

# Causality in Artificial Intelligence

Sargur N. Srihari

University at Buffalo, The State University of New York  
USA

# Topics in Causality

- Importance of Causality in Science
- Importance of Causality in Machine Learning
- Probabilistic Models vs Causal models
- Correlation vs Causality
- Interventions to discover Correlation vs. Causal
- Identifying Causality
- Causal Mechanisms

# Causality

- Causality:
  - Relation between an event (the *cause*) and a second event (the effect), where the second is understood to be a consequence of the first
  - Examples
    - Rain causes Mud
    - Rain causes Umbrella to be opened
    - Smoking causes Cancer
    - Altitude lowers Temperature

# Causality in Philosophy

- Democritus 460-370BC, father of modern science
  - “I would rather discover one causal law than gain the kingdom of Persia”
- Karma in Indian philosophy (Sanatana Dharma)
  - A person's actions causes certain effects in current and future life either positively or negatively

# Causality in Medicine

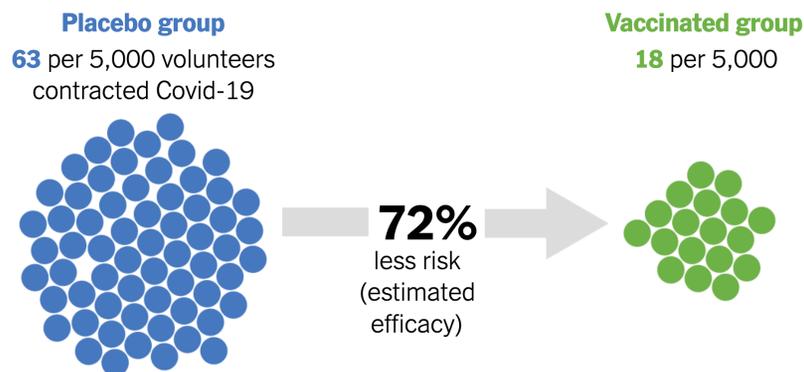
- Medical treatment
  - Possible effects of a medicine
    - Right treatment saves lives
- Vitamin D and Arthritis
  - Correlation versus Causation
  - Need for Randomized Correlation Test
- Vaccination and Covid-19

# Vaccination and Covid-19

- Newspaper report:
  - There is correlation between Pfizer vaccine and fewer Covid-19 deaths
- Causation vs Correlation
  - Correlation may be due to unrelated facts, e.g.,
  - Those countries:
    - Have a younger population
    - Have higher average temperatures

# Randomized Correlation Test

- Test for causality between J&J vaccine and lower incidence of Covid-19
- Set up two groups with and without vaccination
  - Making sure that there are no other intervening hidden variables such as age, temperature of country, etc.
- Then expose both groups to Covid-19

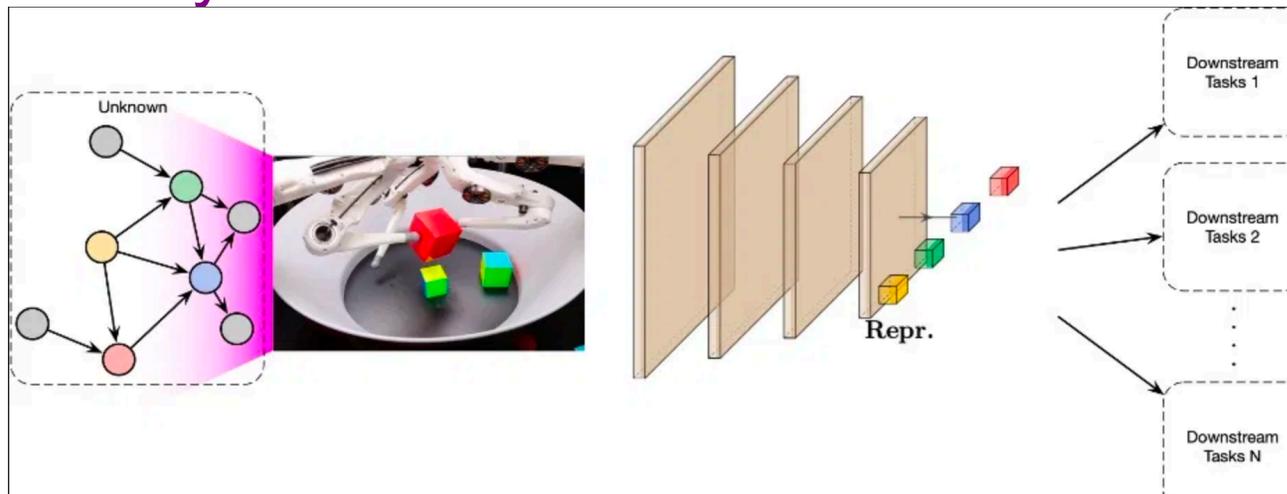


# Determining Causal Relationships

- Best way: Randomized controlled experiments
  - May be too difficult or immoral to perform
- Want to rely on “observational” data to infer causal relationships
- Current statistical methods are good at determining correlation
  - Correlation hints at causality

# Causality and Machine Learning

- Shortcomings of machine learning are:
  - Transferability to other distributions
    - Domain adaptation, life-long learning
  - Explainability of predictions
  - Lack of understanding of cause-effect relationships
    - A hallmark of human cognition
    - Allow interrogating a parsimonious representation
    - Causality will enable ML without much training



# Causality in Computer Vision

- Baseball Player
- Causality:
  - It is the player's arm that is *causing* the bat's movement and not the other way around.
  - The bat is *causing* the sudden change in the ball's direction
- Counterfactual:
  - What would happen if the ball flew a bit higher and didn't hit the bat?

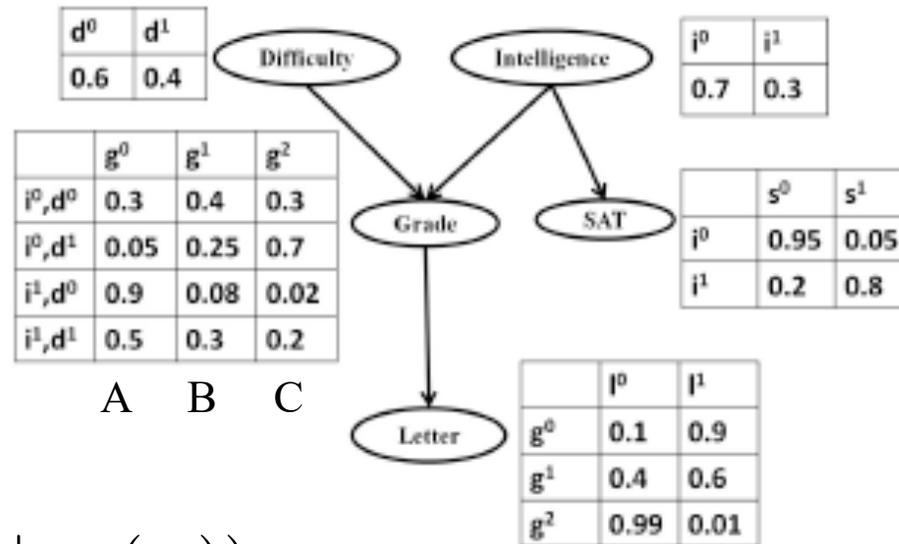


# Probabilistic Models

- For probabilistic queries it does not matter whether model is causal
  - We only care about encoding right distribution

# Joint distribution in student BN

$\mathbf{x} = \{x_1, \dots, x_5\}$   
 $x_1 = D$   
 $x_2 = I$   
 $x_3 = G$   
 $x_4 = S$   
 $x_5 = L$



- CPDs:  $P(x_i \mid pa(x_i))$

- Joint Distribution:

$$P(\mathbf{x}) = P(x_1, x_2, x_3, x_4, x_5)$$

$$P(\mathbf{x}) = \prod_{i=1}^5 P(x_i \mid pa(x_i))$$

$$P(D, I, G, S, L) = P(D)P(I)P(G \mid D, I)P(S \mid I)P(L \mid G)$$

# Extended Student BN

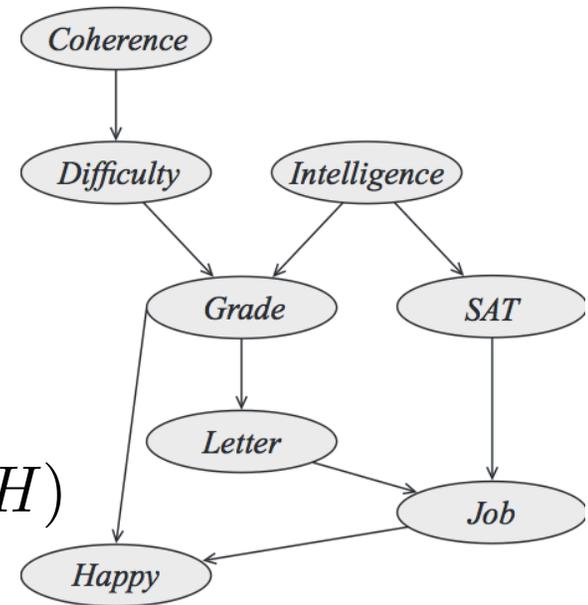
- We wish to infer  $P(J)$

$$P(J) = \sum_H \sum_L \sum_S \sum_G \sum_I \sum_D \sum_C P(C, D, I, G, S, L, J, H)$$

- By chain rule:

$$P(C, D, I, G, S, L, J, H) = P(C)P(D|C)P(I)P(G|I, D)P(S|I)P(L|G)P(J|L)P(H|G, J)$$

– Which is a Sum of Product of factors

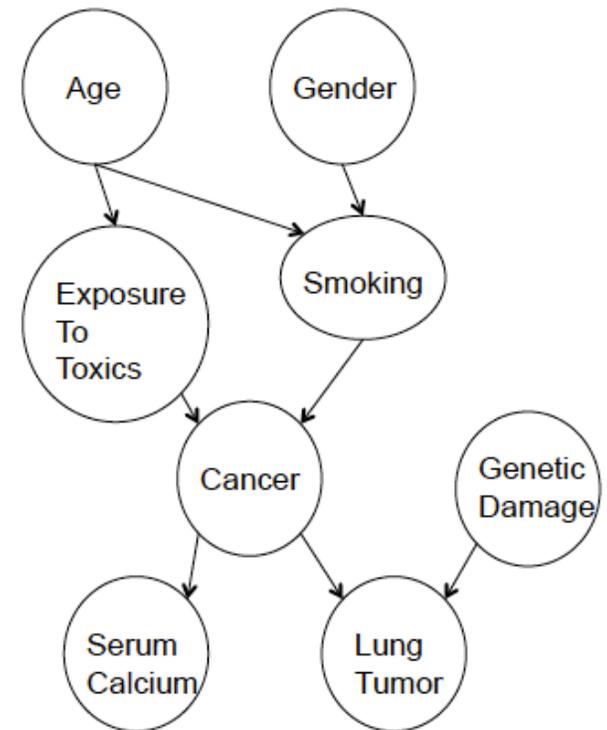


# BN is not necessarily Causal

- Bayesian Net is only an efficient representation of a joint distribution in terms of conditional distributions
- Several different BNs can represent the same distribution— equivalence class
- CW: a good BN should correspond to causality
  - An edge  $X \rightarrow Y$  suggests that  $X$  causes  $Y$  either directly or indirectly
  - Tend to be sparser and more natural

# Value of imposing Causal Semantics

- A type of reasoning in which causal interpretation of the network is critical
  - Will an intervention where prevent smoking in all public spacs is likely to decrease the frequency of lung cancer?
    - To answer needs to understand causal relationship between variables

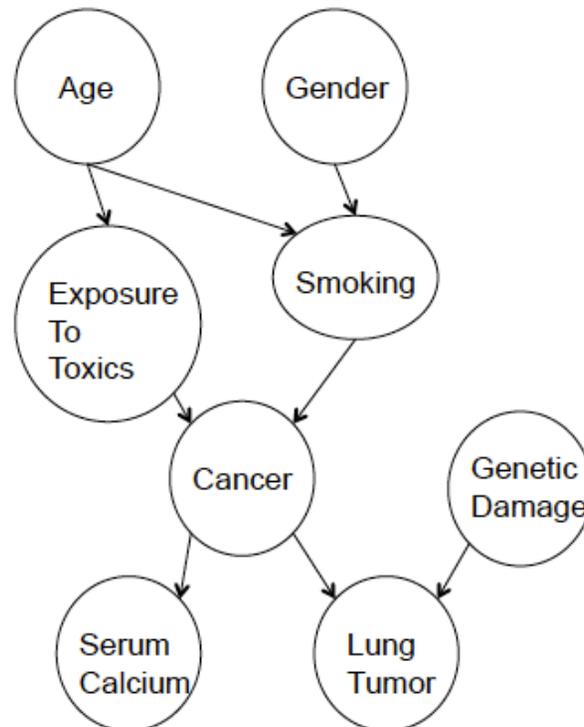


# Interpreting a Causal Model

- Edges have significance
- Interpretation distinguishes between two networks  $X \rightarrow Y$  and  $Y \rightarrow X$  which are equivalent probabilistic models

# Causal Models and Intervention

- Causal models matter when we want an intervention
  - We do not simply observe variable values but can take actions that can manipulate these values



# Effect of Actions

- Actions can affect the world in many ways
- Even a single action can have multiple effects
- Our goal is to isolate the specific issue of understanding causal relationship between variables

# Conditioning and Intervention

- We consider ideal interventions where only a single variable is affected
- Ideal intervention:  $do(Z:=z)$  abbreviated  $do(z)$ 
  - Which forces variable to take the value  $z$  and have no other immediate effect
  - We are interested in queries of the form  $P(Y | do(z))$  or more generally  $P(Y | do(z), X=x)$
  - They are called Intervention queries

# Intervention Queries

- A new type of conditioning of an event of the form  $do(Z := z)$
- Corresponds to case where an agent directly manipulated the world to set  $Z$  to take value  $z$  with probability 1
- We are interested in answering queries of the form  $P(Y|do(z), X=x)$ 
  - Set the variables in  $Z$  to  $z$ , observe values of  $X$  and wish to find the distribution of  $Y$

# Examples of Intervention Queries

## • Diagnosis and Treatment

- If patient takes meds will patient get well?  $P(H \mid do(M=m^1))$ 
  - Where  $H=Health$ . Which is different from  $P(H \mid m^1)$ 
    - » Patients taking meds on their own are healthier

## • Marketing

- If we lower Hamburger price will ketchup sell more?
  - Not a standard observational query but one in which we intervene in the model and thereby change its behavior

## • Policy Making

- If we lower interest rates, withat give rise to inflation?

## • Scientific Discovery

- Does smoking cause cancer?
  - When we formalize it, this is an intervention query “If we force someone to smoke will they be more likely to get cancer?”

# Counterfactual Queries

- A different type of Causal Query
- We already have some information about the true state of the world and want to inquire about the state of the world had we been able to intervene
  - And set the values of certain variables
  - E.g., “Would the US have joined World War II had it not been for the attack on Pearl Harbor?”

# Counterfactual Query Example

- Legal liability case
  - “Did the driver’s intoxicated state cause the accident?”
    - i.e., would accident have happened had driver not been drunk?
- Treatment and Diagnosis
  - “We are faced with a car that does not start, but the lights work. Will replacing the battery make the car start?”

# Answers to queries

- Answer to an intervention query

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

is quite different from the answer to the corresponding probabilistic query

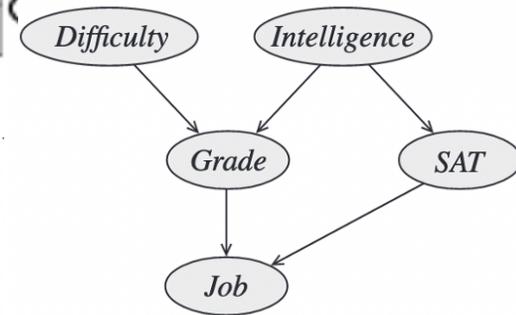
$$P(Y | z, X=x)$$

- Example is given next

# Intervention discovers causality

- Several correlations exist in Student BN. Conditioning on *Grade* of A increases probabilities of high *I*, high *S*, good *J*

	$d^0$	$d^1$	
	0.6	0.4	
	A	B	C
	$g^1$	$g^2$	$g^3$
$i^0, d^0$	0.3	0.4	0.3
$i^0, d^1$	0.05	0.25	0.7
$i^1, d^0$	0.9	0.08	0.02
$i^1, d^1$	0.5	0.3	0.2



$i^0$	$i^1$
0.7	0.3

	$s^0$	$s^1$
$i^0$	0.95	0.05
$i^1$	0.2	0.8

$P(i^1|g^1)$   $P(s^1|g^1)$ ,  $P(j^1|g^1)$   
 can be obtained from  $P(I,D,G,S,J)$   
 with  $P(i^1|g^1) > P(i^1|g^2)$ ,  
 $P(s^1|g^1) > P(s^1|g^2)$  and  $P(j^1|g^1) > P(j^1|g^2)$

- If we set *Grade* to A, we expect  $P(J | do(g^1)) = P(J | g^1)$ 
  - i.e.,  $P(j^1 | do(g^1)) = P(j^1 | g^1)$  and  $P(j^0 | do(g^1)) = P(j^0 | g^1)$
- What about probabilities associated with *Intelligence* and *SAT*?
  - Our beliefs don't change  $P(i^1 | do(g^1)) = P(i^1)$  and  $P(s^1 | do(g^1)) = P(s^1)$
- There is strong correlation between all three variables (*G*, *I*, *S*) and *Job*
  - But the relationship between *Grade* and *Job* is causal
  - Manipulating *Grade* does not change *Intelligence* and *SAT*
- Three types of correlation: (*G*, *J*): *G* causes *J* (causality),  
 (*G*, *I*): *I* causes *G* (opposite causality)  
 (*G*, *S*): joint dependence on *I*

# Goal of Causal Framework

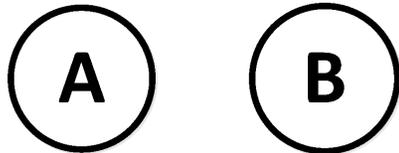
- Provide a rigorous basis for answering queries involving intervention
- Allow us to distinguish between the three cases
  - Causality, Opposite causality, Joint dependence
- Can be used to answer both intervention and counterfactual queries
  - Counterfactual requires much finer-grained information, which is difficult to acquire in practice

# Correlation and Causation

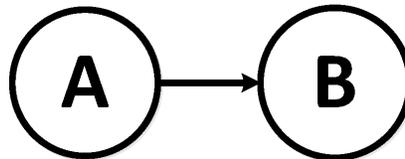
- Correlation between  $X$  and  $Y$  can arise in multiple settings:
  - When  $X$  causes  $Y$
  - When  $Y$  causes  $X$
  - When  $X$  and  $Y$  are both effects of a single cause
- Leads to question of identifiability of causal models
  - If we observe two correlated variables, what can we infer about the causal relationship between them?

# Relationships between events

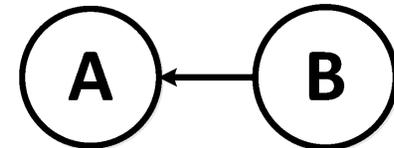
A | B



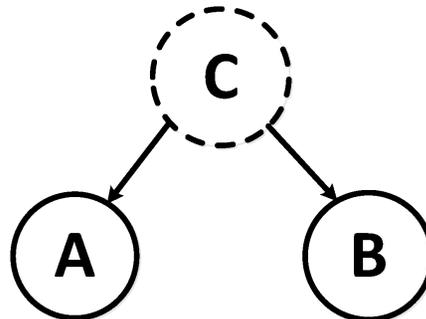
A Causes B



B Causes A



Common Causes for A and B, which do not cause each other



A = ice cream sales  
B = drowning deaths  
C = hot summer

Correlation is a broader concept than causation

# Difficulty due to Latent Factors

- Identifying causality is complicated by many reasons for correlation between  $X$  and  $Y$ 
  - When a variable  $W$  causally affects  $X$  and  $Y$  we observe a correlation between them
  - If we know existence of  $W$  and can observe it,
    - We can disentangle correlation of  $X$  and  $Y$  induced by  $W$  and compute residual correlation between  $X$  and  $Y$  attributed to a direct causal relationship
  - There are many unobserved latent variables
    - Latent variables induce correlations that do not correspond to causal relations
    - Forms a *confounding factor* in determining causality

# Example of Latent Variable

- Correlation between PeptAid and stomach ulcer
  - Because taking PeptAid precedes ulcer we may conclude that PeptAid causes ulcers
- Alternative explanation: correlation due to a latent common cause— pre-ulcer discomfort
  - Even if we account for this there may be many other latent variables, such as stress

# Confounding Factors

- Confounding factors that induce noncausal correlation between variables
  - Latent Variables
  - Selection Bias
  - Mixture of different populations

# Example of Selection Bias

- In alumni survey
  - Negative correlation between participation in athletics ( $A$ ) and GPA ( $G$ )
  - Can we conclude  $A$  reduces  $G$ , or that high  $G$  tend not to participate in  $A$
  - Alternative explanation
    - Survey respondents ( $S = s^1$ ) are unrepresentative
    - Students with good  $A$  and  $G$  respond and not otherwise
    - We have causal links  $A \rightarrow S \leftarrow G$
    - A case of intercausal reasoning  $P(a^1 | s^1) > P(a^1 | g^1, s^1)$
  - Without accounting for bias in population we may falsely explain correlation as causation

# Example of Mixture of Populations

- Girls develop verbal abilities earlier
- Boys tend to be taller at same age
- No known correlation between height and verbal ability
- If we measure across all children, we see negative correlation between verbal ability and height
- Here we have a latent variable of Gender which confounds common cause

# Example of two time series

- Data over past century:
  - Average height of adult world population in that year ( $H$ ) increasing due to nutrition)
  - Total size of polar caps in that year ( $S$ ) decreasing due to global warming
  - Negative correlation between  $H$  and  $S$
- Cannot conclude that size of polar caps causally influences population height
- Also a case of a latent variable: Time

# Challenge of Causal Modeling

- Correlation between a pair of variables  $X$  and  $Y$  may be a consequence of multiple mechanisms
  - Some are causal, others not
- To answer a causal query regarding an intervention at  $X$ , we need to disentangle the different mechanisms to isolate the correlation due to the causal effect of  $X$  on  $Y$

# Causal Model Definition

- A causal model takes same form as a BN
  - A DAG over random variables
  - Each variable  $X$  is governed by a causal mechanism
  - Value of  $X$  is a stochastic function of parent values
    - Difference is in interpretation of edges

# Causal Mechanism (CPD)

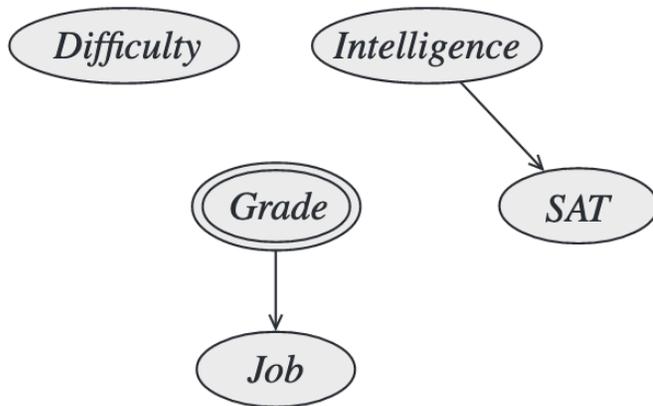
- Has the same form as a standard CPD
- For node  $X$  and its parents  $U$ , a causal model has a stochastic function from the values of  $U$  to the values of  $X$ 's parents are its direct causes
- Causality flows in the direction of the edges
  - So that  $X$ 's value is actually determined via the stochastic function implied by  $X$ 's CPD

# Causal Model and Intervention Queries

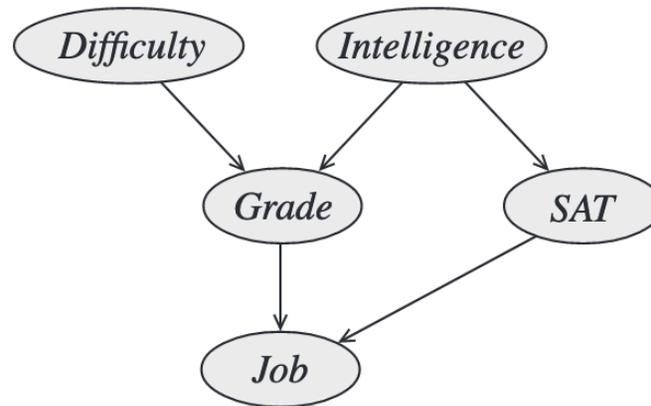
- Forms basis of intervention queries
  - Intervene at  $X$  setting to  $x$  means
    - Replace  $X$ 's CPD with one with  $X = x$  with probability 1
- When we intervene at a variable  $X$ , setting its value to  $x$ , we replace the original causal mechanism with one that dictates it take the value  $x$
- A causal model can answer intervention queries using the appropriate mutilated network

# Example of Mutilated Network

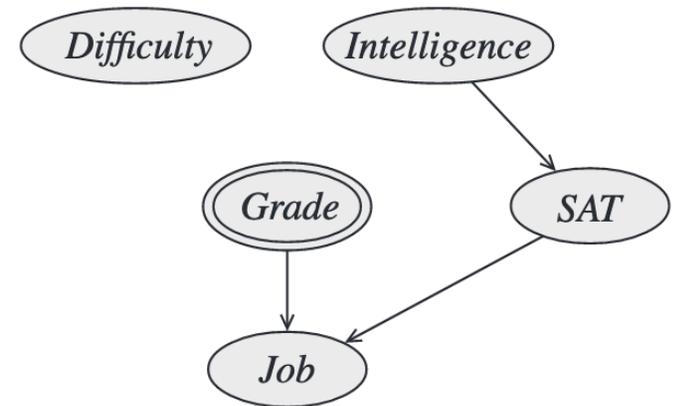
Mutilated student network with intervention at  $G$



Expanded student network with arc  $S \rightarrow J$



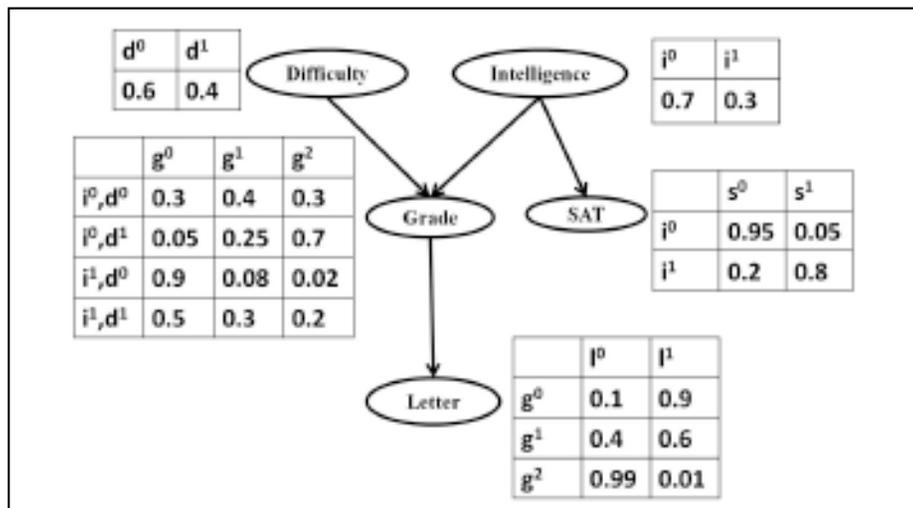
Mutilated expanded network with intervention at  $G$



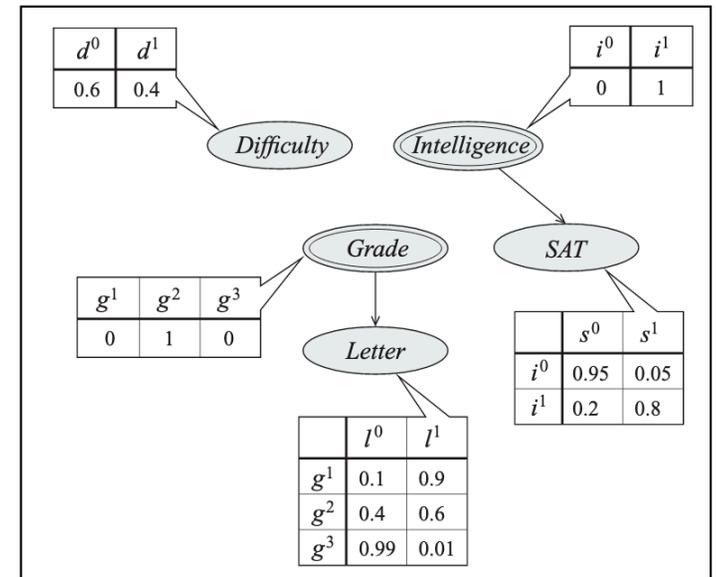
- We eliminate all incoming edges into each variable  $Z_i \in \mathbf{Z}$  and set its value to be  $z_i$  with probability 1

# Mutilated Network and Causal Model

- We eliminate all incoming edges into each variable  $Z_i \in \mathbf{Z}$  and set its value to be  $z_i$  with probability 1



$C_{(\mathbf{Z} = \mathbf{z})}$

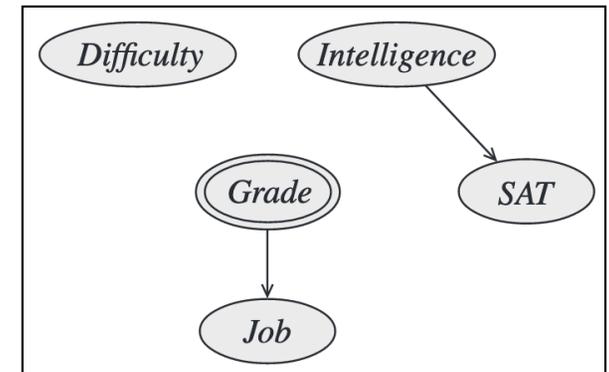
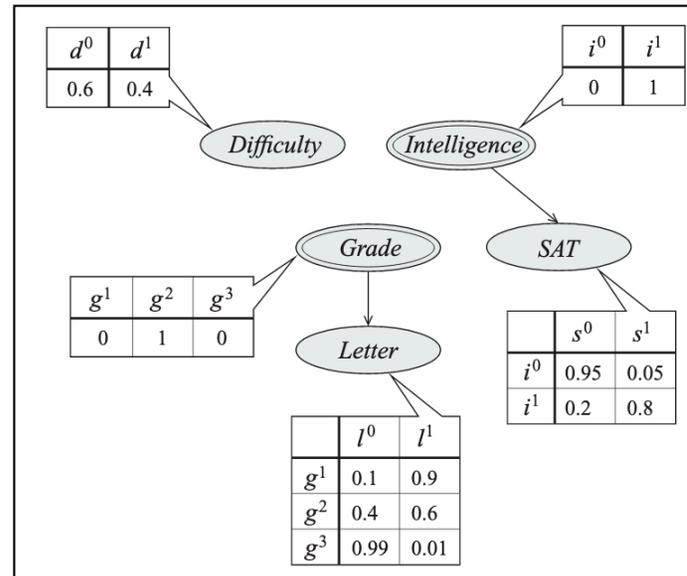
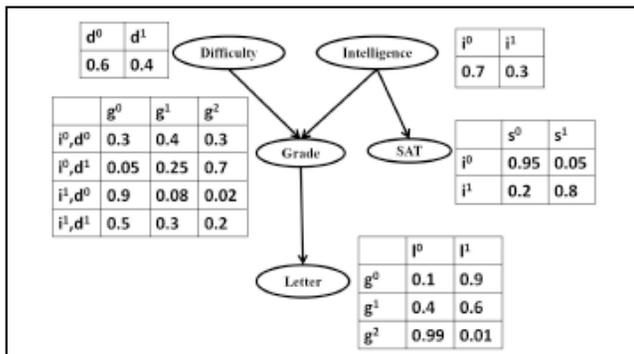


- A causal model  $C$  over  $\mathbf{X}$  is a BN that can answer:
  - Probabilistic queries
  - Intervention queries  $P(\mathbf{Y} | do(\mathbf{z}), \mathbf{X} = \mathbf{x})$  as follows
 
$$P_C(\mathbf{Y} | do(\mathbf{z}), \mathbf{x}) = P_{C(\mathbf{Z} = \mathbf{z})}(\mathbf{Y} | \mathbf{x})$$

# Inference with Mutilated Network

$$C_{(G=g^1)}$$

$C$



- Prior over SAT remains unchanged

$$P(S | do(G=g^2)) = P_{C(G=g^2)}(S) = P_C(S)$$

- Letter is

$$P_C(L | G = g^2)$$

- In general, answer to an intervention query differs from answer to some observational query

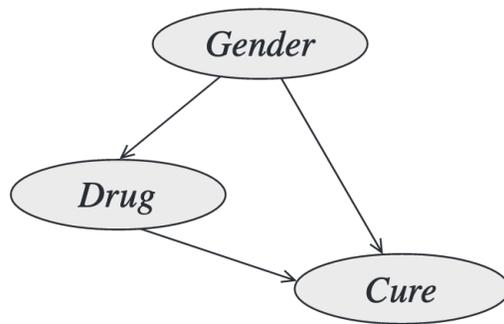
# Simpson's Paradox

- In US House of Reps, Civil Rights Act, 1964
  - 61% of Ds voted for, 80%, of Rs voted for
  - Might conclude: being R, rather than D, was factor in voting for Act
  - Picture changes if we include additional factor
    - Whether legislator is from North/South state

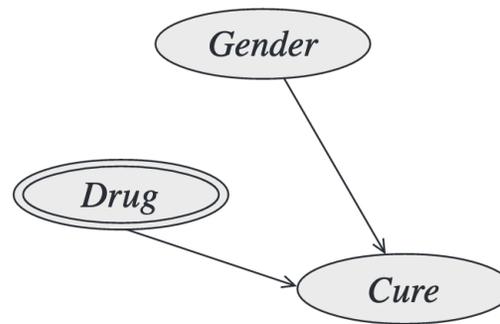
- If include extra factor, situation reverses, in both North *and* South.
- The break down:
  - North:** D (94%), R (85%)
    - D (145/154)
    - R (138/162)
  - South:** D (7%), R (0%)
    - D (7/94), R (0/10)
  - Overall:**
    - D (152/248, 61%),
    - R(138/172, 80%)
- Paradox is due to an unbalanced dataset

# Causal frame resolves Simpson's paradox

General:	$P(c^1   d^1) > P(c^1   d^0)$	$57.5\% > 50\%$
Conditional:	$P(c^1   d^1, G = male) < P(c^1   d^0, G = male)$	$70\% < 80\%$
	$P(c^1   d^1, G = female) < P(c^1   d^0, G = female)$	$20\% < 40\%$



Causal network for Simpson's paradox



Mutilated network for answering  $P(c^1 | do(d^1))$

The appropriate query we need to answer in determining whether to prescribe the drug is not  $P(c^1 | d^1)$ , but rather  $P(c^1 | do(d^1))$

# Structural Causal Identifiability

- Intervention queries can be answered with a fully specified causal model
- Cannot fully specify causal model due to presence of latent variables

# Determining $P(Y | do(X))$

- Consider a pair  $X, Y$  with observed correlation

- Our task is to determine  $P(Y | do(X))$

- Assume that  $X$  temporally precedes  $Y$

- Therefore  $Y$  cannot cause  $X$

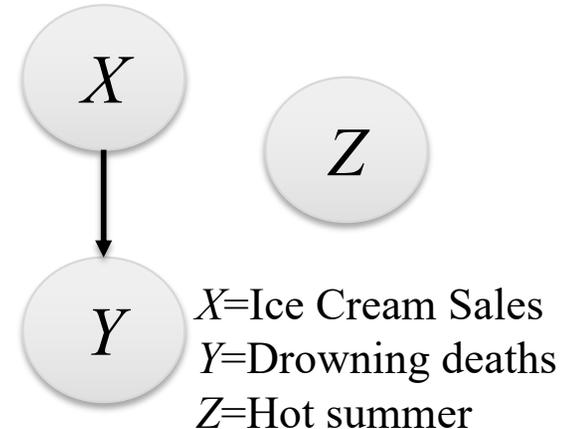
- Three cases:

1. If some correlation is due to a hidden cause, cannot determine how much effect perturbing  $X$  will have on  $Y$

2. If all correlation is due to causal link then  $P(Y|do(X))=P(Y|X)$

3. If all of the correlation is due to the hidden common cause then  $P(Y|do(X)) = P(Y)$

- In general any distribution between the two is possible



# Identifiable intervention queries

- Latent variables are inevitable
  - Sometimes, causal questions answered using only observed correlations
- Which intervention queries are *identifiable*?
  - i.e., answered using only conditional probabilities involving observable variables
    - Probabilities over observed variables can be estimated from data or elicited from an expert.
    - Thus, if we can reduce the answer to a query to an expression involving only such probabilities, we may be able to provide a robust and accurate answer to it

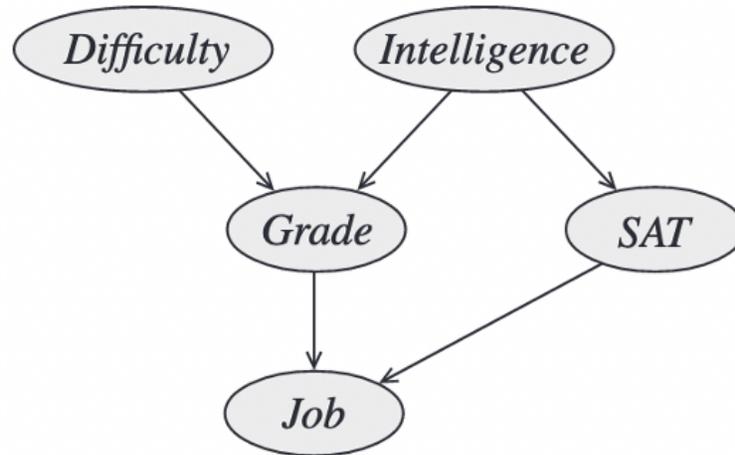
# Augmented Causal Model

- Replace one query with an equivalent one
- May be able to replace intervention queries with purely observational ones
- An *augmented causal model*  $\mathcal{G}^\dagger$  encodes possible effects of intervention in the model
  - View intervention as a new decision variable  $\hat{Z}$  that determines whether we intervene at  $Z$  and if so what value it is

# Query Simplification Rules

- Let  $C$  be a Causal model over graph  $G$ . Then
  1.  $P(Y | do(Z=z), X=x, W=w) = P(Y | do(Z=z), X=x)$   
if  $W$  is d-separated from  $Y$  given  $Z, X$  in the graph  $G_Z^\dagger$
  2.  $P(Y | do(Z=z), X=x, W=w) = P(Y | do(Z=z), X=x, W=w)$   
if  $Y$  is d-separated from  $\hat{X}$  given  $X, Z, W$  in the graph  $G_Z^\dagger$
  3.  $P(Y | do(Z=z), do(X=x), W=w) = P(Y | do(Z=z), W=w)$   
if  $Y$  is d-separated from  $\hat{X}$  given  $Z, W$  in the graph  $G_Z^\dagger$

# Example of Augmented Model

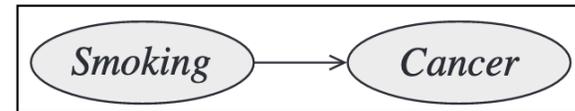


- Query of the form  $P(S \mid do(G))$  taking  $Y=S, X=G$  and  $Z, W = \emptyset$
- Augmented model  $\mathcal{G}$  contains new decision parent  $\hat{G}$

# Does Smoking cause Cancer?

- 1960s data: two competing models

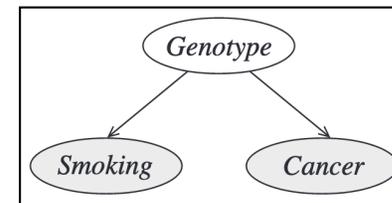
- Model 1: *Surgeon General*



- Correlation is causal

- If we ban smoking, cancer rate will be same as for non-smokers

- Model 2: *Tobacco Industry*



- No causality

- An unobserved genotype produces craving for nicotine

- Both models express same distributions over  $S, C$

- No way to distinguish them based on observed data
- Models give different answers to interventional queries

- Model 1:  $P(c^1|do(S:=s^1)) = P(c^1|s^1)$  →

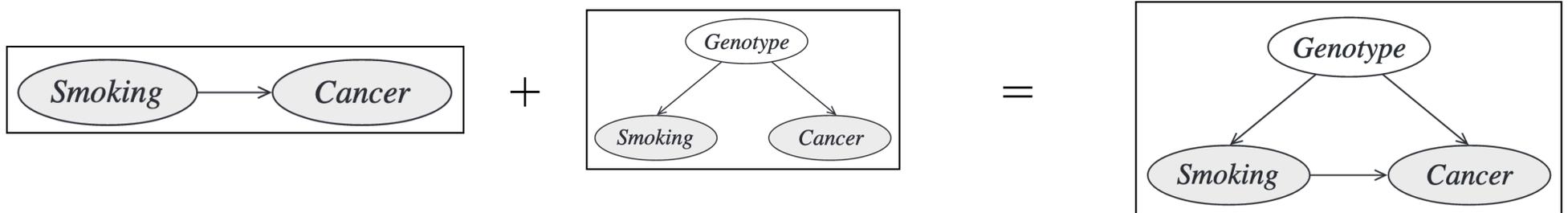
If we force people to smoke their probability of getting cancer is the same as probability conditioned on smoking

- Model 2:  $P(c^1|do(S:=s^1)) = P(c^1)$  →

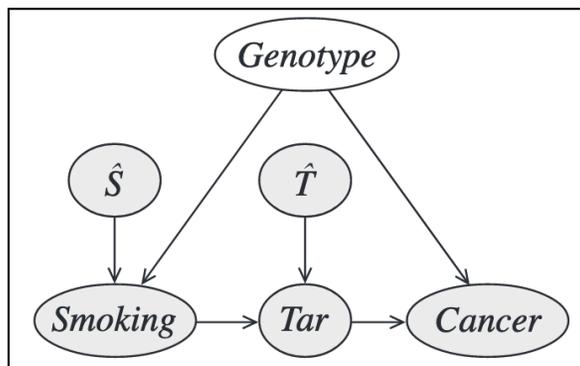
Making population smoke or stop smoking would have no effect on rate of cancer cases

# Solution for Causality of Smoking to Cancer

- Combine the two models into a single one



- But cannot parameterize based on observations
  - Since both original models can perfectly explain data
- If we refine the model by introducing additional assumptions, we can provide such estimates

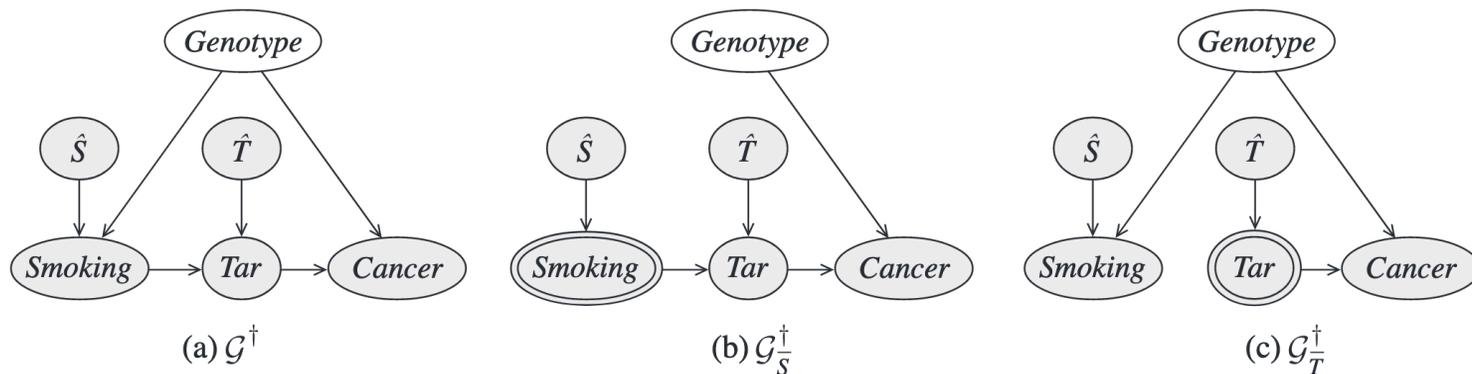


Effect of smoking on cancer is not a direct one, but occurs via tar deposits in the lungs

If we can measure tar amounts we can now determine probability of intervention query  $P(c^1|do(s^1))$  using observed correlations alone

# Answer to cancer intervention query

- Goal: answer to  $P(c^1 | do(s^1))$ , an interventional query
  - Mutilated networks for: (b) intervention at  $S$ , (c) intervention at  $T$



$$P(C | do(s^1)) = \sum_t P(C | do(s^1), t) P(t | do(s^1))$$

$$P(t | do(s)) = P(t | s)$$

$$P(C | do(s^1), do(t)) = P(C | do(t))$$

$$P(C | do(s^1), t) = P(C | do(s^1), do(t))$$

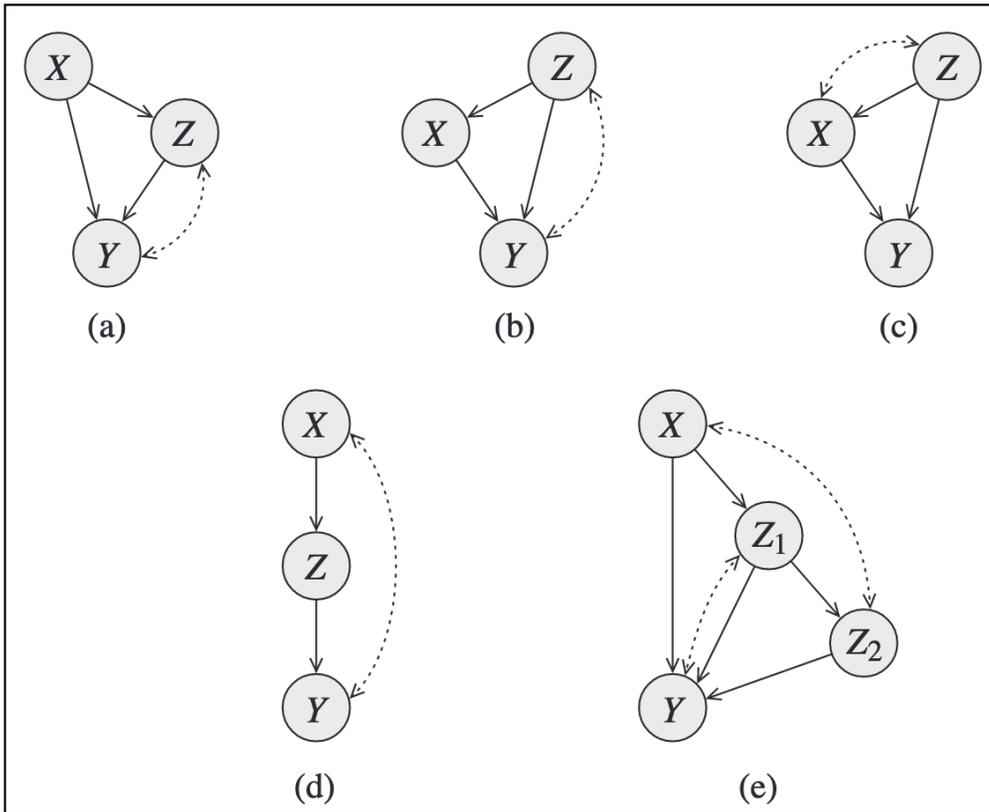
$$\begin{aligned} P(C | do(t)) &= \sum_{s'} P(C | do(t), s') P(s' | do(t)) \\ &= \sum_{s'} P(C | t, s') P(s' | do(t)) \\ &= \sum_{s'} P(C | t, s') P(s'). \end{aligned}$$

$$\begin{aligned} P(c | do(s^1)) &= \sum_t P(c | do(s^1), t) P(t | do(s^1)) \\ &= \sum_t P(c | do(s^1), t) P(t | s^1) \\ &= \sum_t P(t | s^1) \sum_{s'} P(c | t, s') P(s'). \end{aligned}$$

# Identifiable graphs

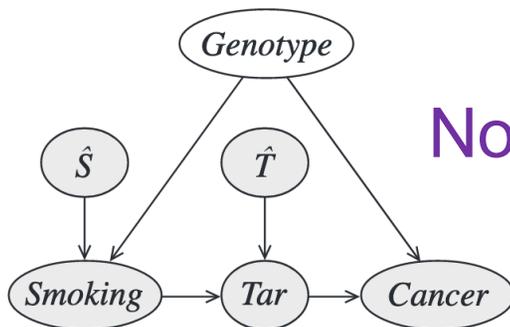
- In an identifiable graph Intervention queries are answerable using only conditional probabilities involving observable variables
- Identifiable graphs are where  $P(Y | do(X))$  is answerable
- Unidentifiable graphs are where  $P(Y | do(X))$  is unanswerable

# Models where $P(Y | do(X))$ is identifiable



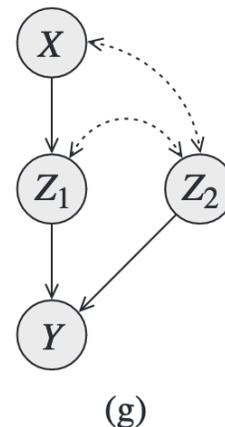
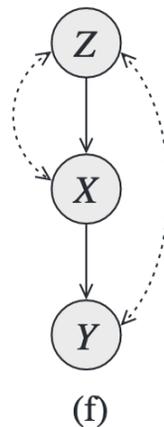
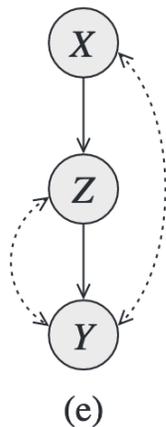
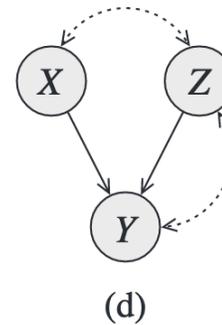
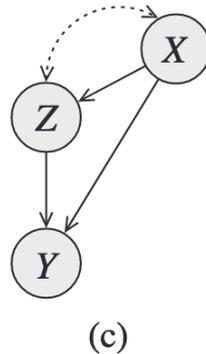
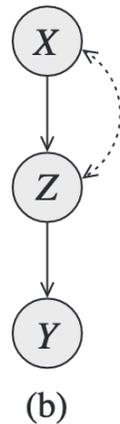
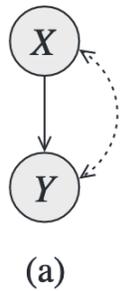
The bidirected dashed arrows indicate cases where a latent variable affects both of the linked variables

These graphs are maximal, i.e., adding edges makes them unidentifiable



Note: Smoking model has the same form as (d)

# Models where $P(Y | do(X))$ is unidentifiable



Removal of edges can only help in identifiability  
Conversely, adding edges can only reduce identifiability  
Unidentifiable queries are quite common in practice

# Three-layer Causal Hierarchy

Level (Symbol)	Typical Activity	Typical Questions	Examples
1. Association $P(y x)$	Seeing	What is? How would seeing $X$ change my belief in $Y$ ?	What does a symptom tell me about a disease? What does a survey tell us about the election results?
2. Intervention $P(y do(x), z)$	Doing Intervening	What if? What if I do $X$ ?	What if I take aspirin, will my headache be cured? What if we ban cigarettes?
3. Counterfactuals $P(y_x x', y')$	Imagining, Retrospection	Why? Was it $X$ that caused $Y$ ? What if I had acted differently?	Was it the aspirin that stopped my headache? Would Kennedy be alive had Os- wald not shot him? What if I had not been smoking the past 2 years?