

(1)

Question:

To build: (1) Bayesian Networks - Strictly Statistical  
 (2) Causality - an intervention on top of a Bayesian network

We'll think of BNs causally (though you don't have to)

- We want to model a joint distribution.
- We want to assume some independencies to simplify it

To model  $P(X_1, \dots, X_n)$  on  $n$  variables

To assume rules:  $\{ P(X_i | X_2, X_3) = P(X_i | X_2), \dots \}$   
 of the form:

Solution: Chain rule for conditional independence

$$P(X_1, \dots, X_n) = P(X_1 | X_2, \dots, X_n) P(X_2 | X_3, \dots, X_n) \dots P(X_n) \\ = \prod_{i=1}^n P(X_i | X_{i+1}, \dots, X_n)$$

example

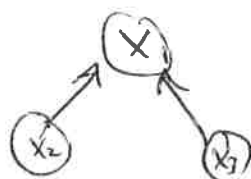
$P(X_1, X_2, X_3)$  where  $P(X_1 | X_2, X_3) = P(X_1 | X_2)$

$$P(X_1, X_2, X_3) = P(X_1 | X_2, X_3) P(X_2 | X_3) P(X_3) \\ = P(X_1 | X_2) P(X_2 | X_3) P(X_3)$$

consider a visual for this.

- (1) Each term is a node
- (2) Each variable conditioned on is a parent.

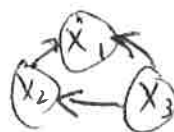
$$P(X_1 | X_2, X_3) \Rightarrow$$



\* No Causal Interpretation \*

\* Yet \*

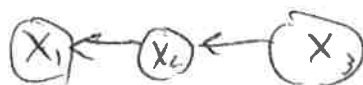
$$P(X_1 | X_2, X_3) P(X_2 | X_3) P(X_3) \Rightarrow$$



use the assumption now...

(2)

$$P(X_1 | X_2) P(X_2 | X_3) P(X_3) \Rightarrow$$



\* Our independence assumption implies we cut an edge! \*

A couple notes: (1) Complete graph  $\Rightarrow$  no assumptions.

(2) Variable ordering is arbitrary! So any complete graph can represent any joint distribution.

(Statistical representation, not a causal representation. Causality comes from extra structures: interventions)

(3)

$P(X_1, X_2, X_3)$  where  $P(X_1 | X_2, X_3) = P(X_1 | X_2)$  can be represented by



but not



"V-structure" or "collider".

cutting the

$1 \leftrightarrow 3$  edge above implies

$$P(X_1 | X_2, X_3) = P(X_1 | X_2)$$

but this graph doesn't imply that. Instead:

$$P(X_1, X_3) = P(X_1) P(X_3)$$

$\Rightarrow$  Graphs with the same skeleton can represent the same joint distributions only if they have the same V-structures

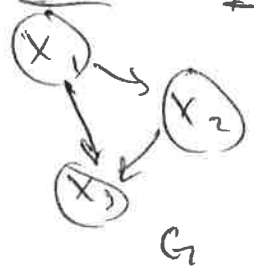
- What we're talking about up to this point is models for joint distributions. These summarize the statistical properties of a system.
- If we only care about statistics, we can use any graph that is "observationally equivalent"; i.e. it has the same colliders and skeleton, that can express the dependencies in the joint distribution.

———— END Bayesian Networks, BEGIN Causality ————

Key assumption: One of these graphs is the right one to model the set of possible interventions.

Intervention: In real life, disrupt the system to fix the value of a node, independently from what the value would have been if you hadn't disrupted it.

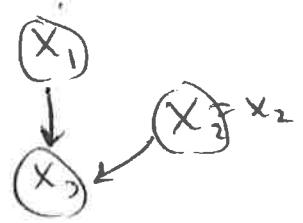
model for observed system



now we intervene on  $X_2$



model for intervention



Steps: (1) Cut all edges in  $G$  going into the node you're intervening on, producing a new graph  $G_{\bar{X}_2}$

notation:  $\left\{ \begin{array}{l} \bar{X}_2 \Rightarrow \text{causes of } X_2 \text{ deleted} \\ \underline{X}_2 \Rightarrow \text{effects of } X_2 \text{ deleted} \end{array} \right.$

(2) Set  $X_2 = x_2$ , the value you choose for your intervention

(3) The rest happens as normal.

observational system

interventional system (4)

$$P(x_3, x_2, x_1) =$$

$$P(x_3 | x_2, x_1) P(x_2 | x_1) P(x_1)$$

$$\Rightarrow P(x_3, x_1 | x_2) = \frac{P(x_3 | x_2, x_1) P(x_2 | x_1) P(x_1)}{P(x_2)}$$

$$P(x_3, x_1 | do(x_2 = x_2)) =$$

$$P(x_3 | x_1, x_2 = x_2) P(x_1)$$

$\Rightarrow$  different from conditioning.

\* but still estimable from the old joint + marginals!!!\*

Note

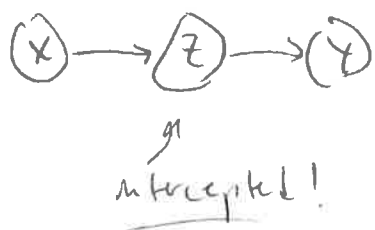
$$P(x_3, x_1 | do(x_2 = x_2)) = \frac{P(x_3, x_2, x_1)}{P(x_2 | x_1)} \quad \text{---} \quad \cancel{P(x_1, \dots, x_n)}$$

In general,

$$P(x_1, \dots, x_n | do(x_i = x_i)) = \frac{P(x_1, \dots, x_n)}{P(x_i | Pa(x_i))}$$

Bayesian Network together with this notion of intervention  
are called "Causal Bayesian Networks".

Does X cause Y? What does cause mean?



yes.  
 (Indirectly,  
 through  
 mechanism Z  
 X, Y correlated)

(Randomized control!)

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

$$P(Y|do(X=HT)) \neq P(Y|do(X=...))$$

assume stability



no  
 ("confounding"  
 lemonade example)  
 correlated X & Y

see (18)



no  
 (not even correlated)

Regression Example for which

Direct vs. Indirect Effects, mediation

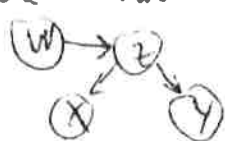
Direct effect: X is parent of Y

Indirect effect: X upstream from Y. (← there is a directed path from X to Y, X & Pa(Y))

- Are there real "direct" effects?  
mediators, suppressing mediators.

- We often (legitimately) ignore mechanisms.

- Okay as long as you don't suppress confounders. Then you have to introduce arcs, e.g.

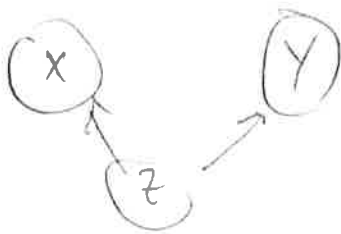


→



→





example: what does confounding look like?

(f.s)

$$Y = \beta_{Yz} z + \epsilon_Y$$

$$X = \beta_{Xz} z + \epsilon_X$$

$$z = \epsilon_z$$

~~$$Y = \epsilon_z \beta_{Yz} + \epsilon_Y$$~~

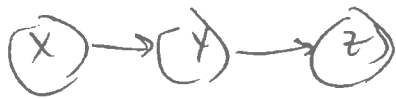
$$\begin{aligned} \beta_{YX} &= \frac{\text{cov}(X, Y)}{\sigma_X^2} = \frac{\text{cov}(\beta_{Yz} z, \beta_{Xz} z)}{\sigma_X^2} \\ &= \beta_{Yz} \beta_{Xz} \frac{\sigma_z^2}{\sigma_X^2} \end{aligned}$$

nonzero (unconditional) regression coefficient, with no directed causal path from X to Y! (or vice versa)

- Fix  $z$ , then  $\sigma_z^2 = 0$  - and  $\beta_{YX} \rightarrow 0$   
 $\rightarrow$  controlling ~~breaks~~ confounding
- set either  $\beta_{Yz}$  or  $\beta_{Xz}$  (or both) to 0, and you break confounding  
 $\Rightarrow$  due to common causes.

How do paths ~~lead to~~ dependence?

(7)



$$P(X, Z) = \sum_Y P(Z|Y) \underbrace{P(Y|X)}_{\text{coupled}} P(X)$$



$$P(X, Y) = \sum_Z P(Z|X, Y) \underbrace{P(X|Z)}_{\text{marginalized}} P(X) P(Y)$$

$$= P(X, Y)$$



$$P(X, Y) = \sum_{A, B, Z} \cancel{P(A|X)} P(X) \cancel{P(A|X)} P(Z|A, B) P(B|Y) P(Y)$$

$$= \sum_{A, B} P(A|X) P(A|X) P(B|Y) P(Y)$$

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



$$P(X, Y) = \sum_Z P(X|Z) \underbrace{P(Z) P(Y|Z)}_{\text{coupled}}$$

Regressions just give conditional expectations (when properly specified, and using MSE loss), so if  $Y \perp X | Z$  (it d-separates  $X$  and  $Y$ ), then

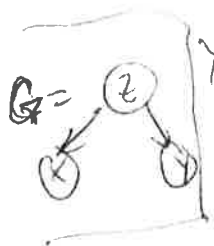
$$Y(X, z) = \sum_y y P(Y|X, z) = \sum_y y \frac{P(Y|z) P(X|z) P(z)}{P(X, z)}$$

$$= \sum_y y P(Y|z) \frac{P(X, z)}{P(z)} \frac{P(z)}{P(X, z)} = \sum_y y P(Y|z) P(Y)$$

the regression estimate is independent of  $X$  when you include  $Z$ ! (but not in general, since  $P(Y|X, z) = P(Y|z) \not\Rightarrow P(Y|X) = P(Y)$  i.e. conditional independence does not imply independence)

We can look at this in the context of confounding, and see

$$P(X, Y, Z) = P(Y|Z) P(X|Z) P(Z) \quad (\text{from factorization implied by } G)$$



$$Y(X, z) = \sum_y y \frac{P(X, Y, z)}{P(X, z) P(z)} = \sum_y y \frac{P(Y|z) P(X|z) P(z)}{P(X, z) P(z)}$$

$$= \sum_y y P(Y|z) \frac{P(X, z)}{P(z)} \frac{P(z)}{P(X, z)} = \sum_y y P(Y|z)$$

independent of  $X$ ! , but without  $Z$ ,

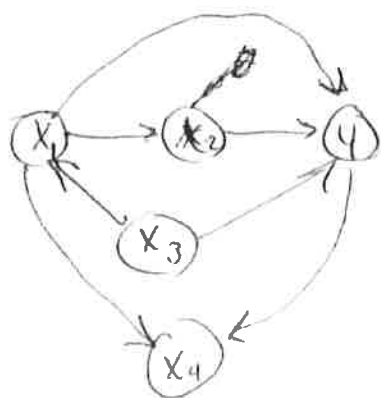
$$Y(X) = \sum_y y \frac{P(Y|X)}{P(X)} = \sum_{z, y} y \frac{P(Y, z|X)}{P(X)} = \sum_{z, y} y \frac{P(X, Y, z)}{P(X)}$$

$$= \sum_{z, y} y \frac{P(Y|z) P(X|z) P(z)}{P(X)}$$

so  $X$  and  $Y$  are coupled by the sum over  $z$ ! (wait factor in general)



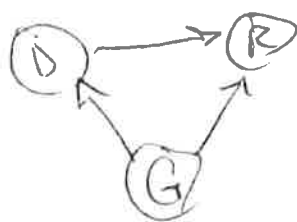
Wanted: Total effect  
 How can we block all the bad paths? (8)  
 which paths should we block? which are good?



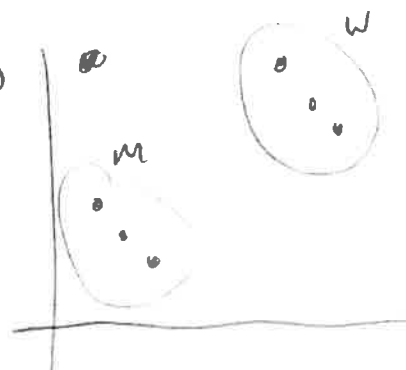
bad  $X_1 \leftarrow X_3 \rightarrow Y$  spurious not downstream!  
 good  $X_1 \rightarrow Y$  direct effect  
 good  $X_1 \rightarrow X_2 \rightarrow Y$  mediated effect.

Knowing when to condition resolves Paradoxes!  
 A drug ~~study~~ study (Simpson's Paradox)  
 Examples ~~Simpson's Paradox~~ collider bias (selection effect... admissions example)

- Positive effect overall
- Negative effect in ~~pop.~~ men / women
- Condition or not?!



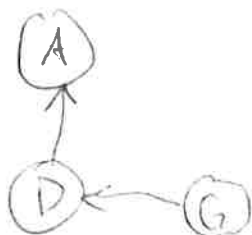
$P(R)$



Pose

~~Conclusion~~

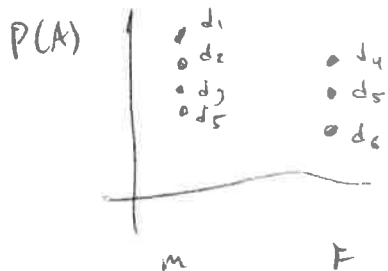
- Women comply more, have naturally higher recovery rates
- Men comply less, have naturally lower rates.



$$P(A|G=F) < P(A|G=M)$$

is there ~~guarantee~~  $\rightarrow$  discrimination?

• Because of mediation, ~~there~~ there is no discrimination.



You really need no spurious paths to establish causation.

Ignorability

- $\rightarrow$  Guaranteed in randomized trial
- $\rightarrow$  Guaranteed when they're "blocked"

$$(Y^0, Y^1) \perp D$$

$$= \begin{pmatrix} \text{assignment} \\ \text{ignorable} \end{pmatrix} = \text{an RCT}$$

$$(Y^0, Y^1) \perp D \mid S$$

$$= \text{conditionally ignorable} = \text{as good as RCT.}$$

$S \Rightarrow Z$ , improve precision of effect measurement.

weaker version:

(regression example)  $\rightarrow$  controlling breaks dependence b/t  $D$  and  $Y^0, Y^1$  gives unbiased estimator for  $\delta$

$$Y^0 \perp D, \text{ but not } (Y^1, Y^0) \perp D : \text{ still identify}$$

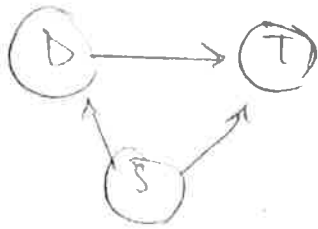
ATT

vs use an ATE

Matching + Exposure modeling : Next time!

(10)

$$P(D|S) = ?$$



same  $P(D|S) \Rightarrow$  effectively fixed  $S$ !