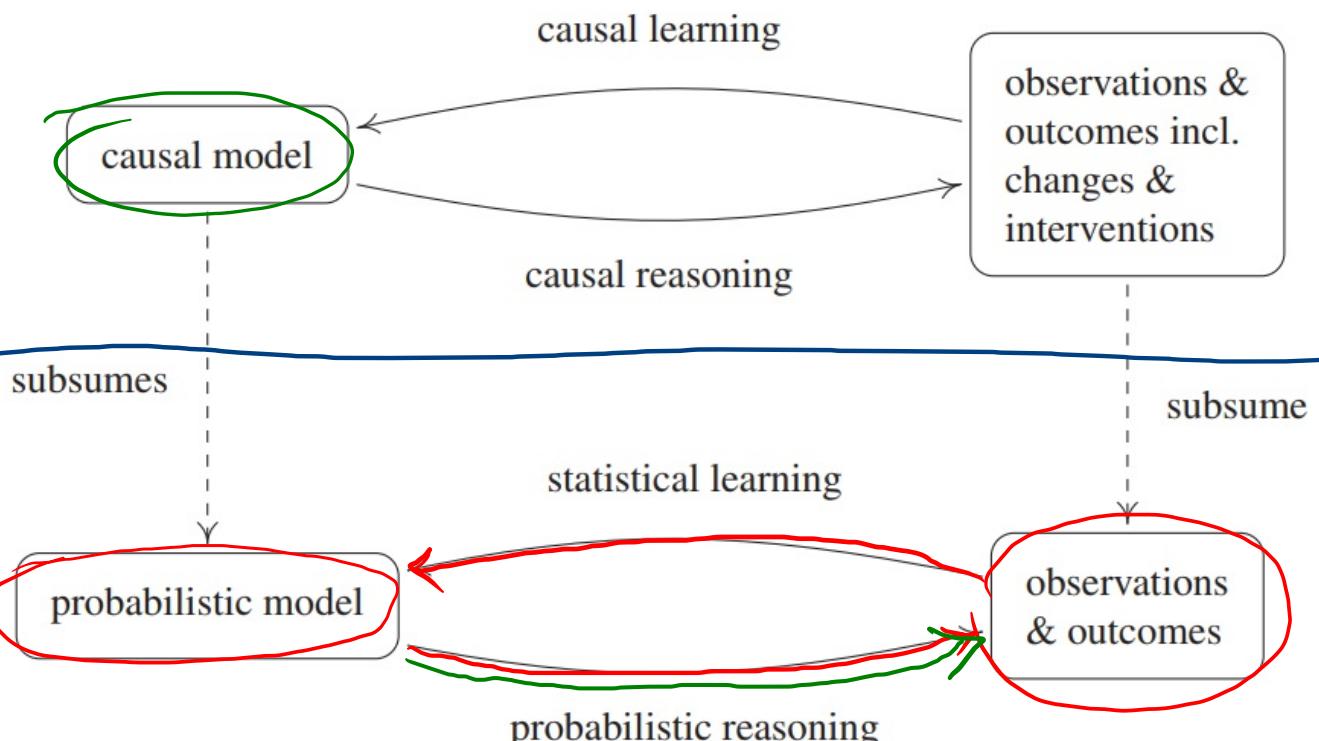


Causality

Dominik Janzing & Bernhard Schölkopf
Max Planck Institute for Intelligent Systems
Tübingen, Germany

<http://ei.is.tuebingen.mpg.de>



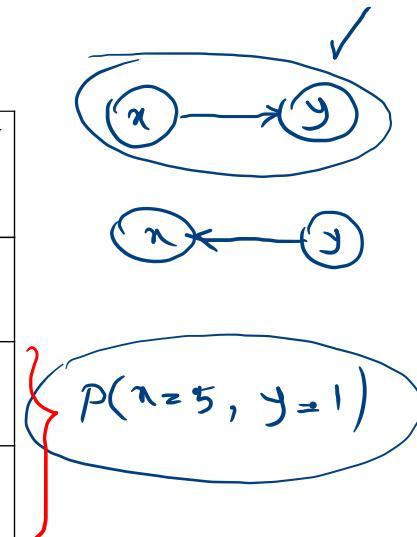


$$\checkmark \quad \underline{P(x,y)} \longrightarrow P(y|x) \quad P(x)$$

A Modeling Taxonomy

model	<u>predict in IID setting</u>	<u>predict under changing distributions / interventions</u>	<u>answer counter-factual questions</u>	<u>obtain physical insight</u>	<u>automatically learn from data</u>
mechanistic model ↓	Y	Y	Y	Y	?
structural causal model	Y	Y	Y	N	Y??
causal graphical model	Y	Y	N	N	Y?
statistical model ~~~~~	Y	N	N	N	Y

SCM →



“Correlation does not tell us anything about causality”

- Better to talk of dependence than correlation
- Most statisticians would agree that causality does tell us something about dependence
- But dependence does tell us something about causality too:

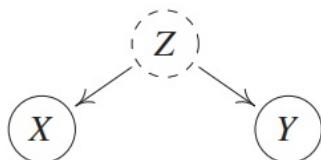


Principle 1.1 (Reichenbach's common cause principle) *If two random variables X and Y are statistically dependent ($X \not\perp\!\!\!\perp Y$), then there exists a third variable Z that causally influences both. (As a special case, Z may coincide with either X or Y .) Furthermore, this variable Z screens X and Y from each other in the sense that given Z , they become independent, $X \perp\!\!\!\perp Y | Z$.*

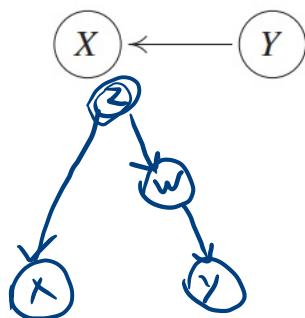


by permission of the
University of Pittsburgh.
All rights reserved.

X Y



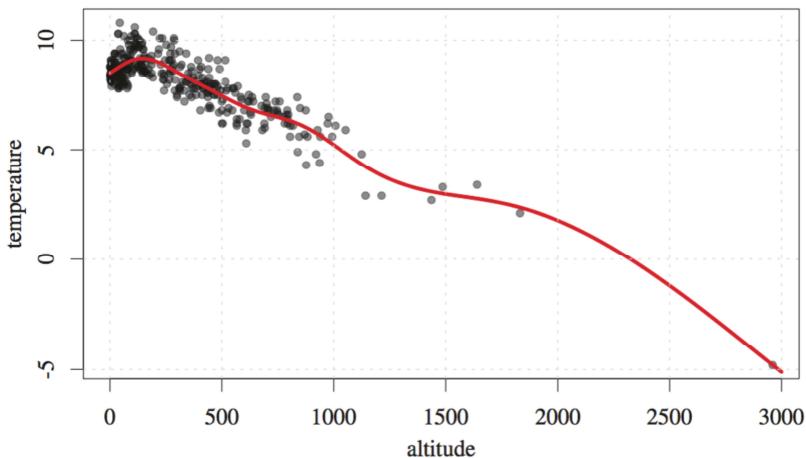
x y



$$P(x,y) = P(x) P(y)$$

$$E[xy] = 0$$

What is cause and what is effect?

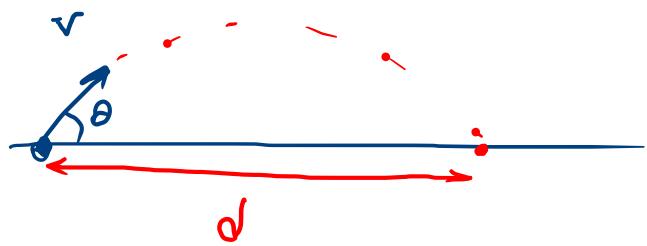


$$P(T, A) = P(A) P(T|A)$$

Two red arrows point from the terms $P(A)$ and $P(T|A)$ in the equation above to the corresponding nodes in the causal diagram.

$$\begin{aligned} p(a, t) &= p(a|t) p(t) && T \rightarrow A \\ &= p(t|a) p(a) && A \rightarrow T \end{aligned}$$





$$p(v,d) = \frac{p(v)}{M_1} \frac{p(d|v)}{M_2}$$

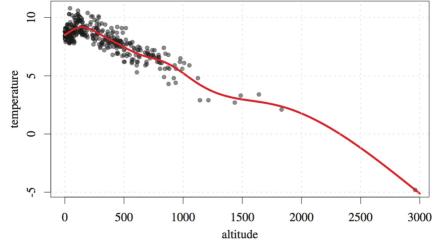
$$p(v) \approx N(10, 2)$$

$$p(v) \approx N(40, 5)$$

$$p(v)$$

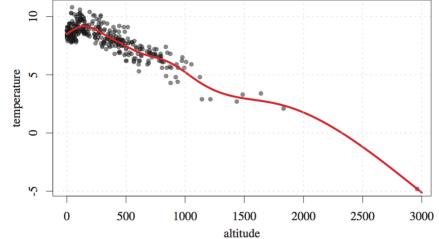
$$p(x_1, \dots, x_n) = \prod_{i=1}^n p(x_i | \text{Pa}(x_i))$$

Autonomous/invariant mechanisms



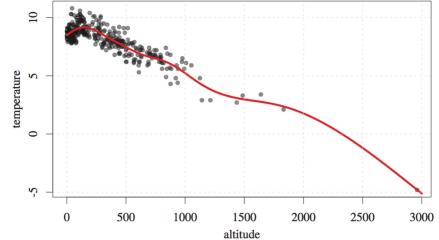
- intervention on a : raise the city, find that t changes
- hypothetical intervention on a : still expect that t changes, since we can think of a physical mechanism $p(t|a)$ that is independent of $p(a)$
- we expect that $p(t|a)$ is invariant across, say, different countries in a similar climate zone

Independence of cause & mechanism



- the conditional density $p(t|a)$ (viewed as a function of t and a) provides no information about the marginal density function $p(a)$
- this also applies if we only have a single density

Independence of noise terms



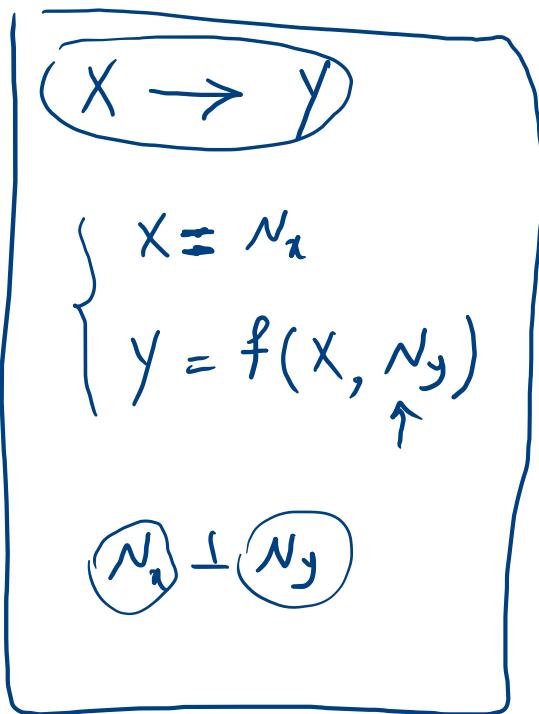
- view the distribution as entailed by a structural causal model (SCM)

$$\overbrace{\quad\quad\quad}^{\text{A} \rightarrow T} \left\{ \begin{array}{l} A := N_A, \\ (T) := f_T(A, N_T), \end{array} \right.$$

where $N_T \perp\!\!\!\perp N_A$

- this allows identification of the causal graph under suitable restrictions on the functional form of f_T

SCM



$$X \sim P(x)$$

$$Y \rightarrow X$$

$$Y = N_y$$

$$X = g(Y, N_x)$$

$$N_x \perp N_y$$

$$P(x, y)$$

$$x \rightarrow y$$

$$P(x, y) = P(x) \underbrace{P(y|x)}$$

$$y = f(x, N_y)$$

Pearl's do-notation

$$x \quad y$$
$$P(y|x=n)$$

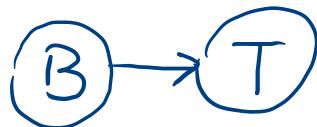
$$P(y|do\ x=n)$$

- Motivation: goal of causality is to infer the effect of interventions

- distribution of Y given that X is set to x :

$$p(Y|do\ X=x) \text{ or } p(Y|do\ x)$$

- don't confuse it with $P(Y|x)$
- can be computed from p and G



$$\underline{P(T|B=1)}$$

$$(B=0, \quad T=20)$$

$$P(T|do\ B=1)$$



y: يشعر بالعمر
x: شرب ١٠ كوب قهوة

Difference between seeing and doing

$$p(y|x)$$

probability that someone gets 100 years old given that we know that he/she drinks 10 cups of coffee per day

$$p(y|do\ x)$$

probability that some randomly chosen person gets 100 years old after he/she has been forced to drink 10 cups of coffee per day



MAX-PLANCK-GESSELLSCHAFT

The Principle of Independent Mechanisms

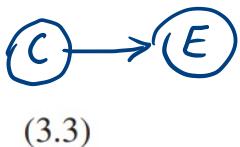
•

Principle 2.1 (Independent mechanisms) *The causal generative process of a system's variables is composed of autonomous modules that do not inform or influence each other.*

In the probabilistic case, this means that the conditional distribution of each variable given its causes (i.e., its mechanism) does not inform or influence the other conditional distributions. In case we have only two variables, this reduces to an independence between the cause distribution and the mechanism producing the effect distribution.

Example 3.2 (Cause-effect interventions) Suppose that the distribution $P_{C,E}$ is entailed by an SCM \mathfrak{C}

$$\begin{aligned} C &:= N_C \\ E &:= \underline{4 \cdot C + N_E}, \end{aligned}$$



with $N_C, N_E \stackrel{\text{iid}}{\sim} \mathcal{N}(0, 1)$, and graph $C \rightarrow E$. Then,

$$P(C|E=2) = \frac{P(C, E=2)}{P(E=2)} = \frac{\overbrace{P(C)}^{\mathcal{N}(0,1)} \overbrace{P(E|C)}^{\mathcal{N}(4c,1)}}{P(E=2)} = \mathcal{N}(2|0, 17)$$

$$P(C, E) = P(C) \underbrace{P(E|C)}_{\mathcal{N}(4c,1)} = \mathcal{N}(0,1) \mathcal{N}(4c,1)$$

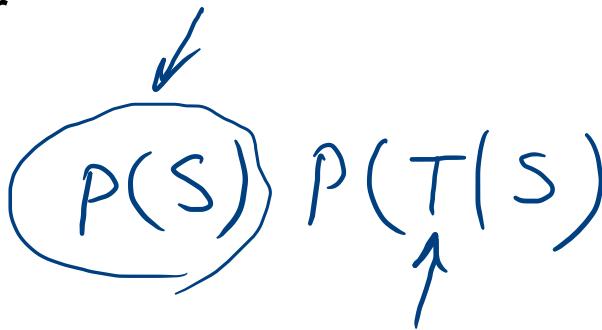
$$P(C | \cancel{d o} E=2) = P(C)$$

$$\cancel{P(C) \quad P(E|C)}$$

Example: Smoking and Teeth Color

S: سیگار باش

T: دندان خاکستری



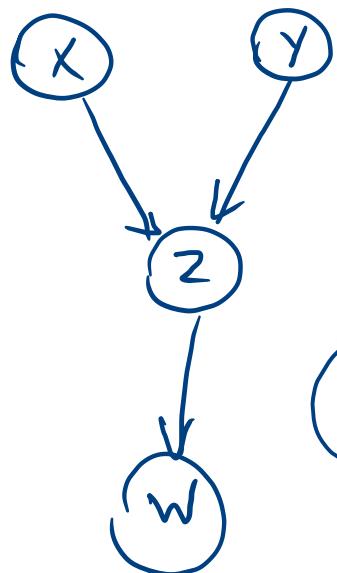
$$P(S | T=1)$$

$$P(S | \text{do } T=1)$$



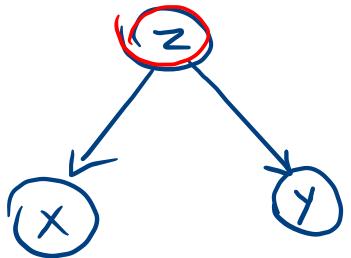
$$P(C \mid \text{do } E=e) = P(C)$$

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



$$p(x, y \mid \text{do } w = \omega) = p(x, y)$$

$$p(z \mid \text{do } y = y) = p(z \mid y)$$



$$P(x, y, z) = P(z) \cdot P(x|z) \cdot P(y|z)$$

The term $P(y|x)$ is crossed out with a large red X.

$$P(x | \text{do } y=y) = P(x)$$

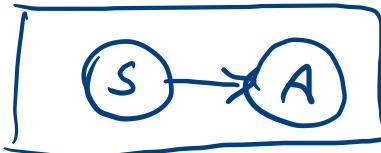
$$P(y | \text{do } x=x) = P(y)$$

$$P(x) = \sum_z P(x, z)$$

$$P(x|y) \neq P(x)$$

S :

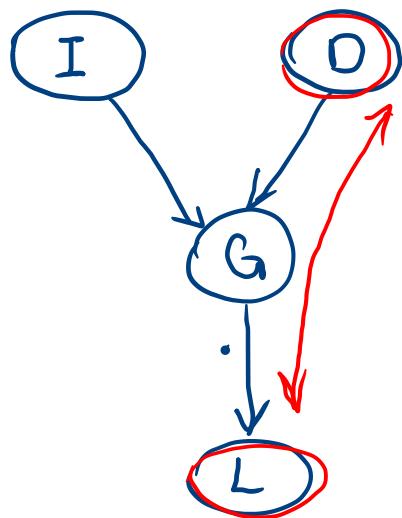
A :



$$P(A|S)$$

$$P(A|\text{do } S=1)$$

Observational Dist
Interventional ~



$$P(D | \text{do } L = l) = P(D)$$

(i, d, g, \cancel{l})
 l'

$$P(x_1, \dots, x_n)$$

Causal Graph

SCM : Structural Causal Model

$$x_1 = u_1$$

$$x_2 = f_2(x_1, u_2)$$

$$x_3 = f_3(x_1, x_2, u_3)$$

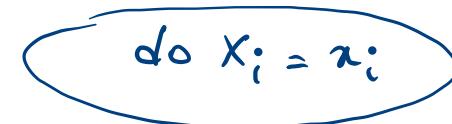
:

$$u_1 \perp u_2 \perp \dots \perp u_n$$

Counterfactual Quest.

$$\prod_{i=1}^n p(x_i | \text{Pa}(x_i))$$

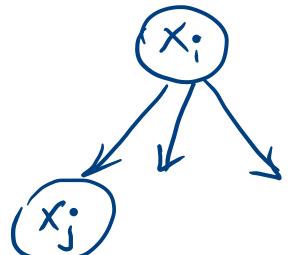
Computing $p(X_k | do x_i)$



summation over x_i yields

$$p(X_1, \dots, X_{i-1}, X_{i+1}, \dots, X_n | do x_i) = \prod_{j \neq i} p(X_j | PA_j(x_i)).$$

- distribution of X_j with $j \neq i$ is given by dropping $p(X_i | PA_i)$ and substituting x_i into PA_j to get $PA_j(x_i)$.
- obtain $p(X_k | do x_i)$ by marginalization



$$P(x_j | x_i) \rightarrow P(x_j | z_i)$$



UNIVERSITÄT WIEN VIENNA UNIVERSITY

Computing $p(X_k|do\ x_i)$

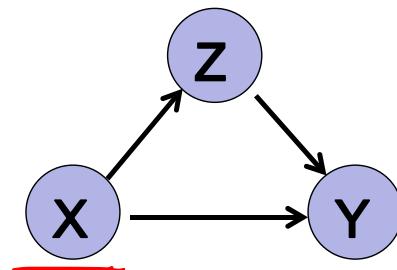
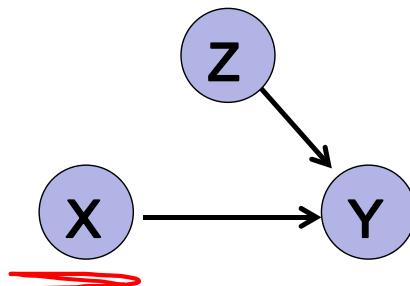
summation over x_i yields



$$p(X_1, \dots, X_{i-1}, X_{i+1}, \dots, X_n | do\ x_i) = \prod_{j \neq i} p(X_j | PA_j(x_i)).$$

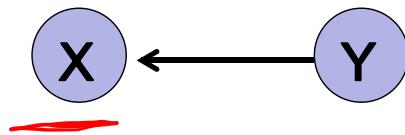
- distribution of X_j with $j \neq i$ is given by dropping $p(X_i | PA_i)$ and substituting x_i into PA_j to get $PA_j(x_i)$.
- obtain $p(X_k | do\ x_i)$ by marginalization

Examples for $p(.|do x) = p(.|\underline{x})$

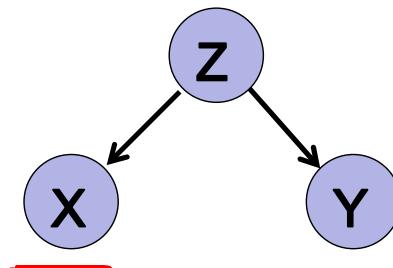


Examples for $p(.|do x) \neq p(.|x)$

- $p(Y|do x) = P(Y) \neq P(Y|x)$



- $p(Y|do x) = P(Y) \neq P(Y|x)$



Counterfactuals

Identifiability problem

e.g. Tian & Pearl (2002)

- given the causal DAG G and two nodes X_i, X_j
- which nodes need to be observed to compute $p(X_i|do\,x_j)$?



MAX-PLANCK-GESELLSCHAFT

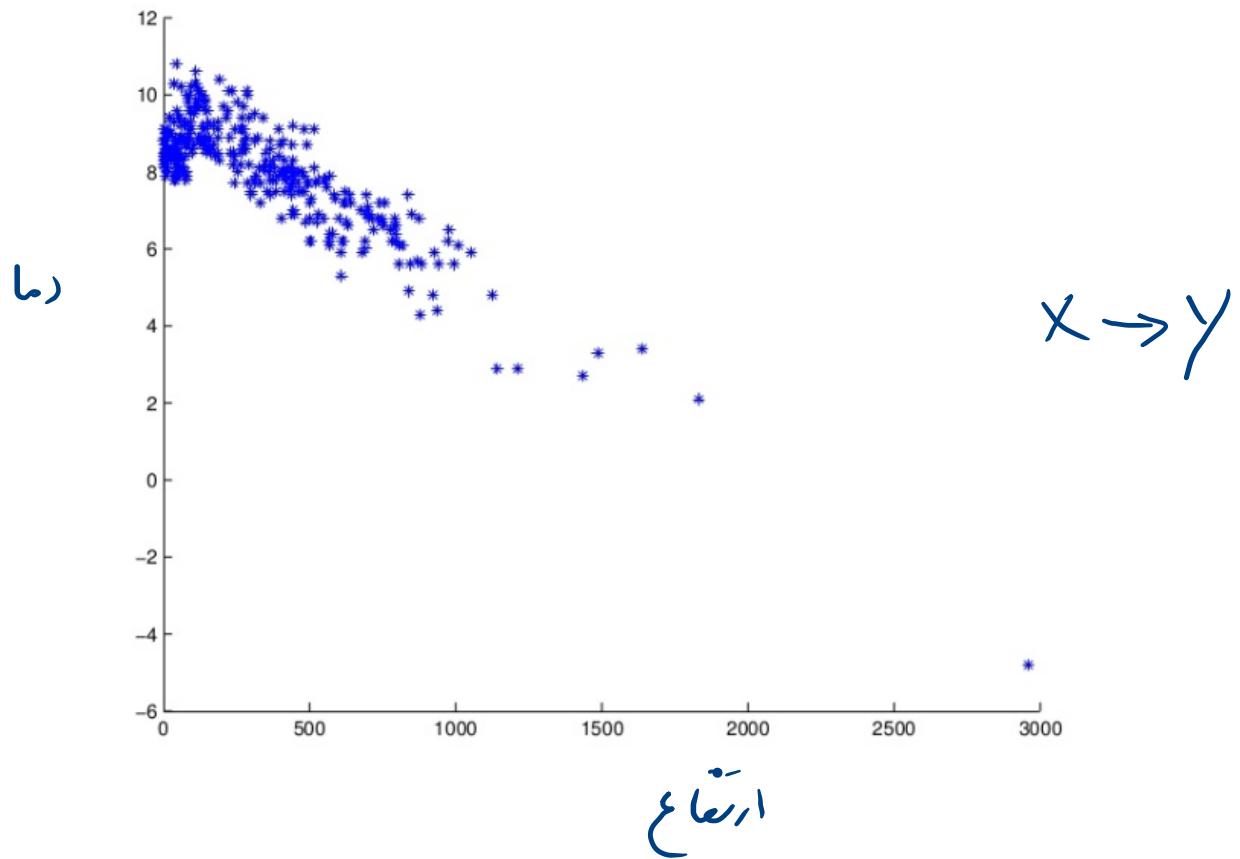
Causal Discovery

Inferring the DAG

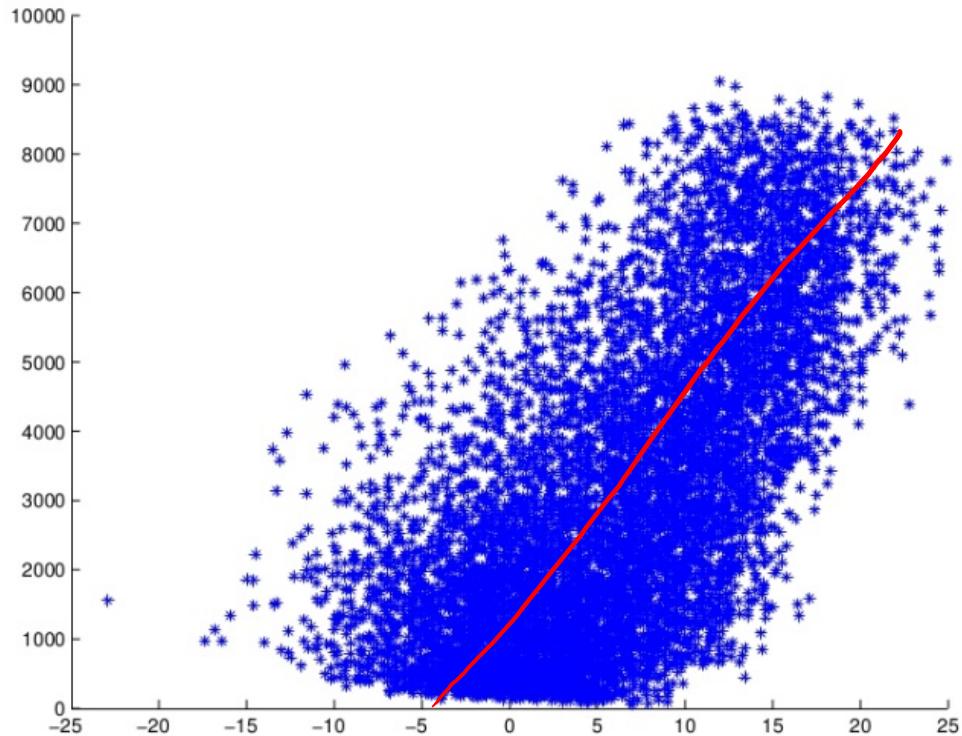
- Key postulate: Causal Markov condition
- Essential mathematical concept: d-separation
(describes the conditional independences required by a causal DAG)



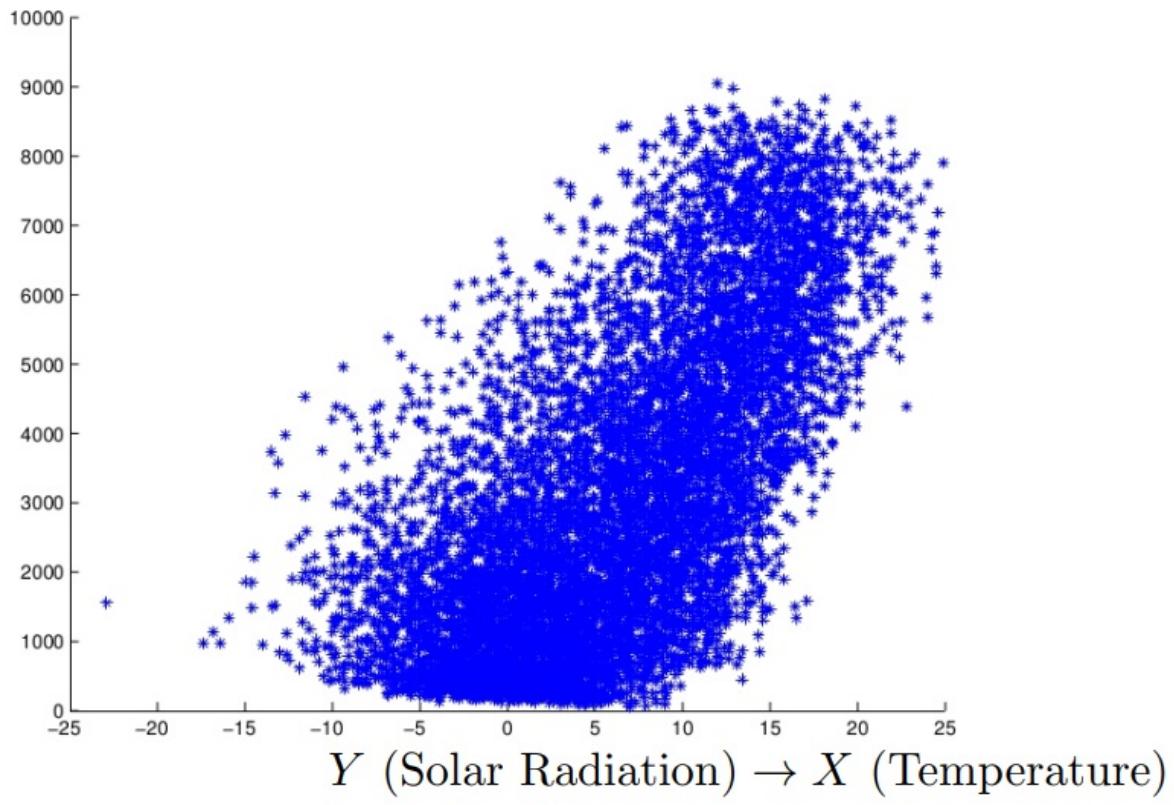
What's the cause and what's the effect?



What's the cause and what's the effect?



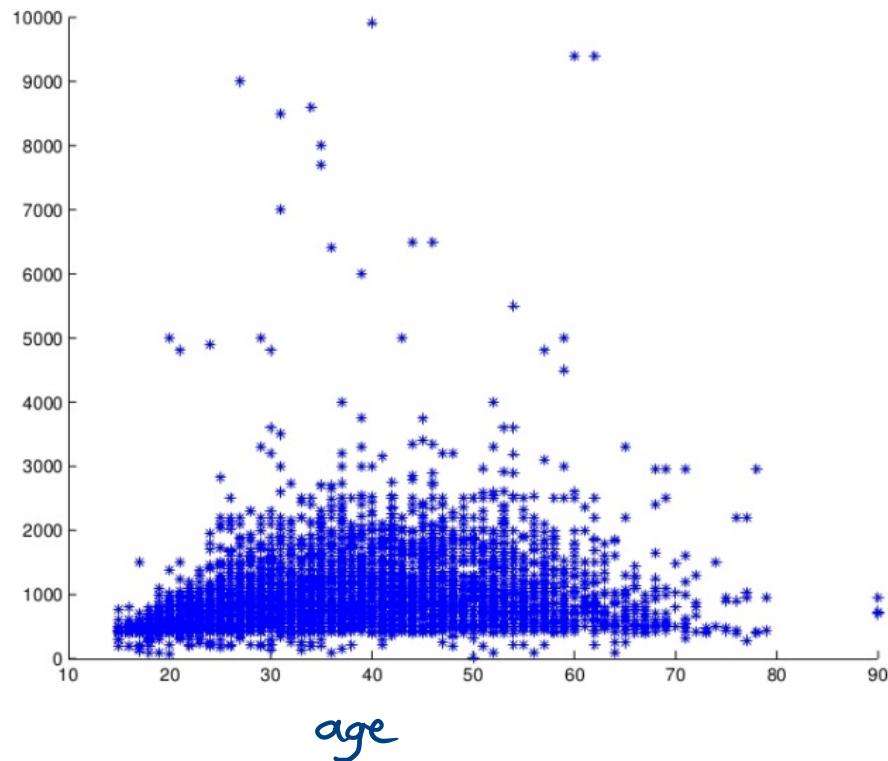
What's the cause and what's the effect?



What's the cause and what's the effect?

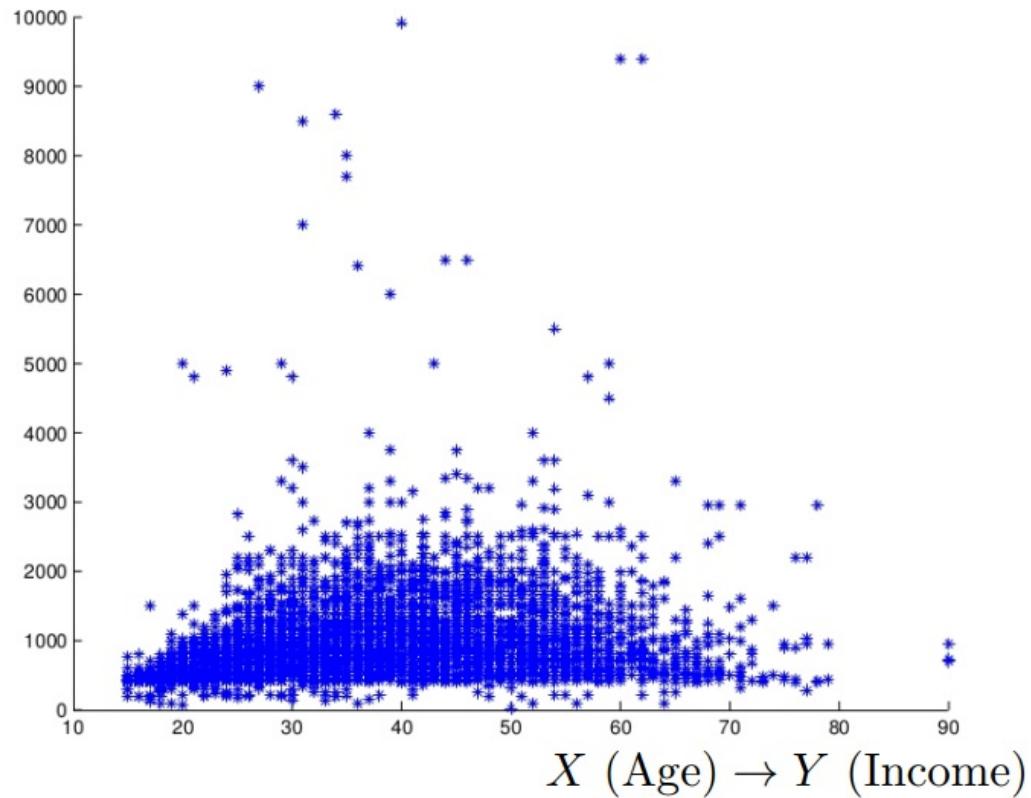
income

age → income



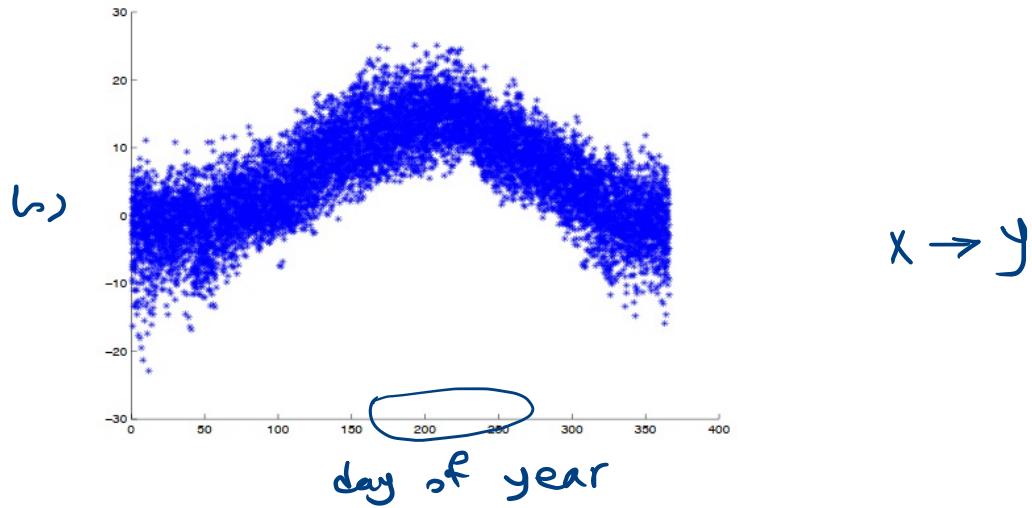
STANFORD UNIVERSITY

What's the cause and what's the effect?



MINERVA LEAPS

What's the cause and what's the effect?



Causal faithfulness

Spirites, Glymour, Scheines



p is called faithful relative to G if only those independences hold true that are implied by the Markov condition, i.e.,

$$(X \perp\!\!\!\perp Y | Z)_G \Leftrightarrow (X \perp\!\!\!\perp Y | Z)_p$$

Recall: Markov condition reads

$$\rightarrow (X \perp\!\!\!\perp Y | Z)_G \Rightarrow (X \perp\!\!\!\perp Y | Z)_p$$



d-separation



MAX-PLANCK-GESELLSCHAFT

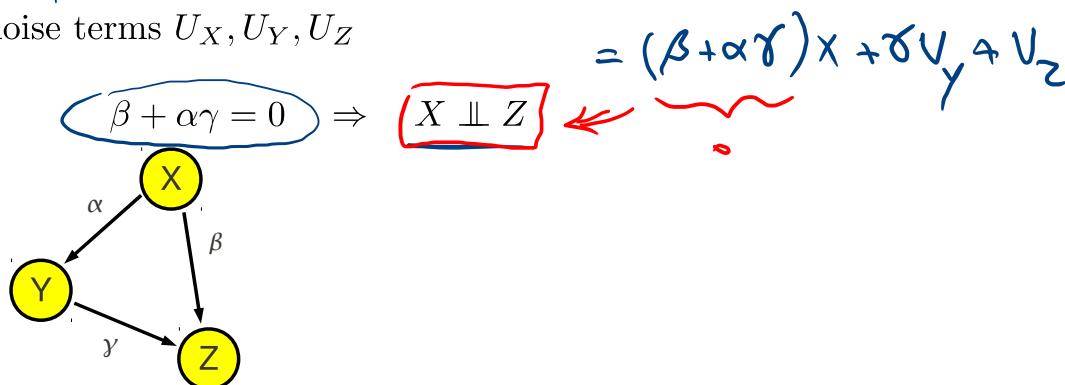
Examples of unfaithful distributions (1)

Cancellation of direct and indirect influence in linear models

SCM

$$\left\{ \begin{array}{l} X = U_X \\ Y = \alpha X + U_Y \\ Z = \beta X + \gamma Y + U_Z \end{array} \right. = \beta X + \gamma \underline{\alpha X} + \gamma U_Y + U_Z$$

with independent noise terms U_X, U_Y, U_Z



Markov equivalence class

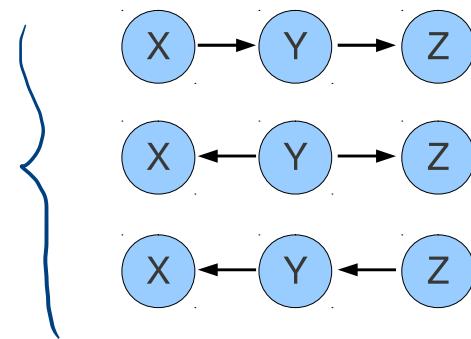
Theorem (Verma and Pearl, 1990): two DAGs are Markov equivalent iff they have the same skeleton and the same *v*-structures.

skeleton: corresponding undirected graph

v-structure: substructure $X \rightarrow Y \leftarrow Z$ with no edge between X and Z



Markov equivalent DAGs



A blue oval contains the conditional independence statement:

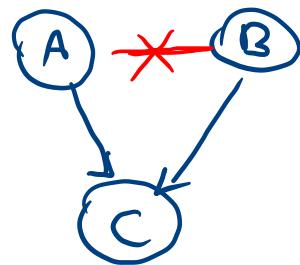
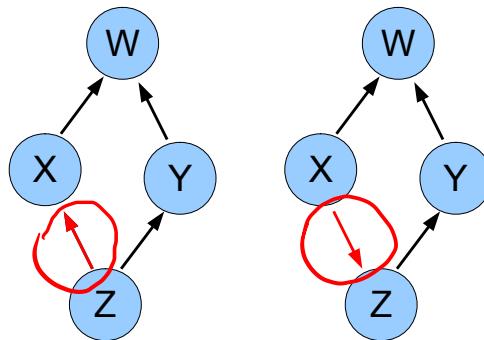
$$X \perp\!\!\! \perp Z \mid Y$$

same skeleton, no *v*-structure

$$X \perp\!\!\! \perp Z \mid Y$$



Markov equivalent DAGs



same skeleton, same v-structure at W



Efficient construction of skeleton

PC algorithm by Spirtes & Glymour (1991)

iteration over size of Sepset

1. remove all edges $X - Y$ with $X \perp\!\!\!\perp Y$
2. remove all edges $X - Y$ for which there is a neighbor $Z \neq Y$ of X with $X \perp\!\!\!\perp Y | Z$
3. remove all edges $X - Y$ for which there are two neighbors $Z_1, Z_2 \neq Y$ of X with $X \perp\!\!\!\perp Y | Z_1, Z_2$
4. ...

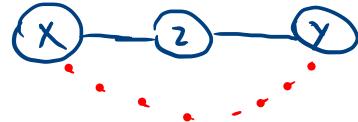


Advantages

- many edges can be removed already for small sets
- testing all sets S_{XY} containing the adjacencies of X is sufficient
- depending on sparseness, algorithm only requires independence tests with small conditioning tests
- polynomial for graphs of bounded degree



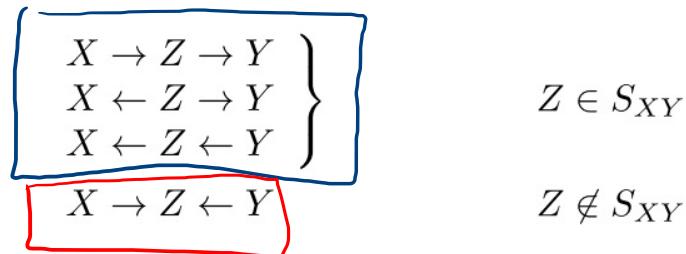
Find v-structures



$X \perp Y | Z \checkmark$

- given $X - \cancel{Z} - Y$ with X and Y non-adjacent
- given S_{XY} with $X \perp Y | S_{XY}$

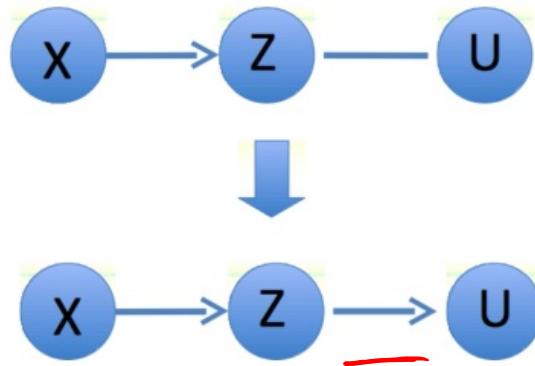
a priori, there are 4 possible orientations:



Orientation rule: create v-structure if $Z \notin S_{XY}$

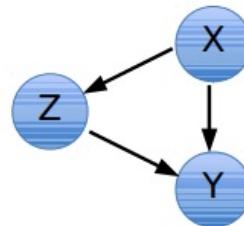
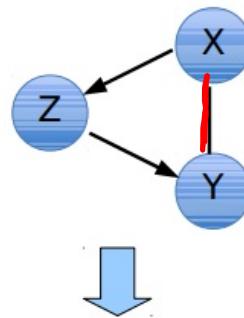


Direct further edges (Rule 1)



(otherwise we get a new v-structure)

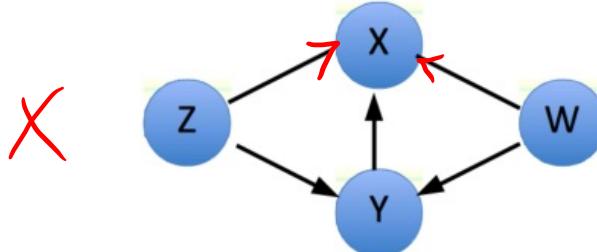
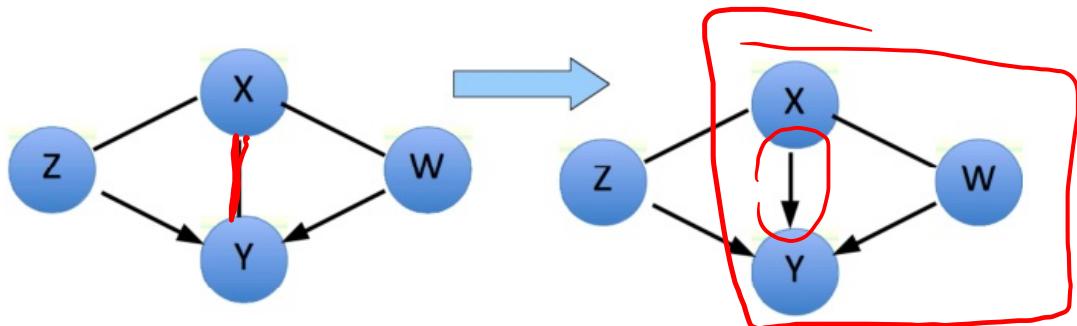
Direct further edges (Rule 2)



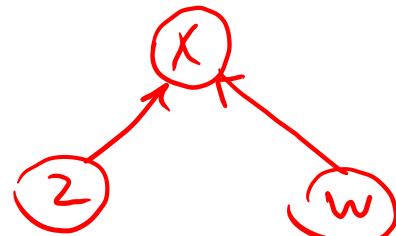
(otherwise one gets a cycle)



Direct further edges (Rule 3)



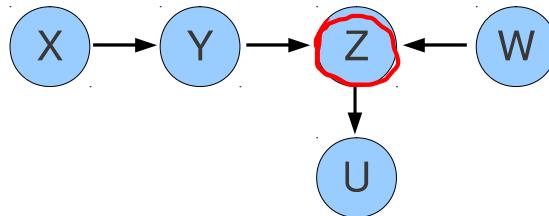
could not be completed
without creating a cycle
or a new v-structure



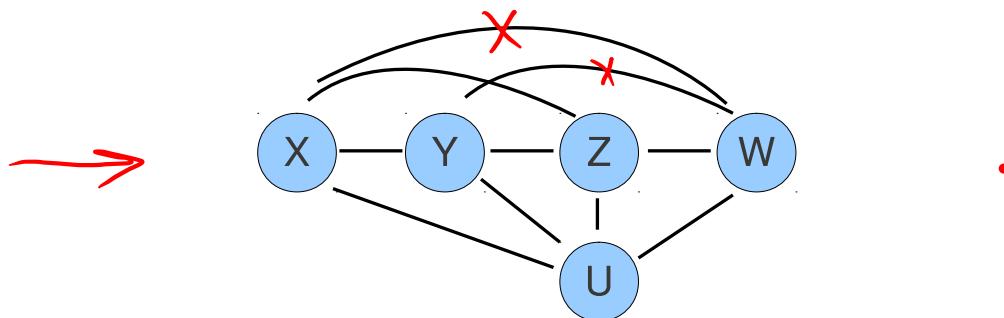
Examples

(taken from Spirtes et al, 2010)

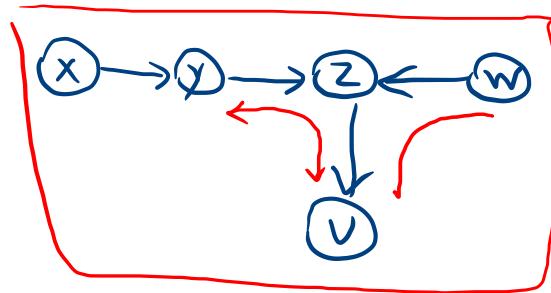
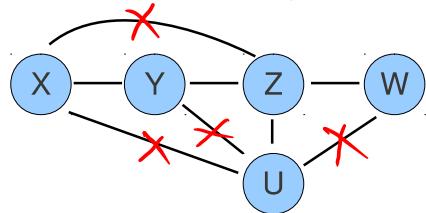
true DAG



start with fully connected undirected graph

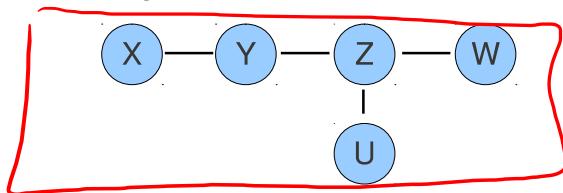


remove all edges $X - Y$ with $X \perp\!\!\!\perp Y | \emptyset$



$$X \perp\!\!\!\perp W \quad Y \perp\!\!\!\perp W$$

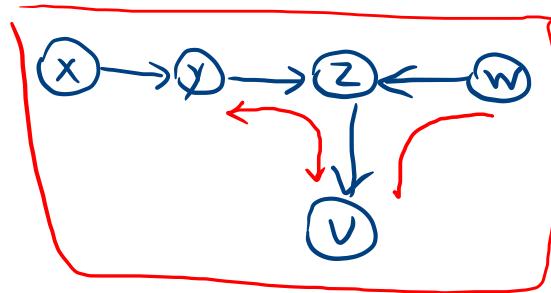
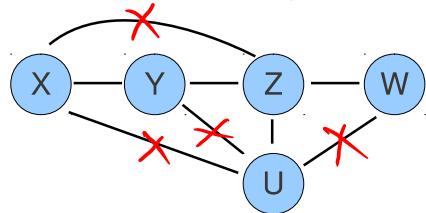
remove all edges having Sepset of size 1



$$X \perp\!\!\!\perp Z | Y \quad X \perp\!\!\!\perp U | Y \quad Y \perp\!\!\!\perp U | Z \quad W \perp\!\!\!\perp U | Z$$

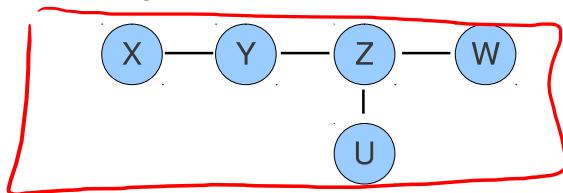


remove all edges $X - Y$ with $X \perp\!\!\!\perp Y | \emptyset$



$$X \perp\!\!\!\perp W \quad Y \perp\!\!\!\perp W$$

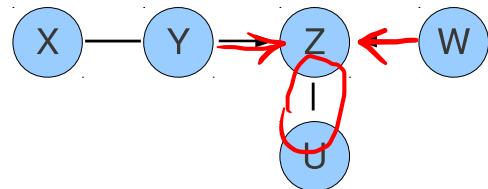
remove all edges having Sepset of size 1



$$X \perp\!\!\!\perp Z | Y \quad X \perp\!\!\!\perp U | Y \quad Y \perp\!\!\!\perp U | Z \quad W \perp\!\!\!\perp U | Z$$

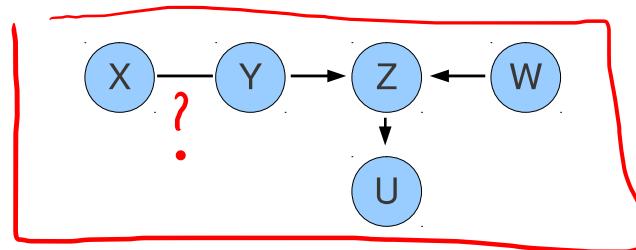


find v-structure



$$Z \notin S_{YW}$$

orient further edges (no further v-structure)



edge $X - Y$ remains undirected