(X|Y) = P(X) \quad \text{if } P(Y) \neq 0$$

**Expected Value** - sum over the probability distribution

$$\mathbb{E}[X] = \sum_x x P(X=x)$$

↑  
random variable  
values RV can take

for continuous variables  $\sum \rightarrow \int$

$$\mathbb{E}[X] = \int x P(x) dx$$

### Variance

$$\text{var}(x) = \mathbb{E}[(X - \mathbb{E}[X])^2] = \sum_x (x - \mathbb{E}[X])^2 p_x(x)$$

↑  
how spread out are the values!

### Covariance

$$\sigma_{xy} = \mathbb{E}[(X - \mathbb{E}[X])(Y - \mathbb{E}[Y])]$$

degree to which  $X$  and  $Y$  covary

when normalised, we get a coefficient:

$$r_{xy} = \frac{\sigma_{xy}}{\sigma_x \sigma_y} \quad \Leftarrow \text{PEARSON CORRELATION}$$

$r_{xy} \in [-1, 1]$  but  $r_{xy} = 0$ , ~~means~~  $X$  and  $Y$  when are independent

but  $r_{xy}$  does not imply independence.

$$P(Y|X) = P(Y)$$

$R^2$  coefficient ~ a measure for goodness of a fit

$$R^2 = 1 - \frac{\sum_i (y_i - \hat{y}_i)^2}{\sum_i (y_i - \bar{y})^2}$$

prediction  
retained point  
mean

if fit is perfect, then  $R^2=1$  and  $\hat{R}^2=0$   
implies fit is no better than baseline  $\bar{y}$

### Lecture 3

In linear regression modeling, optimal solution is reached when

in:

$$y = \alpha + \beta x$$

### Multi regression

$y$  on vers  $x_1$  &  $x_2$

$$y = \alpha + \beta_1 x_1 + \beta_2 x_2$$

$\beta > 0$  - positive correlation

$\beta < 0$  - negative correlation

$\beta = 0$  - no linear correlation

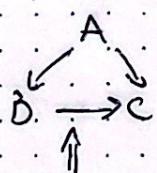
### V-fold

- ① split the data into  $V$  blocks and in each of  $V$  iterations, one block will be used for testing while rest for training

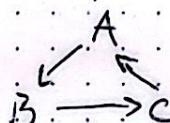


### Directed Graphs

#### Acyclic



#### Cyclic



#### DAGs

### STRUCTURAL CAUSAL MODELS

- 2 sets  $U \& V$  and set of functions of

$V$ -exogenous variable:

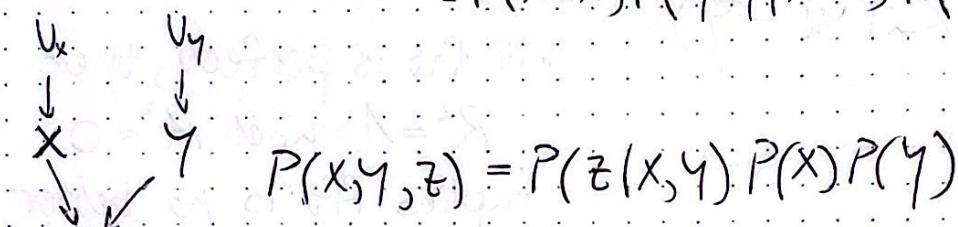
↳ external to the model

$V$ -endogenous:

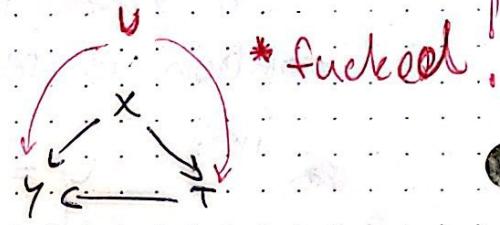
↳ descendant of at least one exogenous variable

## PRODUCT DECOMPOSITION RULE

$$X \rightarrow Y \rightarrow Z \quad P(X=x, Y=y, Z=z) = \\ = P(X=x) P(Y=y | X=x) P(Z=z | Y=y)$$



If  $U_x \perp\!\!\!\perp U_y \perp\!\!\!\perp U_z$  - the outer noises are independent of each other, then we can drop them from our model  
are not confounders



## Lecture 4

# Potential Outcomes Framework [Rubin-Neyman]

$$Y_{\text{observed}}^{(i)} = t^{(i)} Y_1^{(i)} + (1 - t^{(i)}) Y_0^{(i)} = \begin{cases} Y_0^{(i)} & \text{if } t^{(i)} = 0 \\ Y_1^{(i)} & \text{if } t^{(i)} = 1 \end{cases}$$

but in the real world one patient ~~can~~ cannot undergo both scenarios → that's why we introduce **counterfactuals**.

$$Y_{CF}^{(i)} = (1 - t^{(i)}) Y_1^{(i)} + t^{(i)} Y_0^{(i)} \quad \text{X missing outcome}$$

now ideally we want to introduce a measure

$$\tau^{(i)} = Y_1^{(i)} - Y_0^{(i)}$$

telling us what is the **treatment effect** for that individual person.

$$\bar{\tau} = \hat{E}[\tau^{(i)}] = \hat{E}[Y_1^{(i)} - Y_0^{(i)}] = \frac{1}{N} \sum_{i=0}^N (Y_1^{(i)} - Y_0^{(i)})$$

## Regression Adjustment

fitting  $\hat{Y} = \beta_x X + \beta_T T + \varepsilon$



$$\begin{aligned} T &= \begin{pmatrix} 1 & 0 \\ 1 & 0 \\ \vdots & \vdots \\ 0 & 1 \\ 0 & 1 \end{pmatrix} & X &= \begin{pmatrix} 1 & 0 \\ 0 & 1 \\ \vdots & \vdots \\ 1 & 0 \\ 0 & 1 \end{pmatrix} & \rightarrow & \begin{pmatrix} Y_1^{(1)} \\ Y_1^{(2)} \\ \vdots \\ Y_1^{(N-1)} \\ Y_1^{(N)} \end{pmatrix} &= \begin{pmatrix} \beta_{T=0} + \beta_{x=y} \\ \beta_{T=0} + \beta_{x=0} \\ \vdots \\ \beta_{T=1} + \beta_{x=y} \\ \beta_{T=1} + \beta_{x=0} \end{pmatrix} \\ \text{Control} & \text{Drug} & \text{Young} & \text{Old} & & & \end{aligned}$$

Assumptions: Overlap and collinearity

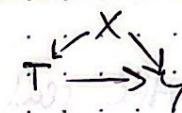
$$\bar{\tau} = \hat{E}[\tau^{(i)}] = \hat{E}[Y_1^{(i)} - Y_0^{(i)}] = \frac{1}{N} \sum_{i=0}^N (Y_1^{(i)} - Y_0^{(i)})$$

since we computed the  $\beta$  we can easily compute the counterfactuals

Another way of computing the counterfactuals is with **MATCHING**, where we try to create clone/ twin for each individual  $\Rightarrow$  try to find most similar unit from sample

Lecture 5

## Adjustment Formula



$$\mathbb{E}[Y_1 - Y_0 | X] = \mathbb{E}[Y_1 | X] - \mathbb{E}[Y_0 | X] =$$

$$\text{CATE} = \mathbb{E}[Y_1 | T=1, X] - \mathbb{E}[Y_0 | T=0, X] =$$

$$\text{ATE} = \mathbb{E}[Y_1 | T=1, X] - \mathbb{E}[Y_1 | T=0, X]$$

$$\mathbb{E}[Y_1 - Y_0] = \mathbb{E}_x [\mathbb{E}[Y_1 - Y_0 | X]] = \mathbb{E}_x [\mathbb{E}[Y | T=1, X] - \mathbb{E}[Y | T=0, X]]$$

adjustment formula

Fit model for  $Q(T, X) = \mathbb{E}[Y | T, X]$

- under linear assumptions:

$$\mathbb{E}[Y | T, X] = \alpha_0 + \beta_X X + \beta_T T + \varepsilon, \quad \mathbb{E}[\varepsilon] = 0$$

$$\text{ATE} = \mathbb{E}_x [\mathbb{E}[Y | T=1, X] - \mathbb{E}[Y | T=0, X]] =$$

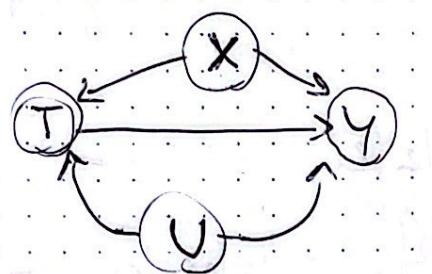
$$= (\cancel{\alpha_0} + \beta_X X + \beta_T \cdot 1) - (\cancel{\alpha_0} + \beta_X X + \beta_T \cdot 0) =$$

$$= \beta_T \quad \text{which is correct! for } \begin{array}{c} X \\ \leftarrow \rightarrow \\ T \end{array}$$

if we condition on  $X$ , then  $Y$  only depends on  $T$ !

## are 6 Instrumental variable

- U - unobserved confounder  $\rightarrow$  conditioning on X alone would not result in a randomised treatment



We don't have access  
to U

What to do? Introduce new **INSTRUMENTAL** variable!

### Example:

T - smoking during pregnancy

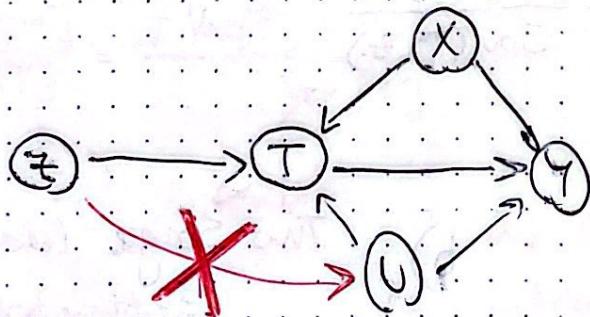
Y - birth weight of child

X - parity / mom's age / weights

U - smth (no access)

Randomise Z (intention-to-treat): either receive encouragement to stop smoking ( $Z=1$ ) or receive usual care ( $Z=0$ )

$$E(Y|Z=1) - E(Y|Z=0) \quad \text{effect } Z \text{ on } Y$$



want to use this  
to get effect T on Y

any effect of Z on Y is  $(Y|Z=1, T) = (Y|Z=0, T)$   
only through T

We want ATE (average treatment effect):

$$E[(Y|t=1) - (Y|t=0)]$$

but we can't get this from our graph directly cuz of  $U$ !

However, we can calculate:

$$T = \frac{E[(Y|z=1) - (Y|z=0)]}{E[(T|z=1) - (T|z=0)]}$$

which is exactly ETA!

\*  $T$  was allocated randomly to subjects.

## ※1 The Wald Estimator (for binary variables)

$$\hat{T} = \frac{\frac{1}{n_{z=1}} \sum_{i \in z=1} (y_{i1} - \bar{y}_{z=0})}{\frac{1}{n_{z=1}} \sum_{i \in z=1} (\bar{T}^{(1)} - \bar{T}^{(0)})}$$

## ※2 IV Estimator

Linear case:  $\hat{T} = \frac{\text{Cov}(Y, t)}{\text{Cov}(T, z)}$

a)  $\hat{T} = \frac{\hat{\text{Cov}}(Y, t)}{\hat{\text{Cov}}(T, z)}$

b) Two-Stage Least-Square Estimator

$$\text{cov}(Y, Z) = \mathbb{E}[YZ] - \mathbb{E}[Y]\mathbb{E}[Z] = \dots =$$

$$= J \text{Cov}(T, Z) + \int_U \text{Cov}(U, Z) =$$

$$= J \text{Cov}(T, Z)$$

\* but since there  
is no covariation  
between  $U$  and  $Z$   
this is 0!

$$Y = JT + \int_U U$$

b) 1 Estimate  $\mathbb{E}[T|Z]$  to obtain  $\hat{T}$  in subspace  $T$

2 Estimate  $\mathbb{E}[Y|\hat{T}]$ , to obtain  $\hat{J}$ , which is fitted coefficient in front of  $\hat{T}$  in this regression

$$Y = \hat{J} \hat{T} + \int_U U$$

Lecture 7. In the case stated above  $T \xrightarrow{J} Y$  naive regression leads to bias!  $\hat{J} = \frac{\text{cov}(T, Y)}{\text{Var}(T)} = \frac{\text{cov}(T, Y)}{\text{Var}(T) + \int_U \text{cov}(T, U)}$

If we perform naive regression:

do these

$$\frac{\text{cov}(T, Y)}{\text{Var}(T)} = \frac{J \text{Var}(T) + \int_U \text{cov}(T, U)}{\text{Var}(T)} = \frac{J \text{Var}(T) + \int_U \cdot g_U \text{Var}(U)}{\text{Var}(T)} =$$

$$= J + \left( \frac{\int_U}{g_U} \right) - \text{biased term}$$

causal term

as a solution we can introduce INSTRUMENTAL variable

Matching is performed with the help of some balancing score  $\Rightarrow$  helps find good matches.

• **Sensitivity Analysis** - no guarantee that matching leads to balance on variables we didn't match for  
 $\hookrightarrow$  try to observe the hidden bias

Q5:

- #1 Does the conclusion change from statistically significant or not?
- #2 Does it change the direction of effect?

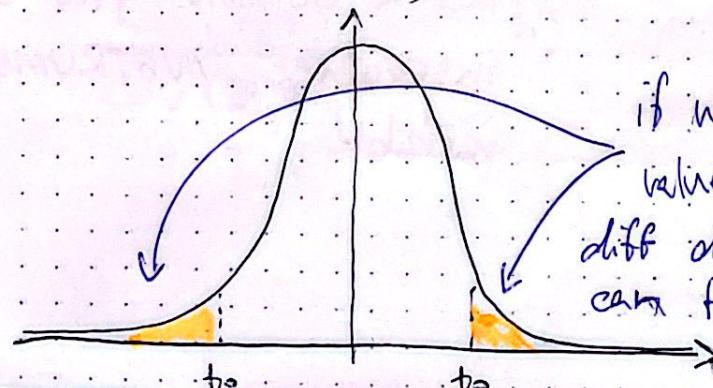
Types of Sensitivity Analysis:

- ① Quick Sensitivity Check
- ② Super Learning other potential confounders
  - $\hookrightarrow$  through cross validation we can see if we are sensitive to certain confounders
- ③ Deriving bounds on the causal statistical estimation

## Hypothesis testing

$$\frac{\text{ATE}}{\text{GATE}} = \frac{\text{signal}}{\text{noise}} \sim t\text{-distributed}$$

$$p\text{-value} = \Pr(|\frac{\text{signal}}{\text{noise}}| > t_0 | H_0)$$



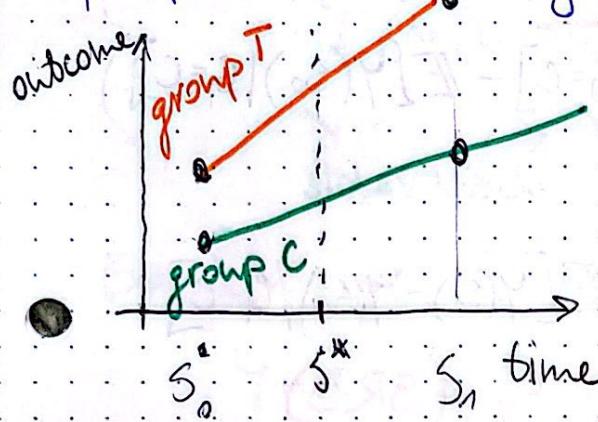
if we get here, then this value probs belongs to diff distribution and we can falsify  $H_0$

## Lecture 8

### DID - Difference in Difference

↳ treatment effect on outcome is estimated as the difference in changes over time between two groups

Group T - receives treatment  
Group C - control group



Both groups receive treatment only after time S\*

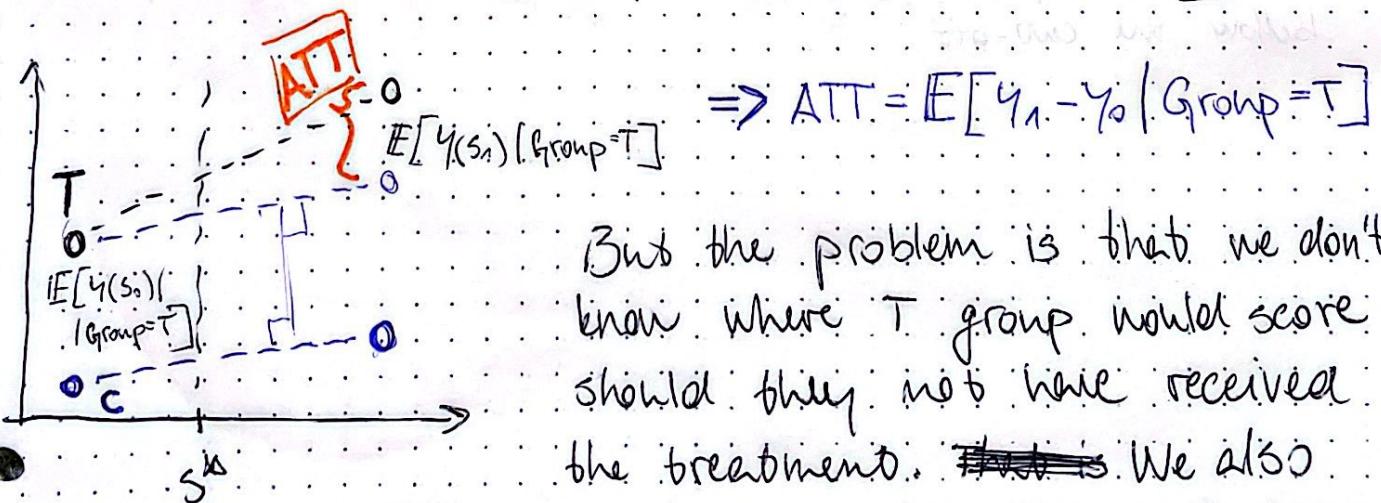
We only have measures at S<sub>0</sub> & S<sub>1</sub> and only know S\* occurred sometime in between

### Average Treatment Effect

$$\text{ATE} \Rightarrow E[Y_1 - Y_0] = E_x [E[Y|T=1, X] - E[Y|T=0, X]]$$

### Average Treatment Effect of the Treated

$$\begin{aligned} \text{ATT} \Rightarrow E[Y_1 - Y_0 | T=1] &= E_x [E[(Y_1 - Y_0) | T=1, X]] = \\ &= E_x [E[Y_1 | T=1, X] - E[Y_0 | T=0, X]] = \\ &= E_x [E[Y | T=1, X] - E[Y | T=0, X]] \end{aligned}$$



$$\Rightarrow \text{ATT} = E[Y_1 - Y_0 | \text{Group} = T]$$

But the problem is that we don't know where T group would score should they not have received the treatment. ~~This~~ We also need to consider that T & C can be fundamentally different, meaning that

there is a set difference between the two.  
Now we can approximate that there would be a movement across the parallel line from that of  $c$  applied to  $T$ . The difference between our actual and theoretical points is ATT.

$$\text{ATT} = \mathbb{E}[Y_1 - Y_0 | G=T] = (\underbrace{\mathbb{E}[Y(s_1) | G=T]}_{\text{observable}} - \mathbb{E}[Y(s_0) | G=T]) - (\mathbb{E}[Y(s_1) | G=c] - \mathbb{E}[Y(s_0) | G=c])$$

$$\mathbb{E}[Y(s_1) | G=T] = \mathbb{E}[Y(s_0) | G=T] + \mathbb{E}[Y(s_1) - Y(s_0) | G=c]$$

## Sharp Regression Discontinuity (SRD)

How can we identify causal effect despite positivity violations?

$W$  determines  $T$

$$T(W) = I\{W \geq c\} = \begin{cases} 1 & W \geq c \\ 0 & W < c \end{cases}, \quad \text{c - cutoff value}$$

which violates positivity

AIM: estimate causal effect:  $\bar{Y}_{SRD} = \mathbb{E}[Y_1 - Y_0 | W=c]$

The SRD looks at discontinuity in outcome at the cut-off

↳ the outcome just above the cut-off minus the outcome just below the cut-off

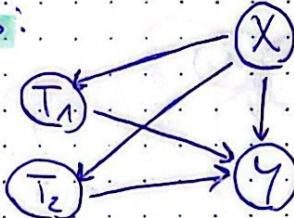
## Conditional Average Treatment Effect (CATE)

$$\mathbb{E}[Y_1 - Y_0 | X=x] = \text{ATE}$$

Suppose we have 2 treatments:

$$I_{12}^a = \left[ \mathbb{E}(Y | (T_1, T_2) = (1, 1), X) - \mathbb{E}(Y | (T_1, T_2) = (0, 1), X) \right]$$

$$- \left[ \mathbb{E}(Y | (T_1, T_2) = (1, 0), X) - \mathbb{E}(Y | (T_1, T_2) = (0, 0), X) \right]$$



Example with a linear model:

$$Y = \alpha_0 + \alpha_1 T_1 + \alpha_2 T_2 + \gamma T_1 T_2$$

$$\mathbb{E}(Y | T_1=1, T_2=1) = \alpha_0 + \alpha_1 + \alpha_2 + \gamma$$

$$\mathbb{E}(Y | T_1=1, T_2=0) = \alpha_0 + \alpha_1 \quad \mathbb{E}(Y | T_1=0, T_2=1) = \alpha_0 + \alpha_2$$

$$\mathbb{E}(Y | T_1=0, T_2=0) = \alpha_0$$

$$\text{ATE}_{T_1} (Y | T_2=1) = \alpha_1 + \gamma \quad \text{ATE}_{T_2} (Y | T_1=1) = \alpha_2 + \gamma$$

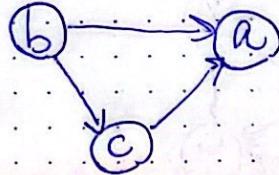
$$\text{ATE}_{T_1} (Y | T_2=0) = \alpha_1 \quad \text{ATE}_{T_2} (Y | T_1=0) = \alpha_2$$

$$I_{12}^a = \gamma = I_{21}^a \rightarrow \text{Symmetrie!}$$

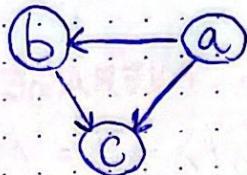
# DAG

\* acyclic!

#1



#1

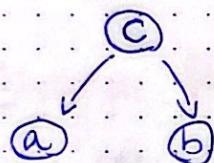


$$P(a) P(b|a) P(c|a, b) = \\ = P(a, b) P(c|a, b) = \underline{P(c|a, b)}$$

we calculate the same probability in different ways

$$P(b) P(c|b) P(a|b, c) = \\ = P(b, c) P(a|b, c) = \underline{P(a, b, c)}$$

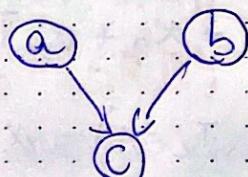
## D-separation



Fork

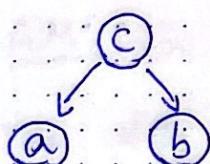


Chain



Collider

## Fork



$$P(a, b, c) = P(c) P(a|c) P(b|c)$$

$$P(a, b) = \sum_c P(a, b, c) = \sum_c P(a|c) P(b|c) P(c) =$$

$$\neq P(a) P(b) \Rightarrow a \not\perp\!\!\! \perp b | \emptyset$$

$$P(a, b | c) = \frac{P(a, b, c)}{P(c)} = \frac{P(a|c) P(b|c) P(c)}{P(c)} = \\ = P(a|c) P(b|c)$$

$\Rightarrow c$  blocks the path from  $a$  to  $b$

$a \perp\!\!\! \perp b | c \Rightarrow a$  and  $b$  are independent if conditioned on  $c$

## Chain



$$P(a, b, c) = P(a) P(b|a) P(c|b)$$

since summing over all b  
 $\Rightarrow P(c|a)$

$$\begin{aligned} P(a, b) &= \sum_b P(a) P(b|a) P(b|c) = P(a) \sum_b P(b|a) P(b|b) = \\ &= P(a) P(c|a) \neq P(a) P(b) \Rightarrow a \not\perp\!\!\! \perp b | \emptyset \end{aligned}$$

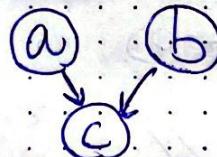
$$P(a, c|b) = \frac{P(a, c, b)}{P(b)} = \frac{P(a) P(b|a) P(c|b)}{P(b)} \xrightarrow{\text{bayes}}$$

$$= \frac{P(a) P(c|b)}{P(b)} \cdot \frac{P(a|b) P(b)}{P(a)} = P(b|a) = \frac{P(a|b) P(b)}{P(a)}$$

$$= P(c|b) P(a|b)$$

$\Rightarrow a \perp\!\!\! \perp b | e$ ; c blocks the path from a to b

## Collider



$$P(a, b, c) = P(a) P(b) P(c|a, b) = 1$$

$$\begin{aligned} P(a, b) &= \sum_c P(a) P(b) P(c|a, b) = P(a) P(b) \sum_c P(c|a, b) = \\ &= P(a) P(b) \Rightarrow a \perp\!\!\! \perp b | \emptyset \end{aligned}$$

$$\begin{aligned} P(a, b|c) &= \frac{P(a, b, c)}{P(c)} = \frac{P(a) P(b) P(c|a, b)}{P(c)} \neq P(a) P(b) \\ &\Rightarrow a \not\perp\!\!\! \perp b | \emptyset \end{aligned}$$

To conclude, only in case of **Collider**  $a \perp\!\!\!\perp b | \emptyset$   
for **Fork** and **Chain** we know that  $a \not\perp\!\!\!\perp b | \emptyset$

Fork	Chain	Collider
$a \not\perp\!\!\!\perp b   \emptyset$	$a \not\perp\!\!\!\perp b   \emptyset$	$a \perp\!\!\!\perp b   \emptyset$
$a \perp\!\!\!\perp b   c$	$a \perp\!\!\!\perp b   c$	$a \not\perp\!\!\!\perp b   c$

Exercise:  $B \rightarrow G \leftarrow F$   $P(B=1) = 0.9$   $P(F=1) = 0.9$

$$P(G=1 | B=1, F=1) = 0.8$$

$$P(G=0 | B=1, F=0) = 0.2$$

$$P(G=0 | B=0, F=1) = 0.2$$

$$P(G=1 | B=0, F=0) = 0.1$$

$$P(F=0) = 1 - P(F=1) = 0.1$$

$$\begin{aligned} P(F=0 | G=0) &= \frac{P(F=0, G=0)}{P(G=0)} = \frac{\sum_b P(F=0, B, G=0)}{P(G=0)} \\ &= \frac{[P(G=0 | F=0, B=1) + P(G=0 | F=0, B=0)]}{P(F=0, B=0)} \\ &= \frac{(0.8 + 0.9) \cdot 0.1 \cdot 0.1}{(0.9 + 0.2 + 0.8 + 0.2)} = // \\ &= \frac{P(G=0 | F=0) P(F=0)}{P(G=0)} \end{aligned}$$

$$\begin{aligned} P(G=0) &= \sum_{B,F} P(G=0 | B, F) = \sum_{B,F} P(G'' | B, F) P(B | F) P(F) = \\ &= \sum_{B,F} P(G=0 | B, F) P(B) P(F) = \\ &\cancel{*1} * P(G=0 | B=0, F=0) P(B=0) P(F=0) + \\ &\cancel{*2} + P(G=0 | B=1, F=0) P(B=1) P(F=0) + \\ &\cancel{*3} + P(G=0 | B=1, F=1) P(B=1) P(F=1) + \\ &\cancel{*4} + P(G=0 | B=0, F=1) P(B=0) P(F=1) = \end{aligned}$$

$$[0.9 \cdot 0.1 \cdot 0.1] + [0.8 \cdot 0.9 \cdot 0.1] + [0.2 \cdot 0.9 \cdot 0.9] + [0.2 \cdot 0.1 \cdot 0.9] = \\ = 0.315$$

$$P(G=0|F=0) = \sum_B P(G=0|F=0, B) P(B) = \\ = P(G=0|F=0, B=0) P(B=0) + P(G=0|F=0, B=1) P(B=1) = \\ = 0.9 \cdot 0.1 + 0.8 \cdot 0.9 = 0.81$$

---


$$= \frac{P(G=0|F=0) P(F=0)}{P(G=0)} = \frac{0.315 \cdot 0.1}{0.81} = 0.0389$$


---

path  $p$  is blocked by a set of nodes  $Z$   
iff:

1)  $p$  contains chain  $A \rightarrow B \rightarrow C$  or a fork

$A \leftarrow B \rightarrow C$  such that  $B$  is in  $Z$

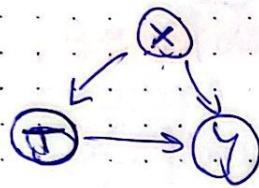
2)  $p$  contains collider  $A \rightarrow B \leftarrow C$  such that

collision node  $B$  is not in  $Z$  and no descendant  
of  $B$  is in  $Z$

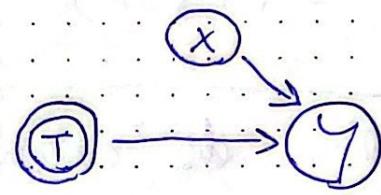
Observation (conditioning)vs Intervention

Learn to distinguish between: var  $T$  taking a value  $t$  naturally and cases where we 'fix'  $T=t$  by eliminating the letter  $\text{do}(T=t)$

$$P(Y=y | T=t)$$

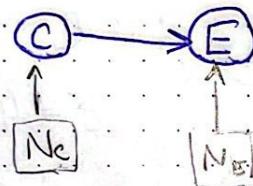


$$P(Y=y | \text{do}(T=t))$$

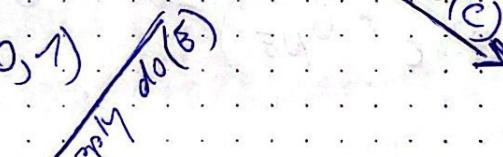
Example:

$$C := N_C$$

$$E := 4 \cdot C + N_E$$



$$N_C, N_E \sim N(0, 1)$$



#graphSurgery

$$P(E | \text{do}(C)) \neq P(E)$$

$\rightarrow = P(E | C)$   
since there are no other confounders

$$P(C | \text{do}(E)) = P(C)$$

$\neq P(C | E)$  hence we get this

Intervention example  $C = N_c \quad E = 4 \cdot C + N_E$

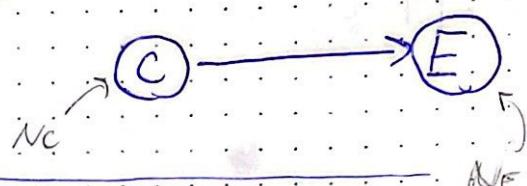
 $N_c, N_E \sim N(0, 1), N_c \perp\!\!\!\perp N_E$

$\text{Var}[\alpha X] = \alpha^2 \text{Var}[X], \quad 4C \sim N(0, 16)$

And since  $4C \perp\!\!\!\perp N_E$ , and the sum of two normally distributed random vars:

$E \sim N(\mu_{ac} + \mu_{NE}, G_{ac}^2 + G_{NE}^2)$

$E \sim N(0, 17)$



Now let's do-calculus

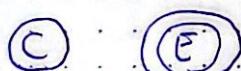


$do(C=2) \Rightarrow P(E|do(C=2)) : C=2$

$E \sim N(\mu_{ac} + \mu_{NE}, G_{ac}^2 + G_{NE}^2)$ 
 $E \sim N(8, 1)$ 
 $E = 4 \cdot C + N_E = 8 + N_E$ 
 $N_E \sim N(0, 1)$

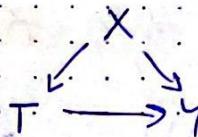
$= P(E|c=2)$

$P(c|do(E=2)) = P(c) = N(0, 1)$



cut we break  
the relation (= bond)

# Adjustment Formula



T - drug usage

X - sex

Y - recovery

Estimate effect of treatment:

$$P(Y=1 | do(T=1)) = P(Y=1 | do(T=0))$$

$P(Y=y | do(T=b))$  for  $\begin{matrix} X \\ \downarrow \\ T \rightarrow Y \end{matrix}$  is equal to

$P_m(Y=y | T=b)$  for modified graph  $\begin{matrix} X \\ \downarrow \\ T \rightarrow Y \end{matrix}$

!modified! (do calculus was applied)

$$P(Y=y | do(T=b)) = P_m(Y=y | T=b) =$$

$$= \sum_x P_m(Y=y | T=b, X=x) \cdot P_m(X=x | T=b) =$$

$$= \sum_x P_m(Y=y | T=b, X=x) P_m(X=x) \quad \text{since } T \text{ and } X \text{ aren't connected}$$

Exercise: Drug      No drug      T-drug

Men    81/87 (93%)    234/270 (87%)    X - sex

Women    192/263 (73%)    56/80 (69%)    Y - outcome

Total:    273/350 (78%)    289/350 (83%)

$$ATE = ACE = P(Y=1 | do(T=1)) - P(Y=1 | do(T=0)) = ?$$

$$P(Y=1 | do(T=1)) = \sum_x P_m(Y=1 | T=1, X=x) P(X=x) =$$

$$= P_m(Y=1 | T=1, X=1) P(X=1) + P_m(Y=1 | T=1, X=0) P(X=0) =$$

$$= 0.93 \cdot (87 + 270) / 700 + 0.73 \cdot \frac{(192 + 56)}{263 + 80} / 700 =$$

$$= 0.832$$

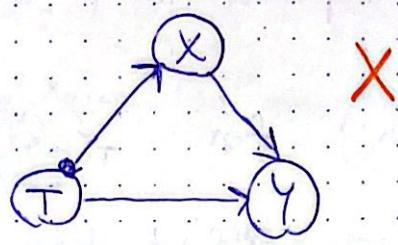
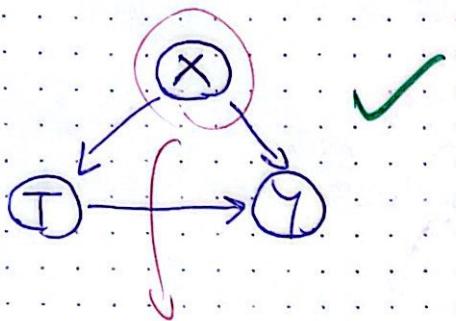
$$P(Y=1 | do(T=0)) = P_m(Y=1 | T=0, X=0) P(X=0) + P_m(Y=1 | T=0, X=1) P(X=1) =$$

$$= 0.69 (263 + 80) / 700 + 0.87 (87 + 270) / 700 = 0.7818$$

$$0.832 - 0.7818 = \boxed{0.0505}$$

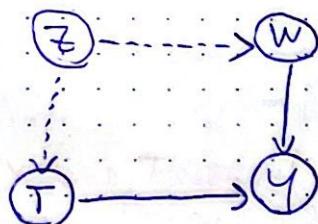
When to use this Backdoor approach?

$$P(Y=y | \text{do}(T=t)) = \sum_x P_m(Y=y | T=t, X=x) p_n(X)$$



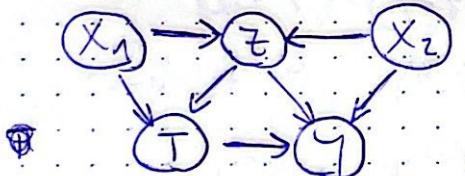
here  $X$  (=confounder) acts as a backdoor

We can use  $W$  as a backdoor!



- 1)  $W$  blocks  $T \leftarrow Z \rightarrow W \rightarrow Y$
- 2)  $W$  leaves directed  $P(T \rightarrow Y)$  uninterrupted
- 3)  $W$  is not a collider or descendants of  $T$

$$P(Y=y | \text{do}(T=t)) = \sum_w P_m(Y=y | T=t, W=w) p_n(W=w)$$



can we condition (use as backdoor) on  $Z$ ? It's a collider, which is a problem.

However, if we condition on  $Z|X_1$  or  $Z|X_2$  or  $Z|X_1X_2$  then we override and it's okay!

## Lecture 11

### Total Law of Probability:

$$P(X=1|Y=1) + P(X=0|Y=1) = 1$$

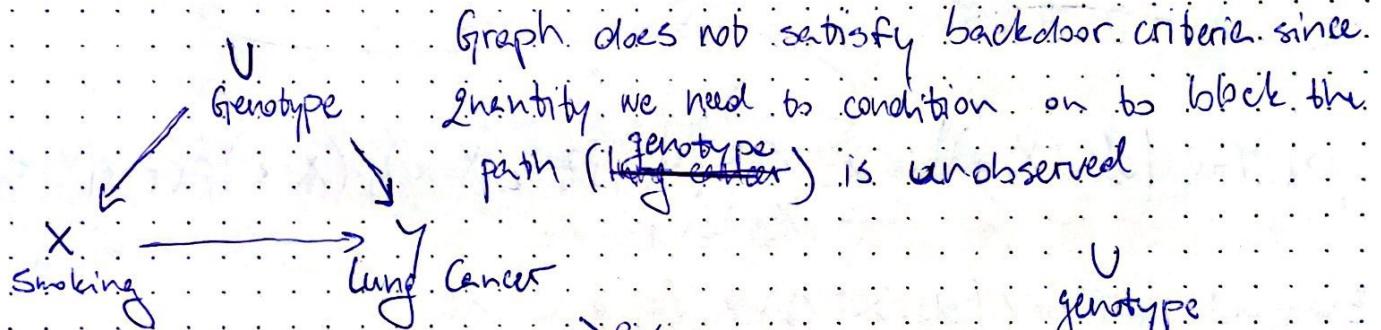
$$P(X=1|Y=1)P(Y=1) + P(X=0|Y=0)P(Y=0) = P(X \neq 1)$$

$$X \rightarrow Y \quad P(X|Y) \neq P(X)$$

$$X \rightarrow Y \rightarrow W \quad P(X|W) \neq P(X)$$

$$P(X, Y, W) = P(X)P(Y|X)P(W|Y)$$

### Pearl's Front-Door Adjustment



Graph still doesn't satisfy the backdoor criteria but we can determine the causal effect

$$P(Y=y | do(X=x))$$

1)  $X \rightarrow Z$  is identifiable, since no back path from  $X$  and  $Z$ .

$$P(Z=z | do(X=x)) = P(Z=z | X=x)$$

2)  $Z \rightarrow Y$  is identifiable, since backdoor from  $Z$  to  $Y$ :

$$Z \leftarrow X \rightarrow U \rightarrow Y$$

is blocked by conditioning on  $X$ :

$$P(Y=y | \text{do}(z=z)) = \sum_x P(Y=y | z=z, X=x) P(X=x)$$

By letting  $z$  be the value  $Z$  takes when setting  $X=x$ :

$$P(Y | \text{do}(X)) = P(Y | \text{do}(X), z) = P(Y | \text{do}(z=z))$$

$$\begin{aligned} P(Y=y | \text{do}(X=x)) &= \sum_z P(Y=y, z | \text{do}(X=x)) \\ &= \sum_z P(Y=y | z, \text{do}(X=x)) P(z | \text{do}(X=x)) \\ &= \sum_z P(Y=y | \text{do}(z=z)) P(z | \text{do}(X=x)) \end{aligned}$$

$$P(Y=y | \text{do}(X=x)) = \sum_z \sum_{x'} P(Y=y | z=z, X=x') P(X=x') P(z=z | X=x)$$

Example:  $P(Y=1 | \text{do}(X=1)) =$

$$= P(Y=1 | z=0, X=1) P(X=1) P(z=0 | X=1)$$

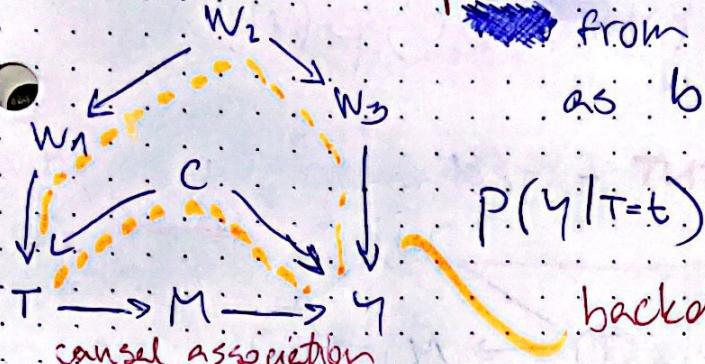
$$+ P(Y=1 | z=0, X=0) P(X=0) P(z=0 | X=1)$$

$$+ P(Y=1 | z=1, X=1) P(X=1) P(z=1 | X=1)$$

$$+ P(Y=1 | z=1, X=0) P(X=0) P(z=1 | X=1) =$$

# Backdoor Adjustment

- we want to block all paths from  $T$  to  $Y$  that can act as backdoors



not a backdoor path!

$$P(Y|T=t)$$

backdoor paths = non-causal association paths

we want to block these and keep only causal association paths

Intervene  
 $do(T=t)$

\* we can remove all incoming flows to  $T$

$$w_1 \leftarrow \textcircled{w_2} \rightarrow w_3 \quad P(Y|do(T=t))$$

$$T \rightarrow M \rightarrow Y$$

we can calculate this by conditioning on  $C$  and  $w_2$  from first graph!

this then blocks the back door 😊

## CRITERION

- A set of variables  $W$  satisfies the backdoor criterion relative to  $T$  and  $Y$  if:
  - $W$  blocks all backdoor paths from  $T$  to  $Y$
  - $W$  does not contain any descendants of  $T$

Then, we can identify causal effect of  $T$  on  $Y$ :

$$P(Y|do(T=t)) = \sum_w P(Y|t, w) p(w)$$

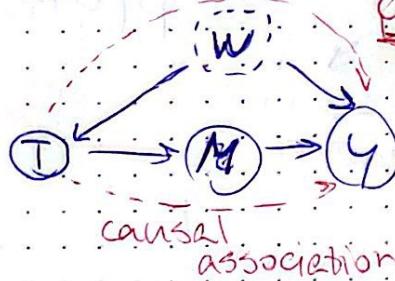
## ADJUSTMENT

for our example that would give us:

$$P(Y | do(T=t)) = \sum_w \sum_t P(Y(t, w, c)) P(w, c)$$

**FRONT DOOR ADJUSTMENT** - used when the confounder is unobservable

1) Identify causal effect of T on M



2) Identify causal effect of M on Y

3) Combine 1) & 2) to get causal effect of T on Y

1  $P(M|do(t)) = P(M|t)$

no backdoor paths  
from T to M

$\hookrightarrow$  cut Y is a collider

backdoor adjustment

2  $P(Y|do(m)) = \sum_t P(Y(t, m)) P(t)$

there is a backdoor path  
from M to Y  
 $M \leftarrow T \leftarrow W \rightarrow Y$

} we can condition  
on T

3  $P(Y|do(t)) = \sum_m p(m|do(t)) P(Y|do(m)) =$

## ADJUSTMENT

$$= \sum_m p(m|t) \sum_{t'} p(Y(t', m) | t) P(t')$$

chaining them together and summing

over middle var ( $= m$ )

## CRITERION

- 1) M completely mediates effect of T on Y if go through M
- 2) There is no unblocked path from T to M
- 3) All backdoor paths from M to Y are blocked by T

# MARKOV BLANKET AND BOUNDARY

**Markov blanket** of random variable  $Y$

- $Y$  in a random variable set  $S = \{X_1, \dots, X_n\}$  is any subset  $S_1$  of  $S$ , conditioned on which other variables are independent with  $Y$ :  $Y \perp\!\!\!\perp S \setminus S_1 \mid S_1$

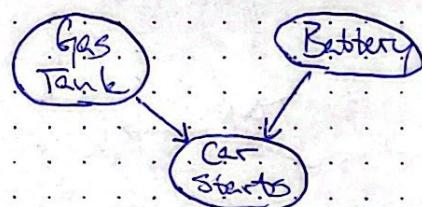
It means  $S_1$  contains all the information needed to infer  $Y$  and the variables  $S \setminus S_1$  are redundant.

In general, **Markov blanket** is  $\Rightarrow$  not unique. Any set  $S$  that contains **Markov blanket** is a **Markov blanket** itself.

**Markov Boundary** is a **Markov blanket** none of whose subsets are **Markov blankets** themselves.

**DISTINGUISHING CAUSAL STRUCTURES:** V-structures

Gas tank  $\perp\!\!\!\perp$  battery



Gas tank  $\not\perp\!\!\!\perp$  battery | car starts = 0

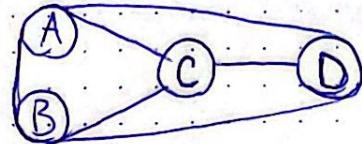
**Markov Equivalence Class (MEC)** - graphs  $G$  and  $G'$  belong to the same equivalence class iff each conditional independence implied by  $G$  is also implied by  $G'$  and vice versa.

$A \perp\!\!\!\perp C \mid B$ :

- 1)  $A \rightarrow B \rightarrow C$
- 2)  $A \rightarrow B \leftarrow C$
- 3)  $A \leftarrow B \leftarrow C$

# Peter-Clark (PC) Algorithm

1) Start with complete graph:

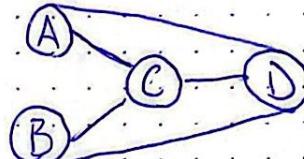


2) Zeroth order Conditional Independence;  $A \perp\!\!\!\perp B$

We perform statistical tests

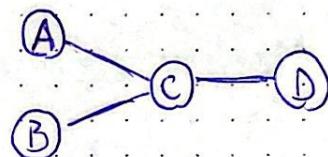
to see if they are independent or not

if not ↗

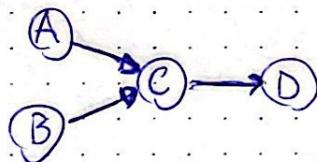


3) First order CI;  $A \perp\!\!\!\perp D | e \wedge B \perp\!\!\!\perp D | e$ :

4) No higher CI order observed :)

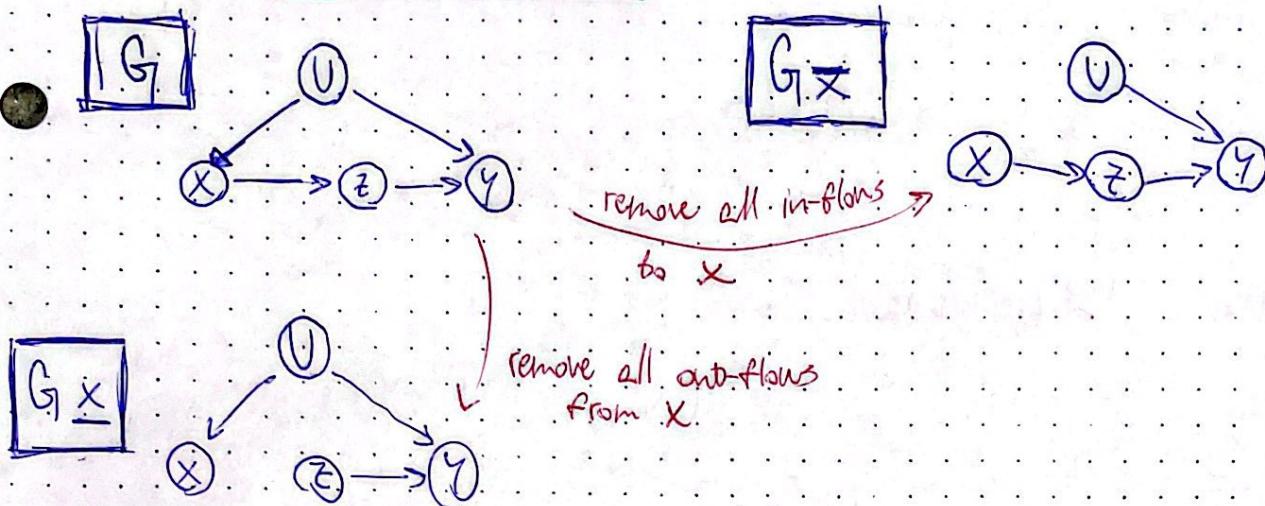


5) Sometimes we can't put the arrows, but here we have a collider!



Lecture 13

## DO-CALCULUS



Rule 1: (insertion/deletion of observations)

$$P(Y | \text{do}(X=x), Z, W) = P(Y | \text{do}(X=x), W) \text{ if } (Y \perp\!\!\!\perp Z) | X, W \text{ in } G \setminus X$$

↳ if  $Y$  and  $Z$  are d-separated by  $X, W$  in a graph where incoming edges in  $X$  have been removed

---

d-separation - two sets of nodes  $X$  and  $Y$  are d-separated by a set of nodes  $Z$  if all of the paths between  $X$  and  $Y$  are blocked by  $Z$ .

---

Rule 2: (action/observation exchange)

$$P(Y | \text{do}(X=x), \text{do}(Z=z), W) = P(Y | \text{do}(X=x), z, W) \text{ if } (Y \perp\!\!\!\perp Z) | X, W$$

↳ if  $Y$  and  $Z$  are d-separated by  $X, W$  in a graph in  $G \setminus Z$  where incoming edges in  $X$  and outgoing edges from  $Z$  have been removed.

Rule 3: Insertion/deletion of actions

$$p(Y \mid do(X=x), do(Z=z), W) = p(Y \mid do(X=x), W) \text{ if } Y \perp\!\!\!\perp Z \mid X, W$$

in  $G_{XZ|W}$

$Z(W)$  is set of  $Z$ -nodes that are not ancestors of any  $W$ -node in  $G_X$

## Optimal Adjustment Sets

$$\underbrace{pa_G(\text{en}_G(X \rightarrow Y)) \setminus (\text{en}_G(X \rightarrow Y) \cup \{X\})}$$

tells us to condition on all parents of  $Y$  excluding those on path  $X \rightarrow Y$

$\Rightarrow$  THEOREM

### LECTURE 14

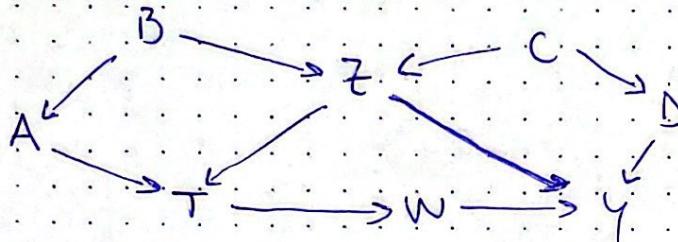
$$do(T=g(Z)) \text{ where: } g(Z) = \begin{cases} 1 & \text{when } Z > z \\ 0 & \text{otherwise} \end{cases}$$

result of such policy is:  $p(Y=y \mid do(T=g(Z))) = ?$

1) introduce "Z specific effect" of T on Y

$$p(Y=y \mid do(T=t), Z=z) = \sum_s p(Y=y \mid T=t, S=s, Z=z) P(S=s \mid Z=z)$$

## Lecture 14



We need to condition on  $Z$ , which is a collider

Q: What's the causal effect of  $T$  on  $Y$ ?

$$P(Y=y | do(T=t)) = \sum_{z,a} P(Y=y | T=t, A=a, Z=z) P(Z=z, A=a)$$

BACKDOOR ADJUSTMENT

need to block the spurious path by condition on  $A$  (parents of  $T$ )

Notes: Collider  $\rightarrow$  closed path  $\xrightarrow{\text{condition}}$   $\Rightarrow$  opened  
Other  $\rightarrow$  opened path  $\xrightarrow{\text{condition}}$   $\Rightarrow$  closed

So here we must obviously condition on  $Z$ , but this is not enough coz it gets OPENED. Hence, must condition on another node along that path  $\Rightarrow (A, B, C, D)$

Q: What's  $c$ -specific effect of  $T$  on  $Y$ ?

$\hookrightarrow$  need to condition on  $C$

$$P(Y=y | do(T=t), C=c) = \sum_z P(Y=y | T=t, C=c, Z=z) P(Z=z | C=c)$$

Q: What's  $z$ -specific effect of  $T$  on  $Y$ ?

$$P(Y=y | do(T=t), Z=z) = \sum_c P(Y=y | T=t, Z=z, C=c) P(C=c | Z=z)$$

Q: What is  $\mathbb{Z}$ -dependent causal effect  $T$  on  $Y$  under strategy:  $g(z) = \begin{cases} 0 & ; z \leq 2 \\ 1 & ; z > 2 \end{cases}; z \in \{1, 2, 3, 4, 5\}$

$$\begin{aligned} P(Y=y | do(T=g(z))) &= \sum_z P(Y=y | do(T=g(z)), z=z) P(z=z) = \\ &= P(Y=y | do(T=0), z=1) P(z=1) \\ &+ P(Y=y | do(T=0), z=2) P(z=2) \\ &+ P(Y=y | do(T=1), z=3) P(z=3) \\ &+ P(Y=y | do(T=1), z=4) P(z=4) \\ &+ P(Y=y | do(T=1), z=5) P(z=5) \end{aligned}$$

|| we can further expand by  
using result from  $\mathbb{Z}$ -specific  
effects (prev. page)

$$\begin{aligned} * (P(Y=y | do(T=t), z=z) &= \sum_{c=c} P(Y=y | T=t, z=z, C=c) \\ &\quad P(C=c | z=z)) \\ \Rightarrow \sum_c P(Y=y | T=0, z=1) P(z=1) P(C=c | z=1) + \\ &+ \sum_c P(Y=y | T=0, z=2) P(z=2) P(C=c | z=2) + \\ &+ \sum_c P(Y=y | T=1, z=3) P(z=3) P(C=c | z=3) + \\ &+ \sum_c P(Y=y | T=1, z=4) P(z=4) P(C=c | z=4) + \\ &+ \sum_c P(Y=y | T=1, z=5) P(z=5) P(C=c | z=5) \end{aligned}$$