

Introduction to Causal Inference

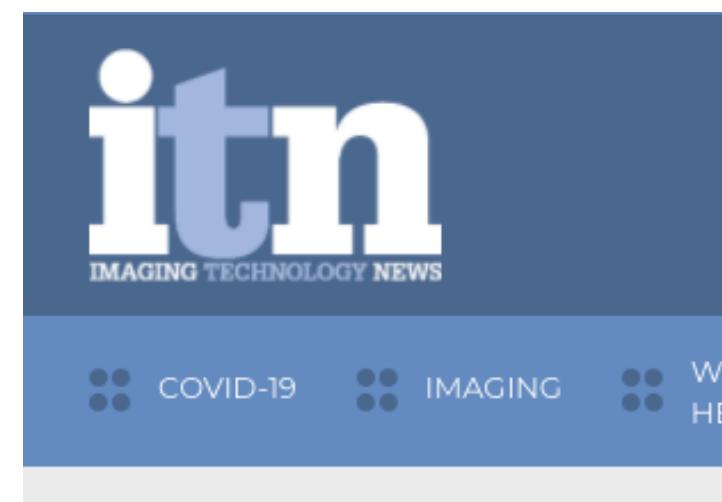
Adèle H. Ribeiro

<https://adele.github.io/> | adele.ribeiro@uni-marburg.de

Data Science in Biomedicine
Faculty of Mathematics and Computer Science
Philipps-Universität Marburg

Tropical Probabilistic AI School
February 2nd, 2024

Current Challenges in AI



nature

Explore content ▾ Journal information ▾ Publish with us ▾ Subscribe

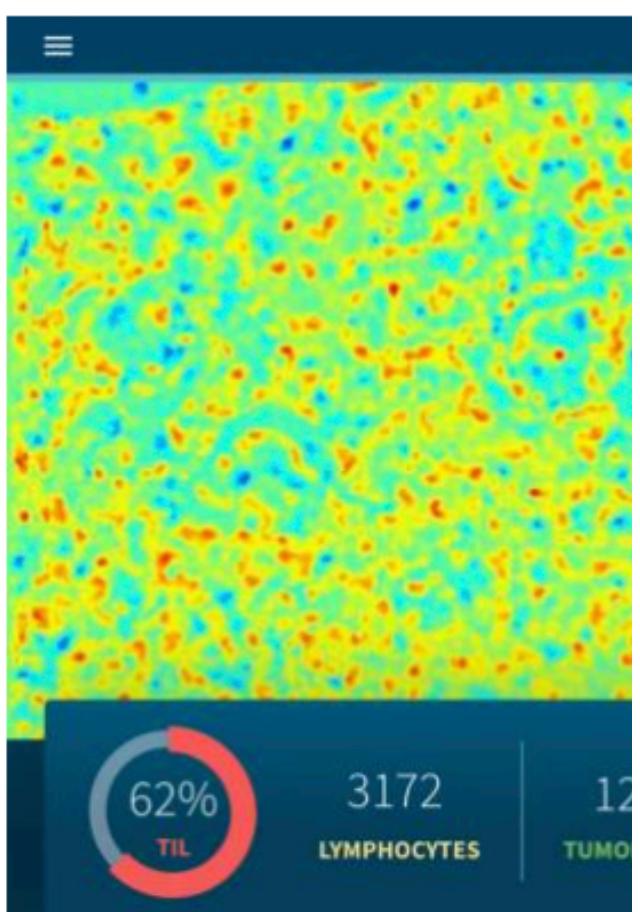
COVID-19 IMAGING WORKING GROUP

nature > outlook > article

NEWS | ARTIFICIAL INTELLIGENCE | MARCH 2023

Making the Role of AI in

Analysis system for the diagnosis of



Detection of tumor-infiltrating lymphocytes (TILs) generate a heatmap showing TILs (red) © of Klauschen/Charité

A fairer way forward for AI in health care

nature

Without careful implementation, AI could exacerbate inequality.

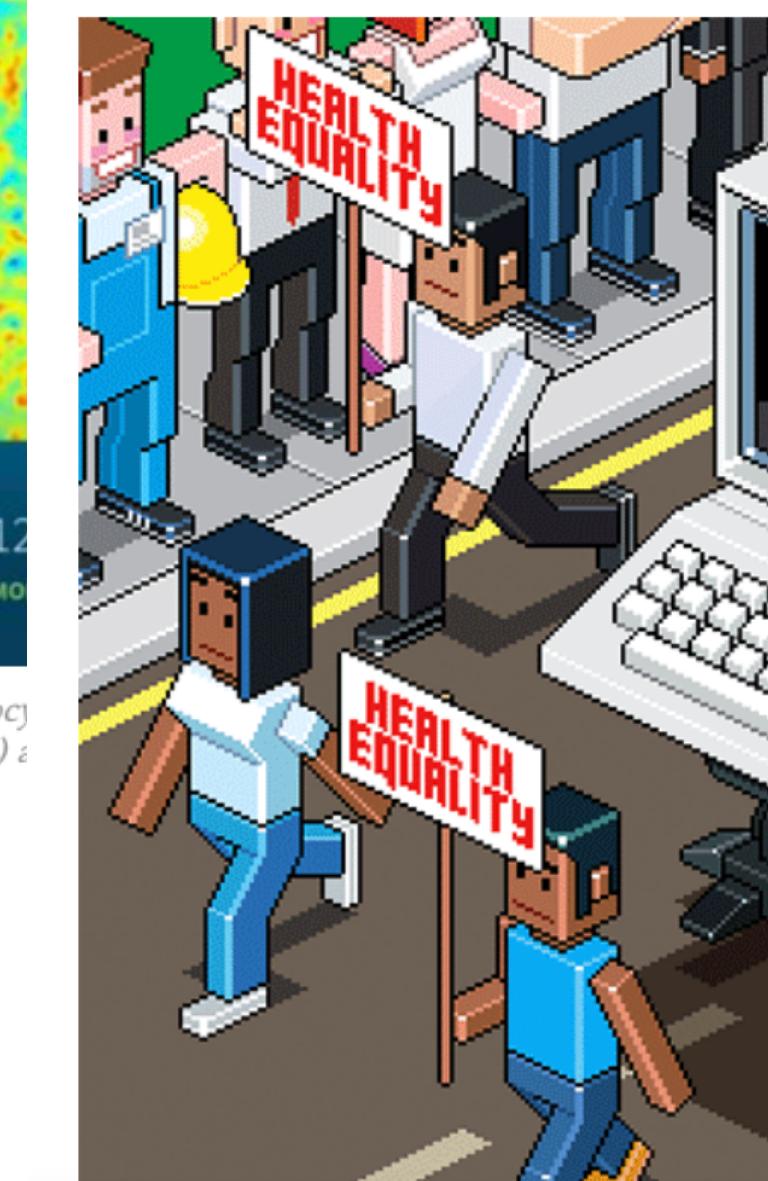
Linda Nordling

Explore content ▾

About the journal ▾

Publish with us ▾

Subscribe



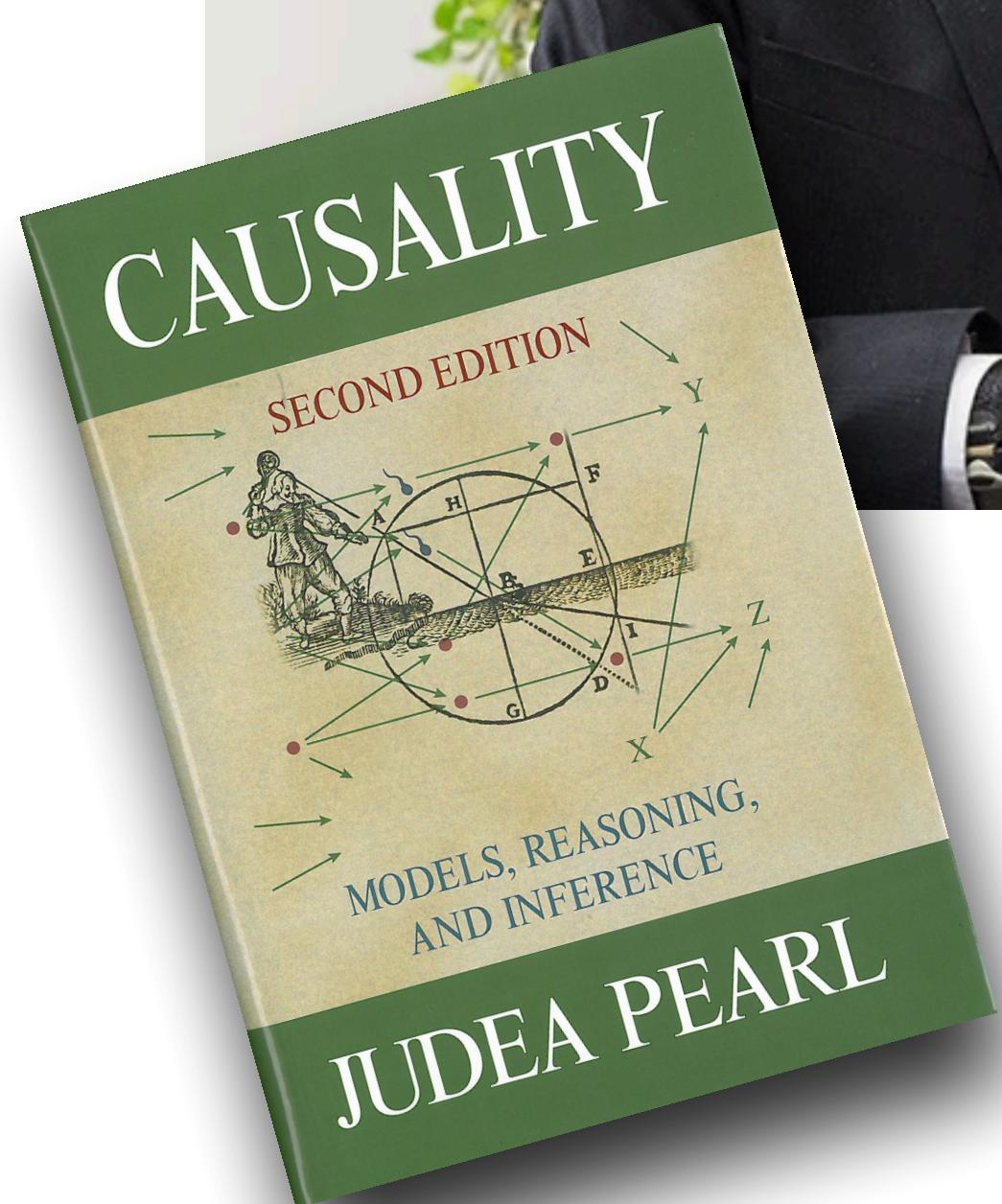
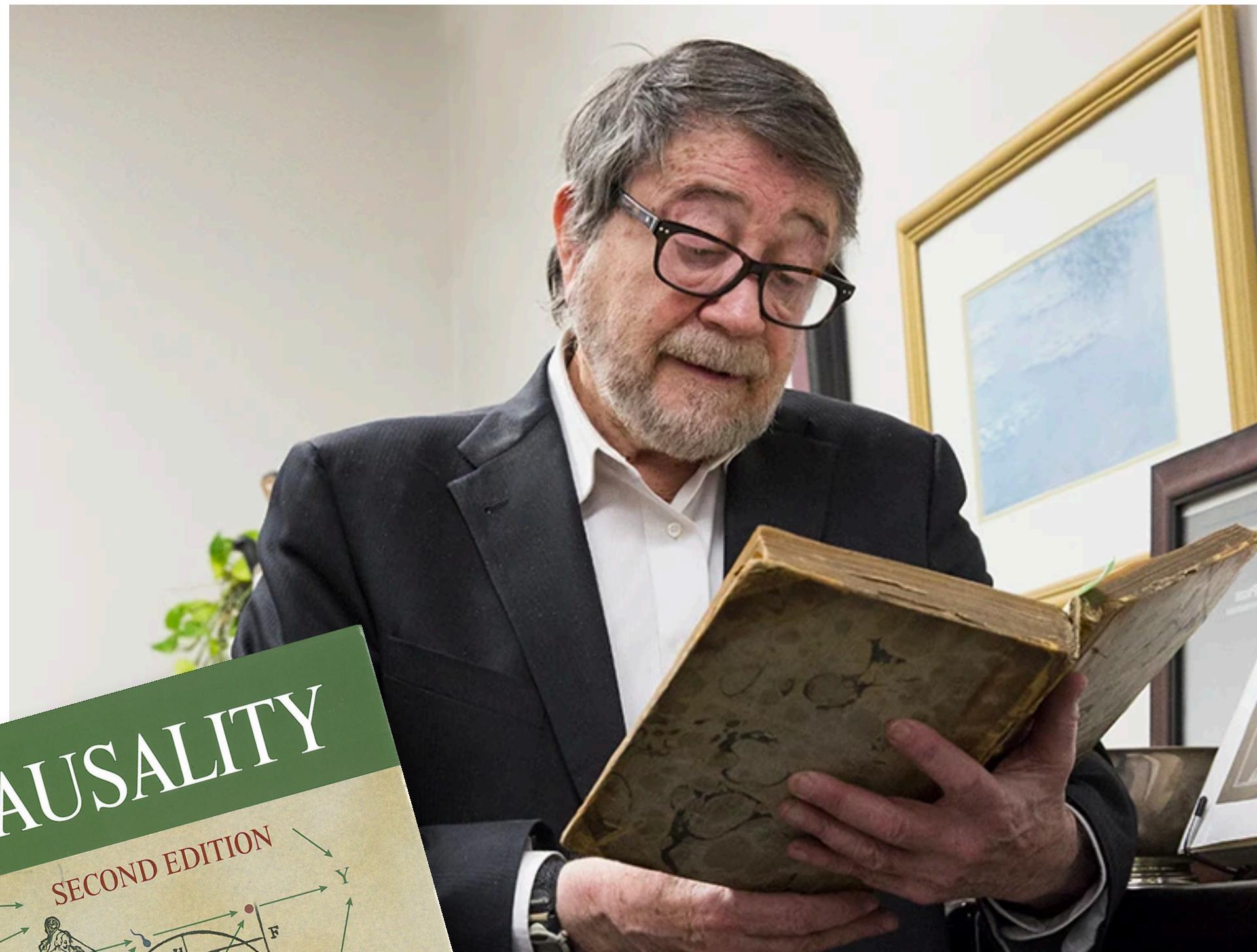
nature > outlook > article

OUTLOOK | 24 February 2023

Why artificial intelligence needs to understand consequences

A machine with a grasp of cause and effect could learn more like a human, through imagination and regret.

Judea Pearl – Causality



Director of the Cognitive Systems Laboratory at the University of California, Los Angeles.

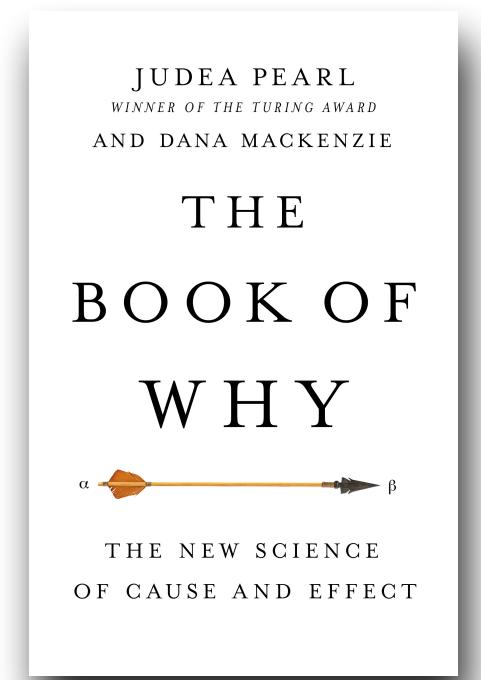
In 2011, he won the A. M. Turing Award (the highest distinction in computer science and a \$250,000 prize)

“for fundamental contributions to artificial intelligence through the development of a calculus for probabilistic and causal reasoning.”

— Association for Computing Machinery (ACM)

“Deep learning has instead given us machines with truly impressive abilities but no intelligence. The difference is profound and lies in the absence of a model of reality.”

— The Book of Why: The New Science of Cause and Effect



Guido W. Imbens & Joshua D. Angrist



Guido W. Imbens

Professor of Applied
Econometrics in
Stanford University

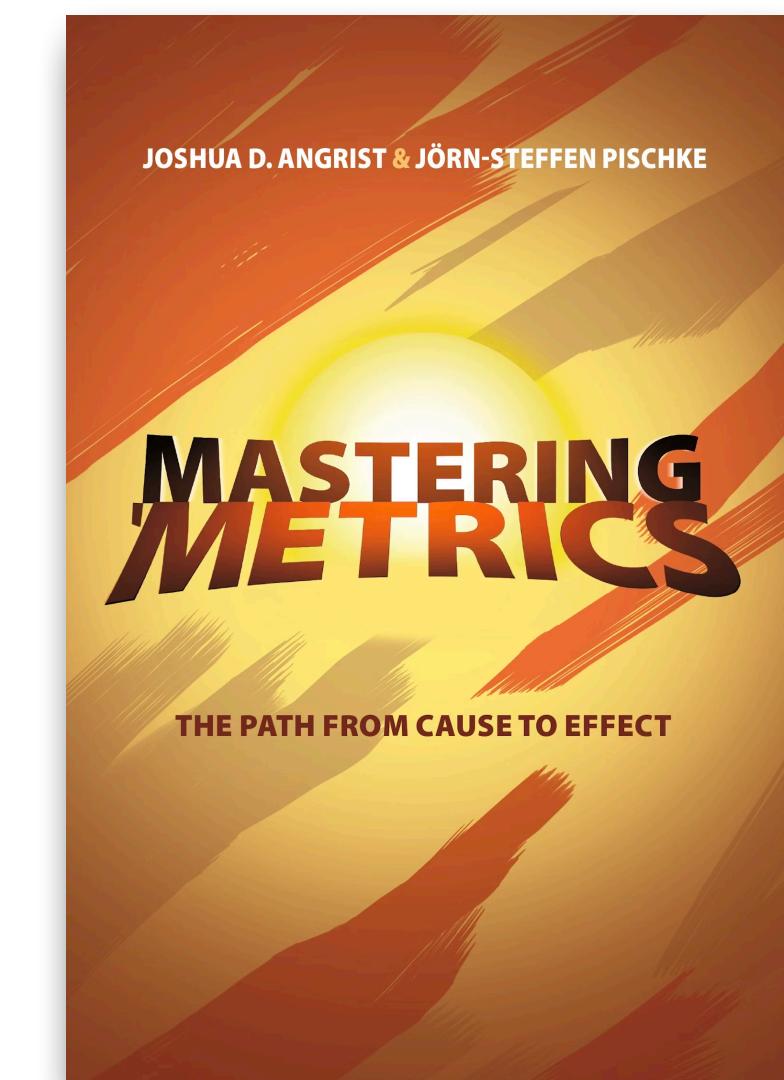
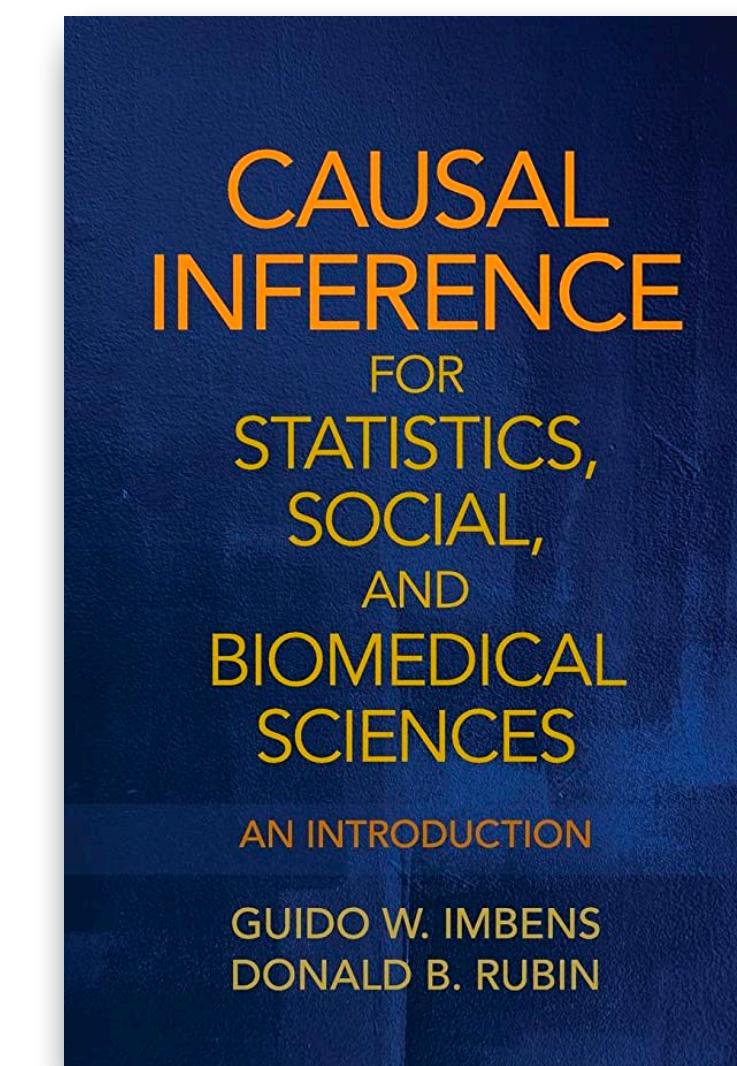


Joshua D. Angrist

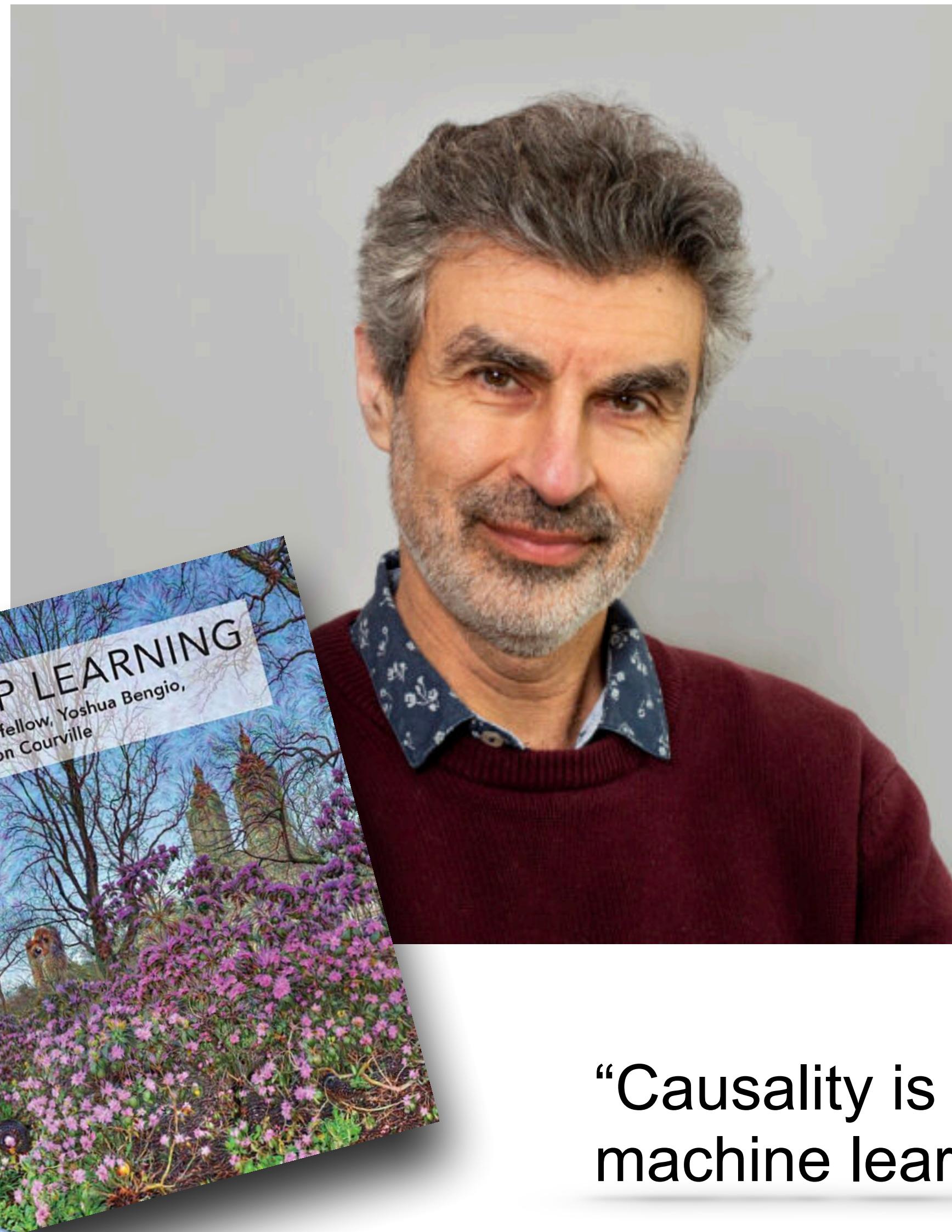
Professor of Economics
at the Massachusetts
Institute of Technology

In 2021, they won the Nobel Prize
in Economics (about \$1 million)

“for their methodological contributions
to the analysis of causal relationships”



Yoshua Bengio – Deep Learning



Professor at the University of Montreal, and the Founder and Scientific Director of Mila – Quebec AI Institute

In 2018, he won the A. M. Turing Award, with Geoffrey Hinton, and Yann LeCun

“for conceptual and engineering breakthroughs that have made deep neural networks a critical component of computing.”

— Association for Computing Machinery (ACM)

“Causality is very important for the next steps of progress of machine learning,” — interview with *IEEE Spectrum*.

Why causality is so important?

Causality allows important capabilities such as

Explainability: provides a better understanding of the underlying mechanisms

- **Causal Discovery**

Causal Effect: can determine the effect of *unrealized* interventions rather than just predicting an outcome (i.e., can distinguish between association and causation)

- **Causal Effect Identification and Estimation**

Fairness: captures and disentangles any mechanisms of discrimination that may be present, including direct, indirect-mediated, and indirect-confounded.

Generalizability: allows the transportability of causal effects across different domains.

Data Fusion: provides language and theory to cohesively combine prior knowledge and data from multiple and heterogeneous studies.

Causal Data Science

Goal is to develop language, criteria, and algorithms for:

- **Data-Fusion:** cohesively combining heterogenous datasets,
- **Causal Inference:** inferring the effects of interventions, and
- **Decision-Making:** making robust and generalizable decisions.



Causal inference and the data-fusion problem

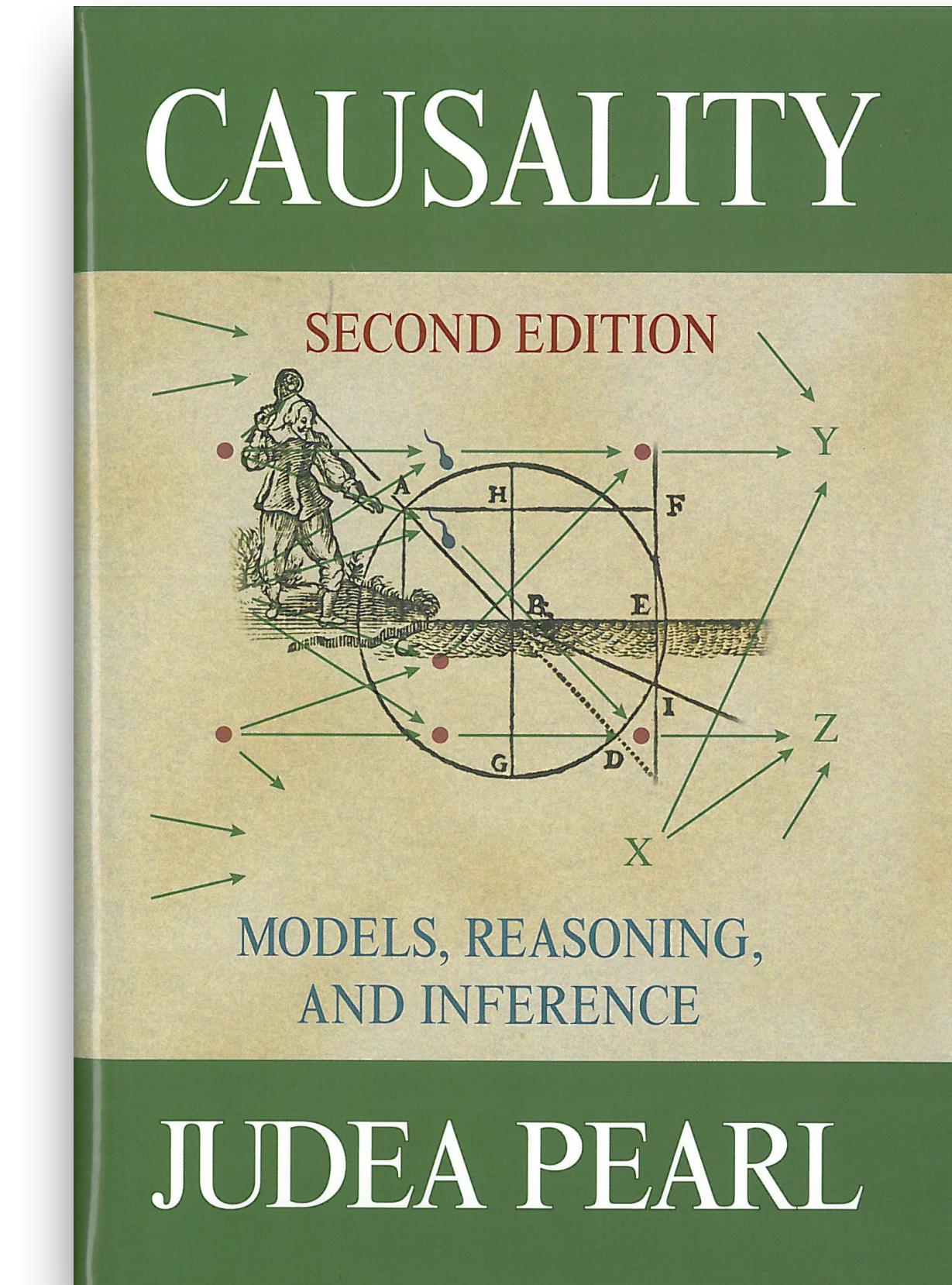
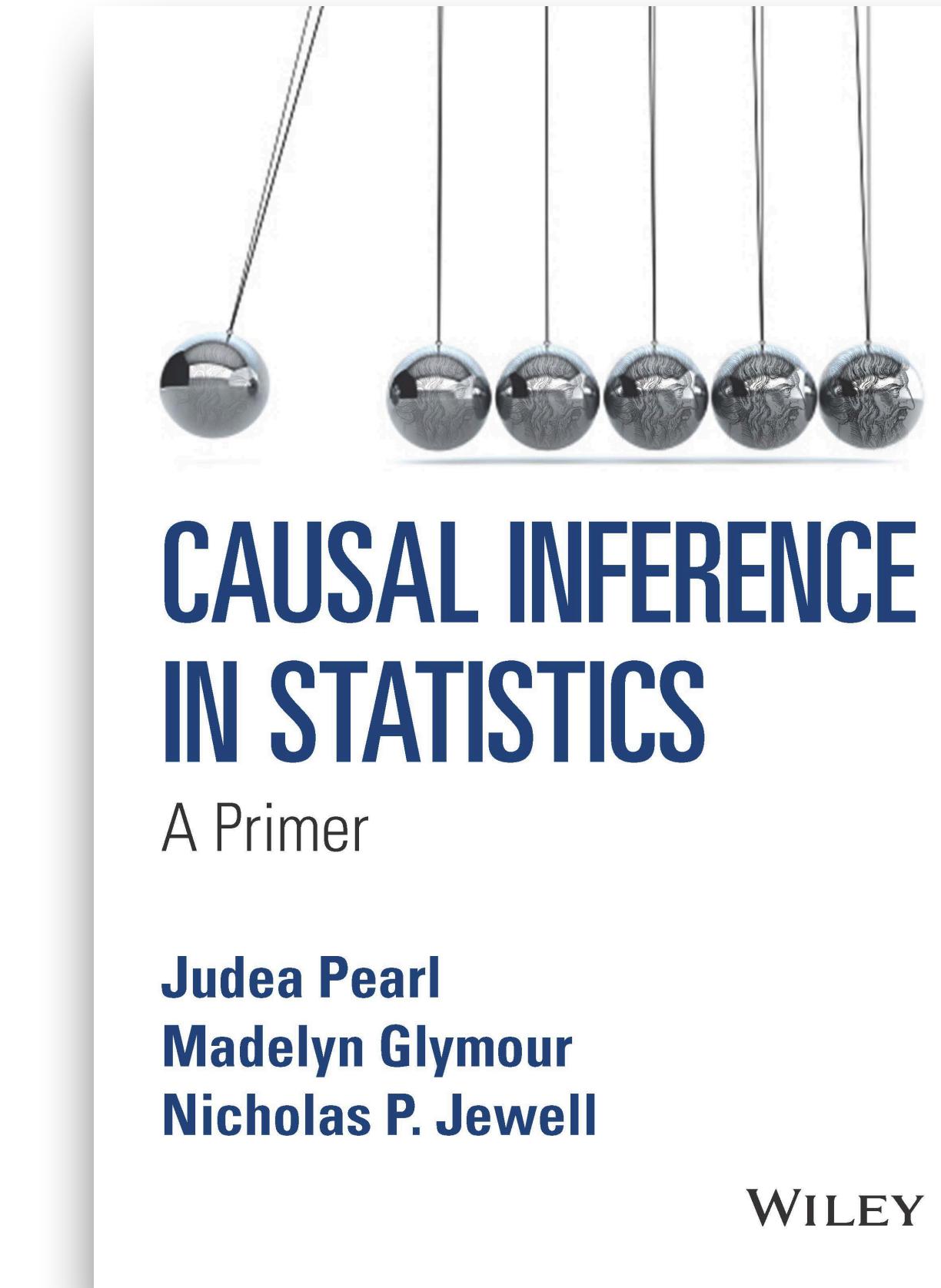
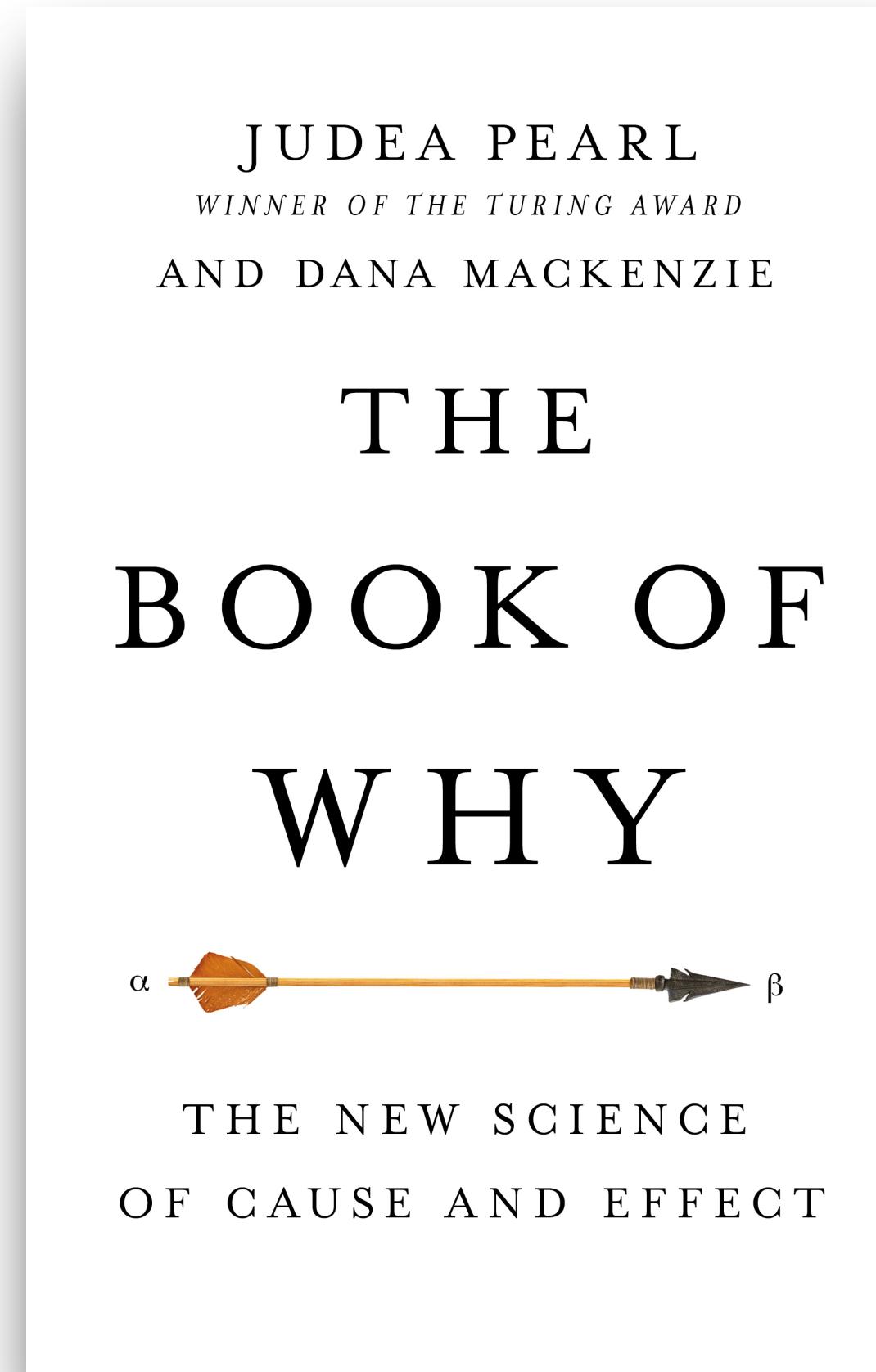
Elias Bareinboim^{a,b,1} and Judea Pearl^a

^aDepartment of Computer Science, University of California, Los Angeles, CA 90095; and ^bDepartment of Computer Science, Purdue University, West Lafayette, IN 47907

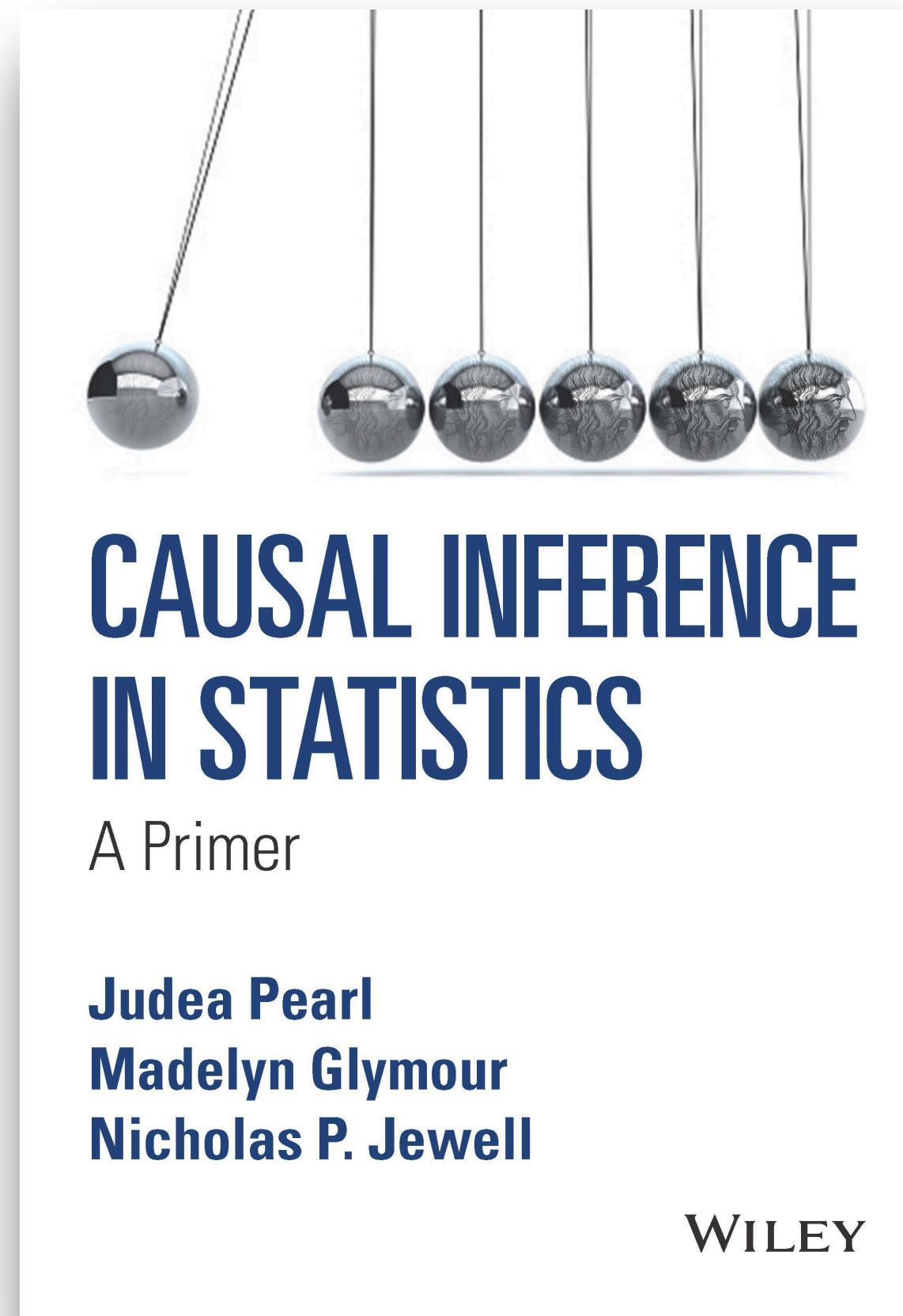
Edited by Richard M. Shiffrin, Indiana University, Bloomington, IN, and approved March 15, 2016 (received for review June 29, 2015)

<http://causalfusion.net>

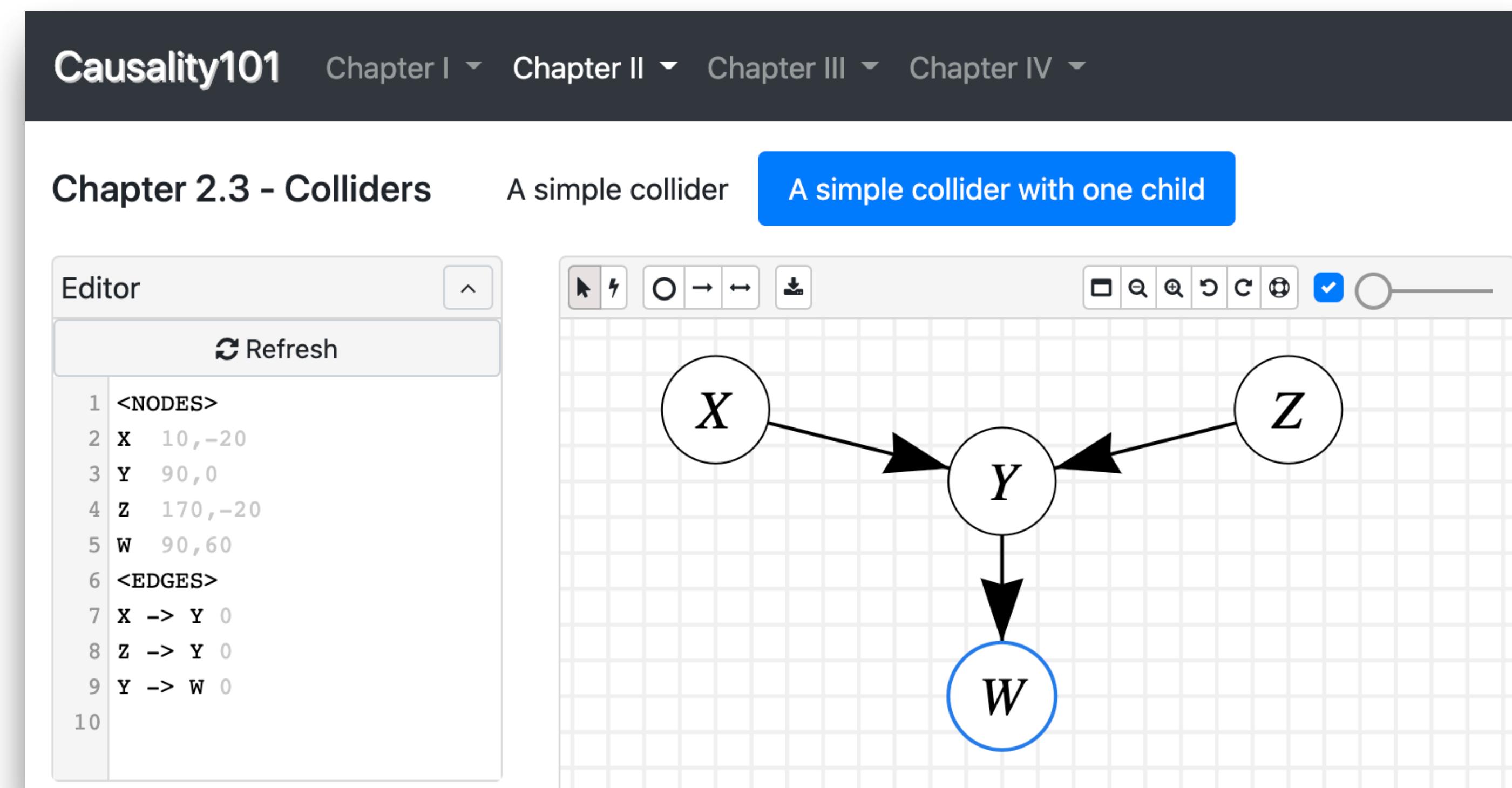
Causality Theory by Judea Pearl



Causality Theory by Judea Pearl



<https://causality101.net/>



Prediction vs Effect of Interventions

Statistical Association vs Causation

Predictive Tasks

Task: Can I guess how severe is a fire by **observing** the number of firefighters?

Yes!

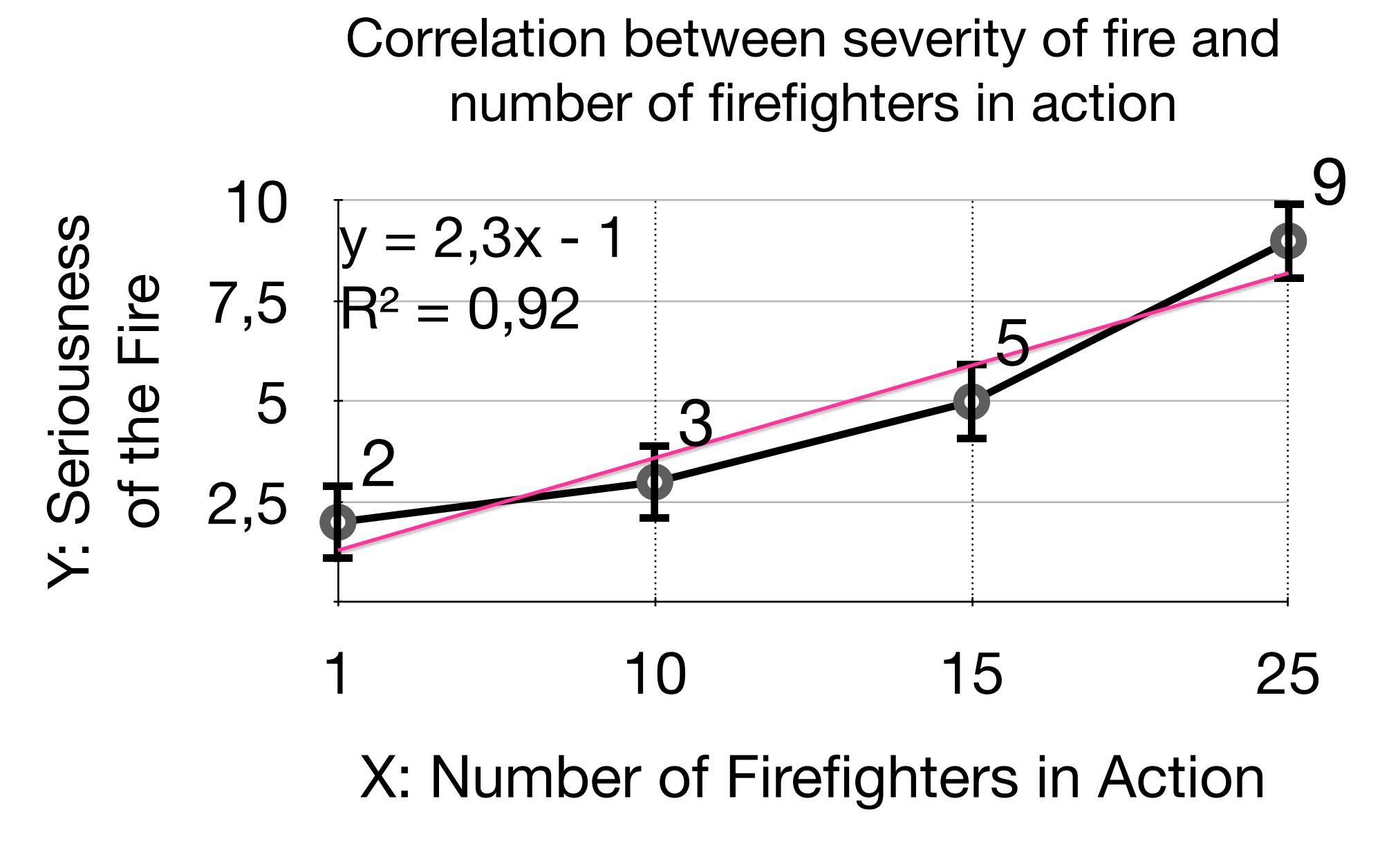
X : Number of firefighters in action
 Y : Severity of the (initial) fire

$\rho_{XY} \neq 0 \implies X \text{ is a good predictor of } Y$

$$P(Y = y | \textcolor{red}{X} = \textcolor{red}{x}) \neq P(Y = y)$$



**Observational
Probability Distribution**



Conclusion: The severity of the fire increases with the number of firefighters.

Prediction \Rightarrow Decision-Making?

Conclusion: The severity of the fire increases with the number of firefighters.

The fewer firefighters, the weaker the fire.



Should we decrease the number of firefighters to reduce the fire?

Causal Effect \neq Statistical Association

The **causal effect** of X on Y is a quantity that tells us how much Y changes after an intervention $do(X = x)$, e.g., $E[Y | do(X = x)]$.



If a different number of firefighters were dispatched, $do(X = x)$, would this change the expected severity of the initial fire Y ?

Causal Effect \equiv Effect of an Intervention

This quantity is derived from the probability of Y obtained after making the intervention $do(X = x)$, which is denoted by $P(Y | do(X = x))$ and is commonly referred as ***Interventional Distribution***.

Causal Effect \equiv Effect of an Intervention

This question can also be answered if knowledge about the underlying reality is available!

X : Number of firefighters in action

Y : (Initial) Severity of the fire

$$\begin{cases} X = f_X(Y, U_X) \\ Y = f_Y(U_Y) \end{cases}$$

**Underlying
Structural Causal Model (SCM)**

Y is not a function of X

In other words, **X is not a cause of Y**

In this case, $\forall x, E[Y | do(X = x)] = E[Y]$.

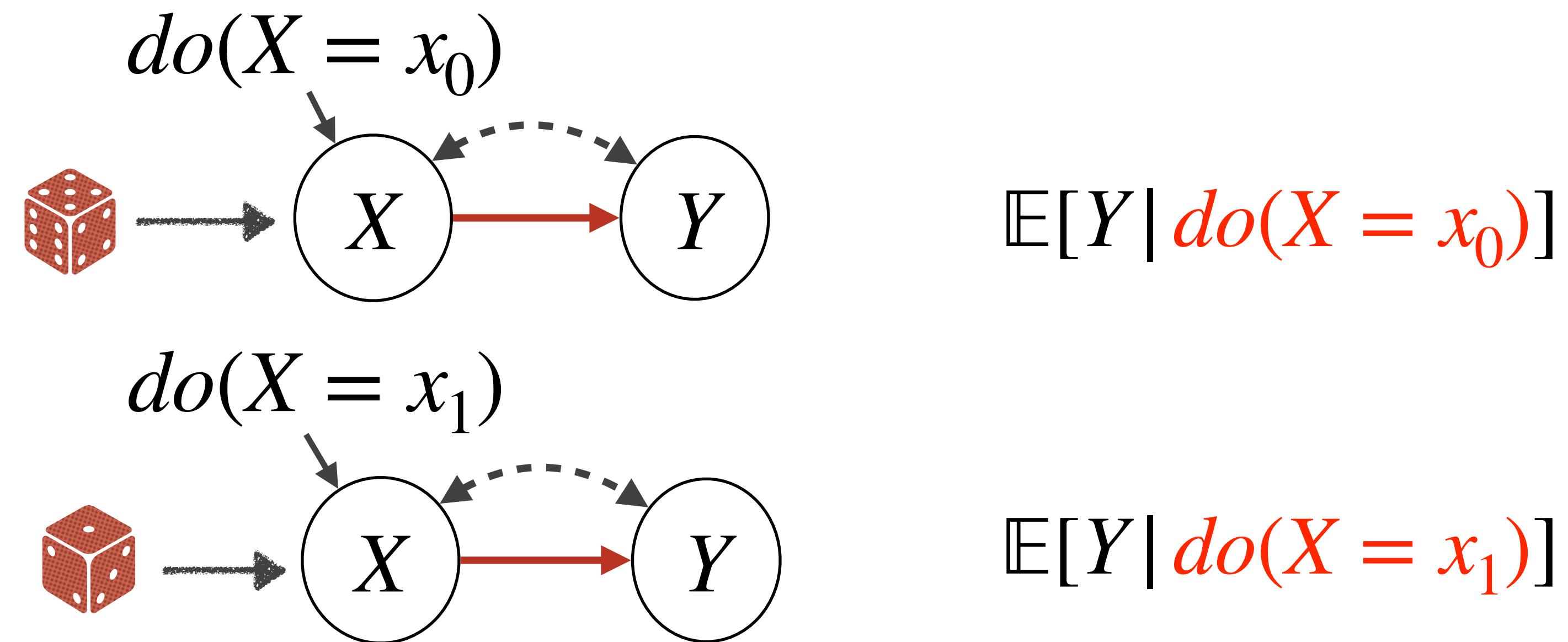
Changing the number of firefighters does not change the (initial) severity of the fire.

The action/intervention on X , $do(X = x)$ is independent of Y ,
i.e., $P(Y = y | do(X = x)) = P(Y = y)$

Randomized Experiments

One way to access the interventional distribution $P(Y | \text{do}(X = x))$ is through a *perfectly realized* Randomized Experiments / Control Trials (e.g. RCT):

Randomization of the
 X 's assignment



Average Causal Effect: $\mathbb{E}[Y | \text{do}(X = x_0)] - \mathbb{E}[Y | \text{do}(X = x_1)]$

What is the causal effect of the number of firefighters X on the severity of the initial fire Y ?

Pearl's Causal Hierarchy (PCH)

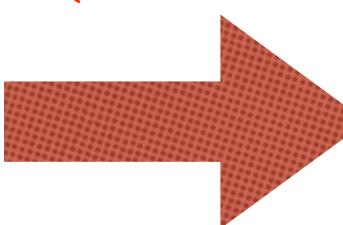
The Three Inferential Layers

What is induced by the SCM?

Observational SCM

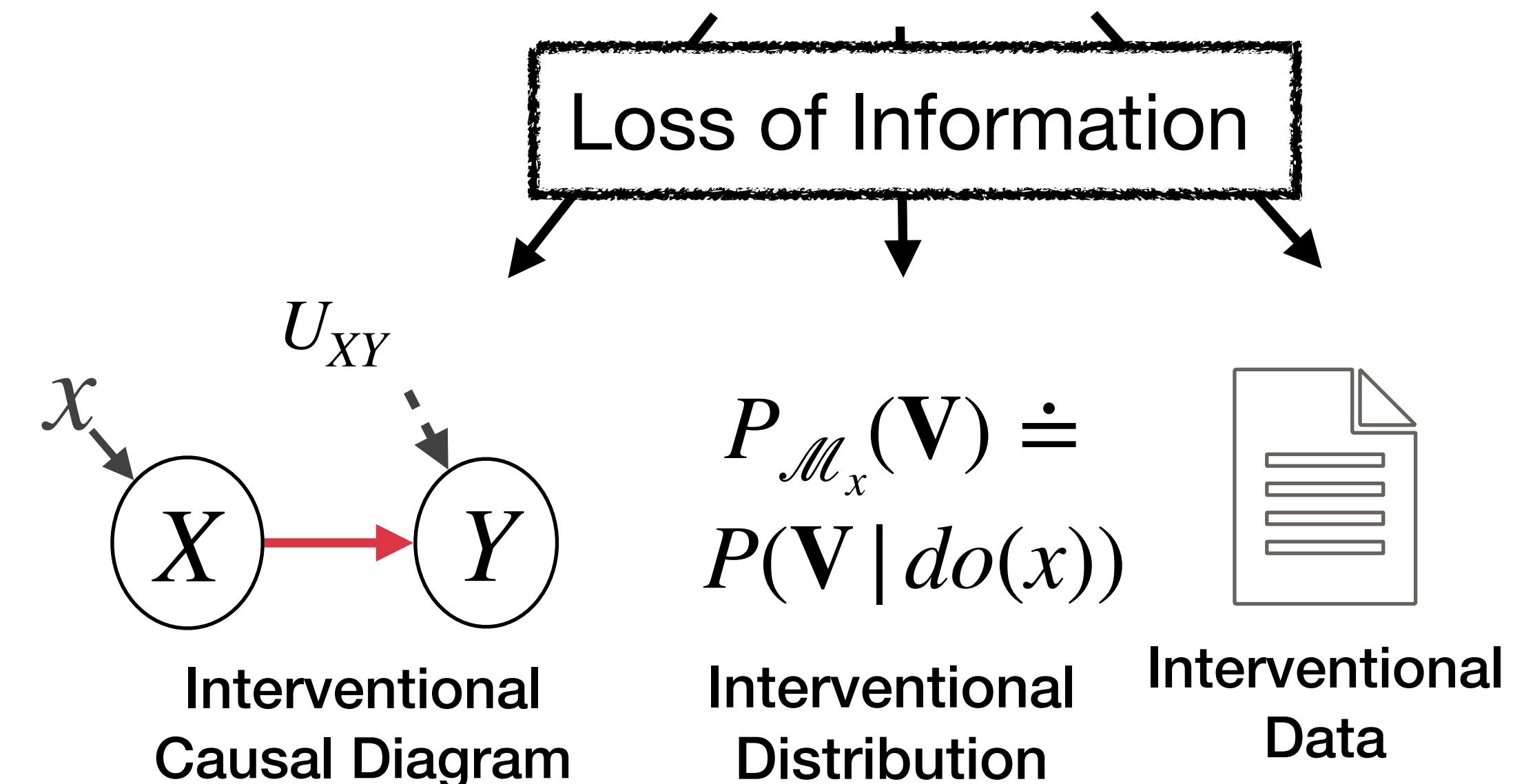
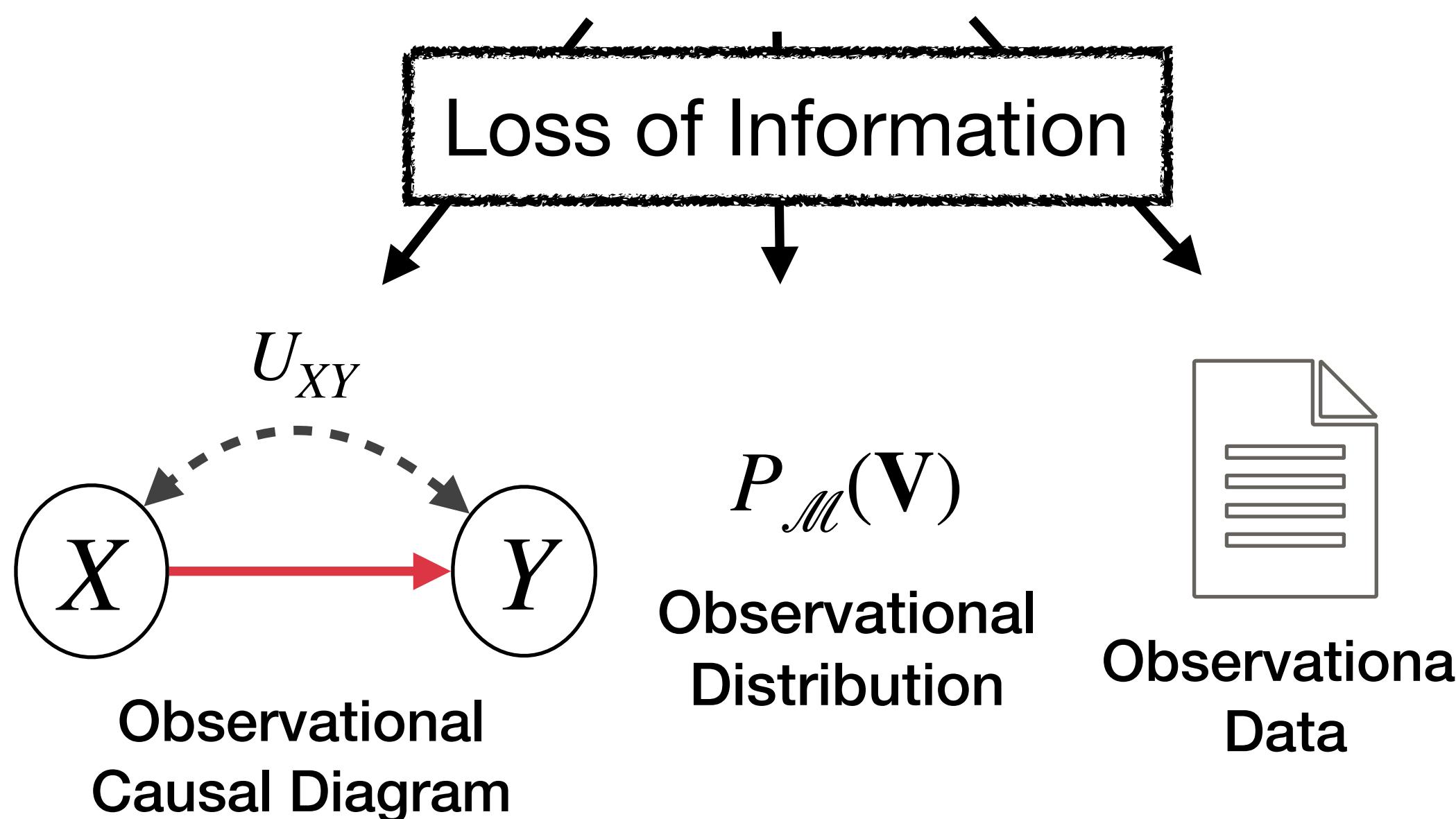
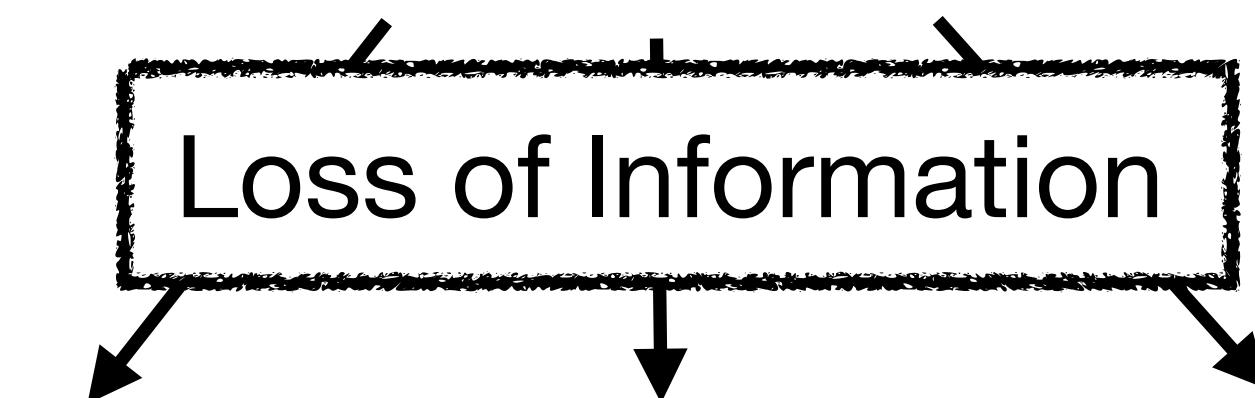
$$\mathcal{M} = \begin{cases} \mathbf{V} = \{X, Y\} \\ \mathbf{U} = \{U_{XY}, U_X, U_Y\} \\ \mathcal{F} = \begin{cases} X = f_X(U_X, U_{XY}) \\ Y = f_Y(X, U_Y, U_{XY}) \end{cases} \\ P(\mathbf{U}) \end{cases}$$

$do(X = x)$



Interventional SCM

$$\mathcal{M}_x = \begin{cases} \mathbf{V} = \{X, Y\} \\ \mathbf{U} = \{U_{XY}, U_X, U_Y\} \\ \mathcal{F} = \begin{cases} X = x \\ Y = f_Y(x, U_Y, U_{XY}) \end{cases} \\ P(\mathbf{U}) \end{cases}$$



Observational

Reality

Structural Causal Model (SCM)

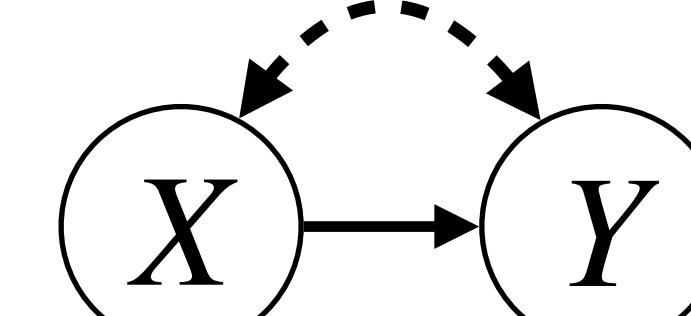
$$\mathcal{M} = \langle \mathbf{V}, \mathbf{U}, \mathcal{F}, P(\mathbf{u}) \rangle$$

$$\mathcal{M} = \begin{cases} \mathbf{V} = \{X, Y\} \\ \mathbf{U} = \{U_{XY}, U_X, U_Y\} \\ \mathcal{F} = \left\{ \begin{array}{l} X \leftarrow f_X(U_X, U_{XY}) \\ Y \leftarrow f_Y(X, U_Y, U_{XY}) \end{array} \right. \\ P(\mathbf{U}) \end{cases}$$

Structural Knowledge

Causal Diagram

G

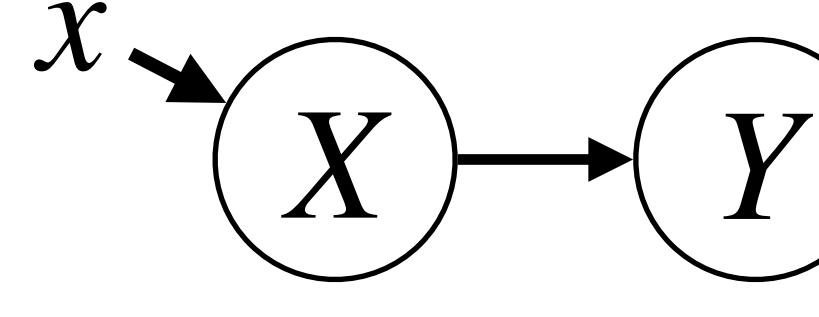


Interventional

Data

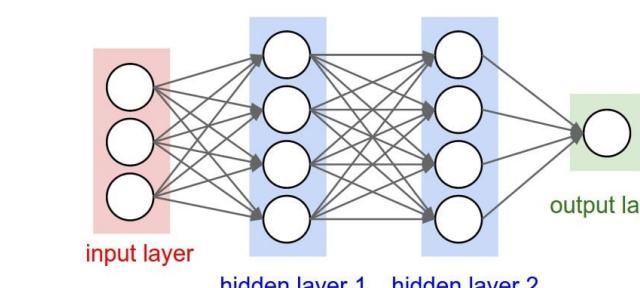
$$\mathcal{M}_x = \begin{cases} \mathbf{V} = \{X, Y\} \\ \mathbf{U} = \{U_{XY}, U_X, U_Y\} \\ \mathcal{F} = \left\{ \begin{array}{l} X \leftarrow x \\ Y \leftarrow f_Y(x, U_Y, U_{XY}) \end{array} \right. \\ P(\mathbf{U}) \end{cases}$$

$G_{\bar{X}}$



$$\hat{P}(Y | do(X = x)) = ?$$

X	Z	Y
---	---	---



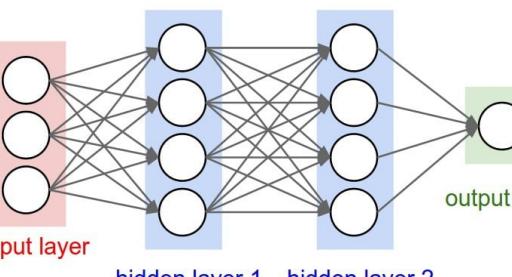
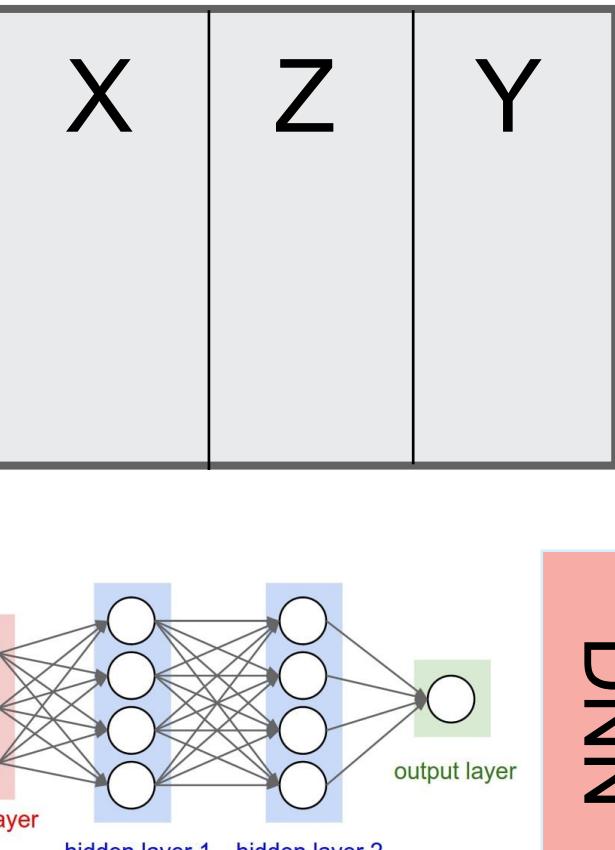
$$\hat{P}(Y | X = x)$$

Seeing

Doing



Observational

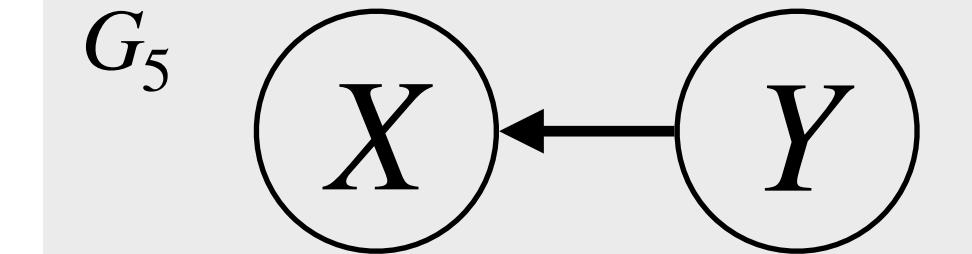
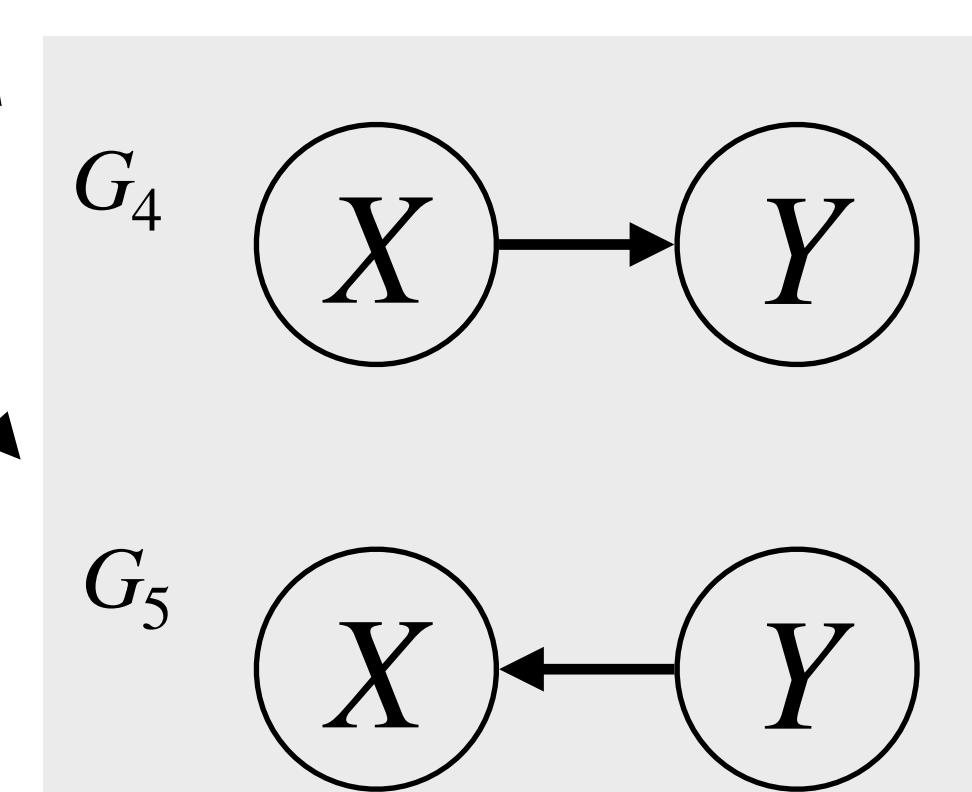
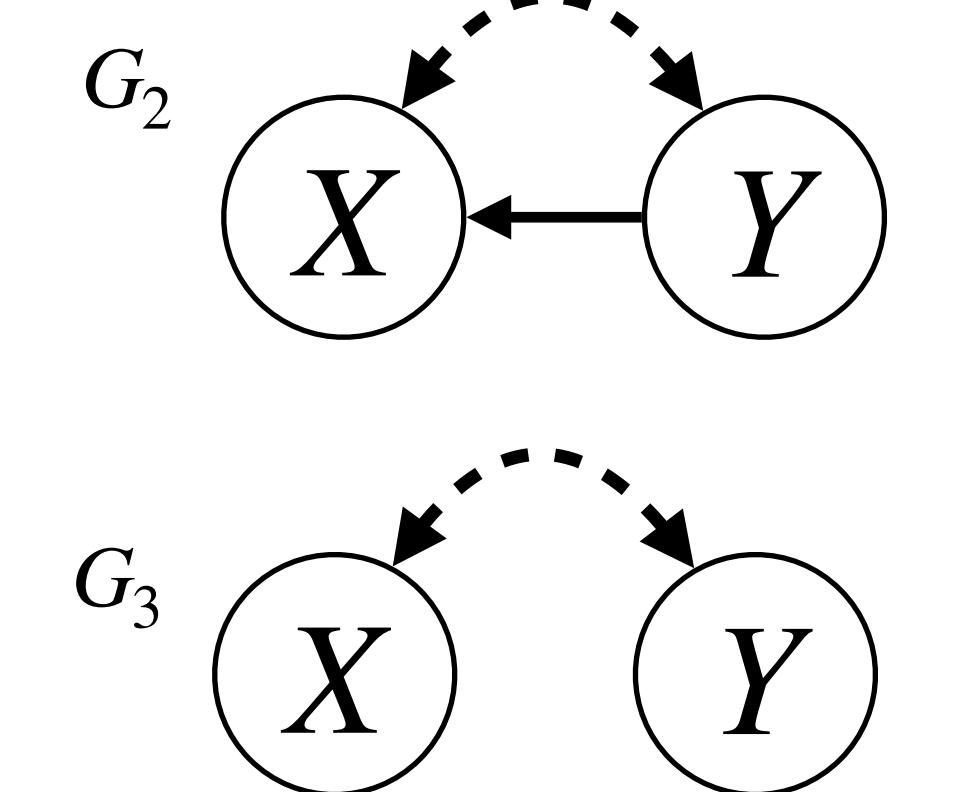
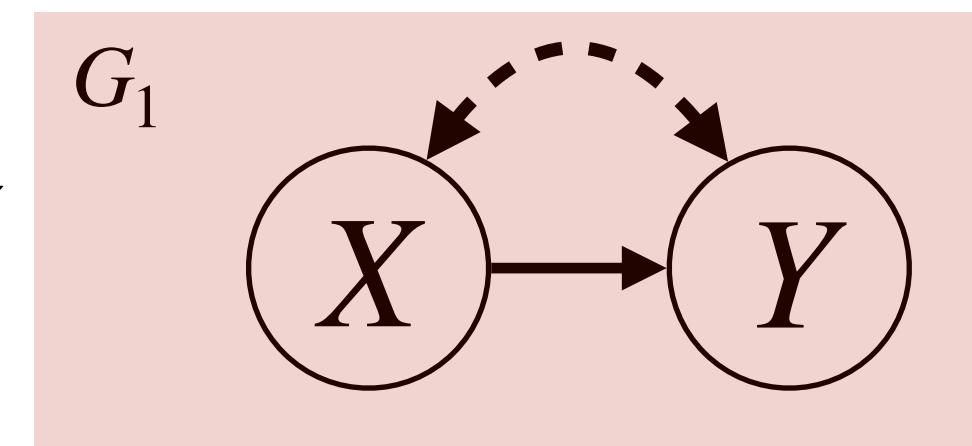


$$P(Y|X=x)$$

Data

Potential Causal Diagrams

Potential SCMs



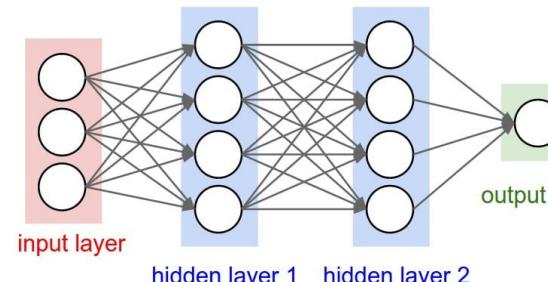
Encoded Knowledge / Assumptions

- $\mathcal{M}_{11} = \langle \mathbf{V}, \mathbf{U}_1, \mathcal{F}_{11}, P_{11}(\mathbf{u}_1) \rangle$
- \vdots
- $\mathcal{M}_{1k_1} = \langle \mathbf{V}, \mathbf{U}_1, \mathcal{F}_{1k_1}, P_{1k_1}(\mathbf{u}_1) \rangle$
- True Model
- $\mathcal{M}_{21} = \langle \mathbf{V}, \mathbf{U}_2, \mathcal{F}_{21}, P_{21}(\mathbf{u}_2) \rangle$
- \vdots
- $\mathcal{M}_{2k_2} = \langle \mathbf{V}, \mathbf{U}_2, \mathcal{F}_{2k_2}, P_{2k_2}(\mathbf{u}_2) \rangle$
- $\mathcal{M}_{31} = \langle \mathbf{V}, \mathbf{U}_3, \mathcal{F}_{31}, P_{31}(\mathbf{u}_3) \rangle$
- \vdots
- $\mathcal{M}_{3k_3} = \langle \mathbf{V}, \mathbf{U}_3, \mathcal{F}_{3k_3}, P_{3k_3}(\mathbf{u}_3) \rangle$
- $\mathcal{M}_{41} = \langle \mathbf{V}, \mathbf{U}_4, \mathcal{F}_{41}, P_{41}(\mathbf{u}_4) \rangle$
- \vdots
- $\mathcal{M}_{4k_4} = \langle \mathbf{V}, \mathbf{U}_4, \mathcal{F}_{4k_4}, P_{4k_4}(\mathbf{u}_4) \rangle$
- $\mathcal{M}_{51} = \langle \mathbf{V}, \mathbf{U}_5, \mathcal{F}_{51}, P_{51}(\mathbf{u}_5) \rangle$
- \vdots
- $\mathcal{M}_{5k_5} = \langle \mathbf{V}, \mathbf{U}_5, \mathcal{F}_{5k_5}, P_{5k_5}(\mathbf{u}_5) \rangle$

Parametrization

Observational

X	Z	Y
---	---	---



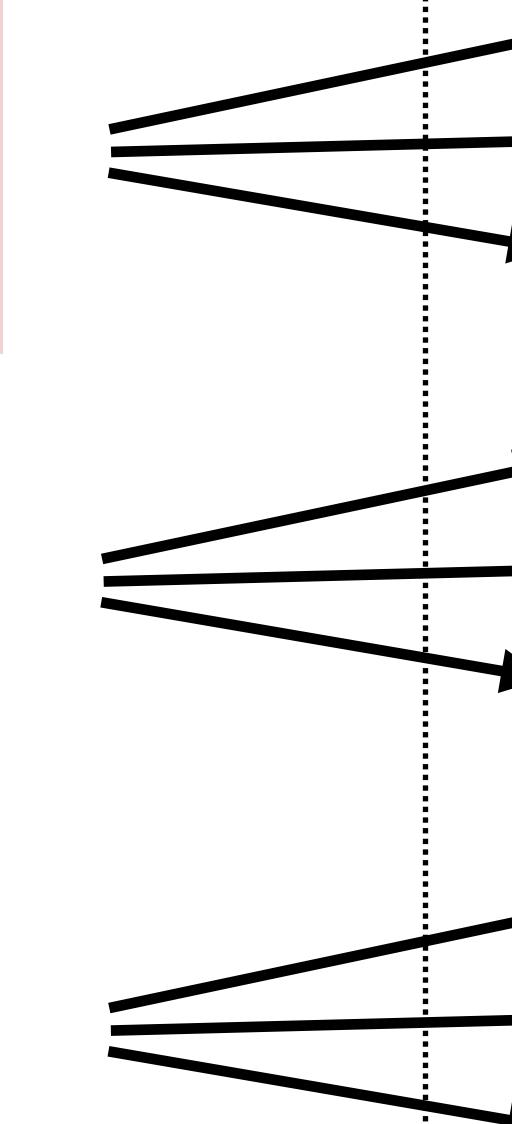
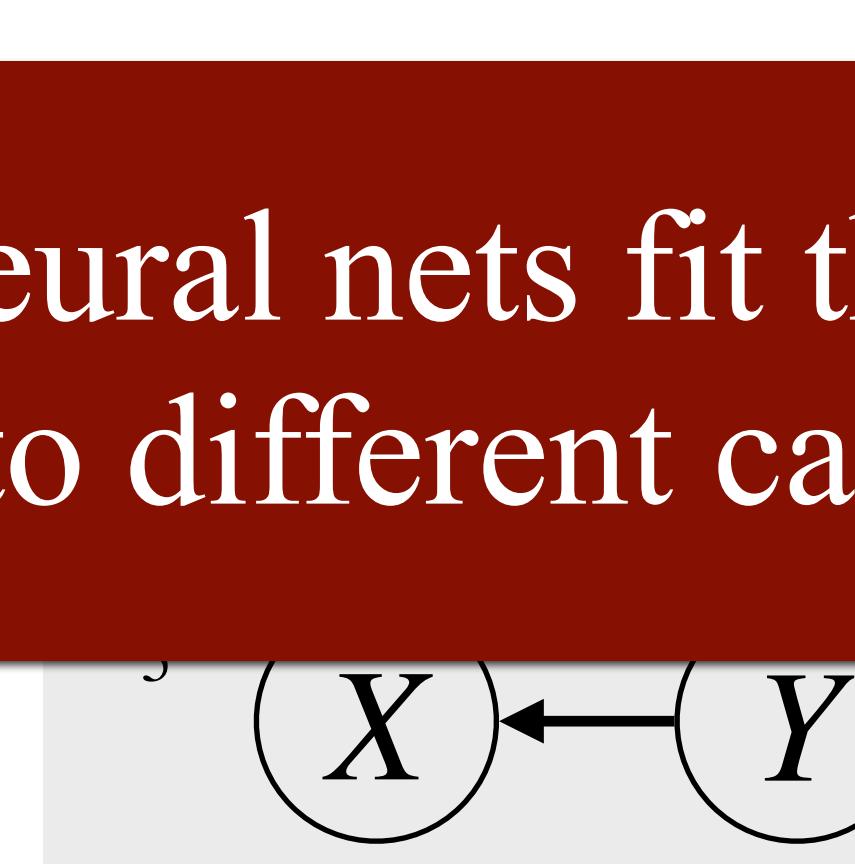
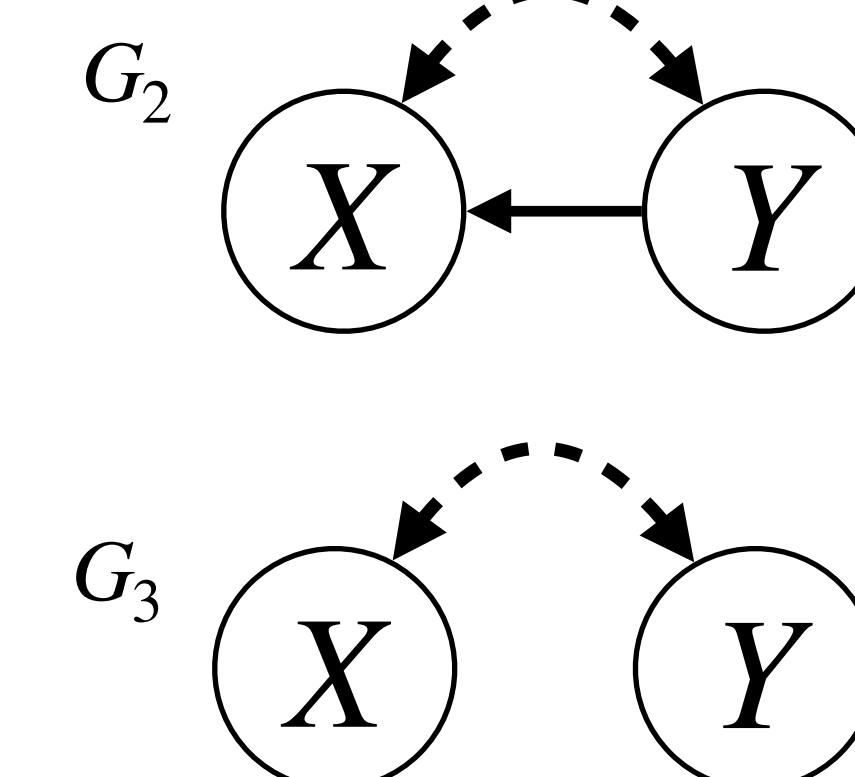
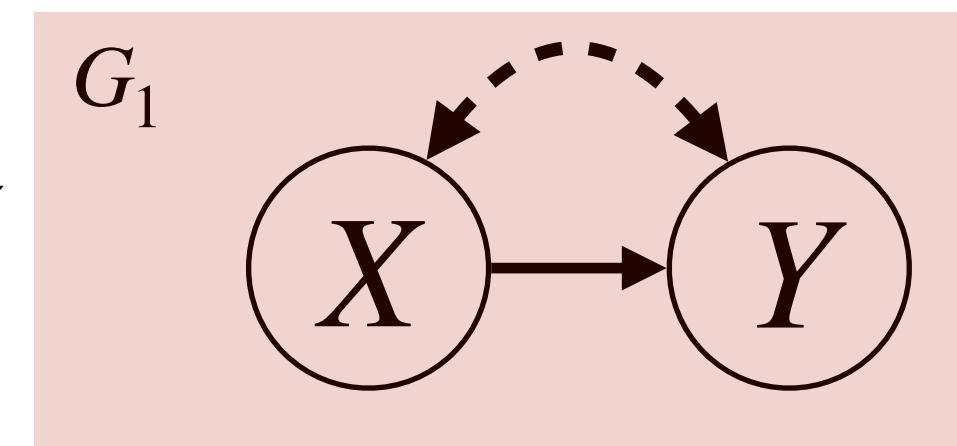
$$P(Y|X = r)$$

Multiple neural nets fit the data equally well,
leading to different causal explanations!

Data

Potential Causal Diagrams

Potential SCMs



$$\mathcal{M}_{11} = \langle \mathbf{V}, \mathbf{U}_1, \mathcal{F}_{11}, P_{11}(\mathbf{u}_1) \rangle$$

⋮

$$\mathcal{M}_{1k_1} = \langle \mathbf{V}, \mathbf{U}_1, \mathcal{F}_{1k_1}, P_{1k_1}(\mathbf{u}_1) \rangle$$

True Model

$$\mathcal{M}_{21} = \langle \mathbf{V}, \mathbf{U}_2, \mathcal{F}_{21}, P_{21}(\mathbf{u}_2) \rangle$$

⋮

$$\mathcal{M}_{2k_2} = \langle \mathbf{V}, \mathbf{U}_2, \mathcal{F}_{2k_2}, P_{2k_2}(\mathbf{u}_2) \rangle$$

$$\mathcal{M}_{31} = \langle \mathbf{V}, \mathbf{U}_3, \mathcal{F}_{31}, P_{31}(\mathbf{u}_3) \rangle$$

⋮

$$\mathcal{M}_{3k_3} = \langle \mathbf{V}, \mathbf{U}_3, \mathcal{F}_{3k_3}, P_{3k_3}(\mathbf{u}_3) \rangle$$

$$\mathcal{M}_{41} = \langle \mathbf{V}, \mathbf{U}_4, \mathcal{F}_{41}, P_{41}(\mathbf{u}_4) \rangle$$

$$\mathcal{M}_{4k_4} = \langle \mathbf{V}, \mathbf{U}_4, \mathcal{F}_{4k_4}, P_{4k_4}(\mathbf{u}_4) \rangle$$

$$\mathcal{M}_{51} = \langle \mathbf{V}, \mathbf{U}_5, \mathcal{F}_{51}, P_{51}(\mathbf{u}_5) \rangle$$

$$\mathcal{M}_{5k_5} = \langle \mathbf{V}, \mathbf{U}_5, \mathcal{F}_{5k_5}, P_{5k_5}(\mathbf{u}_5) \rangle$$

Parametrization

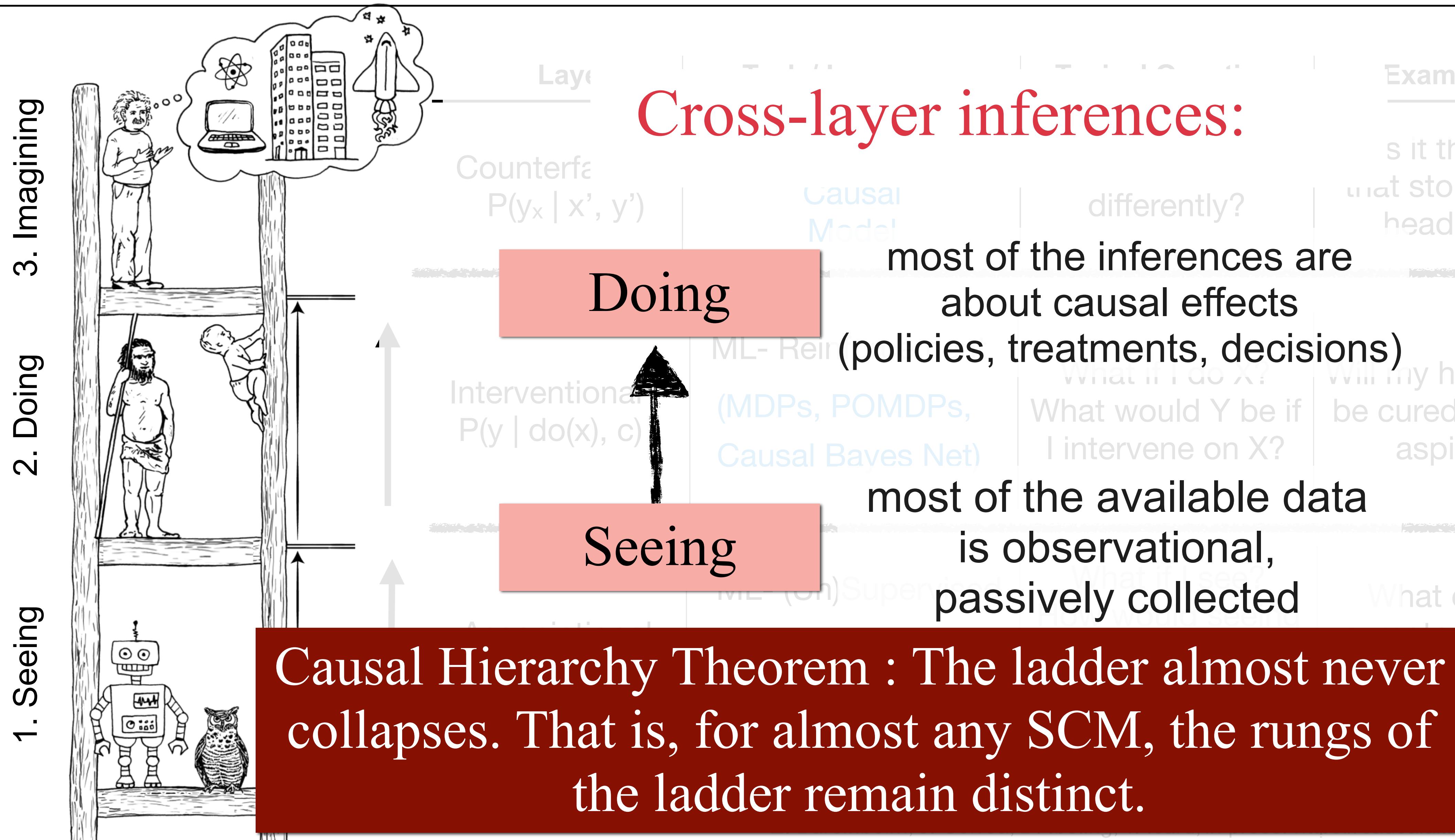
Encoded Knowledge / Assumptions

Ladder of Causation

Layer	Task / Language	Typical Question	Examples
3. Imagining	Counterfactual $P(y_x x', y')$	Structural Causal Model	What if I had acted differently? Was it the aspirin that stopped my headache?
2. Doing	Interventional $P(y \text{do}(x), c)$	ML- Reinforcement (MDPs, POMDPs, Causal Bayes Net)	What if I do X? What would Y be if I intervene on X? Will my headache be cured if I take aspirin?
1. Seeing	Associational $P(y x)$	ML- (Un)Supervised (Decision trees, Deep nets, ...)	What if I see? How would seeing X change my belief in Y? What does a symptom tell us about the disease?

* Book of Why & On Pearl's Hierarchy and the Foundations of Causal Inference, E. Bareinboim, J. Correa, D. Ibeling, T. Icard, in press. <https://causalai.net/r60.pdf> 21

Ladder of Causation



Cross-layer inferences:

most of the inferences are
about causal effects
(policies, treatments, decisions)

most of the available data
is observational,
passively collected

3. Imagining

2. Doing

1. Seeing

Layer

Counterfactual
 $P(y_x | x', y')$

Causal
Model

ML- Reinforcement
Learning

Interventional
 $P(y | \text{do}(x), c)$

Causal Bayes Net

differently?

Examples

Was it the aspirin
that stopped my
headache?

What would Y be if
I intervene on X?
What if I do X?
What will my headache
be cured if I take
aspirin?

What does a
twin tell us
about the
cause?

difference,
et/r60.pdf

Structural Causal Model (SCM)

THE DATA GENERATING MODEL

Structural Causal Model (SCM)

Definition: A structural causal model \mathcal{M} (or, data generating model) is a tuple $\langle \mathbf{V}, \mathbf{U}, \mathcal{F}, P(\mathbf{u}) \rangle$, where

- $\mathbf{V} = \{V_1, \dots, V_n\}$: are endogenous variables
- $\mathbf{U} = \{U_1, \dots, U_m\}$: are exogenous variables
- $\mathcal{F} = \{f_1, \dots, f_n\}$: are functions determining \mathbf{V} , i.e., $v_i \leftarrow f_i(pa_i, u_i)$ where $Pa_i \subseteq \mathbf{V}$ are endogenous causes (parents) of V_i and $U_i \subseteq \mathbf{U}$ are exogenous causes of V_i .
- $P(\mathbf{U})$ is the probability distribution over \mathbf{U} .

Assumption: \mathcal{M} is recursive, i.e., there are no feedback (cyclic) mechanisms.

Structural Equation Model (SEM)

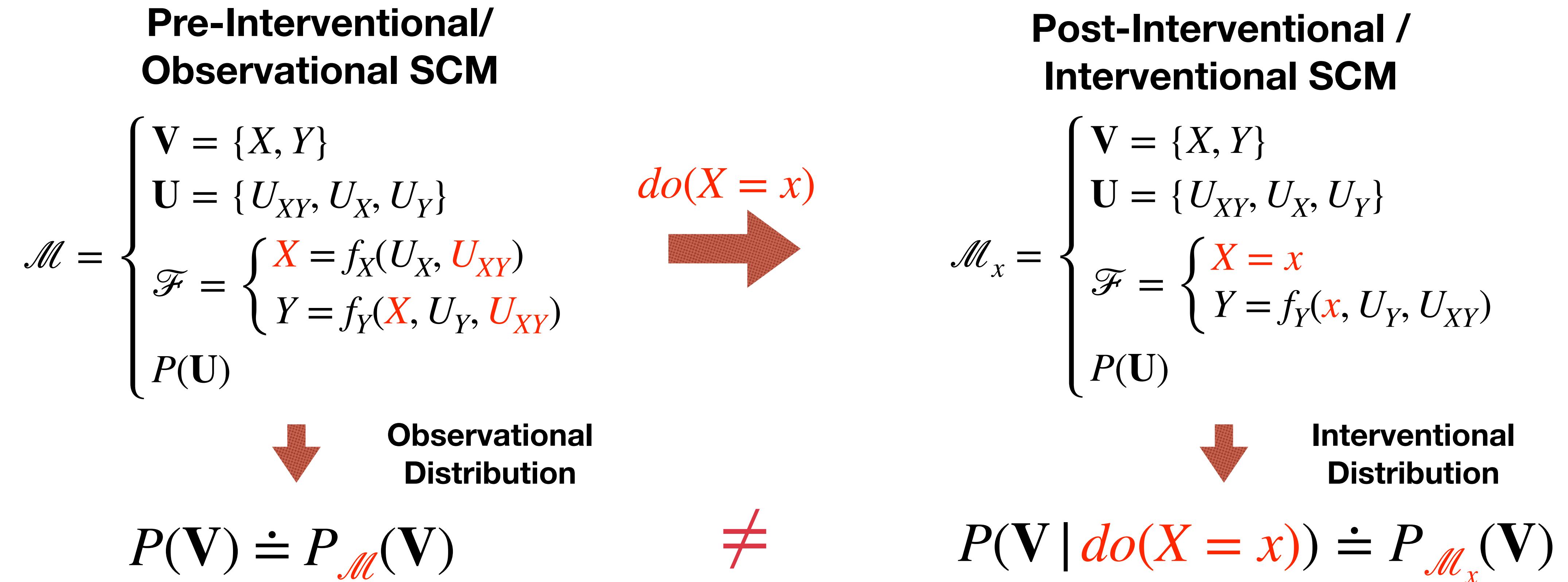
$$\mathcal{M} = \begin{cases} V = \{X, Y, Z\} \\ U = \{\epsilon_X, \epsilon_Y, \epsilon_Z\} \\ \mathcal{F} = \begin{cases} Z = \beta_{Z0} + \epsilon_Z \\ X = \beta_{X0} + \beta_{XZ}Z + \epsilon_X \\ Y = \beta_{Y0} + \beta_{YZ}Z + \beta_{YX}X + \epsilon_Y \end{cases} \\ U \sim \mathcal{N}\left(\mathbf{0}, \Sigma = \begin{bmatrix} \sigma_X & 0 & 0 \\ 0 & \sigma_Y & 0 \\ 0 & 0 & \sigma_Z \end{bmatrix}\right) \end{cases}$$

- **Linear functions**
- **Normal distribution**
- **Markovianity / Causal Sufficiency:**
Error terms in U are independent of each other (diagonal covariance matrix).

Full specification of an SCM requires parametric and distributional assumptions.

Estimation of such models usually requires strong assumptions (e.g., Markovianity).

Effect of Interventions in SCMs



Can we **predict** better the value of Y after
observing que $X = x$?

$P(Y = y | X = x) \neq P(Y = y) \implies X \text{ is } \text{correlated} \text{ to } Y$

Can we **predict** better the value of Y after
making an intervention $do(X = x)$?

$\exists x \text{ s.t. } P_{\mathcal{M}_x}(Y = y) \neq P(Y = y) \implies X \text{ is a cause of } Y$ 26

Effect of Interventions in SCMs

Observational
Distribution

\neq

$$P(\mathbf{V}) \doteq P_{\mathcal{M}}(\mathbf{V}) = \sum_{\mathbf{u}} \prod_{V_i \in \mathbf{V}} P(v_i | pa_i, u_i) P(\mathbf{u})$$

$do(X = x)$

Factorization obtained by
Chain Rule and
conditional independencies
implied by the SCM \mathcal{M} .

Interventional
Distribution

$$P(\mathbf{V} | do(X = x)) \doteq P_{\mathcal{M}_x}(\mathbf{V})$$

Truncated factorization
implied by the SCM \mathcal{M}_x .

$$= \sum_{\mathbf{u}} \prod_{V_i \in \mathbf{V} \setminus \mathbf{X}} P(v_i | pa_i, u_i) P(\mathbf{u}) \Big|_{\mathbf{X}=\mathbf{x}}$$

SCM: Encoder of Functional Knowledge

The knowledge required to fully specify an SCM is usually *unavailable* in practice.

Is it possible to identify the effect of interventions from *observational* data without fully specifying the SCM (i.e., in a non-parametric fashion)?



Yes, with structural knowledge encoded as a causal diagram!

Graphical Causal Model

**The DAG, possibly with latent confounders
(ADMG), induced by an SCM**

Directed Acyclic Graphs

Acyclic Directed Mixed Graphs

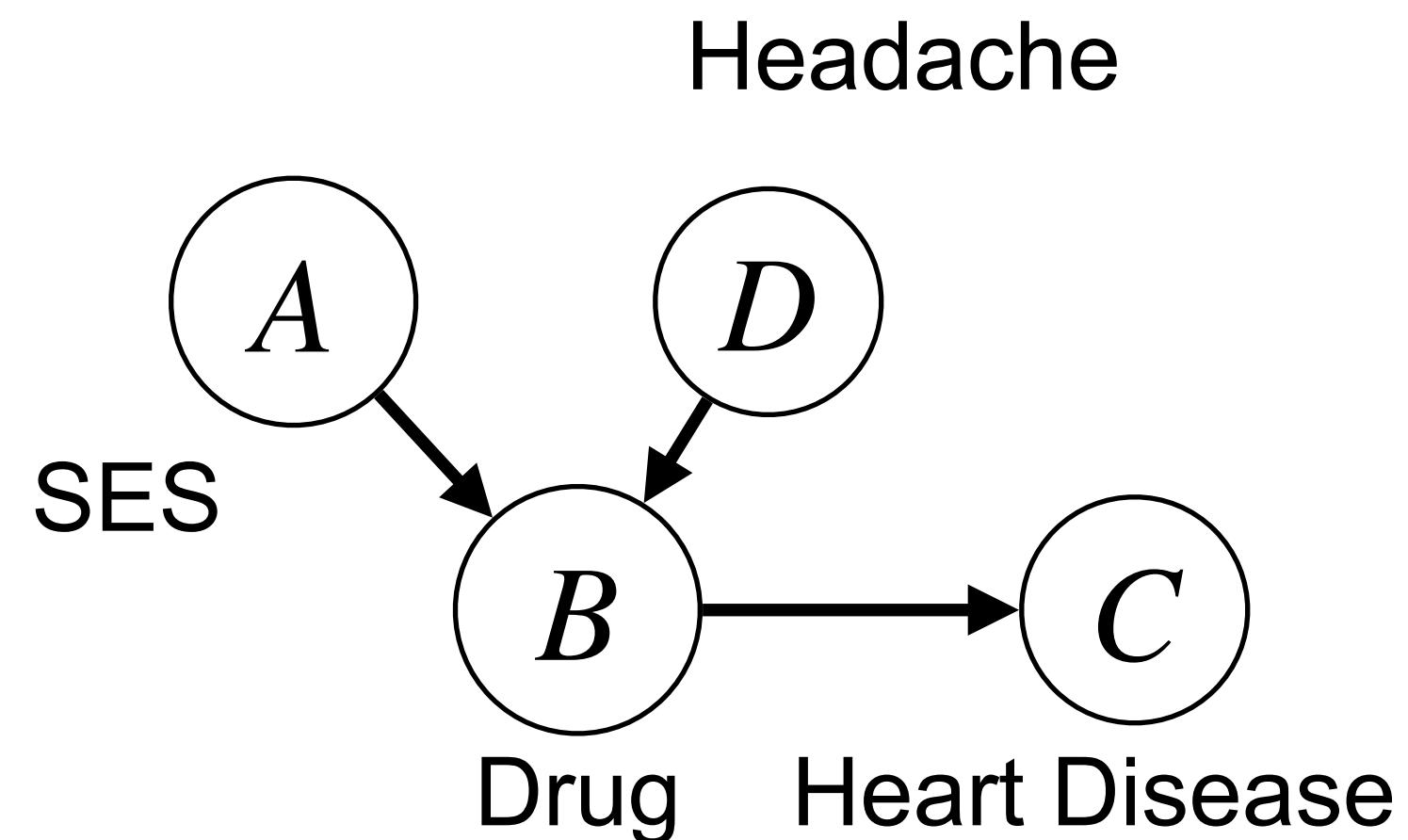
Causal Diagram: Encoder of Structural Knowledge

Structural Causal Model (SCM)

$$\mathcal{M} = \langle \mathbf{V}, \mathbf{U}, \mathcal{F}, P(\mathbf{u}) \rangle$$

$$\mathcal{M} = \left\{ \begin{array}{l} \mathbf{V} = \{A, B, C, D\} \\ \mathbf{U} = \{U_A, U_B, U_C, U_D, U_{CD}\} \\ \mathcal{F} = \left\{ \begin{array}{l} A \leftarrow f_A(U_A) \\ B \leftarrow f_B(A, D, U_B) \\ D \leftarrow f_Z(U_D, U_{CD}) \\ C \leftarrow f_X(B, U_C, U_{CD}) \end{array} \right. \\ P(\mathbf{U}) \end{array} \right.$$

Induced Causal Diagram



An SCM $\mathcal{M} = \langle \mathbf{V}, \mathbf{U}, \mathcal{F}, P(\mathbf{u}) \rangle$ induces a causal diagram such that, **for every** $V_i, V_j \in \mathbf{V}$:

$V_i \rightarrow V_j$, if V_i appears as argument of $f_j \in \mathcal{F}$.

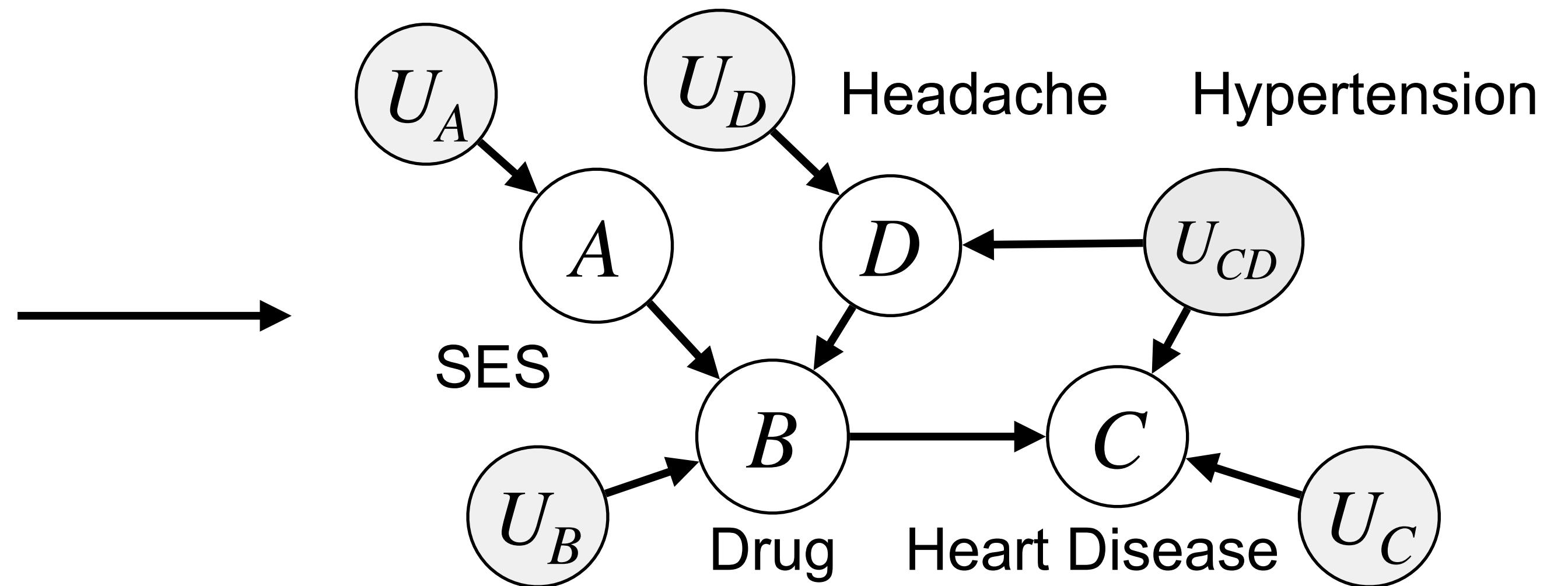
Causal Diagram: Encoder of Structural Knowledge

Structural Causal Model (SCM)

$$\mathcal{M} = \langle \mathbf{V}, \mathbf{U}, \mathcal{F}, P(\mathbf{u}) \rangle$$

$$\mathcal{M} = \begin{cases} \mathbf{V} = \{A, B, C, D\} \\ \mathbf{U} = \{U_A, U_B, U_C, U_D, U_{CD}\} \\ \mathcal{F} = \begin{cases} A \leftarrow f_A(U_A) \\ B \leftarrow f_B(A, D, U_B) \\ D \leftarrow f_Z(U_D, U_{CD}) \\ C \leftarrow f_X(B, U_C, U_{CD}) \end{cases} \\ P(\mathbf{U}) \end{cases}$$

Induced Causal Diagram



An SCM $\mathcal{M} = \langle \mathbf{V}, \mathbf{U}, \mathcal{F}, P(\mathbf{u}) \rangle$ induces a causal diagram such that, **for every** $V_i, V_j \in \mathbf{V}$:

$V_i \rightarrow V_j$, if V_i appears as argument of $f_j \in \mathcal{F}$.

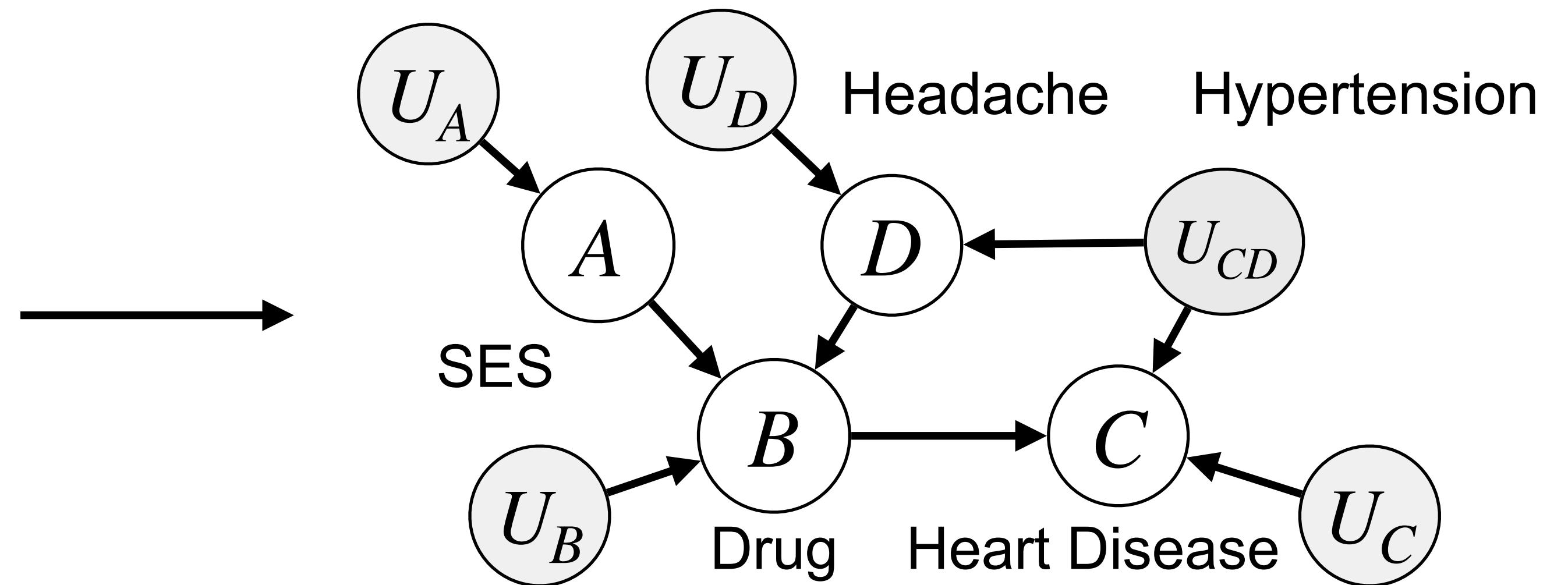
Causal Diagram: Encoder of Structural Knowledge

Structural Causal Model (SCM)

$$\mathcal{M} = \langle \mathbf{V}, \mathbf{U}, \mathcal{F}, P(\mathbf{u}) \rangle$$

$$\mathcal{M} = \begin{cases} \mathbf{V} = \{A, B, C, D\} \\ \mathbf{U} = \{U_A, U_B, U_C, U_D, U_{CD}\} \\ \mathcal{F} = \begin{cases} A \leftarrow f_A(U_A) \\ B \leftarrow f_B(A, D, U_B) \\ D \leftarrow f_Z(U_D, U_{CD}) \\ C \leftarrow f_X(B, U_C, U_{CD}) \end{cases} \\ P(\mathbf{U}) \end{cases}$$

Induced Causal Diagram



An SCM $\mathcal{M} = \langle \mathbf{V}, \mathbf{U}, \mathcal{F}, P(\mathbf{u}) \rangle$ induces a causal diagram such that, **for every** $V_i, V_j \in \mathbf{V}$:

$V_i \rightarrow V_j$, if V_i appears as argument of $f_j \in \mathcal{F}$.

$V_i \leftrightarrow V_j$ if the corresponding $U_i, U_j \in \mathbf{U}$ are correlated or f_i, f_j share some argument $U \in \mathbf{U}$.

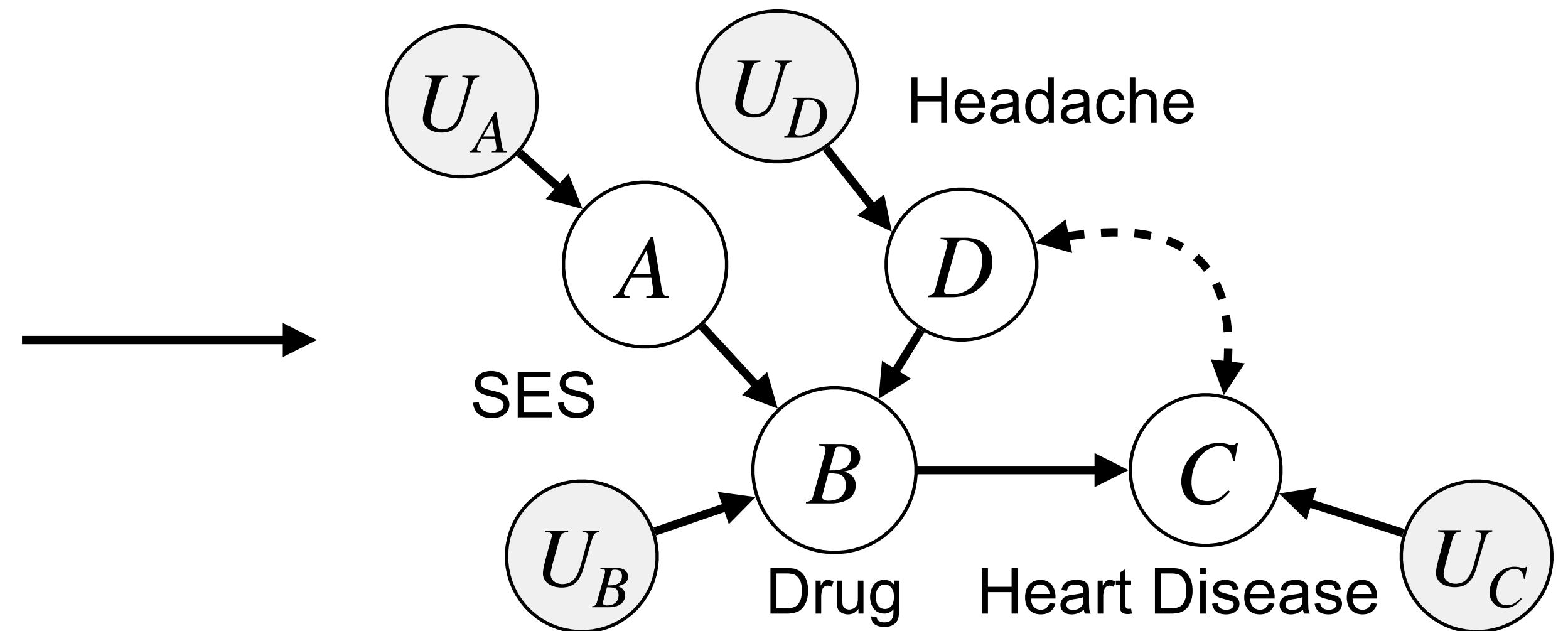
Causal Diagram: Encoder of Structural Knowledge

Structural Causal Model (SCM)

$$\mathcal{M} = \langle \mathbf{V}, \mathbf{U}, \mathcal{F}, P(\mathbf{u}) \rangle$$

$$\mathcal{M} = \begin{cases} \mathbf{V} = \{A, B, C, D\} \\ \mathbf{U} = \{U_A, U_B, U_C, U_D, U_{CD}\} \\ \mathcal{F} = \begin{cases} A \leftarrow f_A(U_A) \\ B \leftarrow f_B(A, D, U_B) \\ D \leftarrow f_Z(U_D, U_{CD}) \\ C \leftarrow f_X(B, U_C, U_{CD}) \end{cases} \\ P(\mathbf{U}) \end{cases}$$

Induced Causal Diagram



An SCM $\mathcal{M} = \langle \mathbf{V}, \mathbf{U}, \mathcal{F}, P(\mathbf{u}) \rangle$ induces a causal diagram such that, **for every** $V_i, V_j \in \mathbf{V}$:

$V_i \rightarrow V_j$, if V_i appears as argument of $f_j \in \mathcal{F}$.

$V_i \leftrightarrow V_j$ if the corresponding $U_i, U_j \in \mathbf{U}$ are correlated or f_i, f_j share some argument $U \in \mathbf{U}$.

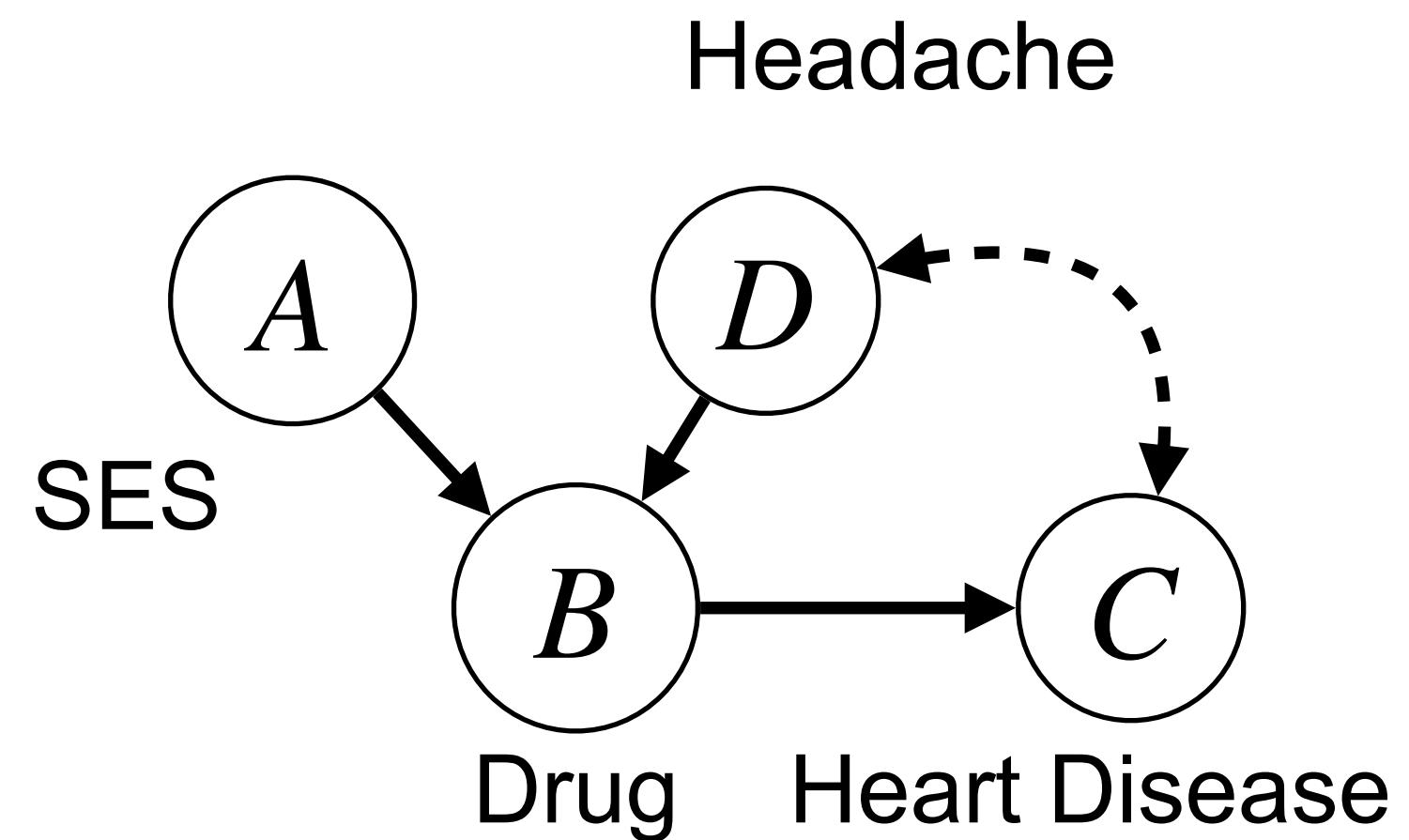
Causal Diagram: Encoder of Structural Knowledge

Structural Causal Model (SCM)

$$\mathcal{M} = \langle \mathbf{V}, \mathbf{U}, \mathcal{F}, P(\mathbf{u}) \rangle$$

$$\mathcal{M} = \left\{ \begin{array}{l} \mathbf{V} = \{A, B, C, D\} \\ \mathbf{U} = \{U_A, U_B, U_C, U_D, U_{CD}\} \\ \mathcal{F} = \left\{ \begin{array}{l} A \leftarrow f_A(U_A) \\ B \leftarrow f_B(A, D, U_B) \\ D \leftarrow f_Z(U_D, U_{CD}) \\ C \leftarrow f_X(B, U_C, U_{CD}) \end{array} \right. \\ P(\mathbf{U}) \end{array} \right.$$

Induced Causal Diagram



An SCM $\mathcal{M} = \langle \mathbf{V}, \mathbf{U}, \mathcal{F}, P(\mathbf{u}) \rangle$ induces a causal diagram such that, **for every** $V_i, V_j \in \mathbf{V}$:

$V_i \rightarrow V_j$, if V_i appears as argument of $f_j \in \mathcal{F}$.

$V_i \leftrightarrow V_j$ if the corresponding $U_i, U_j \in \mathbf{U}$ are correlated or f_i, f_j share some argument $U \in \mathbf{U}$.

D-Separation

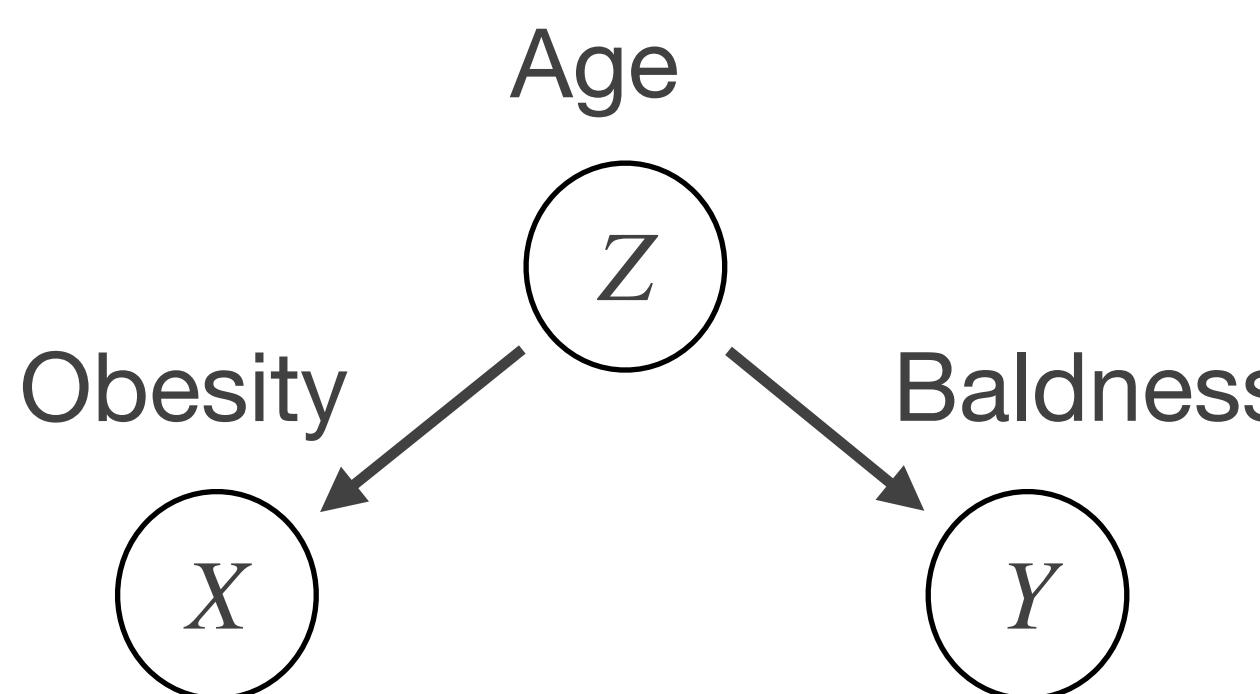
Graphical Tool for DAGs and ADMGs

Directed Acyclic Graphs

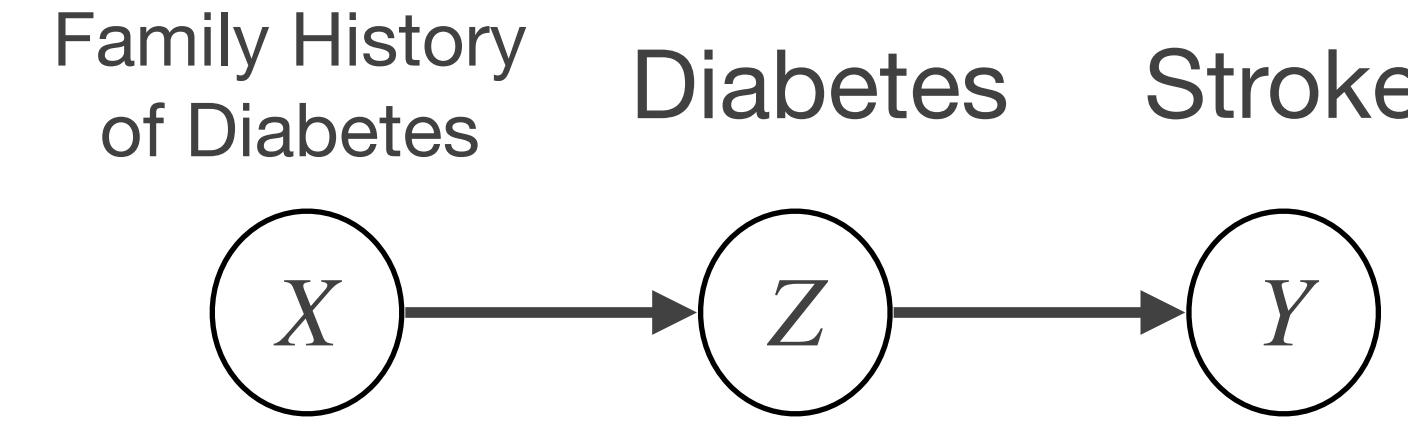
Acyclic Directed Mixed Graphs

Encoding Conditional Independencies

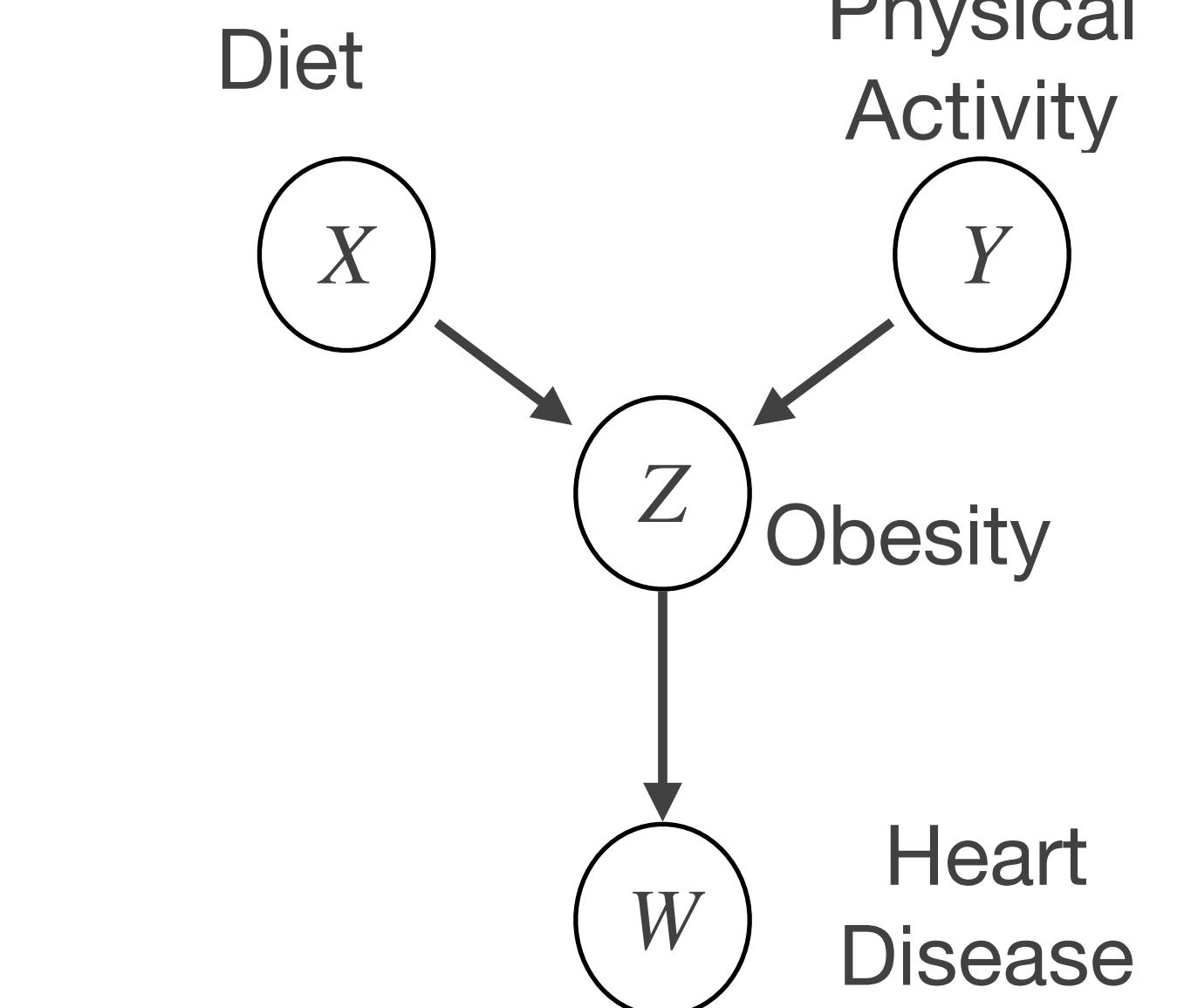
Fork



Chain



V-Structure



In both cases, Z is a non-collider!

$X \perp\!\!\!\perp Y$

$X \perp\!\!\!\perp Y | Z$

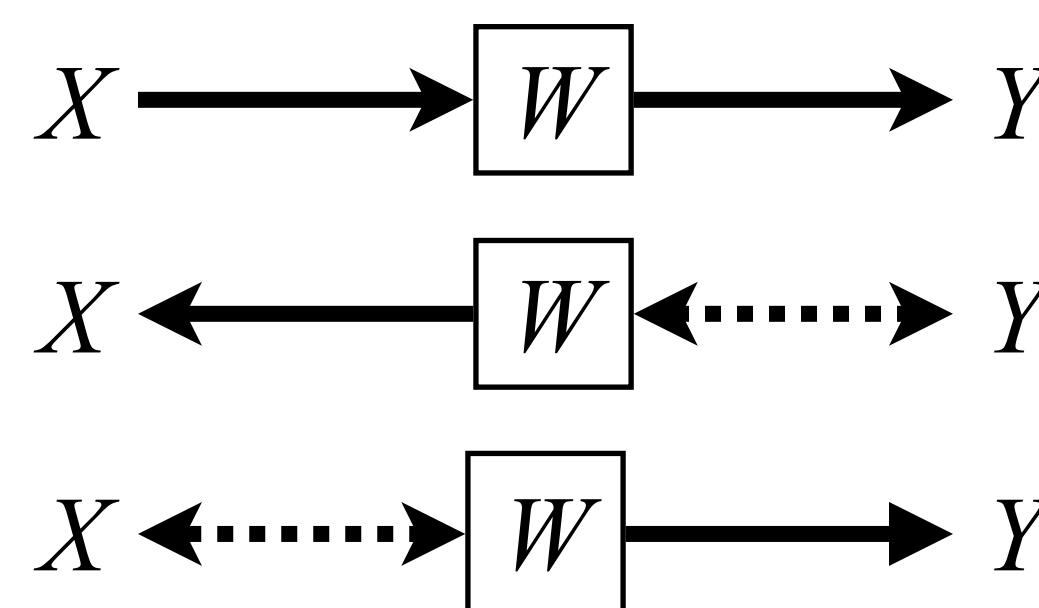
$X \perp\!\!\!\perp Y | W$

Active and Inactive Triplets

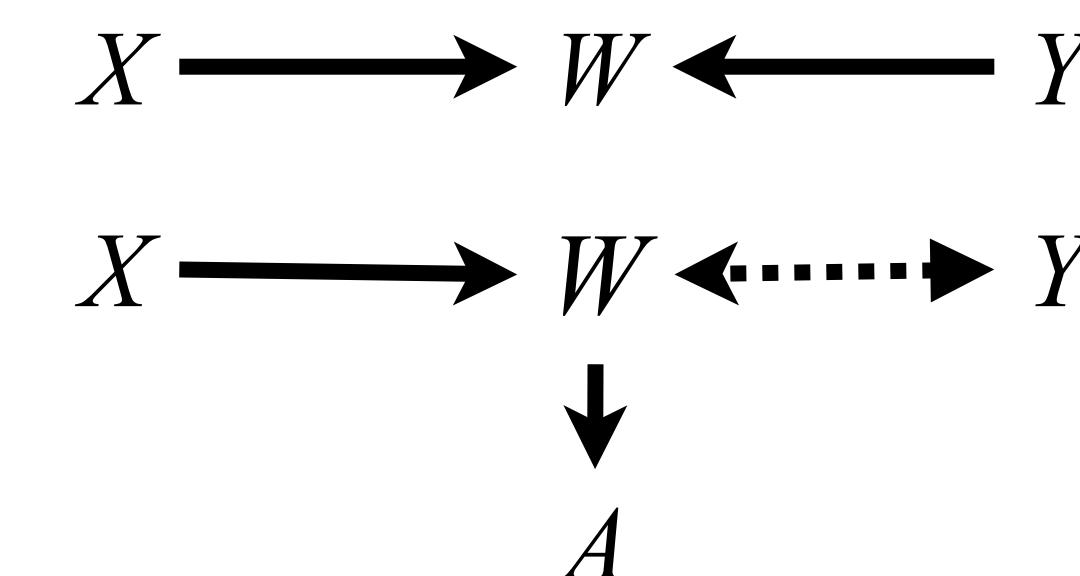
Definition (inactive): A triplet $\langle V_i, V_m, V_j \rangle$ is said to be *inactive* relative to a set Z if the middle node V_m :

1. Is a non-collider and is in Z ; or
2. Is a collider and neither it nor any of its descendants in Z .

W is non-collider
and $W \in Z$



W is (descendant of) a
collider and $W, A \notin Z$



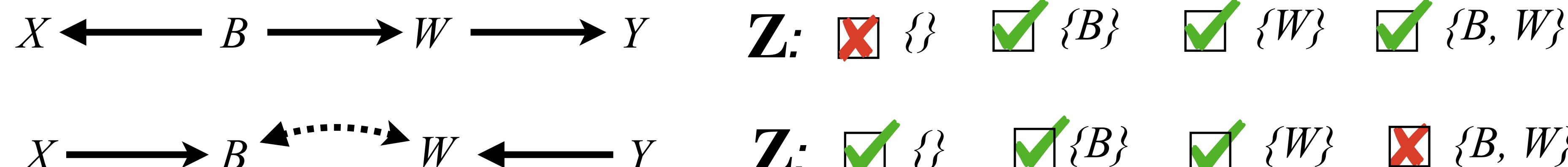
D-Separation

An important tool that allows us to read the conditional (in)dependencies implied by an Acyclic Directed Mixed Graph (ADMG).

Definition (d-separation): A path p in an ADMG G is said to be **d-separated** (or blocked) by a set of variables Z if and only if p contains an inactive triplet in it.

A set Z d-separates X and Y if and only if Z blocks every path between a node in X and a node in Y . We denote that by $(X \perp\!\!\!\perp Y | Z)_G$.

Does Z d-separate X and Y ?

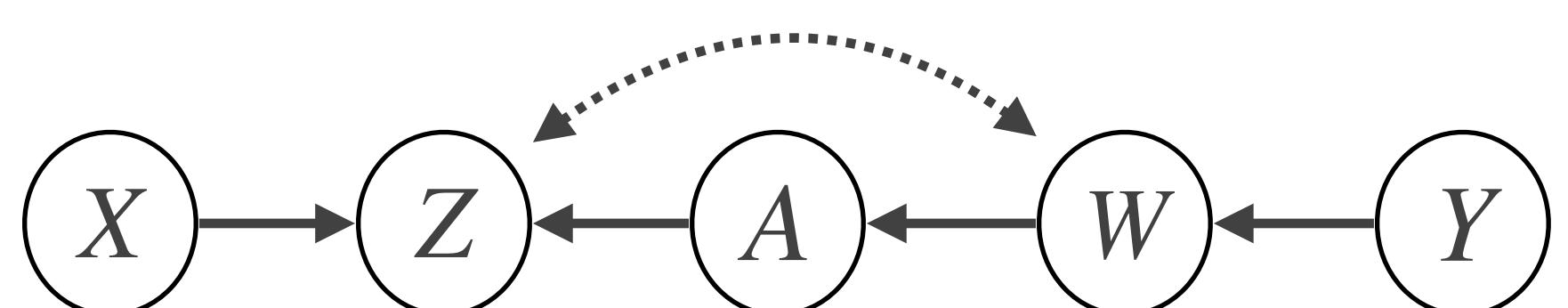


D-Separation

An important tool that allows us to read the conditional (in)dependencies implied by an Acyclic Directed Mixed Graph (ADMG).

Definition (d-separation): A path p in an ADMG G is said to be **d-separated** (or blocked) by a set of variables Z if and only if p contains an inactive triplet in it.

A set Z d-separates X and Y if and only if Z blocks every path between a node in X and a node in Y . We denote that by $(X \perp\!\!\!\perp Y | Z)_G$.



Which set d-separates A and Y ? $\{W\}$, i.e., $(A \perp\!\!\!\perp Y | W)_G$

What about A and X ? $\{\}$, i.e., $(A \perp\!\!\!\perp X)_G$
Z and Y ? None

Global Markov property: $(X \perp\!\!\!\perp Y | Z)_G \Rightarrow (X \perp\!\!\!\perp Y | Z)_P$

D-separations in G imply
conditional independencies in P

Causal Effect Identification from Causal Diagrams

Causal Effect

The **causal effect** of a (set of) treatment variable(s) \mathbf{X} on a (set of) outcome variable(s) \mathbf{Y} is a quantity derived from $P(\mathbf{Y} | do(\mathbf{X}))$ that tells us how much \mathbf{Y} changes due to an intervention $do(\mathbf{X} = \mathbf{x})$.

Examples:

- *Average Treatment Effect (ATE)* for discrete treatments:

$$\mathbb{E}[\mathbf{Y} | do(\mathbf{X} = \mathbf{x}')] - \mathbb{E}[\mathbf{Y} | do(\mathbf{X} = \mathbf{x})],$$

where $\mathbb{E}[\mathbf{Y} | do(\mathbf{X} = \mathbf{x})] = \sum_{\mathbf{y} \in \Omega_{\mathbf{Y}}} \mathbf{y} P(\mathbf{y} | do(\mathbf{x}))$

defined for two treatment levels \mathbf{x}' and \mathbf{x} of \mathbf{X} .

- *Average Treatment Effect (ATE)* for continuous treatments,

$$\frac{\partial \mathbb{E}[Y_i | do(X_j = x_j)]}{\partial x_j}, \text{ for all } Y_i \in \mathbf{Y}, \text{ and } X_j \in \mathbf{X}.$$

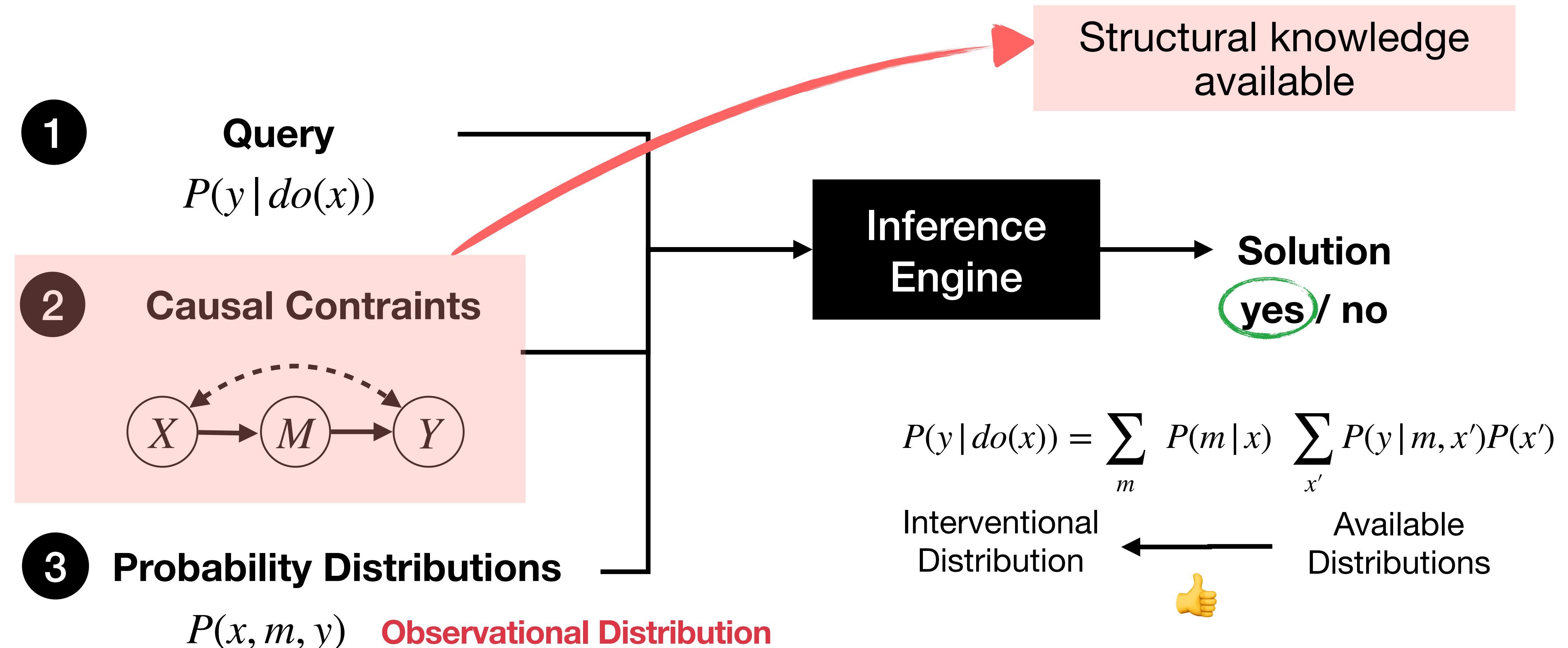
Jacobian of $\mathbb{E}[\mathbf{Y} | do(\mathbf{X} = \mathbf{x})]$, where

$$\mathbb{E}[\mathbf{Y} | do(\mathbf{X} = \mathbf{x})] = \int_{\Omega_{\mathbf{Y}}} \mathbf{y} P(\mathbf{y} | do(\mathbf{x})) d\mathbf{y},$$

and $\Omega_{\mathbf{Y}}$ is the space of all possible values that \mathbf{Y} might take on

The derivative shows the rate of change of \mathbf{Y} w.r.t. $do(\mathbf{X} = \mathbf{x})$

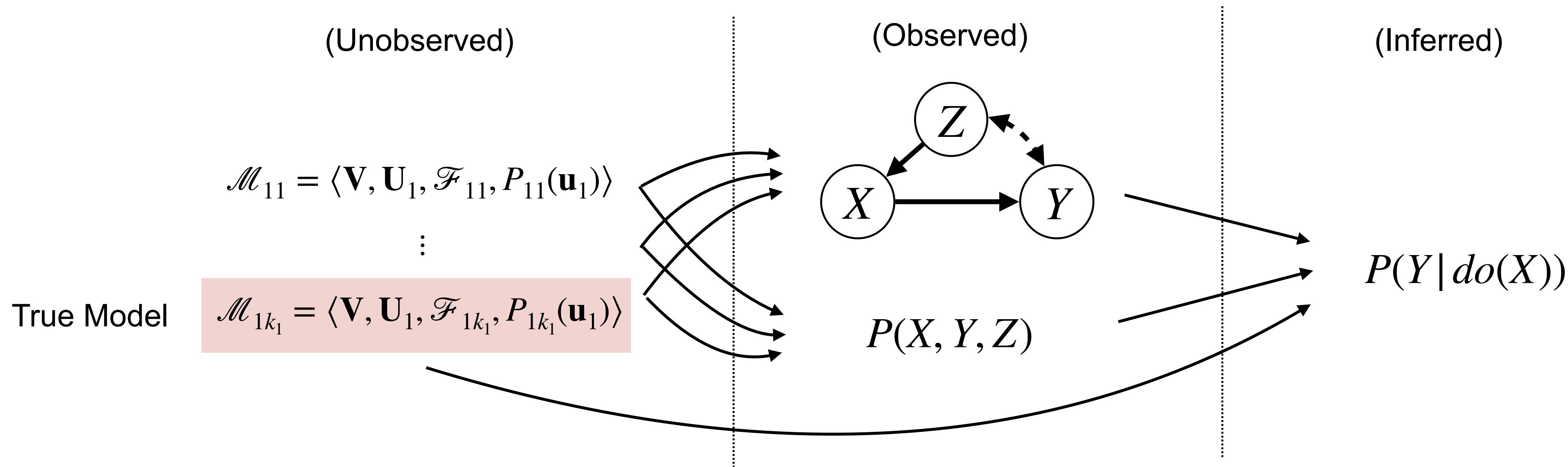
Classical Causal Effect Identification



- Tian, J. and Pearl, J. A General Identification Condition for Causal Effects. In Proceedings of the Eighteenth National Conference on Artificial Intelligence (AAAI 2002), pp. 567–573, Menlo Park, CA, 2002. AAAI Press/MIT Press.

The Effect Identification Problem

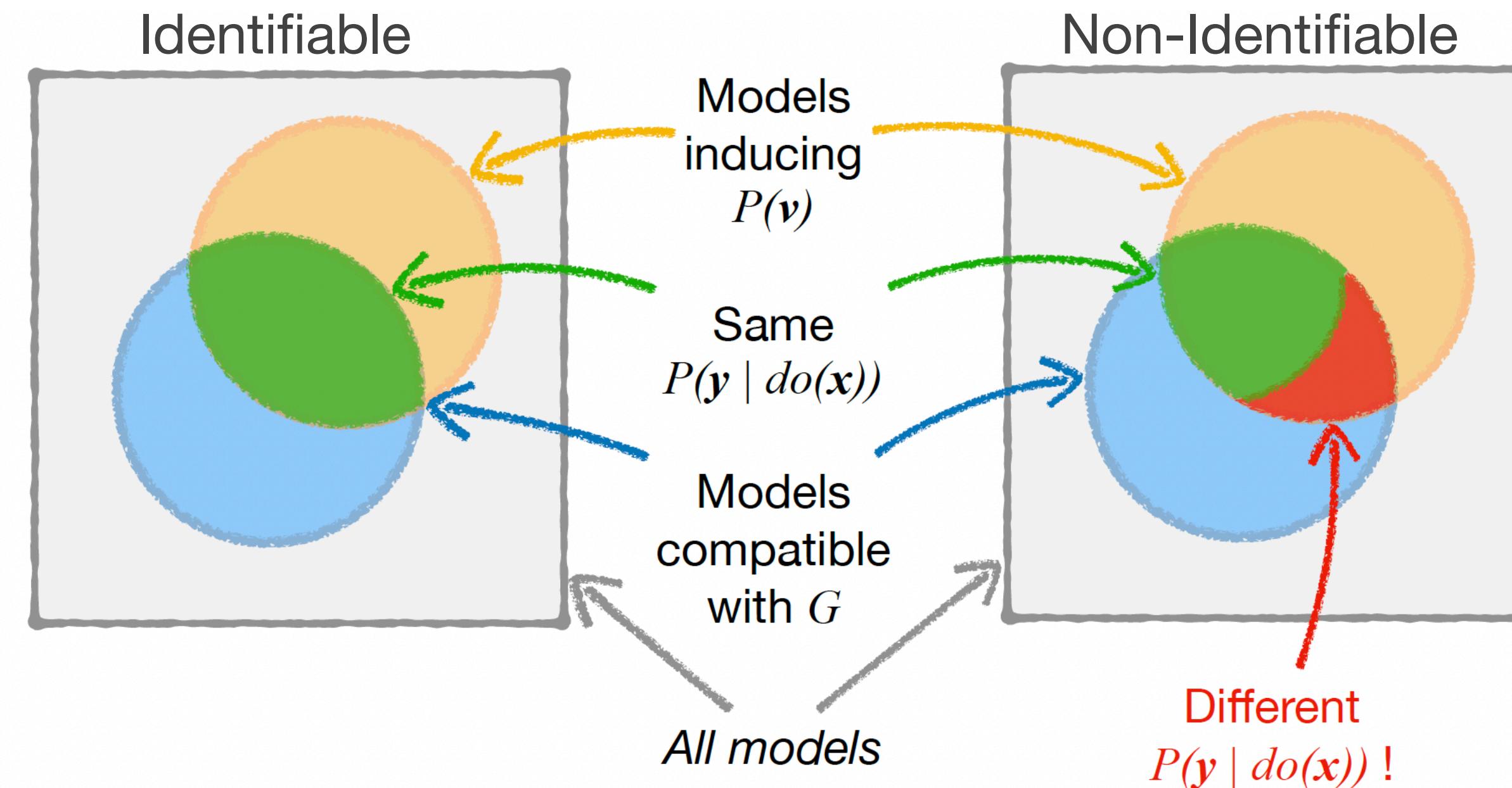
Causal Effect Identifiability: The causal effect of a (set of) treatment variable(s) \mathbf{X} on a (set of) outcome variable(s) \mathbf{Y} is said to be identifiable from a causal diagram G and the probability of the observed variables $P(\mathbf{V})$ if the interventional distribution $P(\mathbf{Y} | do(\mathbf{X}))$ is *uniquely computable*, i.e., if for every pair of SCMs \mathcal{M}_1 and \mathcal{M}_2 that induce G and $P^{\mathcal{M}_1}(\mathbf{V}) = P^{\mathcal{M}_2}(\mathbf{V}) = P(\mathbf{V}) > 0$, $P^{\mathcal{M}_1}(\mathbf{Y} | do(\mathbf{X})) = P^{\mathcal{M}_2}(\mathbf{Y} | do(\mathbf{X})) = P(\mathbf{Y} | do(\mathbf{X}))$.



In words, causal effect identifiability means that, no matter the form of true SCM, for all models \mathcal{M} agreeing with $\langle G, P(\mathbf{V}) \rangle$, they also agree in $P(\mathbf{y} | do(\mathbf{x}))$.

The Effect Identification Problem

Causal Effect Identifiability: The causal effect of a (set of) treatment variable(s) \mathbf{X} on a (set of) outcome variable(s) \mathbf{Y} is said to be identifiable from a causal diagram G and the probability of the observed variables $P(\mathbf{V})$ if the interventional distribution $P(\mathbf{Y} | do(\mathbf{X}))$ is *uniquely computable*, i.e., if for every pair of SCMs \mathcal{M}_1 and \mathcal{M}_2 that induce G and $P^{\mathcal{M}_1}(\mathbf{V}) = P^{\mathcal{M}_2}(\mathbf{V}) = P(\mathbf{V}) > 0$, $P^{\mathcal{M}_1}(\mathbf{Y} | do(\mathbf{X})) = P^{\mathcal{M}_2}(\mathbf{Y} | do(\mathbf{X})) = P(\mathbf{Y} | do(\mathbf{X}))$.



In words, causal effect identifiability means that, no matter the form of true SCM, for all models \mathcal{M} agreeing with $\langle G, P(\mathbf{V}) \rangle$, they also agree in $P(\mathbf{y} | do(\mathbf{x}))$.

Identification in Markovian Models

Truncated Factorization – Markovian: Let G be a causal diagram for the collection \mathbf{P}_* of all interventional distributions $P_{\mathbf{x}}(\mathbf{V})$, for any $\mathbf{X} \subseteq \mathbf{V}$. It follows that $P_{\mathbf{x}}(\mathbf{V})$ factorizes as:

$$P_{\mathbf{x}}(\mathbf{v}) \doteq P(\mathbf{v} \mid do(\mathbf{x})) = \prod_{V_i \in \mathbf{V} \setminus \mathbf{X}} P_{\mathbf{x}}(v_i \mid pa_i) \Big|_{\mathbf{X}=\mathbf{x}}$$

Follows from $P_{\mathbf{x}}(\mathbf{v}) \doteq P(\mathbf{v} \mid do(\mathbf{x}))$
being *Markov* relative to $G_{\overline{\mathbf{X}}}$

$$= \prod_{V_i \in \mathbf{V} \setminus \mathbf{X}} P(v_i \mid pa_i) \Big|_{\mathbf{X}=\mathbf{x}}$$

Markovian SCMs have the modularity
property, i.e., $P_{\mathbf{x}}(v_i \mid pa_i) = P(v_i \mid pa_i)$

Causal Effect of \mathbf{X} on \mathbf{Y} :

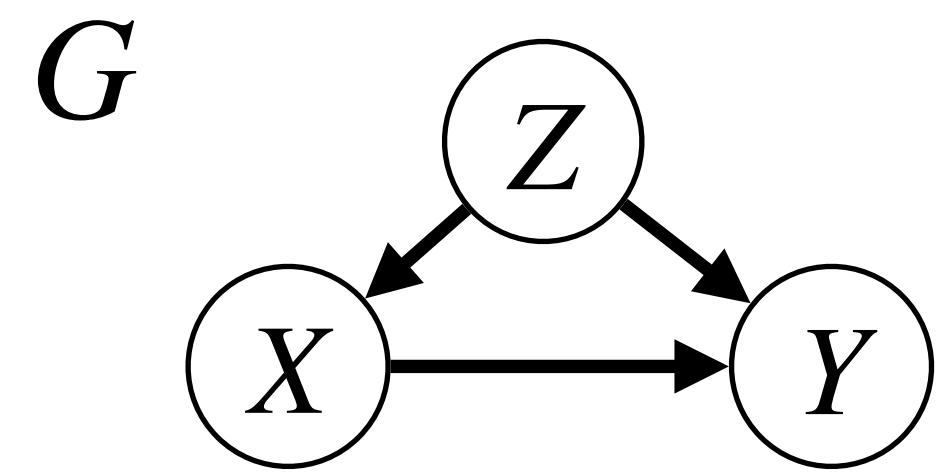
$$P(\mathbf{y} \mid do(\mathbf{x})) = \sum_{\mathbf{V} \setminus (\mathbf{Y} \cup \mathbf{X})} \prod_{V_i \in \mathbf{V} \setminus \mathbf{X}} P(v_i \mid pa_i) \Big|_{\mathbf{X}=\mathbf{x}}$$

- In Markovian Models, the joint interventional distribution (and hence any causal effect) is always identifiable.
- This factorization is a.k.a “manipulation theorem” (Spirtes et al. 1993) or G-computation (Robins 1986, p. 1423).

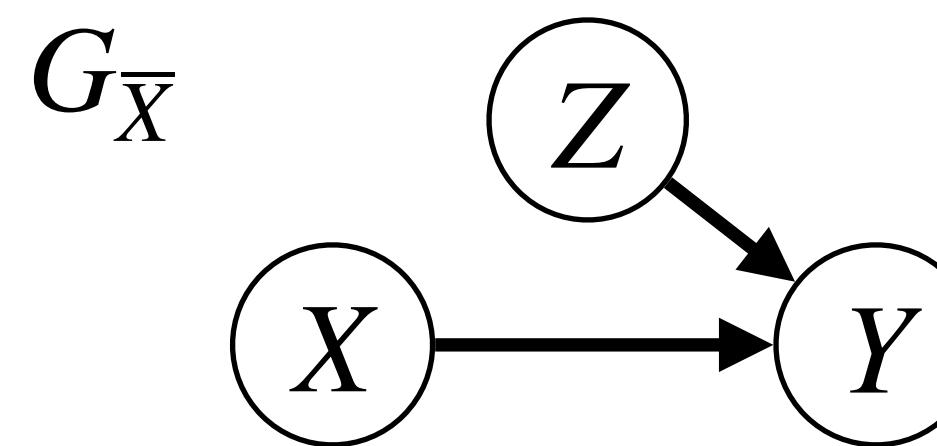
Example: Identifiable Effect

Causal Effect of X on Y :

$$P(y \mid do(x)) = \sum_{V \setminus (Y \cup X)} \prod_{V_i \in V \setminus X} P_x(v_i \mid pa_i) \Big|_{X=x}$$



$do(X = x)$



$$P(x, y, z) = P(z)P(x \mid z)P(y \mid x, z)$$

$$P(y, z \mid do(x)) = P(z)P(y \mid x, z)$$

$$\implies P(y \mid do(x)) = \sum_z P(z)P(y \mid x, z)$$

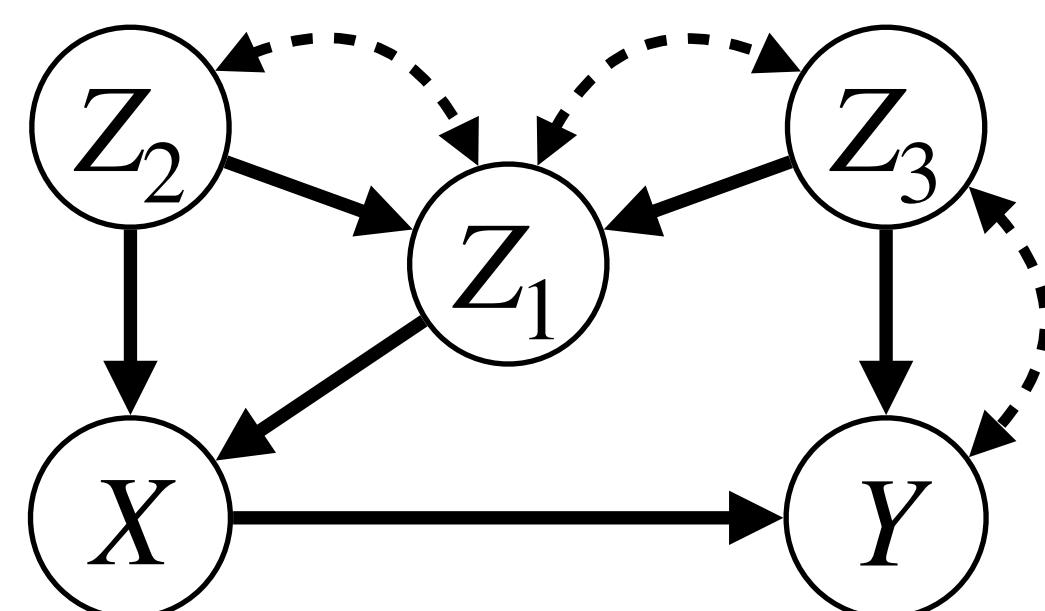
Identification Via Adjustment over Parents

Let G be a causal graph with **all parents observed**.

Then, the effect of \mathbf{X} on \mathbf{Y} is given by:

$$P(\mathbf{y} | do(\mathbf{x})) = \sum_{\mathbf{pa}_{\mathbf{x}}} P\left(\mathbf{y} | \mathbf{x}, \mathbf{pa}_{\mathbf{x}}\right) P\left(\mathbf{pa}_{\mathbf{x}}\right)$$

Proof follows from the truncated factorization for Markovian models.
Try at home!



$$Pa_x = \{Z_1, Z_2\}$$

$$\begin{aligned}\mathbf{X} &= \{X\} \\ \mathbf{Y} &= \{Y\} \\ \mathbf{Pa}_x &= \{Z_1, Z_2\}\end{aligned}$$

$$P(y | do(x)) = \sum_{z_1, z_2} P(y | x, z_1, z_2) P(z_1, z_2)$$

Identification in Semi-Markovian Models

Truncated Factorization – Semi-Markovian: Let G be the causal diagram for the collection \mathbf{P}_* of all interventional distributions $P(\mathbf{V} \mid do(\mathbf{x}))$, for any $\mathbf{X} \subseteq \mathbf{V}$. The interventional distribution $P(\mathbf{V} \mid do(\mathbf{x}))$ factorizes as followings:

$$P(\mathbf{v} \mid do(\mathbf{x})) = \sum_{\mathbf{u}} \prod_{V_i \in \mathbf{V} \setminus \mathbf{X}} P(v_i \mid pa_i, u_i) P(\mathbf{u}) \Big|_{\mathbf{X}=\mathbf{x}}$$

For Semi-Markovian Models, more advanced tools are necessary for evaluating $P_{\mathbf{x}}(\mathbf{v})$ as an expression that depends only on $P(\mathbf{V})$, i.e., a do-free and U-free expression.

Identification via Backdoor Criterion

Let \mathbf{X} be a set of treatment variables and \mathbf{Y} a set of outcome variables in the causal graph G .

If there exists a set \mathbf{Z} such that:

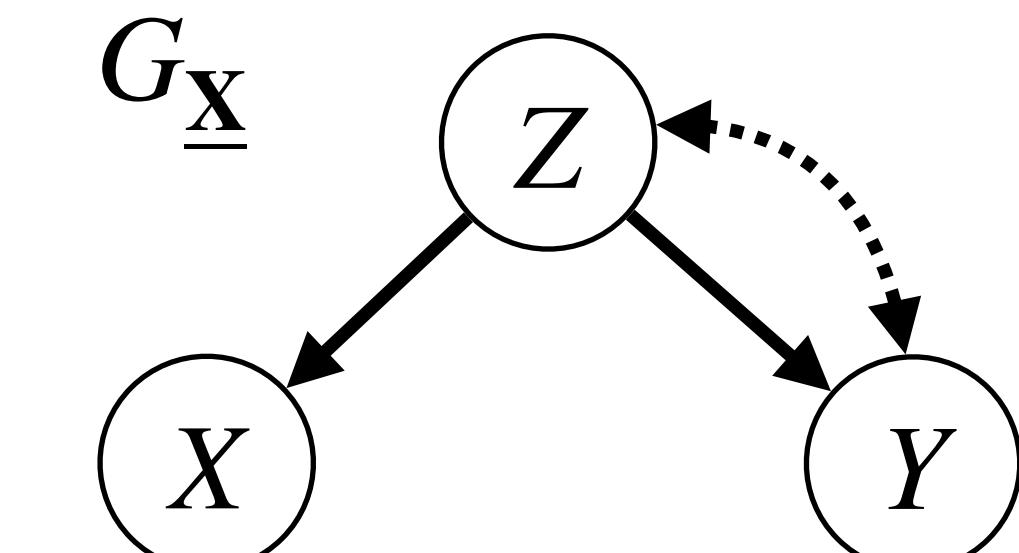
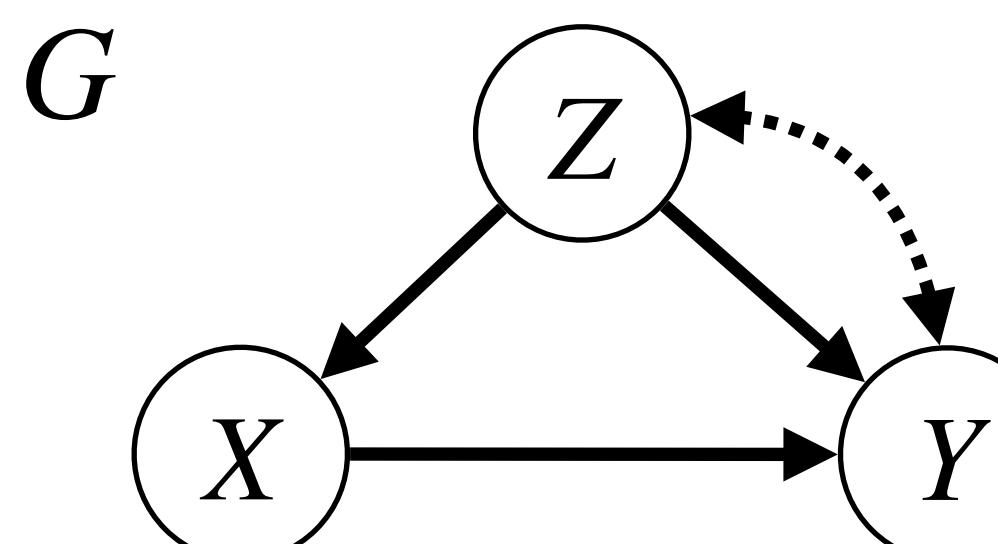
1. \mathbf{Z} d-separates \mathbf{X} and \mathbf{Y} in the graph $G_{\underline{\mathbf{X}}}$, i.e., the graph resulting from cutting the arrows out of \mathbf{X}
2. no node in \mathbf{Z} is a descendant of a variable $X \in \mathbf{X}$ in G (all variables in \mathbf{Z} are pre-treatment)

Then, \mathbf{Z} satisfies the **backdoor criterion** for (\mathbf{X}, \mathbf{Y}) and, then the effect of \mathbf{X} on \mathbf{Y} is given by:

$$P(\mathbf{y} | do(\mathbf{x})) = \sum_{\mathbf{z}} P(\mathbf{y} | \mathbf{x}, \mathbf{z}) P(\mathbf{z})$$

$$\begin{aligned}\mathbf{X} &= \{X\} \\ \mathbf{Y} &= \{Y\} \\ \mathbf{Z} &= \{Z\}\end{aligned}$$

\mathbf{Z} , a set of covariates, admissible for backdoor adjustment



In $G_{\underline{\mathbf{X}}}$, all non-backdoor paths are severed

Identification via Backdoor Criterion

Let \mathbf{X} be a set of treatment variables and \mathbf{Y} a set of outcome variables in the causal graph G .

If there exists a set \mathbf{Z} such that:

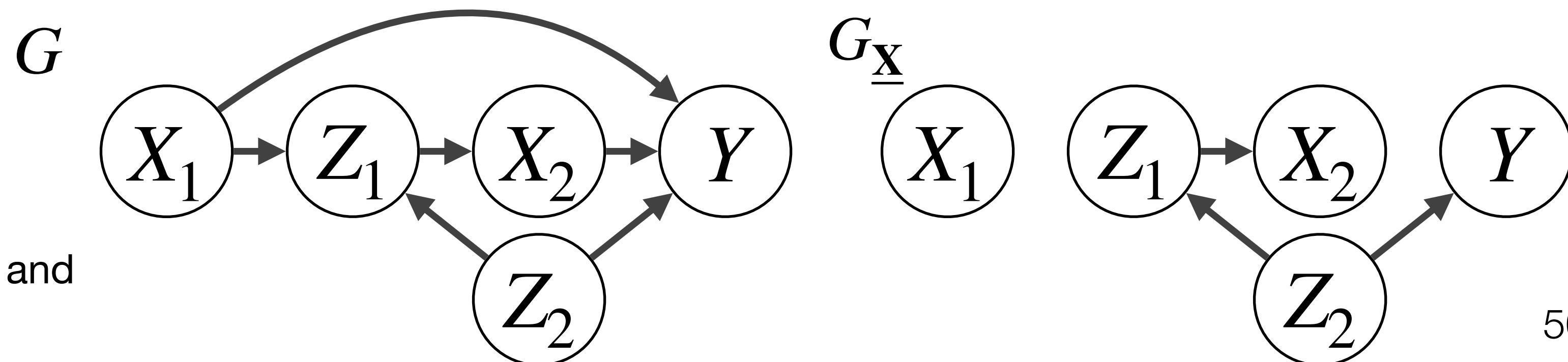
1. \mathbf{Z} d-separates \mathbf{X} and \mathbf{Y} in the graph $G_{\underline{\mathbf{X}}}$, i.e., the graph resulting from cutting the arrows out of \mathbf{X}
2. no node in \mathbf{Z} is a descendant of a variable $X \in \mathbf{X}$ in G (all variables in \mathbf{Z} are pre-treatment)

Then, \mathbf{Z} satisfies the **backdoor criterion** for (\mathbf{X}, \mathbf{Y}) and, then the effect of \mathbf{X} on \mathbf{Y} is given by:

$$P(\mathbf{y} | do(\mathbf{x})) = \sum_{\mathbf{z}} P(\mathbf{y} | \mathbf{x}, \mathbf{z}) P(\mathbf{z})$$

$$\begin{aligned}\mathbf{X} &= \{X_1, X_2\} \\ \mathbf{Y} &= \{Y\} \\ \mathbf{Z} &= \{Z_1, Z_2\}\end{aligned}$$

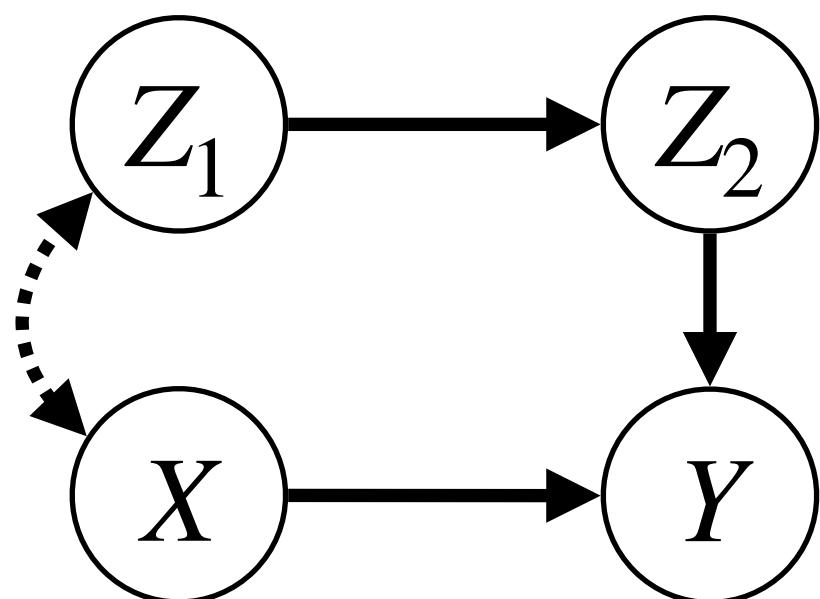
\mathbf{Z} , a set of covariates, admissible for backdoor adjustment



Admissible Sets for BD Adjustment

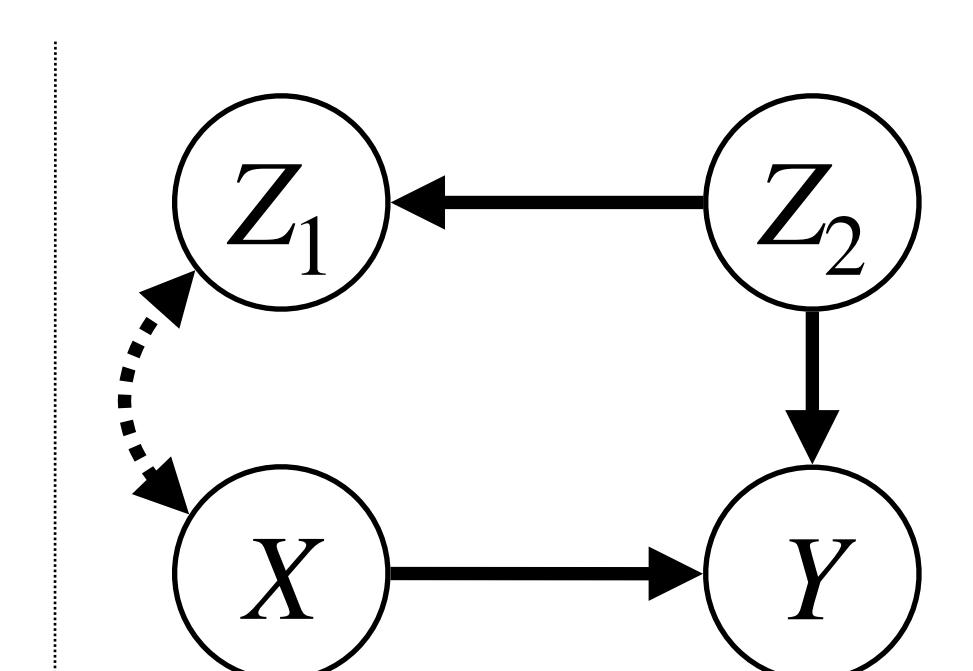
Z satisfies the **backdoor criterion** for or (X, Y) in the causal graph G if:

1. Z d-separates X and Y in the graph $\underline{G_X}$, i.e., the graph resulting from cutting the arrows out of X
2. no node in Z is a descendant of a variable $X \in X$ in G (all variables in Z are pre-treatment)



Minimal BD
Adjustment Sets $\left\{ \begin{array}{l} \{Z_1\}, \\ \{Z_2\}, \\ \{Z_1, Z_2\} \end{array} \right.$

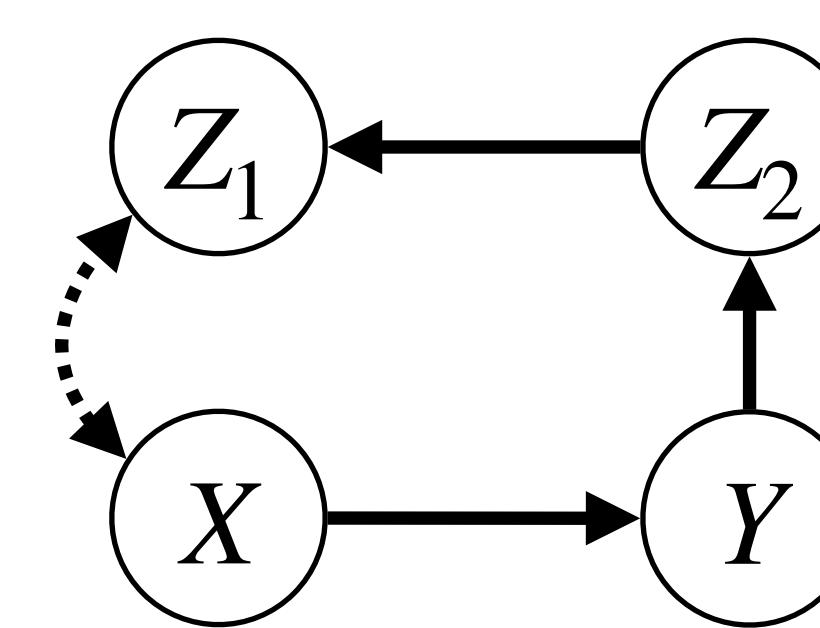
$$P(y|do(x)) = \sum_{z_1} P(y|x, z_1) P(z_1)$$



$\{\}$,

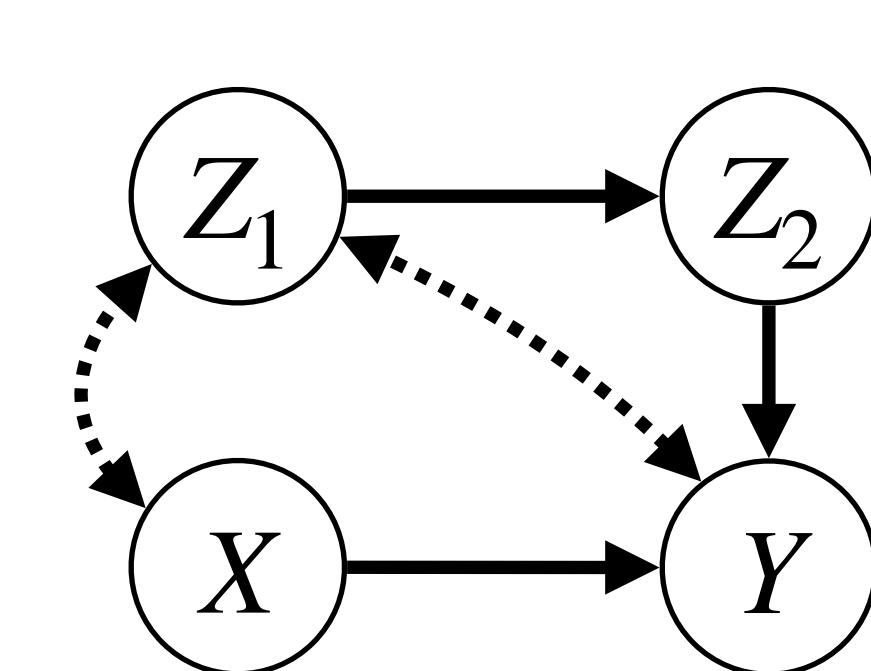
$\{Z_2\}$,
 $\{Z_1, Z_2\}$

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



$\{\}$

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



There is no BD
Adjustment Set!

$P(y|do(x))$ is
non-identifiable

The screenshot shows three parallel causal fusion interfaces, each displaying a causal graph and analysis results.

Graphical Editor:

```

1 <NODES>
2 X -45,-15
3 Y 45,-15
4 Z 0,-60
5
6 <EDGES>
7 X --> Y
8 Z --> X
9 Z --> Y
Populations
  
```

Summary:

Treatment : X
 Outcome : Y
 Adjusted :
 Query : $P(Y|do(X))$

Editor:

Graphical Structural

Confounding Analysis:

- Admissible Sets
- Admissibility Test
- Instrumental Variables
- IV Admissibility Test

Path Analysis:

- D-Separation
- Causal Paths
- Confounding Paths
- Biassing Paths

Do-Calculus Analysis:

- Do-Inspector
- Do-Separation

σ -Calculus Analysis:

- σ -Inspector
- σ -Separation

Compute:

The causal effect of X on Y conditional on Z with do : \equiv (Query: $P(Y|do(X))$ from $P(v)$)

Non-Parametric: **Clear:**

Results:

$P(Y|do(X)) = \sum_Z P(Y|X, Z) P(Z)$

Load:

Estimation:

Derivation:

Remove:

causalfusion.net/app

Fusion^(β)

Summary

Treatment : X
Outcome : Y
Adjusted :
Query : $P(Y|do(X))$

Show More Details

Editor

Graphical Structural Refresh

```
1 <NODES>
2 X -45,-15
3 Y 45,-15
4 Z 0,-60
5
6 <EDGES>
7 X -> Y
8 Z -> X
9 Z -> Y
```

Populations

Confounding Analysis

Admissible Sets
Admissibility Test
Instrumental Variables
IV Admissibility Test

Path Analysis

D-Separation
Causal Paths
Confounding Paths
Biasing Paths

Do-Calculus Analysis

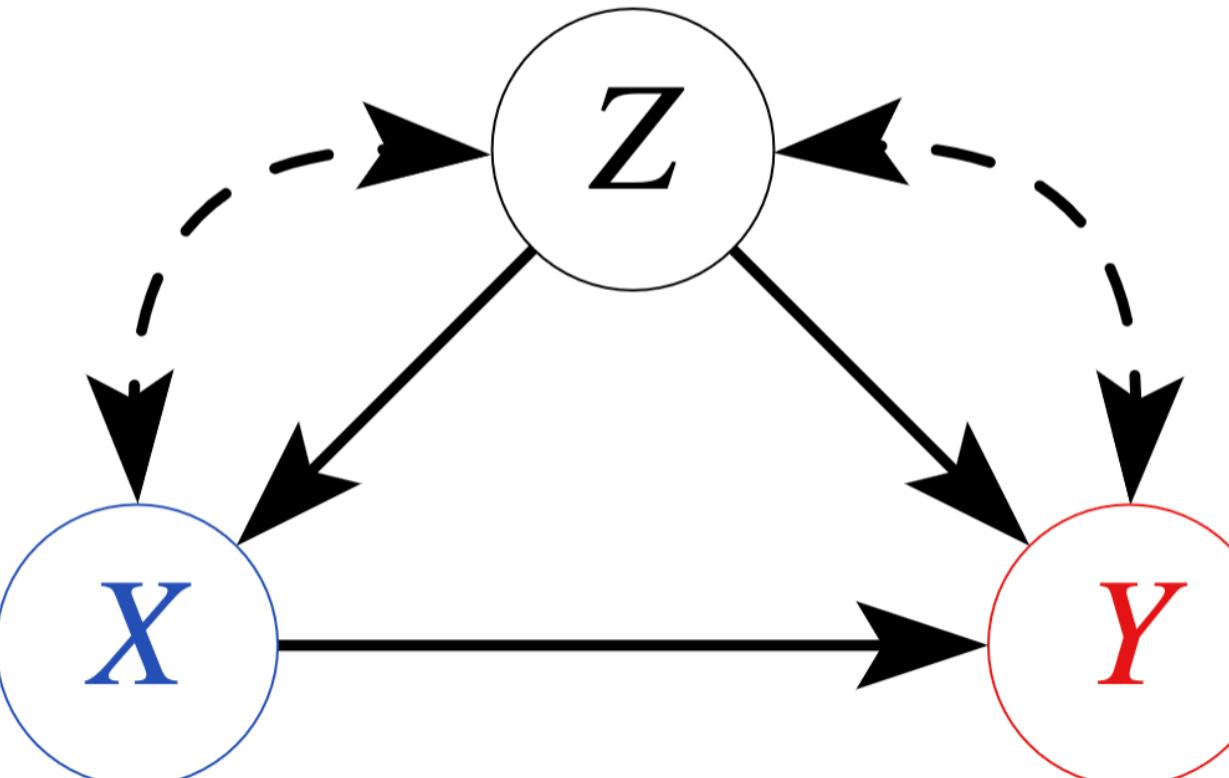
Do-Inspector
Do-Separation

σ-Calculus Analysis

σ-Inspector
σ-Separation

Compute The causal effect of X on Y conditional on with do : \equiv (Query: $P(Y|do(X))$ from $P(v)$) Non-Parametric Clear

1 $P(Y|do(X))$ is not identifiable from $P(X, Y, Z)$.



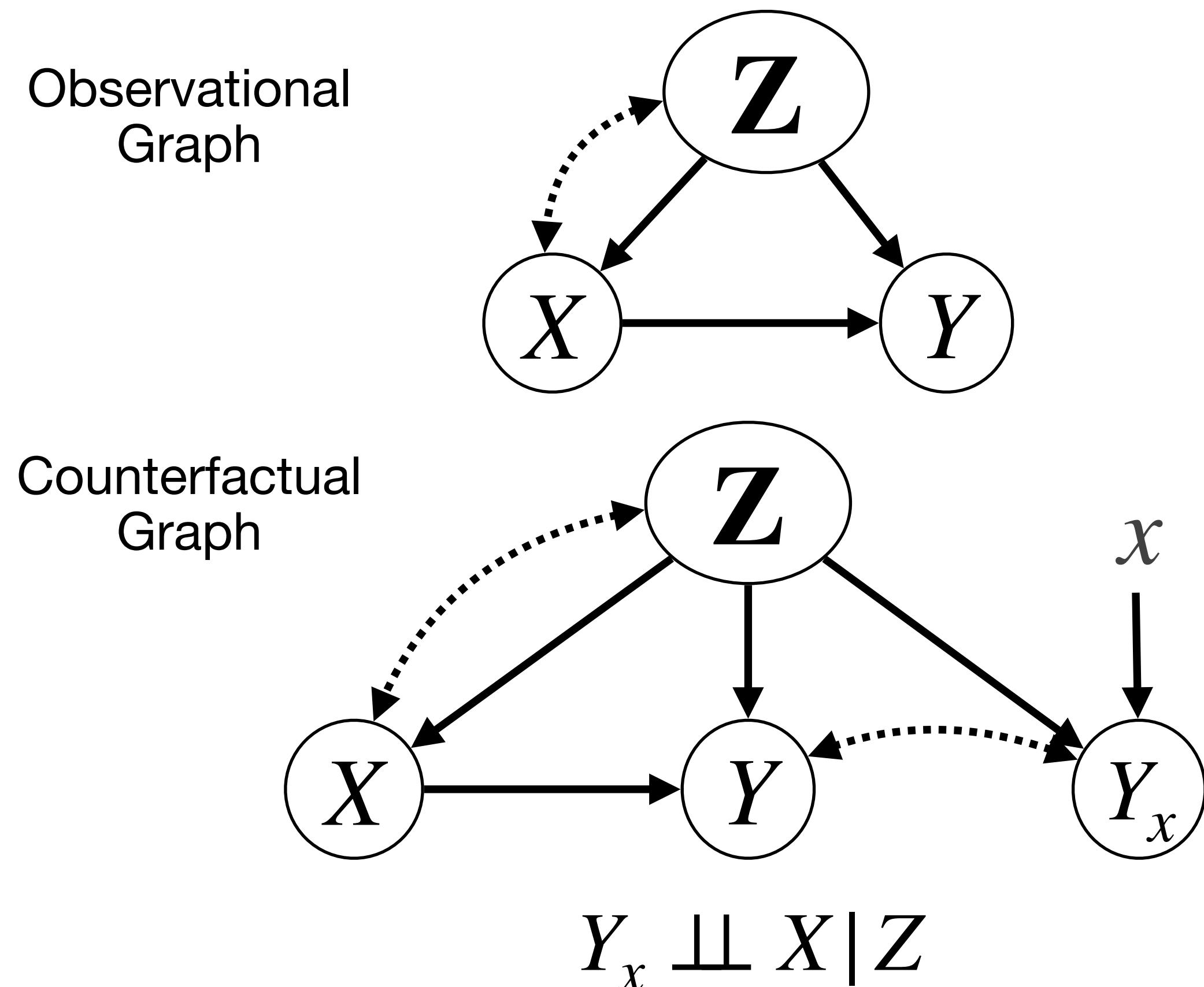
Load Remove

Counterfactual Interpretation of Backdoor

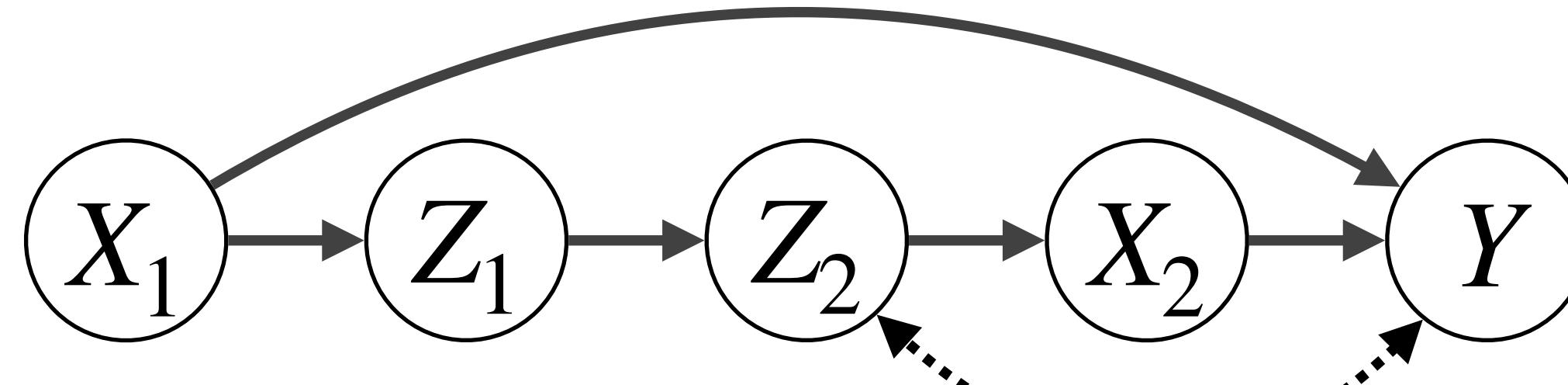
Theorem 4.3.1, Pearl's Primer Book

Theorem: If a set Z satisfies the *backdoor criterion* w.r.t. the ordered pair (X, Y) , then, for all x , it holds that $Y_x \perp\!\!\!\perp X | Z$.

Although the satisfiability of Z to the *backdoor criterion* can be tested given a causal diagram or a PAG, the condition $Y_x \perp\!\!\!\perp X | Z$ is sometimes framed as an assumption, referred to as **(conditional) ignorability, exchangeability or unconfoundedness**.



Is BD complete for covariate adjustment?



There are open backdoor paths between $(\{X_1, X_2\}, \{Y\})$, e.g.. $\langle X_2, Z_2, Y \rangle$ and $\langle X_2, Z_2, Z_1, X_1, Y \rangle$.

However, no set is admissible for **backdoor adjustment** for $(\{X_1, X_2\}, \{Y\})$, as both Z_1 and Z_2 are descendants of $\{X_1, X_2\}$ in G .

$$\text{However, we have } P(y | do(x_1, x_2)) = \sum_z P(y | x_1, x_2, z_1, z_2) P(z_1, z_2)$$

Verify on causalfusion.net

Backdoor Criterion is **sound (sufficient)** but not **complete (necessary)** for covariate adjustment.

There are descendants of \mathbf{X} that may be used (and sometimes needed) for adjustment.

Can we have a sound and complete graphical criterion for covariate adjustment?

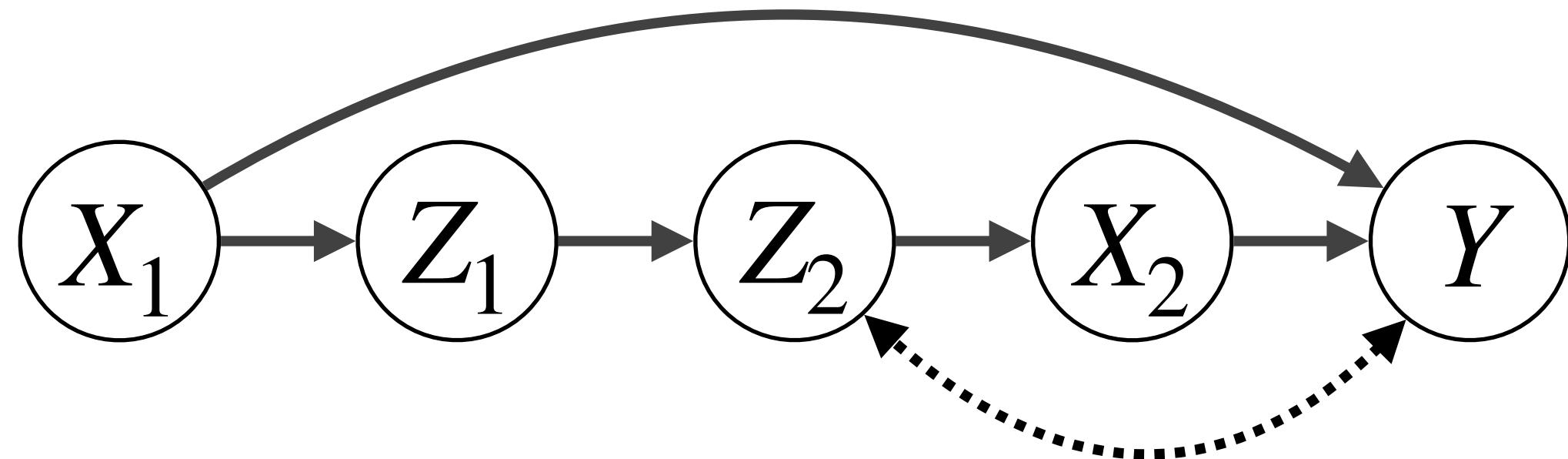
The Adjustment Criterion

Sound and Complete Graphical Criterion for Covariate Adjustment:

- Shpitser. et al (2010) - Shpitser, I., VanderWeele, T., & Robins, J. M. On the validity of covariate adjustment for estimating causal effects. UAI'10: Proceedings of the Twenty-Sixth Conference on Uncertainty in Artificial Intelligence - ([Link](#))

Proper Paths

Definition (Proper Path): A path p between a node in \mathbf{X} and a node in \mathbf{Y} is said to be *proper* if only its first node is in \mathbf{X} .



Causal paths between \mathbf{X} and \mathbf{Y} :

- | | | |
|---|---|------------------------|
| $\langle X_1, Y \rangle$ | } | Proper causal paths |
| $\langle X_2, Y \rangle$ | | Non-proper causal path |
| $\langle X_1, Z_1, Z_2, X_2, Y \rangle$ | | |

$$\mathbf{X} = \{X_1, X_2\}$$

$$\mathbf{Y} = \{Y\}$$

Non-causal paths between \mathbf{X} and \mathbf{Y} :

- | | | |
|---|---|----------------------------|
| $\langle X_2, Z_2, Y \rangle$ | } | Proper non-causal path |
| $\langle X_2, Z_2, Z_1, X_1, Y \rangle$ | | Non-proper non-causal path |

Adjustment Criterion for Causal Diagrams

Let \mathbf{X} be a set of treatment variables and \mathbf{Y} a set of outcome variables in the causal graph G .

$$W \in Desc(W)$$

If there exists a set \mathbf{Z} such that:

1. for every $X \in \mathbf{X}$ and $Y \in \mathbf{Y}$, \mathbf{Z} blocks every **proper non-causal path between X and Y** , and
2. no node in \mathbf{Z} is a descendant of a variable $W \notin \mathbf{X}$ which lies on a proper causal path from \mathbf{X} to \mathbf{Y}

Z cannot contain a variable in a proper causal path or a descendant of such variable.

Then, \mathbf{Z} satisfies the **adjustment criterion** for (\mathbf{X}, \mathbf{Y}) and, then the effect of \mathbf{X} on \mathbf{Y} is given by:

$$P(\mathbf{y} | do(\mathbf{x})) = \sum_{\mathbf{z}} P(\mathbf{y} | \mathbf{x}, \mathbf{z}) P(\mathbf{z})$$

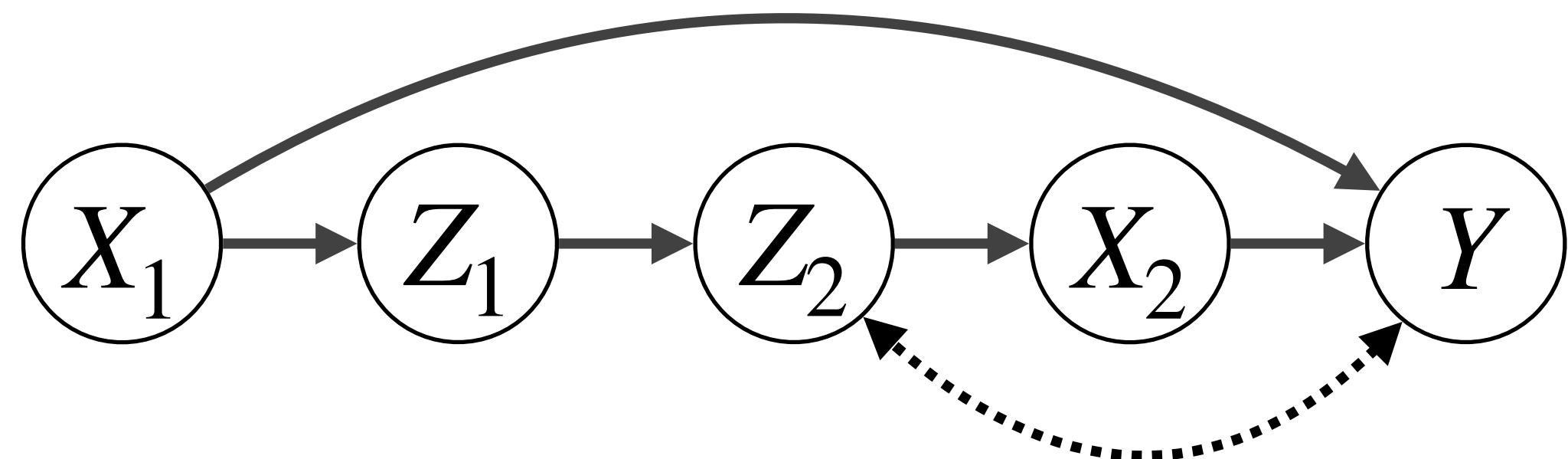
- $\mathbf{Z} = \{\}$ X
- $\mathbf{Z} = \{Z_2\}$ X
- $\mathbf{Z} = \{Z_1, Z_2\}$ ✓

Proper non-causal paths between \mathbf{X} and \mathbf{Y} :

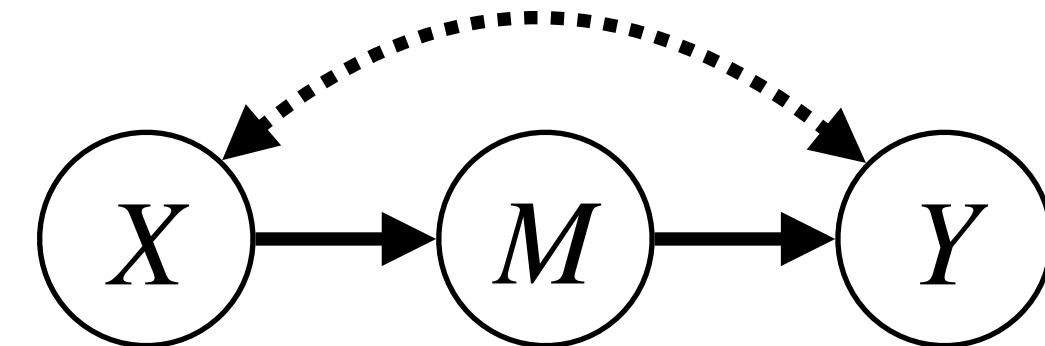
$$\langle X_2, Z_2, Y \rangle$$

Proper causal paths between \mathbf{X} and \mathbf{Y} :

$$\begin{aligned} &\langle X_1, Y \rangle \\ &\langle X_2, Y \rangle \end{aligned}$$

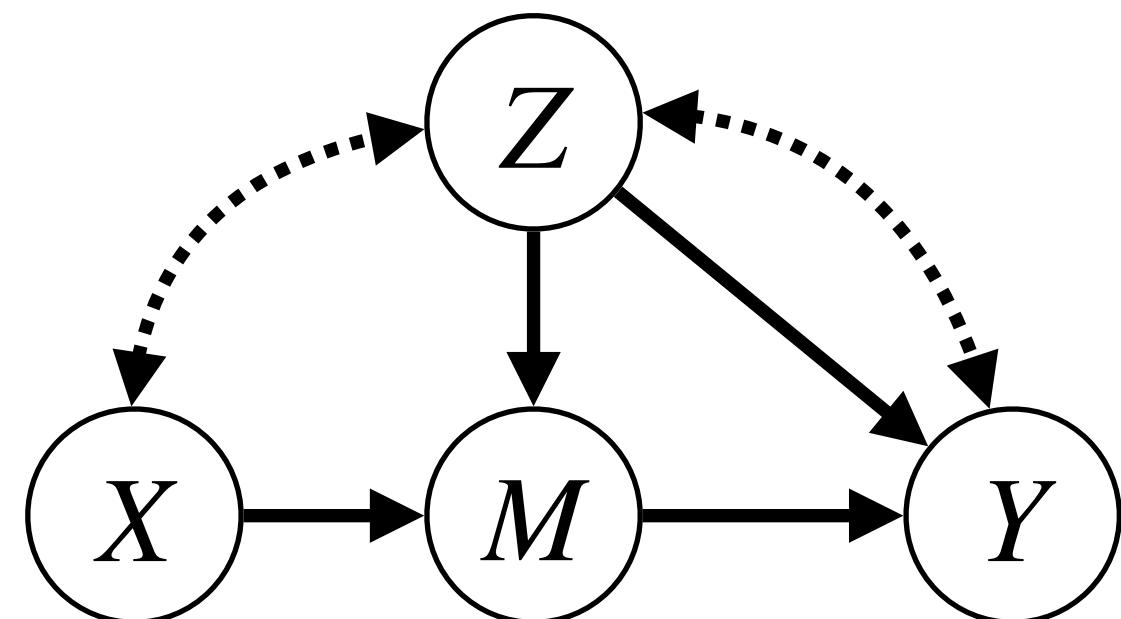


Many Scenarios Beyond Adjustment!

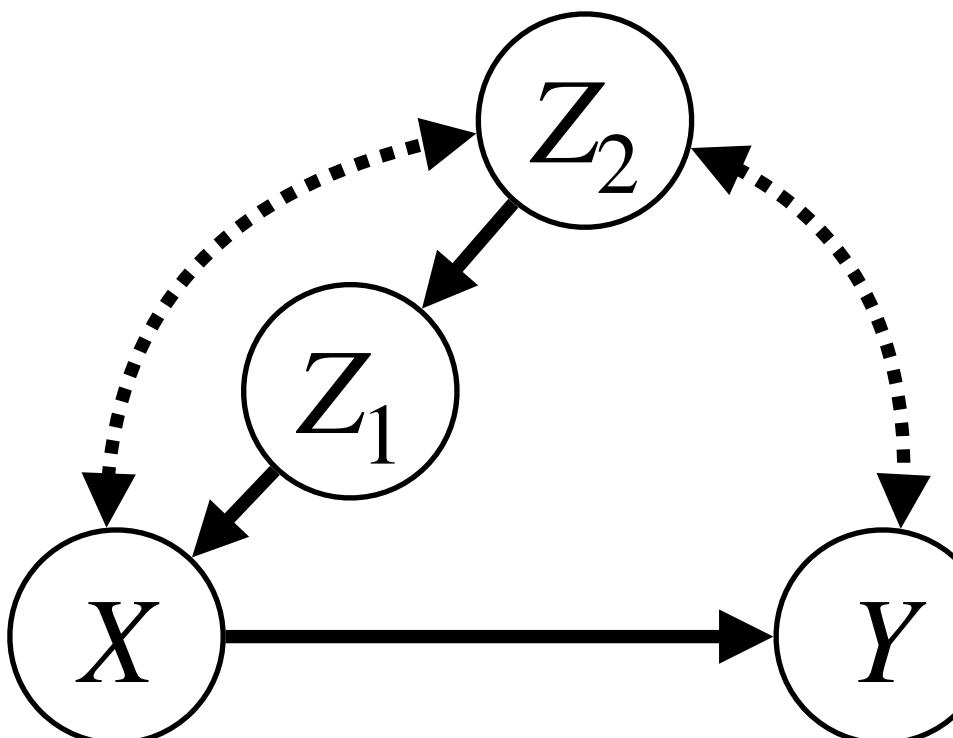


$$P(y | do(x)) = \sum_m P(m | x) \sum_{x'} P(y | m, x') P(x')$$

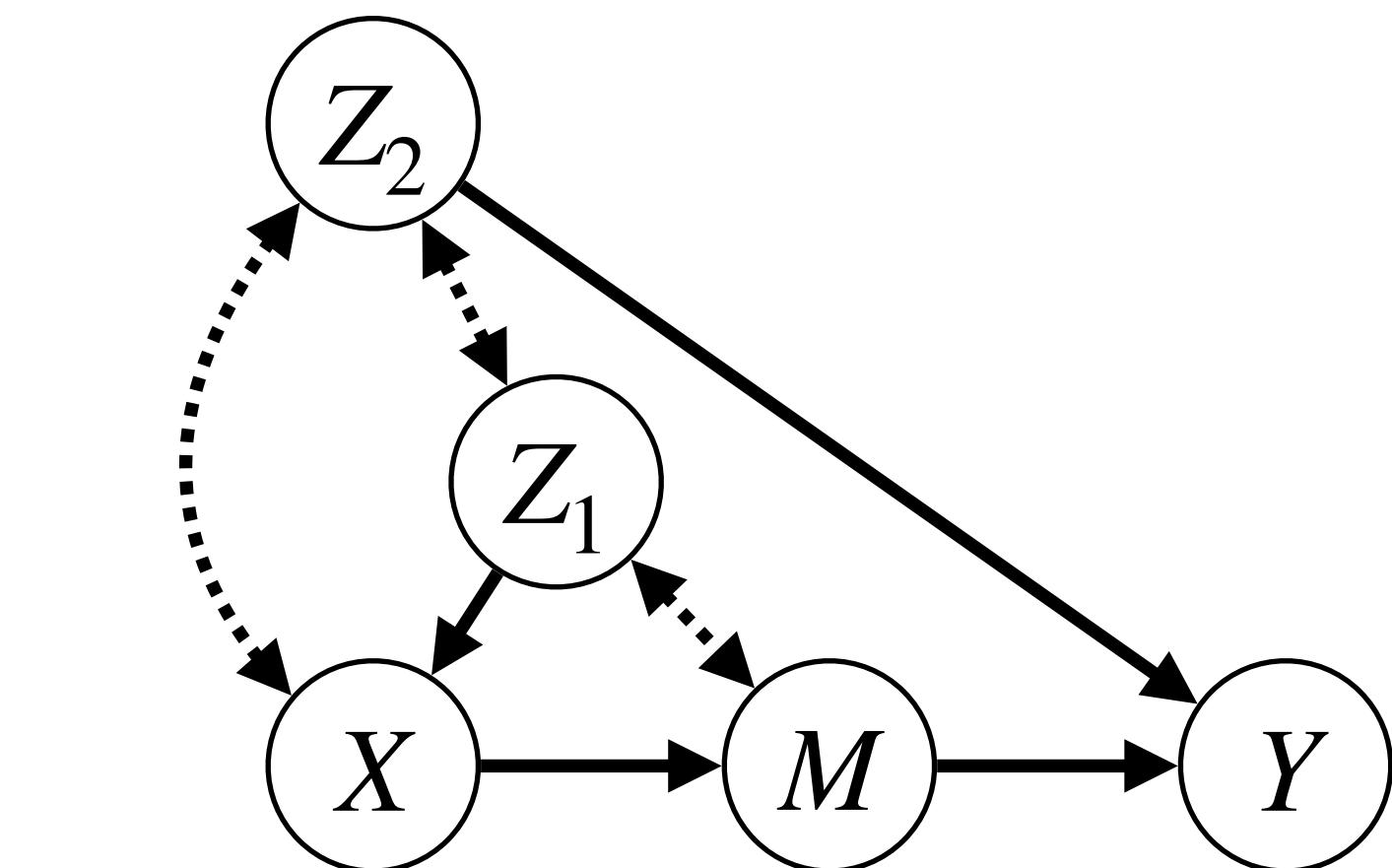
Front-Door



Conditional Front-Door



Napkin



Unnamed

$$P(y | do(x)) = \sum_{m,z} P(m | x, z) + \sum_{x'} P(y | m, x', z) P(x', z)$$

$$P(y | do(x)) = \frac{\sum_{z_2} P(x, y | z_1, z_2) P(z_2)}{\sum_{z_2} P(x | z_1, z_2) P(z_2)}$$

$$P(y | do(x)) = \sum_{z_2, z_3} P(y | x, z_1, z_2, z_3) P(z_2) + \sum_{z_1} P(z_3 | x, z_1) P(z_1)$$

And many others....

Tools for Causal Identification

1. Markovian Models (No Unobserved Confounders)
 - i. Truncated Factorization / G-computation or G-formula
2. Adjustment over Parents (No Unobserved Parents)
3. Non-Markovian Models (Under the Presence of Unobserved Confounders)
 - i. Graphical criteria (Backdoor Adjustment, Generalized Adjustment, Front-door Adjustment)
 - ii. Do-Calculus (a.k.a Causal Calculus)
 - iii. Identify Algorithm (a.k.a. ID algorithm)

Pearl, J. (2000). *Causality: Models, Reasoning, and Inference*. Cambridge University Press, New York. <http://dx.doi.org/10.1017/CBO9780511803161>

Jin Tian. Studies in causal reasoning and learning. PhD thesis, University of California, Los Angeles, 2002.

Questions?



Coding Exercises

Causality Tutorial:

- **Code:** <https://github.com/adele/Causality-Tutorial/blob/main/main.Rmd>
- **HTML Output:** <https://github.com/adele/Causality-Tutorial/blob/main/main.html>

Check Part I:

1. Causal Modeling
2. Causal Effect Identification from Causal Diagrams

**What if domain knowledge does not allow
you construct a causal diagram?**



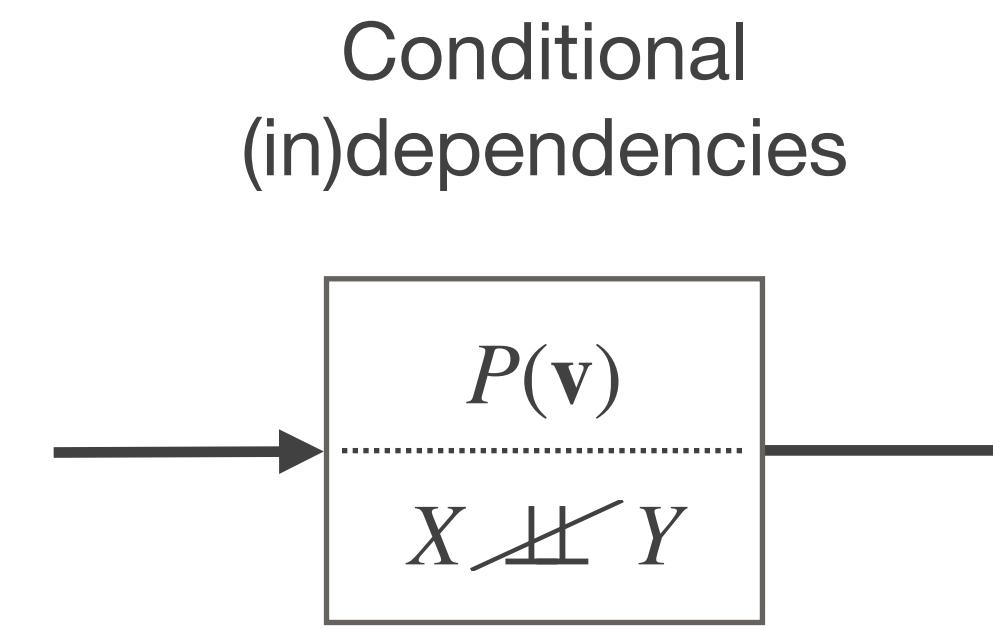
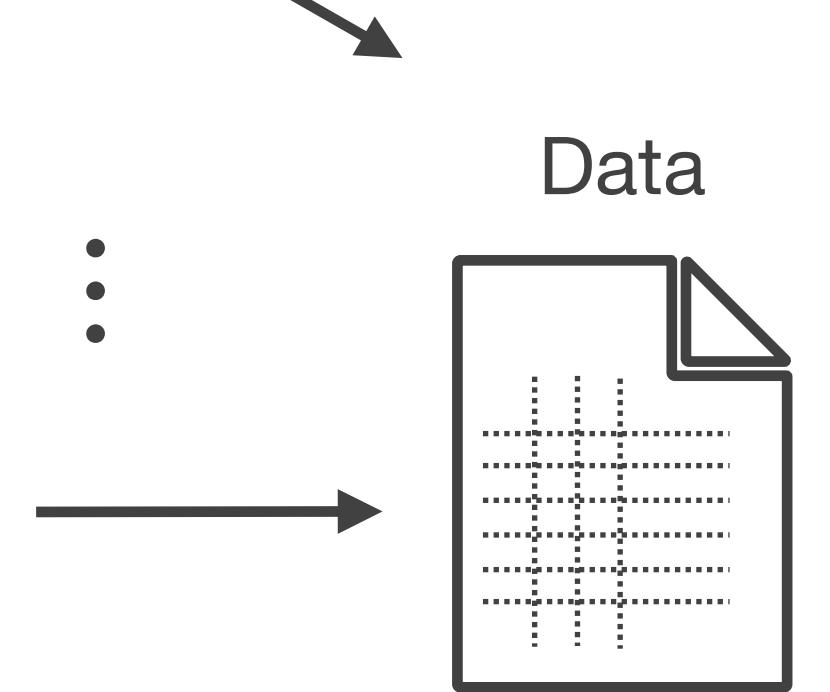
Markov Equivalence Class

$$\mathcal{M}_1 = \begin{cases} \mathbf{V} = \{X, Y\} \\ \mathbf{U} = \{U_x, U_Y\} \\ \mathcal{F} = \left\{ f_X(U_X) \atop f_Y(X, U_Y) \right. \\ P(\mathbf{U}) \end{cases}$$

⋮

$$\mathcal{M}_{N-1} = \begin{cases} \mathbf{V} = \{X, Y\} \\ \mathbf{U} = \{U_x, U_Y, U_{X,Y}\} \\ \mathcal{F} = \left\{ f_X(Y, U_X, U_{X,Y}) \atop f_Y(U_Y, U_{X,Y}) \right. \\ P(\mathbf{U}) \end{cases}$$

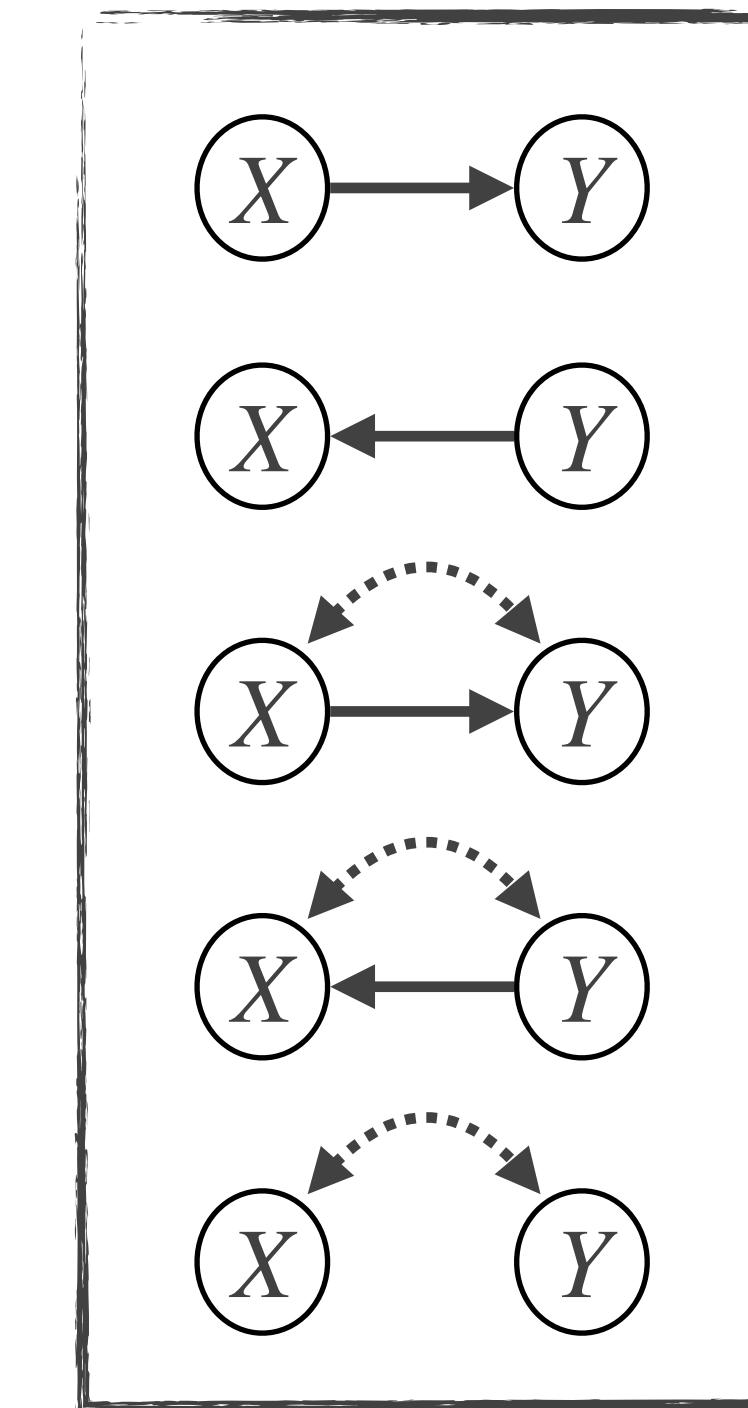
$$\mathcal{M}_N = \begin{cases} \mathbf{V} = \{X, Y\} \\ \mathbf{U} = \{U_x, U_Y\} \\ \mathcal{F} = \left\{ f_X(U_X) \atop f_Y(U_Y) \right. \\ P(\mathbf{U}) \end{cases}$$



$$P(x, y) = \sum_{u_x, u_y} P(x|y)P(y)P(u_x, u_y)$$

$$P(x, y) = \sum_{u_x, u_y} P(y|x)P(x)P(u_x, u_y)$$

Markov Equivalence Class
(class of models implying the same set of conditional independencies)



Correlation does not imply causation!

Super-Exponential Growth

The space of DAGs grows super-exponentially with the number n of variables, as shown by the following recurrence relation (Robinson, 1973):

$$|DAG(n)| = \sum_{i=1}^n \binom{n}{i} 2^{i(n-i)} |DAG(n-1)|$$

Inference through enumeration
is not a good idea!

n	$ DAG(n) $
2	3
3	27
4	729
5	59,049
6	1.4349×10^7
7	1.0460×10^{10}
8	2.2877×10^{13}

Super-Exponential Growth

The space of ADMGs also grows super-exponentially with the number n of variables, and it is much bigger than the space of DAGs:

$$|ADMG(n)| = |DAG(n)| \times 2^{n(n-1)/2}$$

$$|ADMG(n)| \gg |DAG(n)|$$

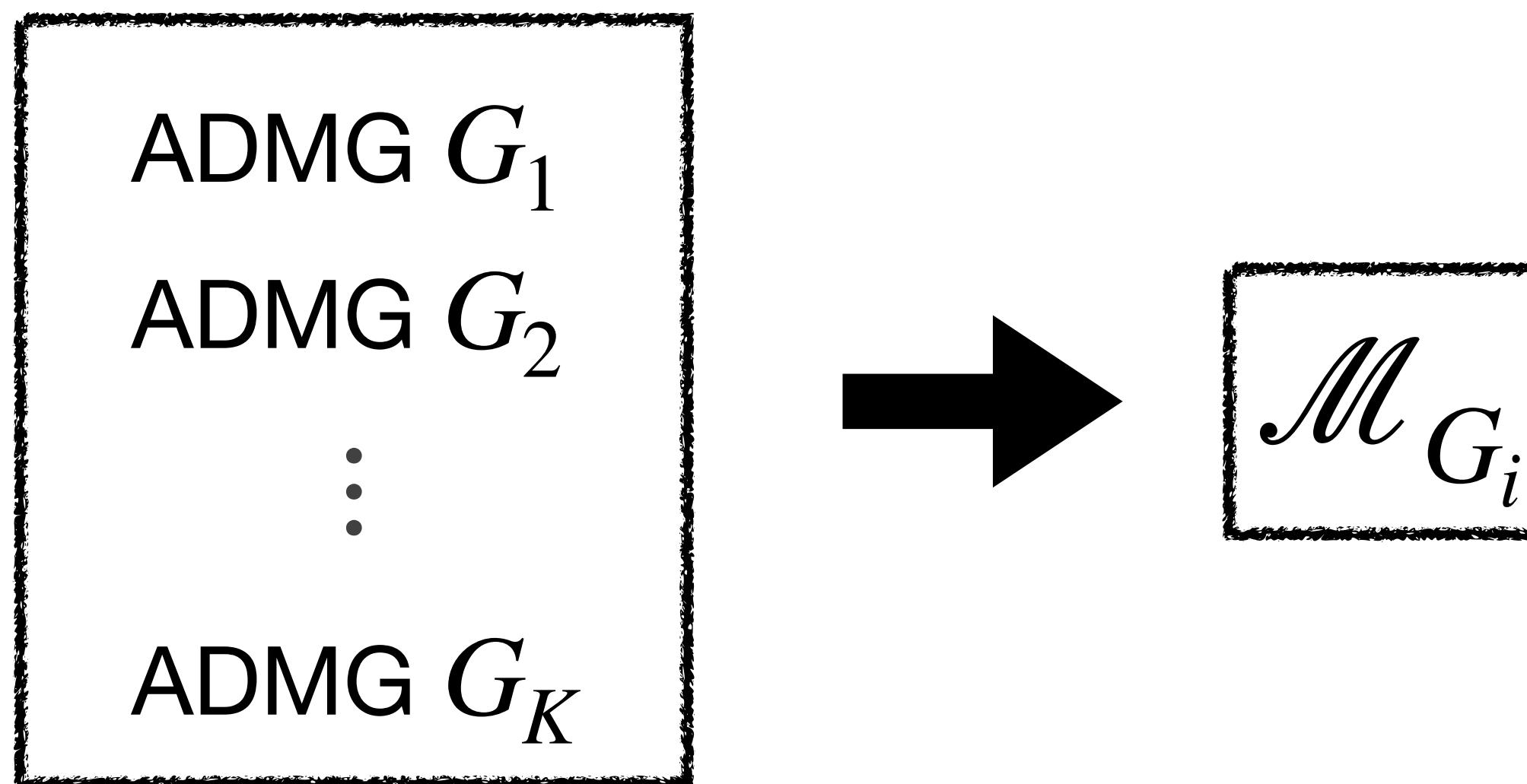
n	$ DAG(n) $	$ ADMG(n) $
2	3	6
3	27	216
4	729	46,656
5	59,049	6.0457×10^7
6	1.4349×10^7	4.7019×10^{11}
7	1.0460×10^{10}	2.1936×10^{16}
8	2.2877×10^{13}	6.1410×10^{21}

Ancestral Graphs

Markov Equivalence Class (MEC): Two models over \mathbf{V} , \mathcal{M}_1 and \mathcal{M}_2 , belong to the same MEC if for every disjoint sets $\mathbf{X}, \mathbf{Y}, \mathbf{Z} \subset \mathbf{V}$, $(\mathbf{X} \perp\!\!\!\perp \mathbf{Y} | \mathbf{Z})_{\mathcal{M}_1} \Leftrightarrow (\mathbf{X} \perp\!\!\!\perp \mathbf{Y} | \mathbf{Z})_{\mathcal{M}_2}$

Given any ADMG G over $\mathbf{V} \cup \mathbf{U} \cup \mathbf{S}$, there exists a *maximal ancestral graph* (MAG) \mathcal{M}_G over \mathbf{V} alone such that for any three disjoint sets of variables $\mathbf{X}, \mathbf{Y}, \mathbf{Z} \subseteq \mathbf{V}$,

$$(\mathbf{X} \perp\!\!\!\perp \mathbf{Y} | \mathbf{Z}, \mathbf{S})_G \Leftrightarrow (\mathbf{X} \perp\!\!\!\perp \mathbf{Y} | \mathbf{Z}, \mathbf{S})_{\mathcal{M}_G}$$



\mathcal{M}_G encodes all probabilistic constraints of multiple ADMG G without considering the \mathbf{U} variables!

Space of MAGs and Posets

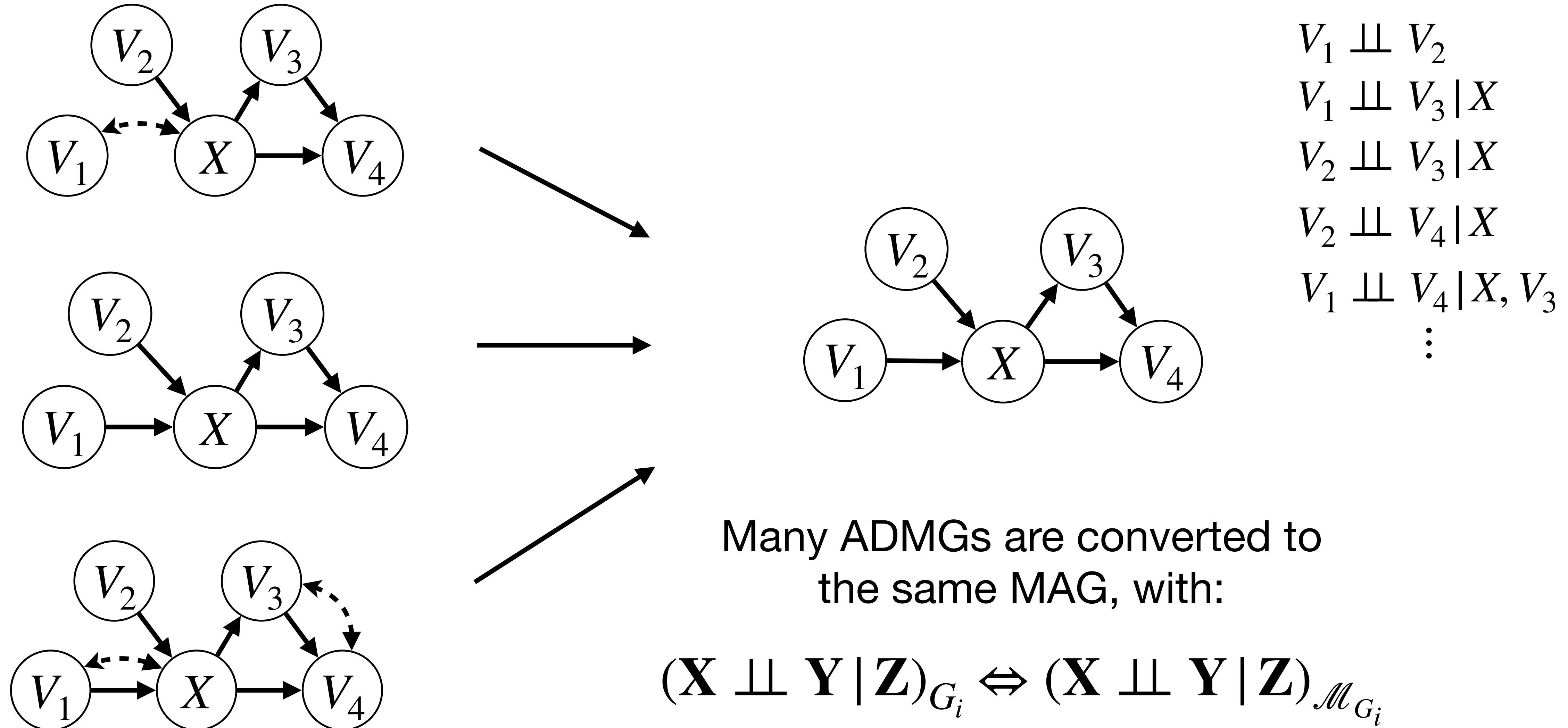
The space of MAGs is the same as the space of partially ordered sets (“posets”).

The recurrence can be found here: <https://oeis.org/A001035>

Growth is still super exponential, but slower than for DAGs and ADMGs.

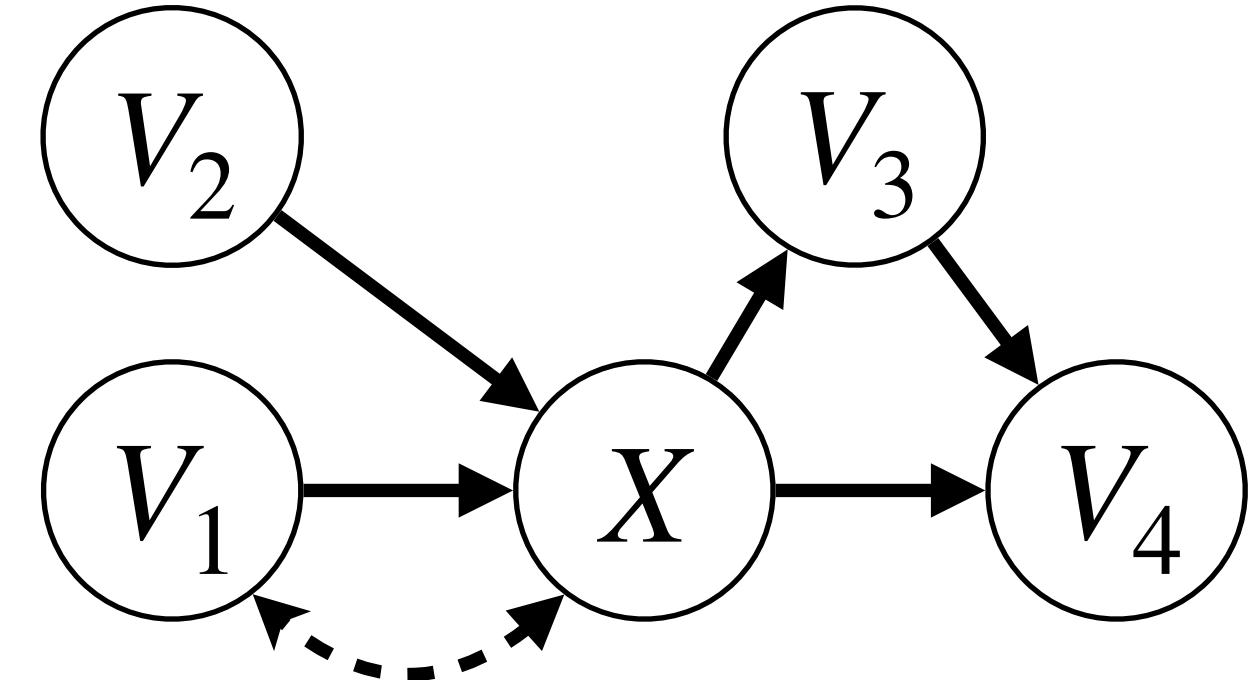
n	$ DAG(n) $	$ ADMG(n) $	$ MAG(n) $
2	3	6	3
3	27	216	19
4	729	46,656	219
5	59,049	6.0457×10^7	4231
6	1.4349×10^7	4.7019×10^{11}	130,023
7	1.0460×10^{10}	2.1936×10^{16}	6,129,859
8	2.2877×10^{13}	6.1410×10^{21}	4.3172×10^8

MAGs: Equivalence Class of ADMGs

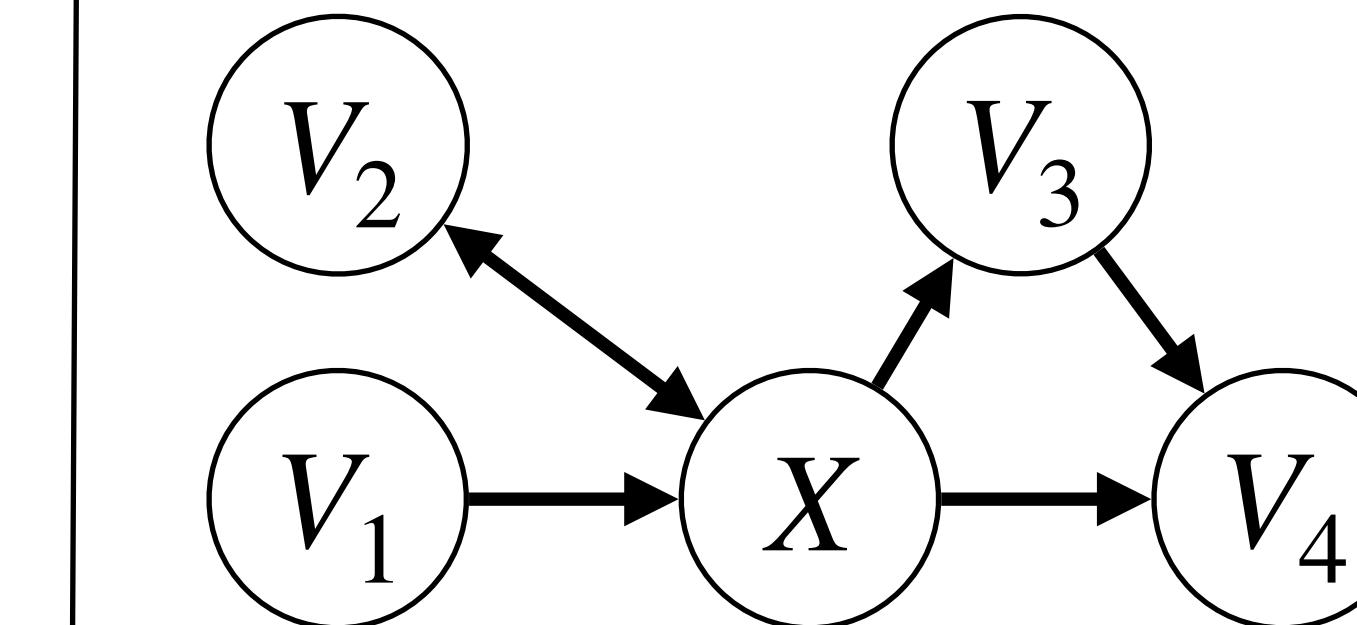
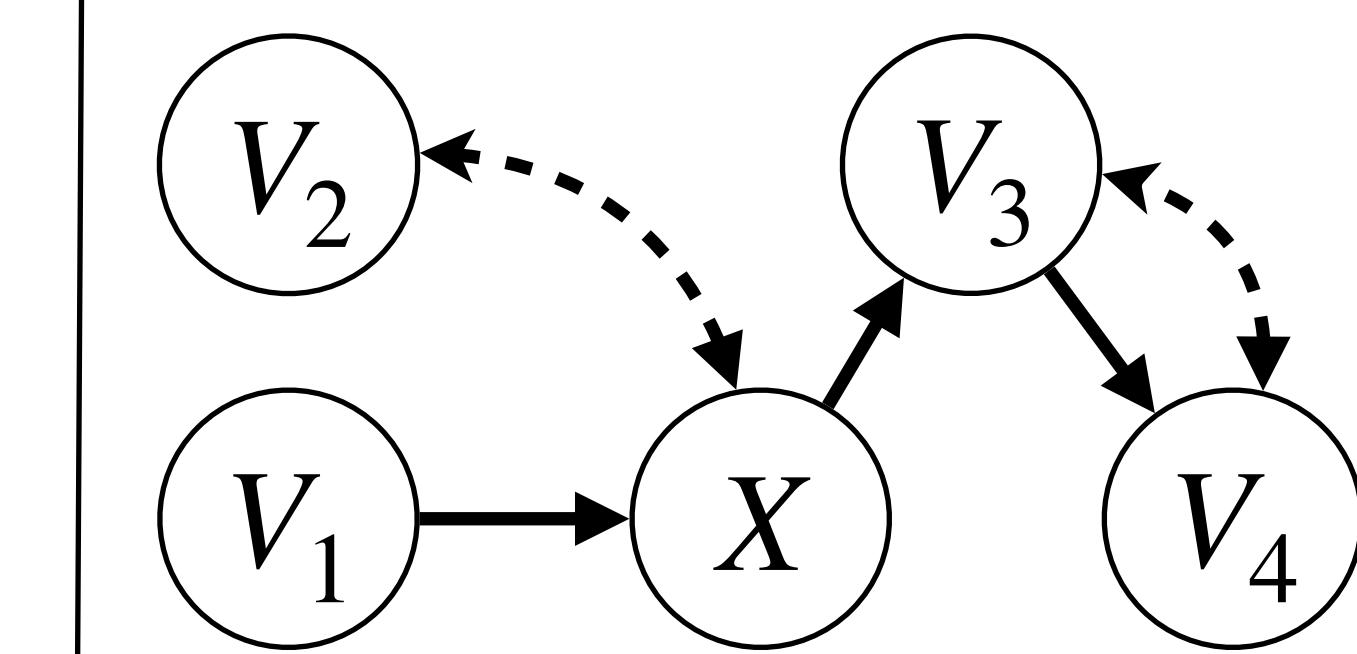
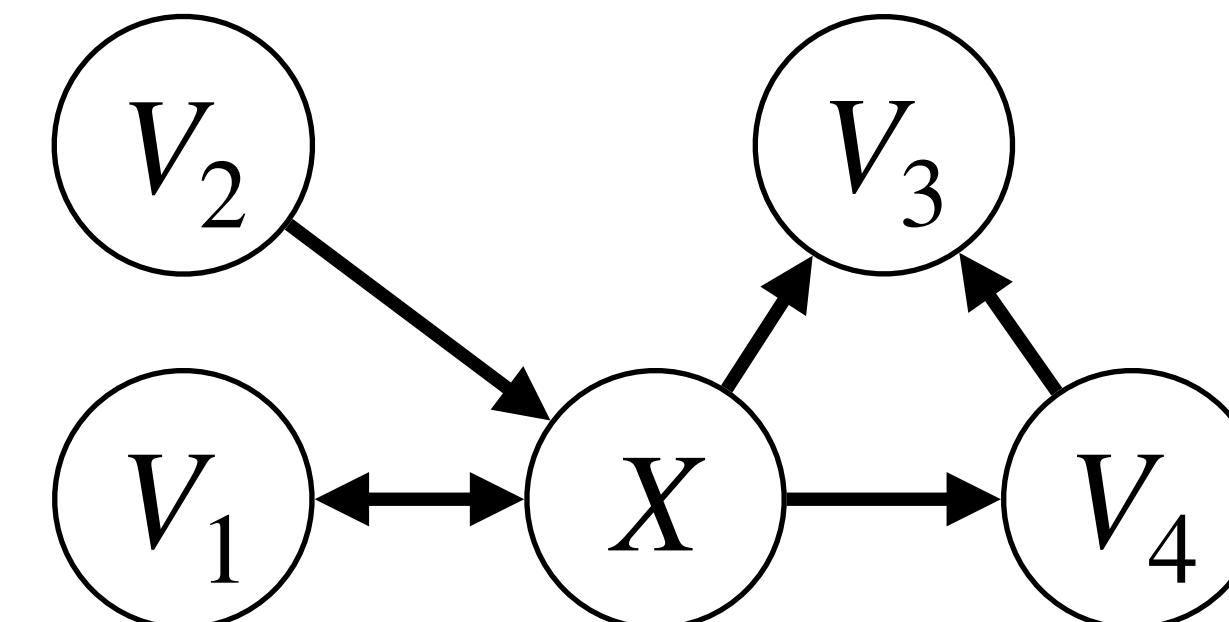
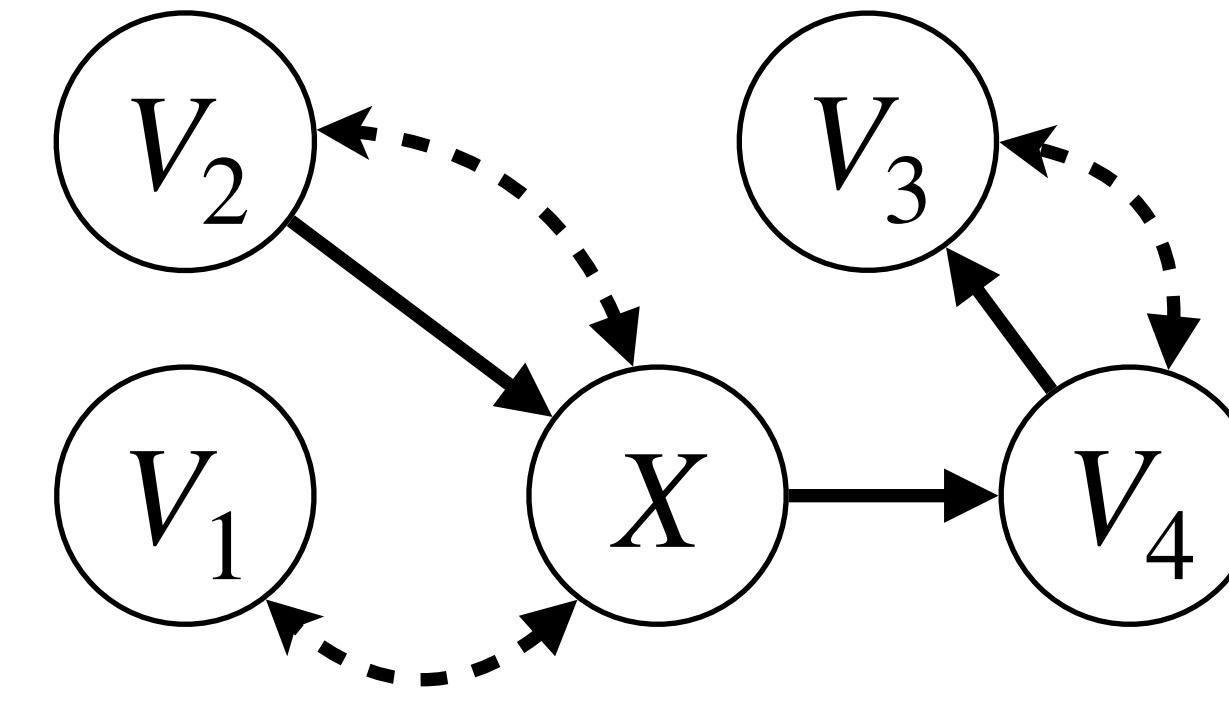
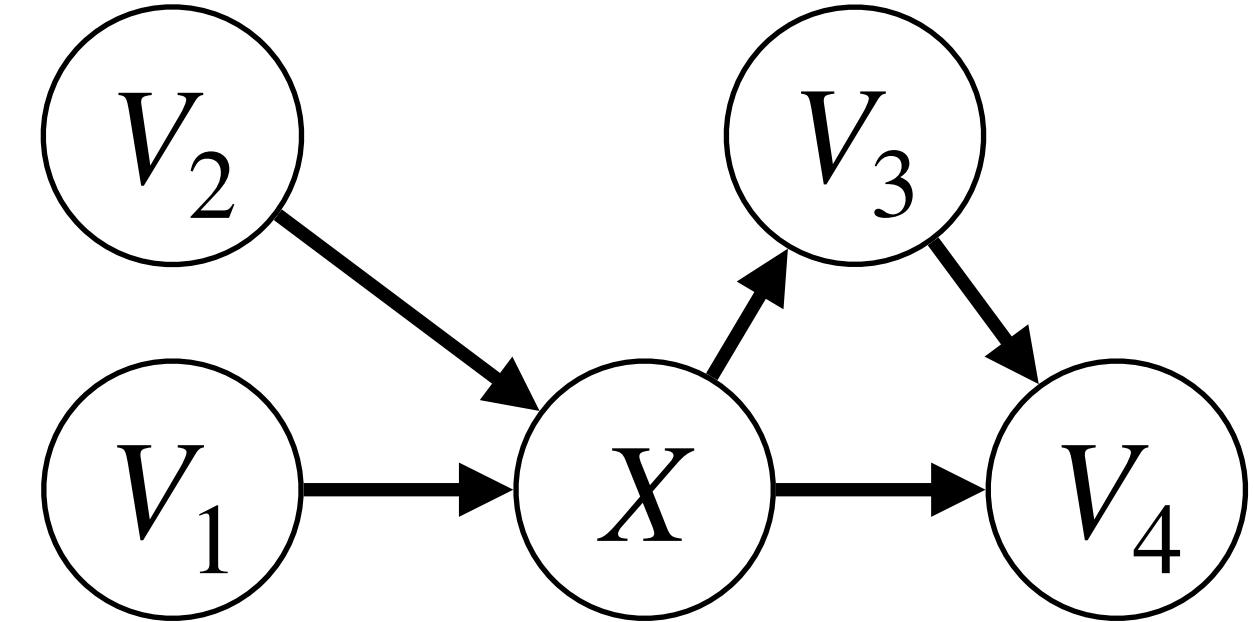


Markov Equivalent MAGs

ADMG



MAG

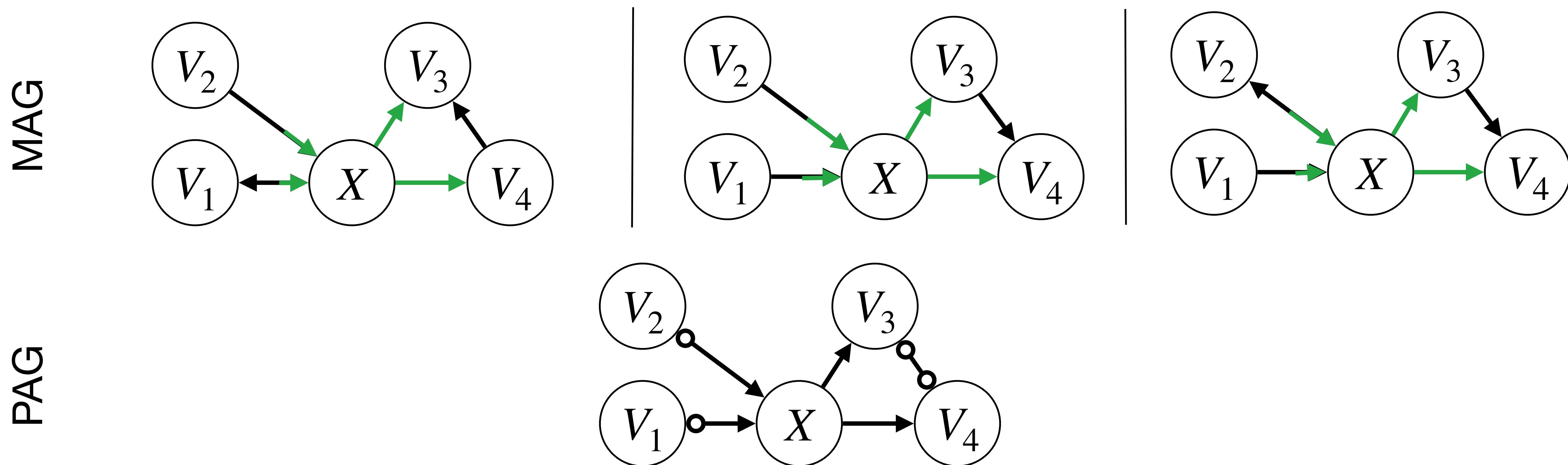


Different MAGs may entail the same set of conditional independencies.

PAGs: Markov Equivalence Class of MAGs

A **partial ancestral graph (PAG)** for $[\mathcal{G}]$ is a graph \mathcal{P} with possibly three kinds of marks (and hence six kinds of edges: $-$, \rightarrow , \leftrightarrow , \circ , $\circ\circ$, $\circ\rightarrow$), such that

- (1) \mathcal{P} has the same adjacencies as any member of \mathcal{G} does; and
- (2) every non-circle mark in \mathcal{P} is an invariant mark in $[\mathcal{G}]$.



Learning the Markov Equivalence Class

Causal Discovery:

Many models are statistically indistinguishable without additional parametric / distributional assumptions.

In non-parametric settings, causal discovery algorithms can only learn a graphical representation of its *Markov equivalence class* (MEC)!

Fast Causal Inference (FCI): Sound and complete causal discovery algorithm, even in the presence of unobserved confounders and selection bias.

Zhang, J. (2008). On the completeness of orientation rules for causal discovery in the presence of latent confounders and selection bias. *Artificial Intelligence*, 172(16):1873–1896. [Link](#)

Learning Structural Invariances

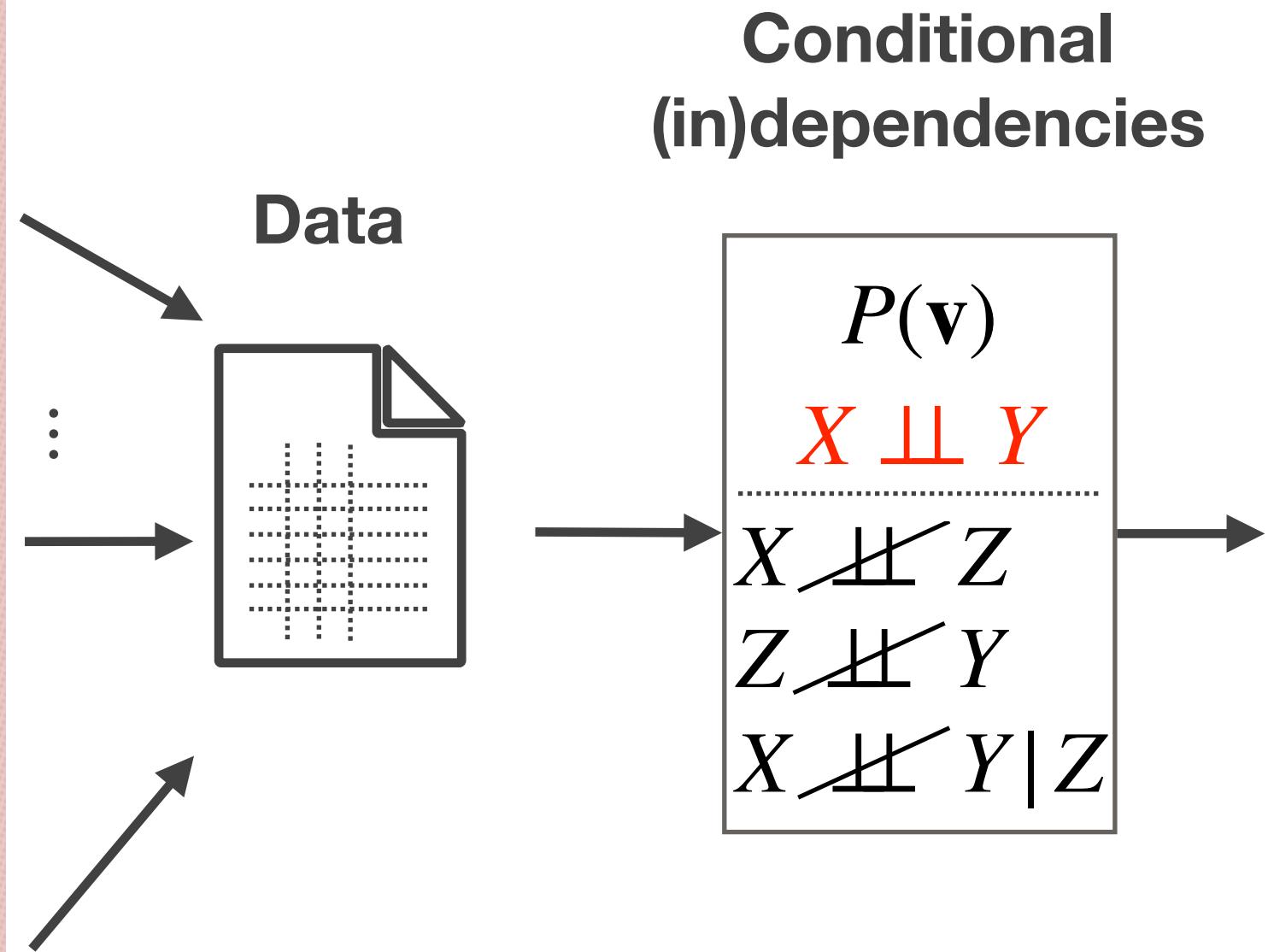
$$\mathcal{M}_1 = \begin{cases} V = \{X, Y, Z\} \\ U = \{U_x, U_y, U_z\} \\ \mathcal{F} = \begin{cases} X \leftarrow f_X(U_x) \\ Z \leftarrow f_Z(X, Y, U_z) \\ Y \leftarrow f_Y(U_y) \end{cases} \\ P(U) \end{cases}$$

⋮

$$\mathcal{M}_{N-1} = \begin{cases} V = \{X, Y, Z\} \\ U = \{U_{xz}, U_{yz}, U_x, U_y, U_z\} \\ \mathcal{F} = \begin{cases} X \leftarrow f_X(U_{xz}, U_x) \\ Z \leftarrow f_Z(Y, U_{xz}, U_z) \\ Y \leftarrow f_Y(U_y) \end{cases} \\ P(U) \end{cases}$$

⋮

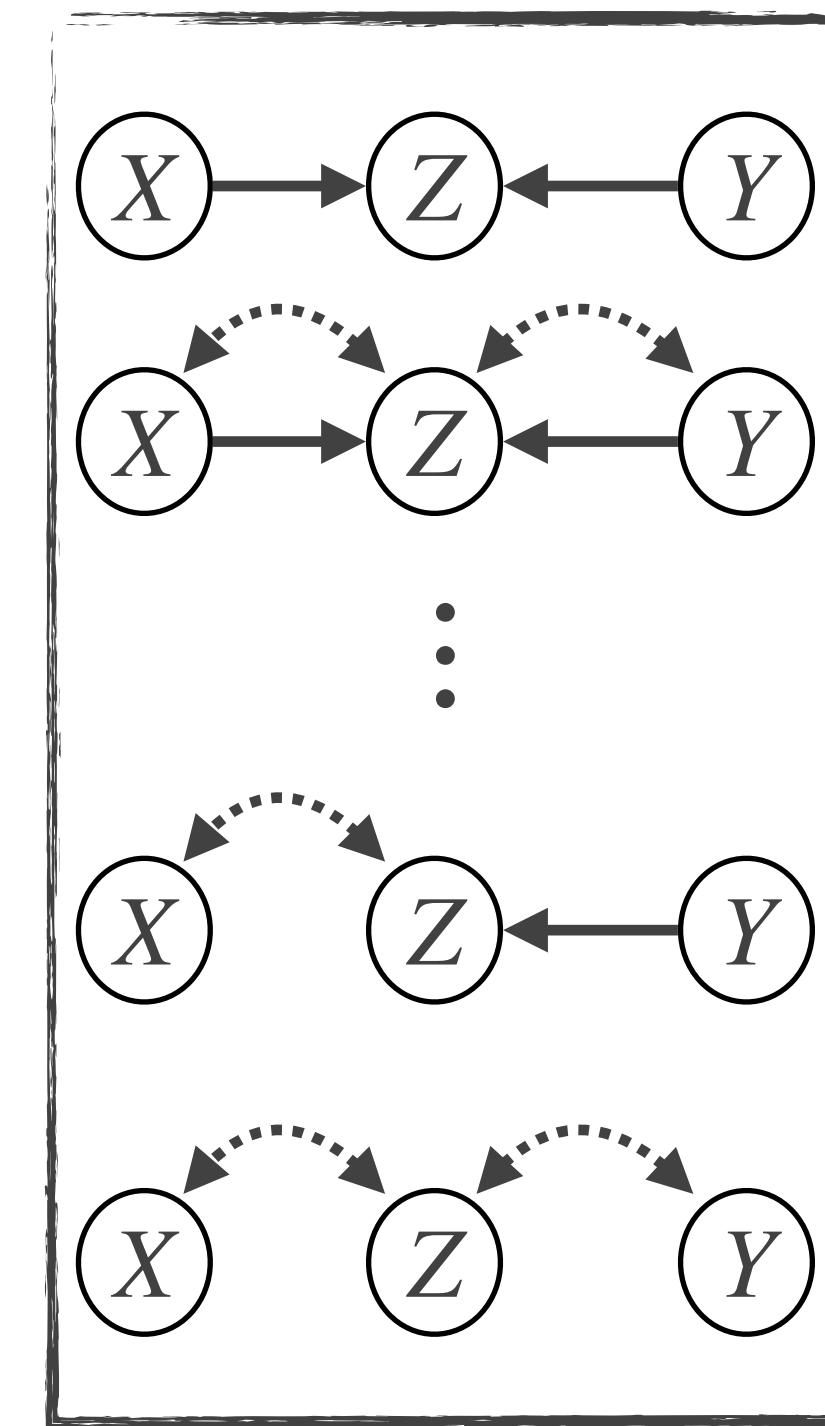
$$\mathcal{M}_N = \begin{cases} V = \{X, Y, Z\} \\ U = \{U_{xz}, U_{yz}, U_x, U_y, U_z\} \\ \mathcal{F} = \begin{cases} X \leftarrow f_X(U_{xz}, U_x) \\ Z \leftarrow f_Z(U_{xz}, U_{yz}, U_z) \\ Y \leftarrow f_Y(U_{yz}, U_y) \end{cases} \\ P(U) \end{cases}$$



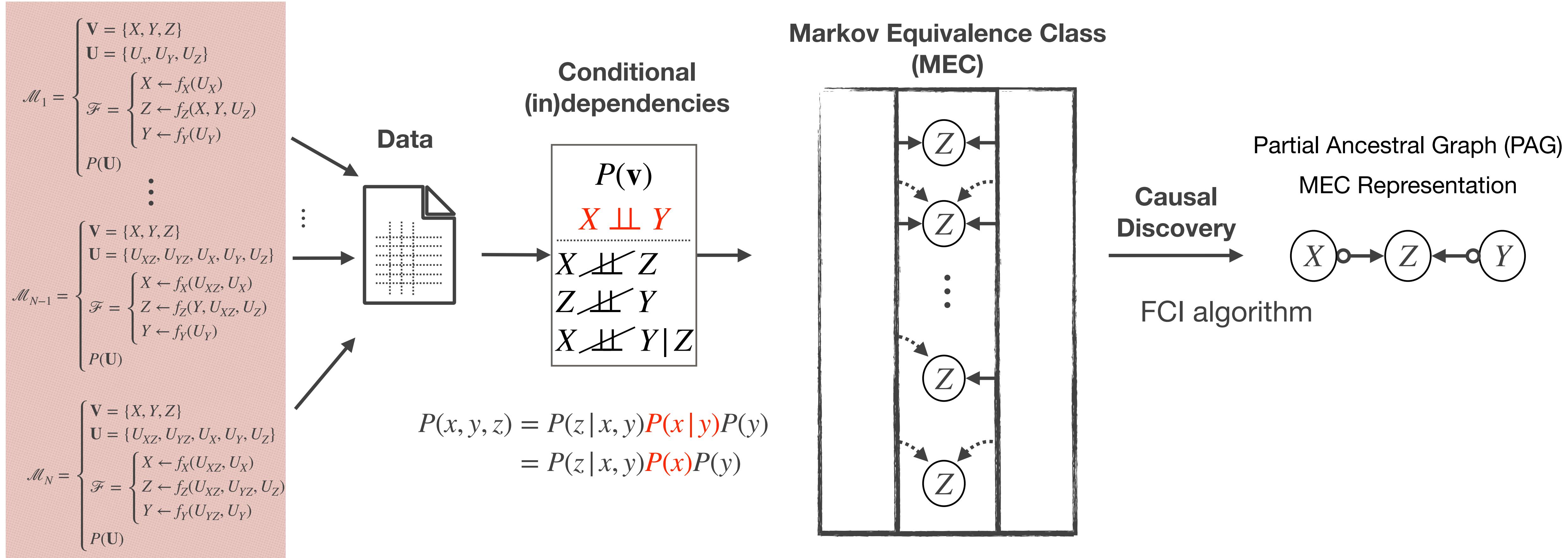
$$P(x, y, z) = P(z|x, y) \color{red} P(x|y) P(y)$$

$$= P(z|x, y) \color{red} P(x) P(y)$$

Markov Equivalence Class (MEC)

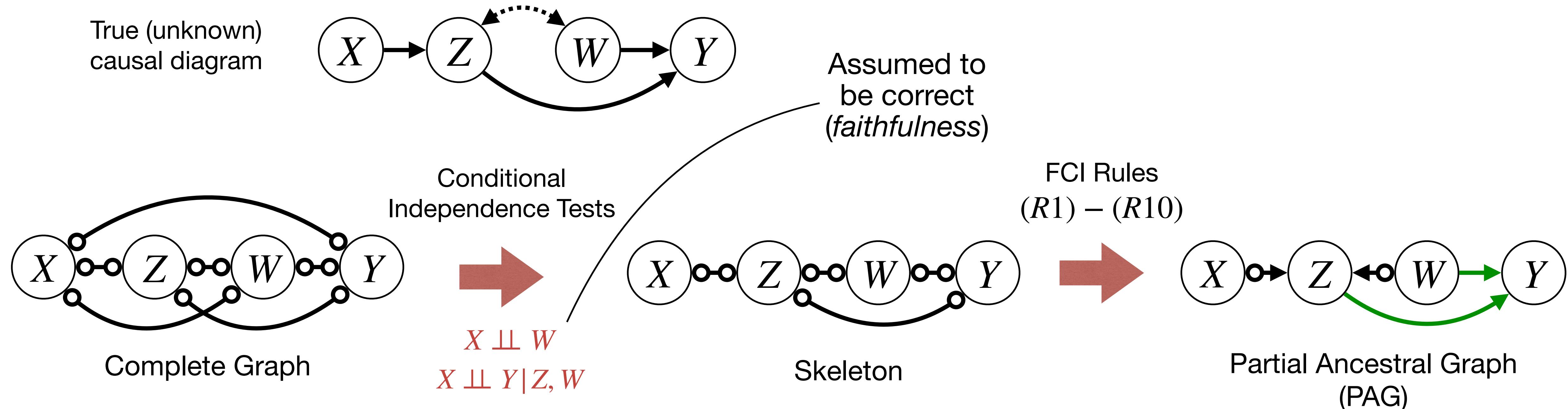


Learning Structural Invariances



Zhang, J. (2008). On the completeness of orientation rules for causal discovery in the presence of latent confounders and selection bias. *Artificial Intelligence*, 172(16):1873–1896. [Link](#)

Fast Causal Inference (FCI) Algorithm



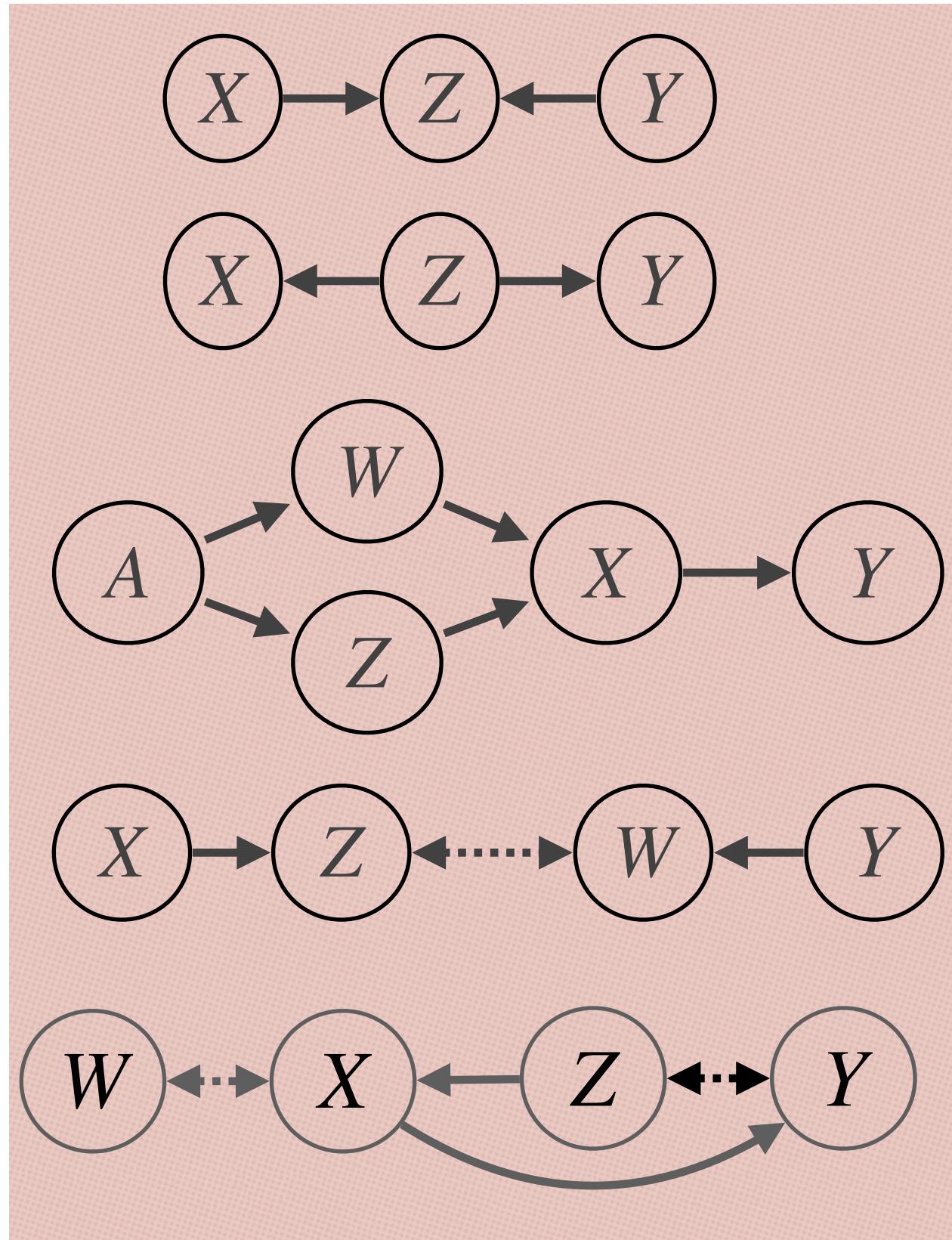
Arrowhead \implies non-ancestrality
Tail \implies ancestrally
Circle \implies non-invariance

$A \longleftrightarrow B$ – spurious association
 $A \dashv B$ – selection bias

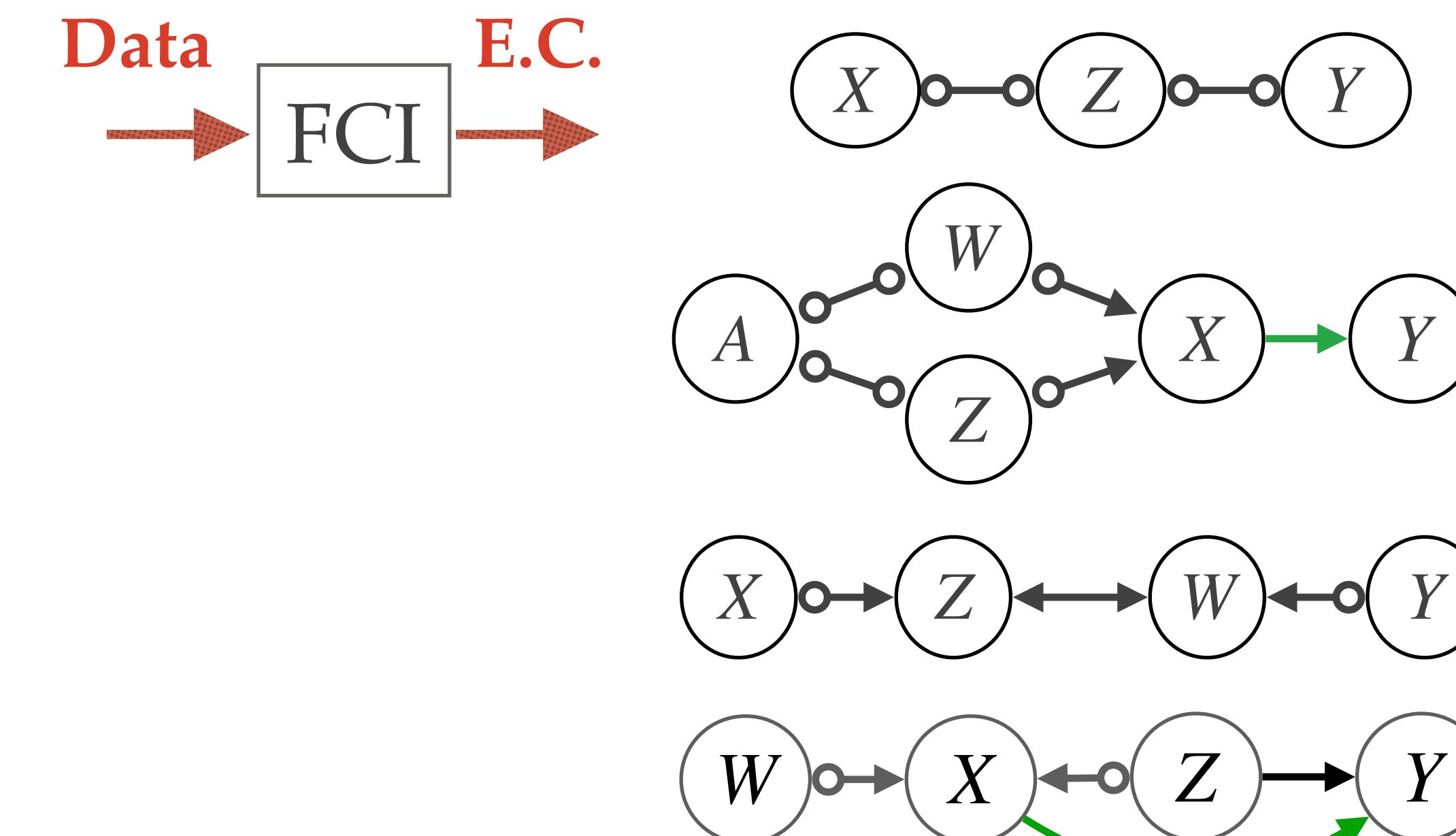
Z is not an ancestor of X or W.
 Z and W are ancestors (and definite causes) of Y.

Fast Causal Inference (FCI) Algorithm

Underlying Causal Diagram



Partial Ancestral Graph



Conditional Independence Tests

Gaussian errors and independent observations: partial correlation test

Fisher, R.A. (1921). *On the "Probable Error" of a Coefficient of Correlation Deduced from a Small Sample.*
R package: <https://cran.r-project.org/web/packages/pcalg/>

Kernel-based non-parametric test:

Zhang, K., Peters, J., Janzing, D., & Schölkopf, B. (2012). *Kernel-based conditional independence test and application in causal discovery.* In: Uncertainty in artificial intelligence. AUAI Press; 2011. p.804–13
R package: <https://cran.r-project.org/web/packages/CondIndTests>

Continuous (conditional Gaussian) or Discrete (Binary, Ordinal, Multinomial) - Linear Regression

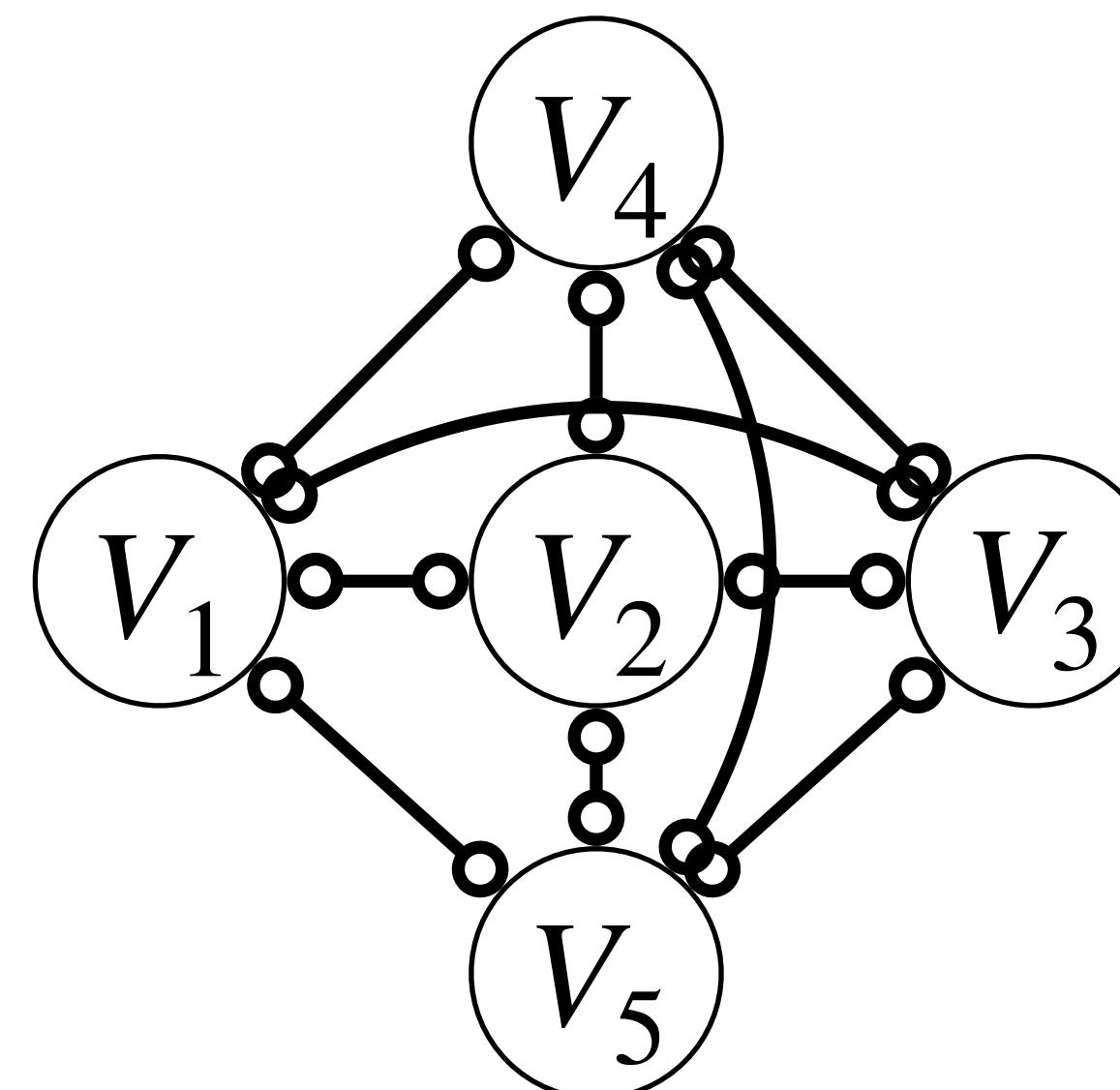
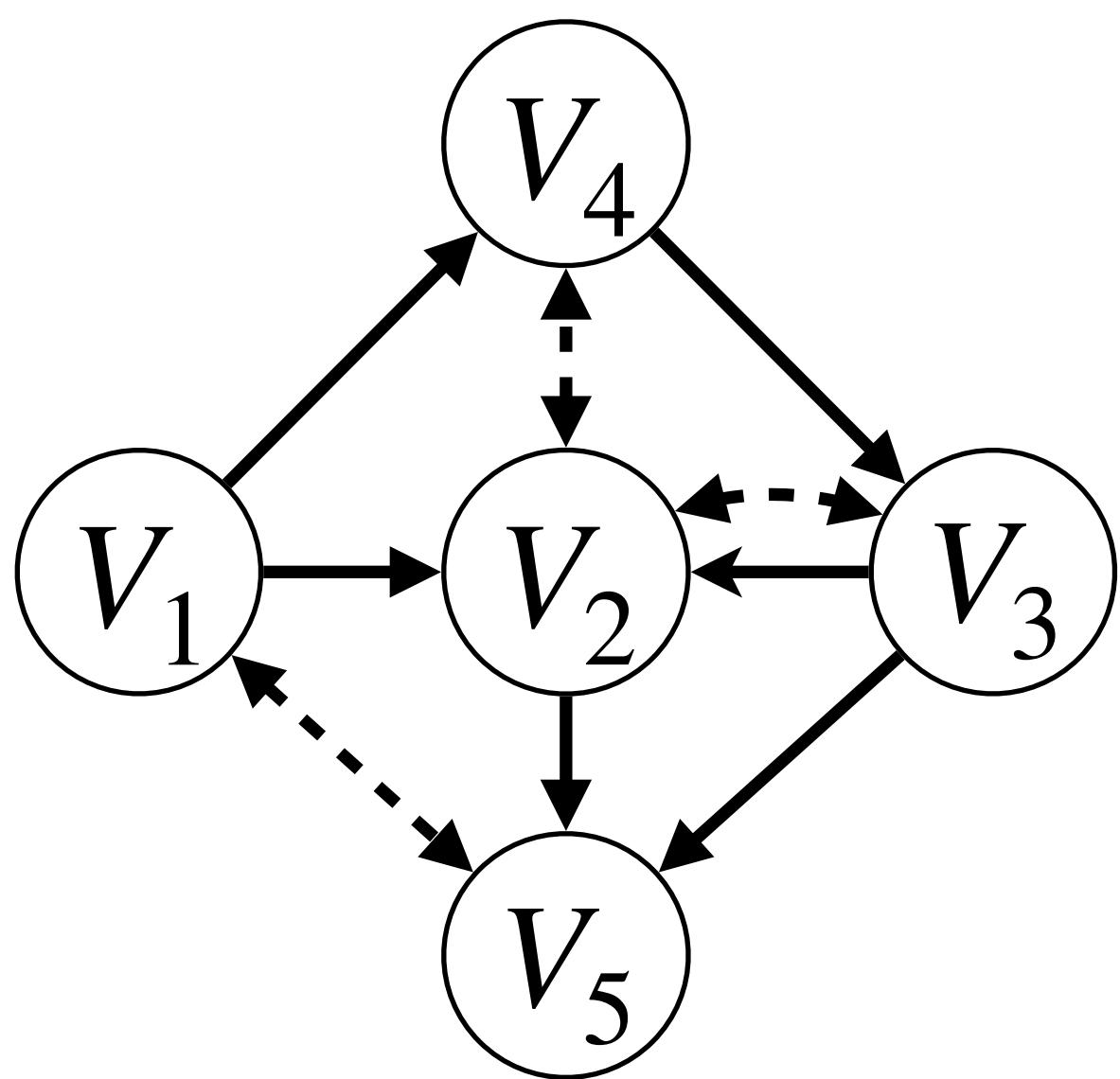
- Tsagris, M., Borboudakis, G., Lagani, V. et al. (2018) Constraint-based causal discovery with mixed data. *Int J Data Sci Anal* 6, 19–30. ([Link](#))
- R package: <https://cran.r-project.org/web/packages/MXM/>

Gaussian errors and correlated observations (family data) :

Ribeiro A.H., Soler J.M.P. (2020). *Learning Genetic and environmental graphical models from family data,* Statistics in Medicine.
R package: <https://github.com/adele/FamilyBasedPGMs>

FCI - Skeleton

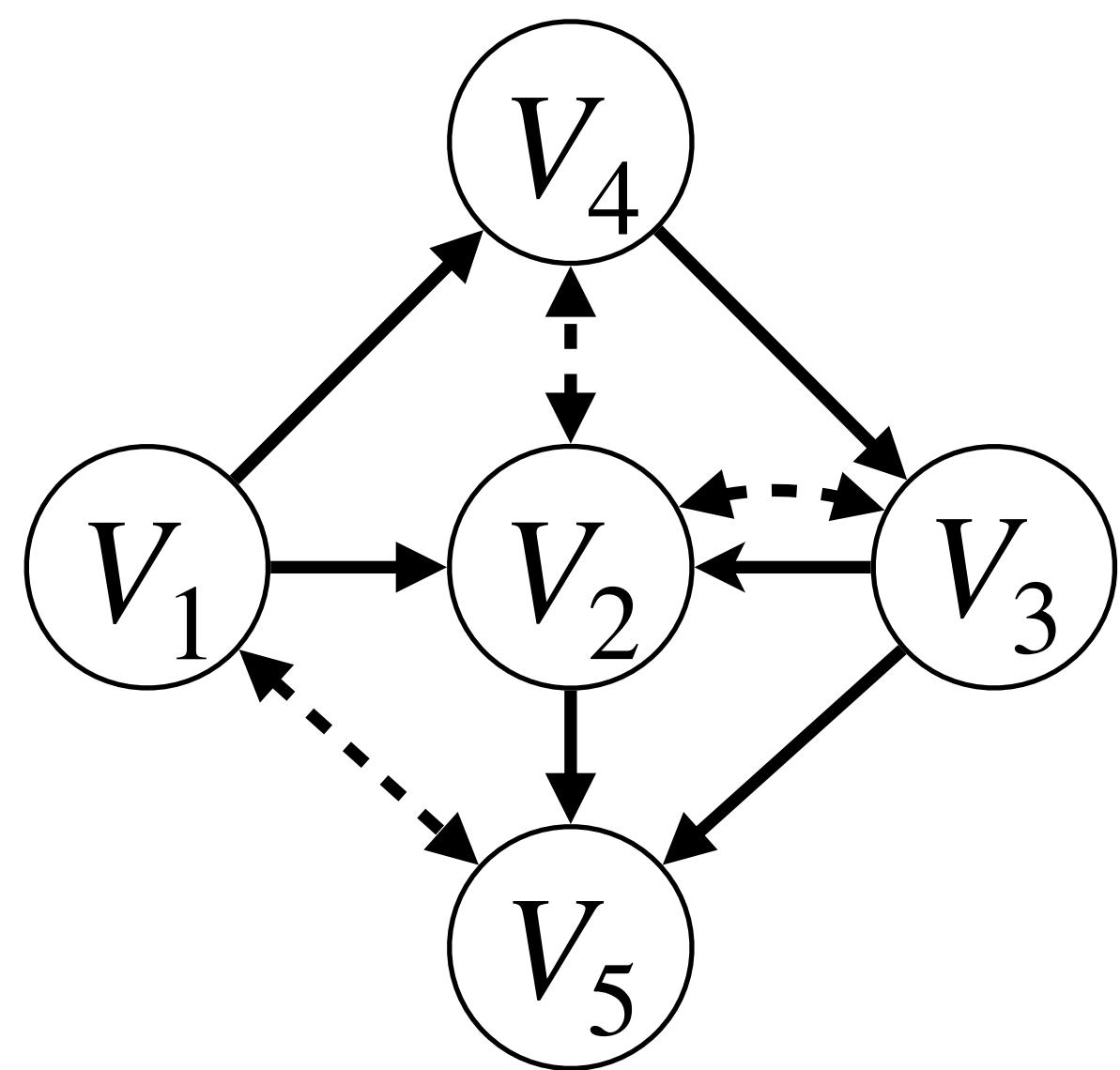
Form a complete graph on the set of variables, in which there is a circle-circle edge between every pair of variables;



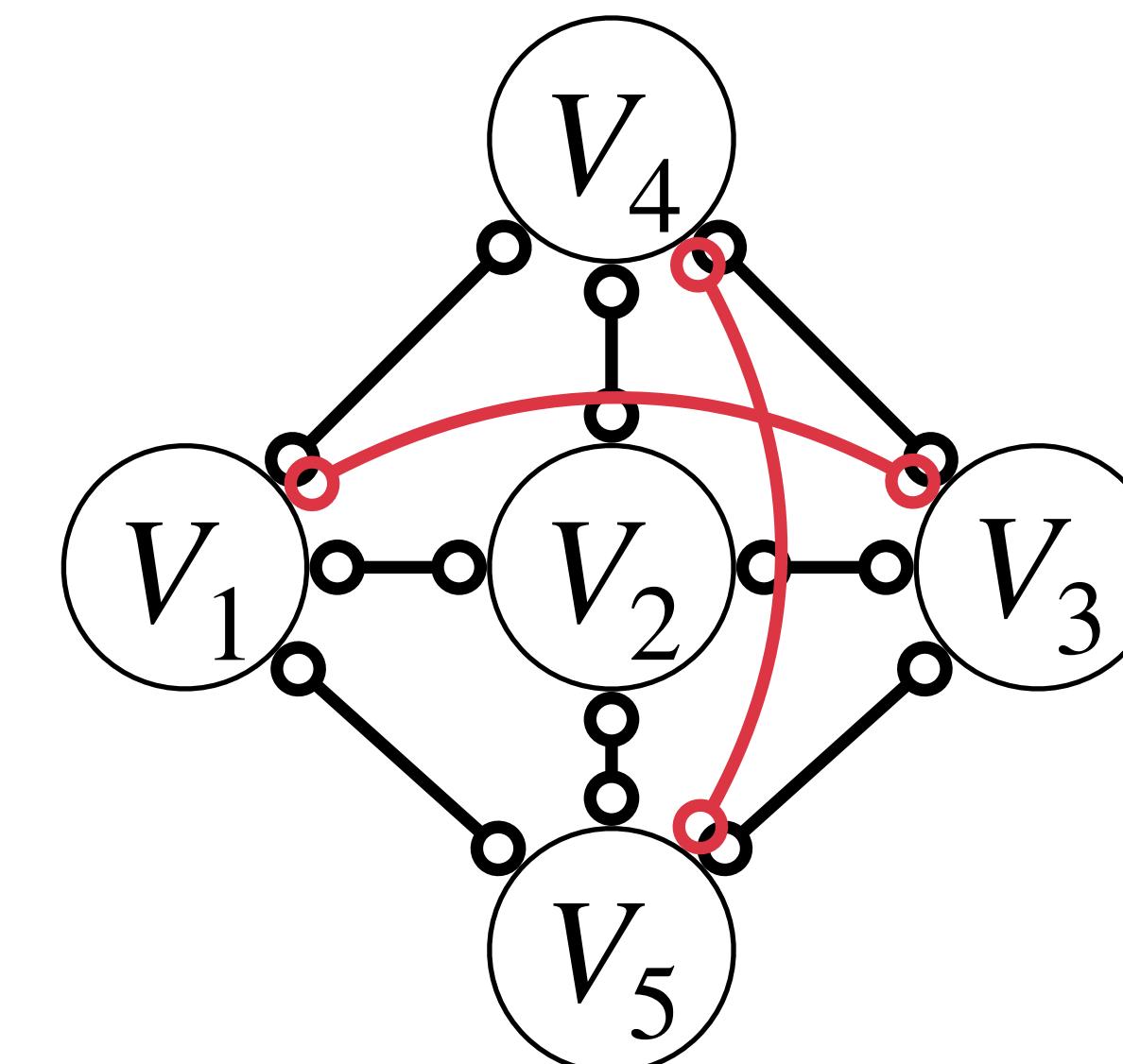
True, unknown ADMG

FCI - Skeleton

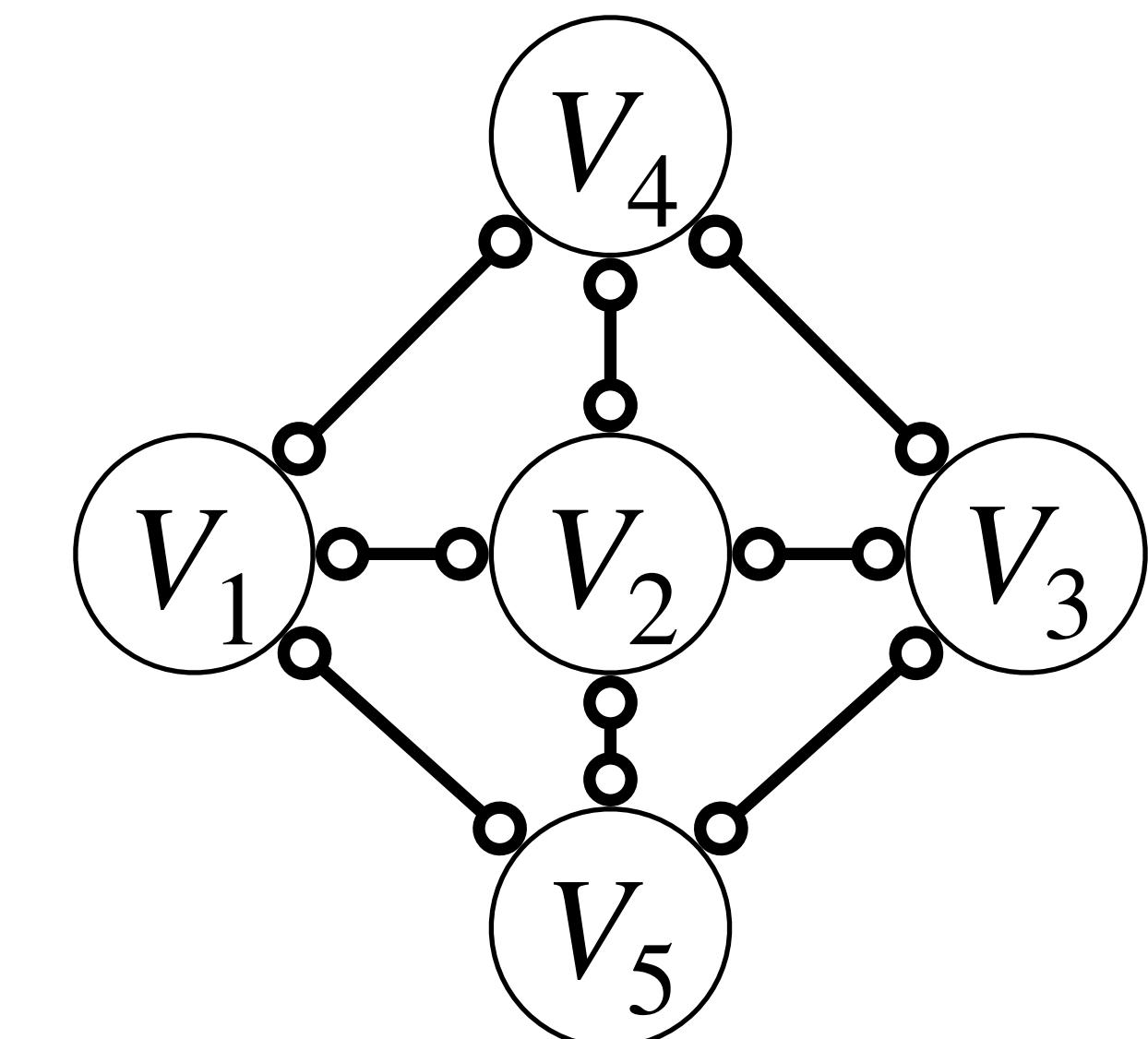
For every pair of variables V_1 and V_2 , if exists a set $S_{1,2}$ such that $V_1 \perp\!\!\!\perp V_2 | S_{1,2}$, then remove the edge between V_1 and V_2 and add $S_{1,2}$ in $\text{Sepset}(V_1, V_2)$.



True, unknown ADMG

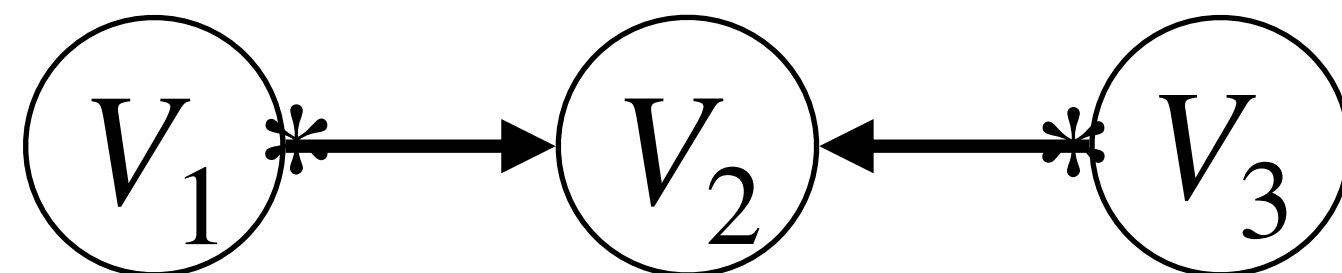


$V_1 \perp\!\!\!\perp V_3 | V_4$ and $V_4 \perp\!\!\!\perp V_5 | V_1, V_2, V_3$

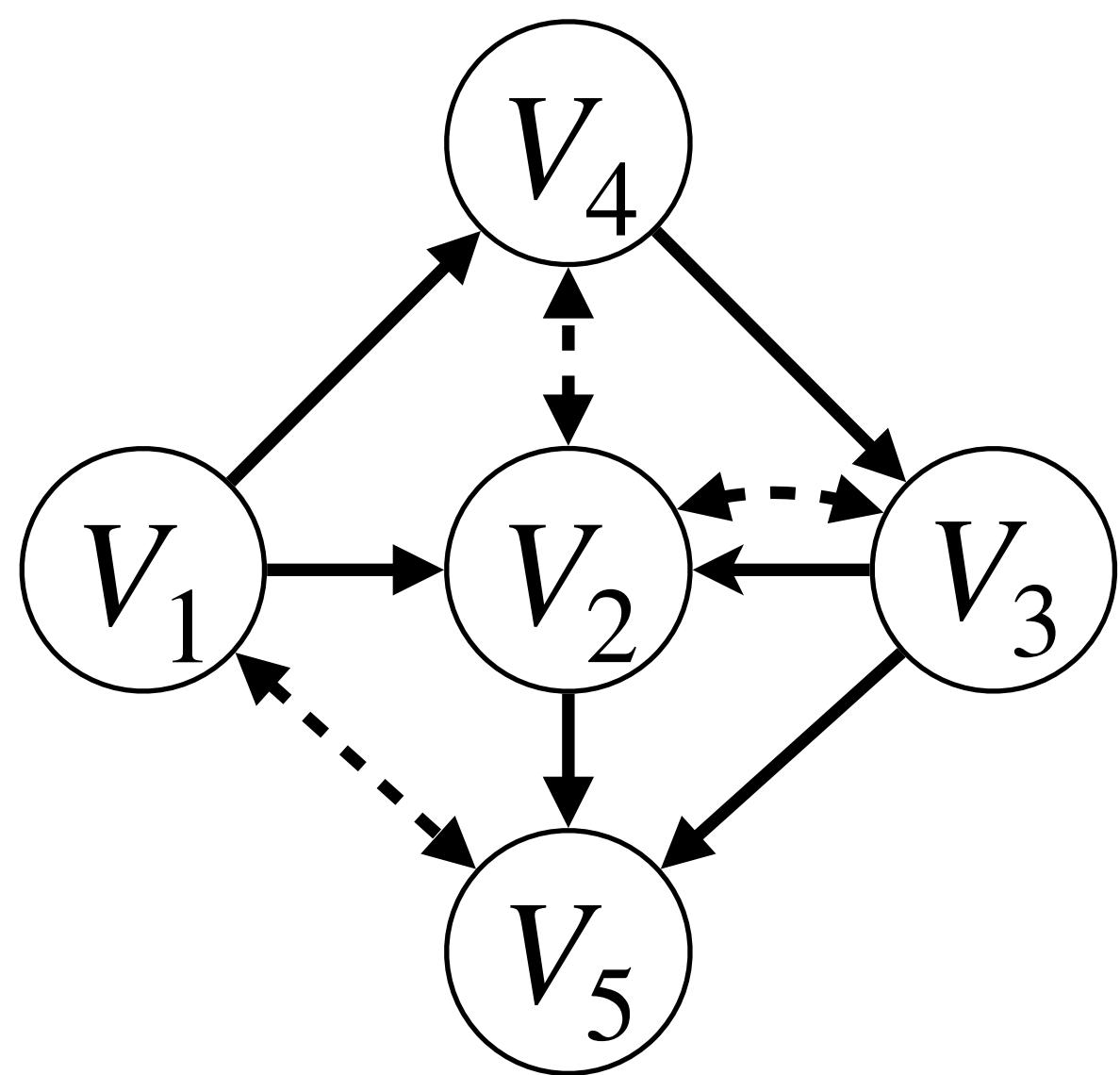


FCI - Orienting the Colliders

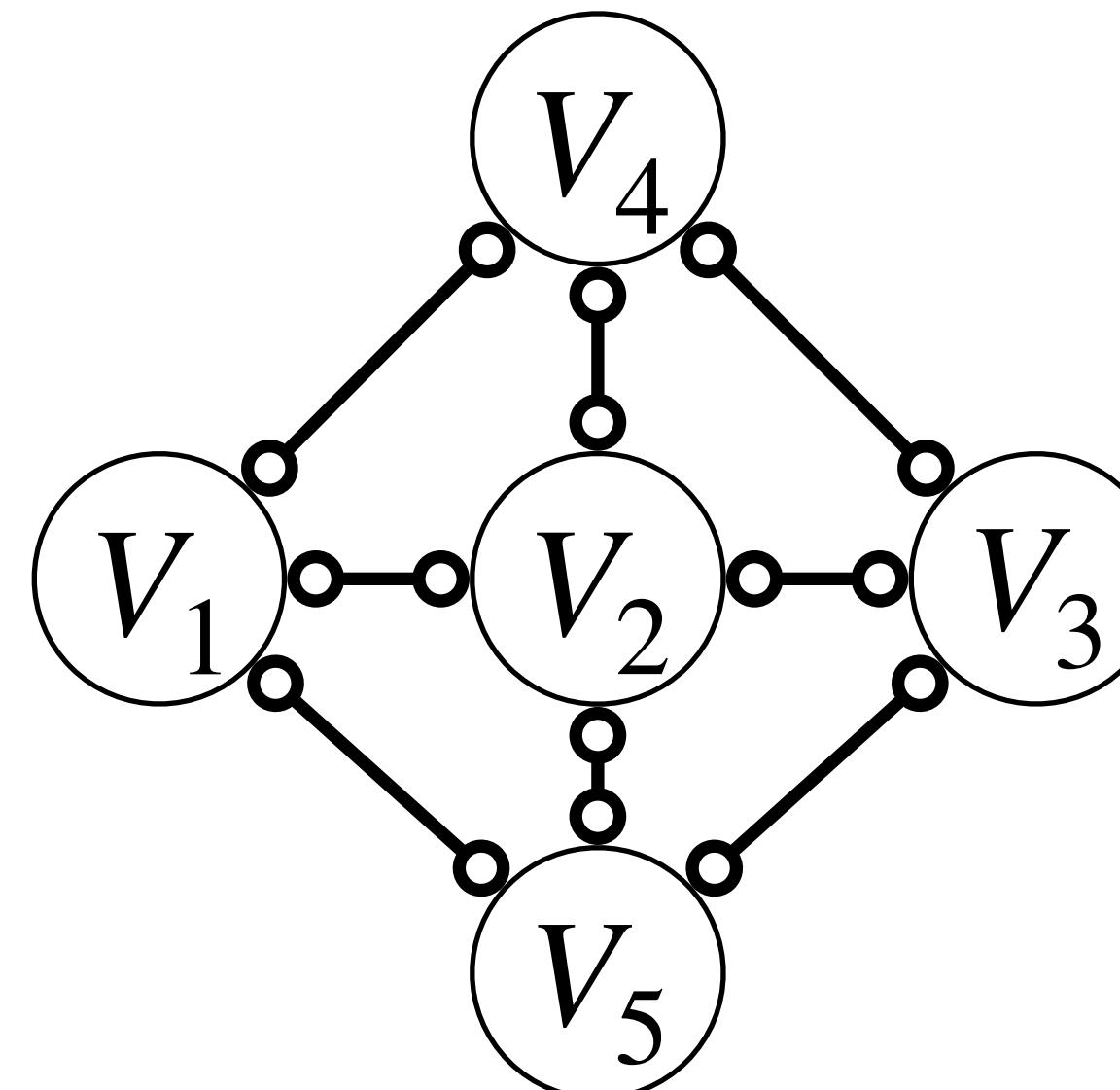
R0: If $\langle V_1, V_2, V_3 \rangle$ is unshielded and $V_2 \notin \text{Sepset}(V_1, V_3)$, then



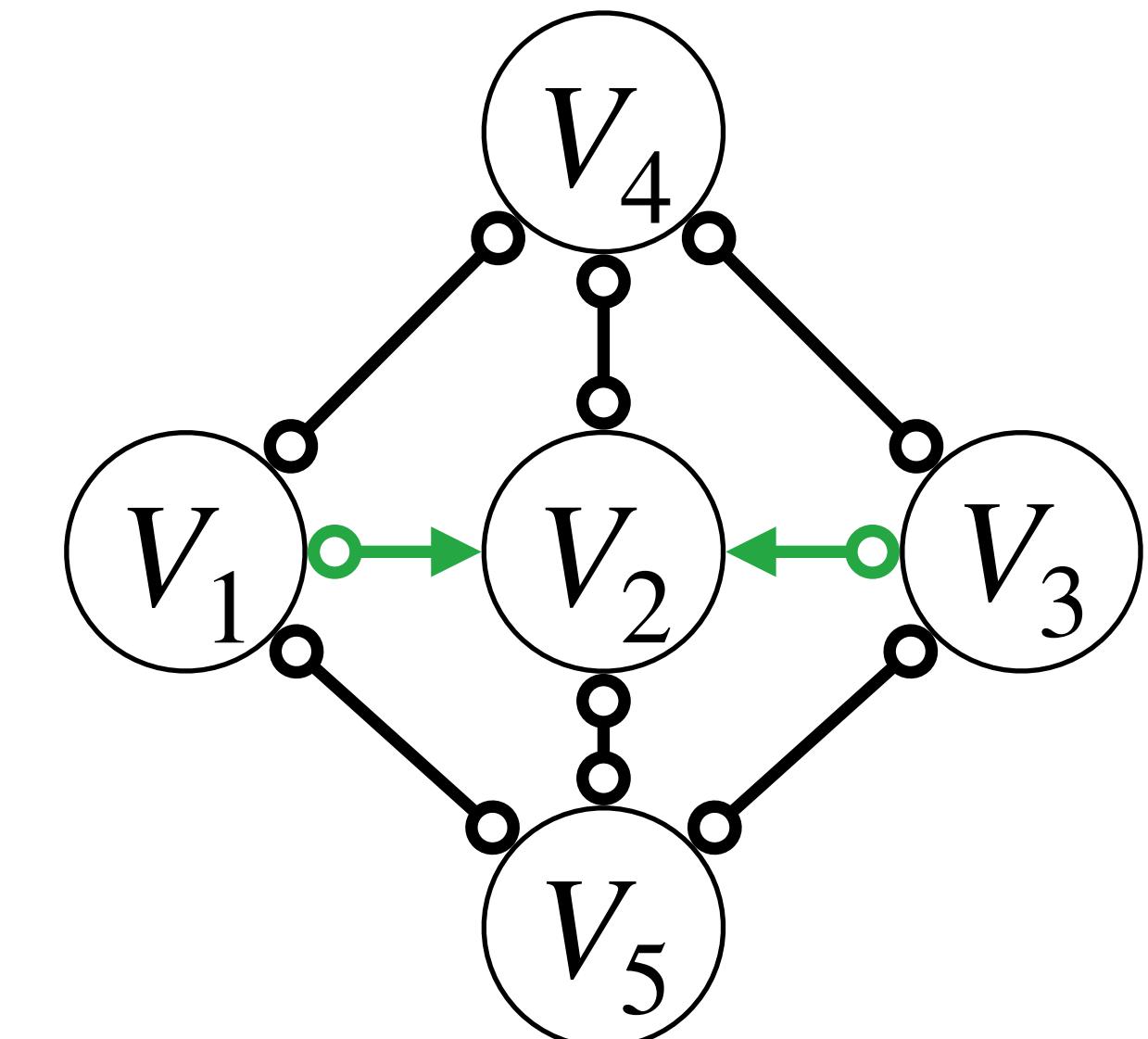
That is the only way for the path between V_1 and V_3 to be blocked when not conditioning on V_2



True, unknown ADMG

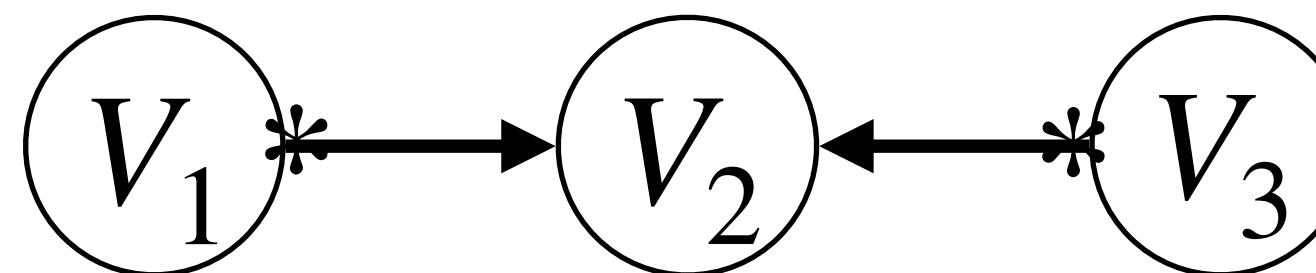


$V_1 \perp\!\!\!\perp V_3 | V_4$ and $V_1 \not\perp\!\!\!\perp V_3 | V_4, V_2$

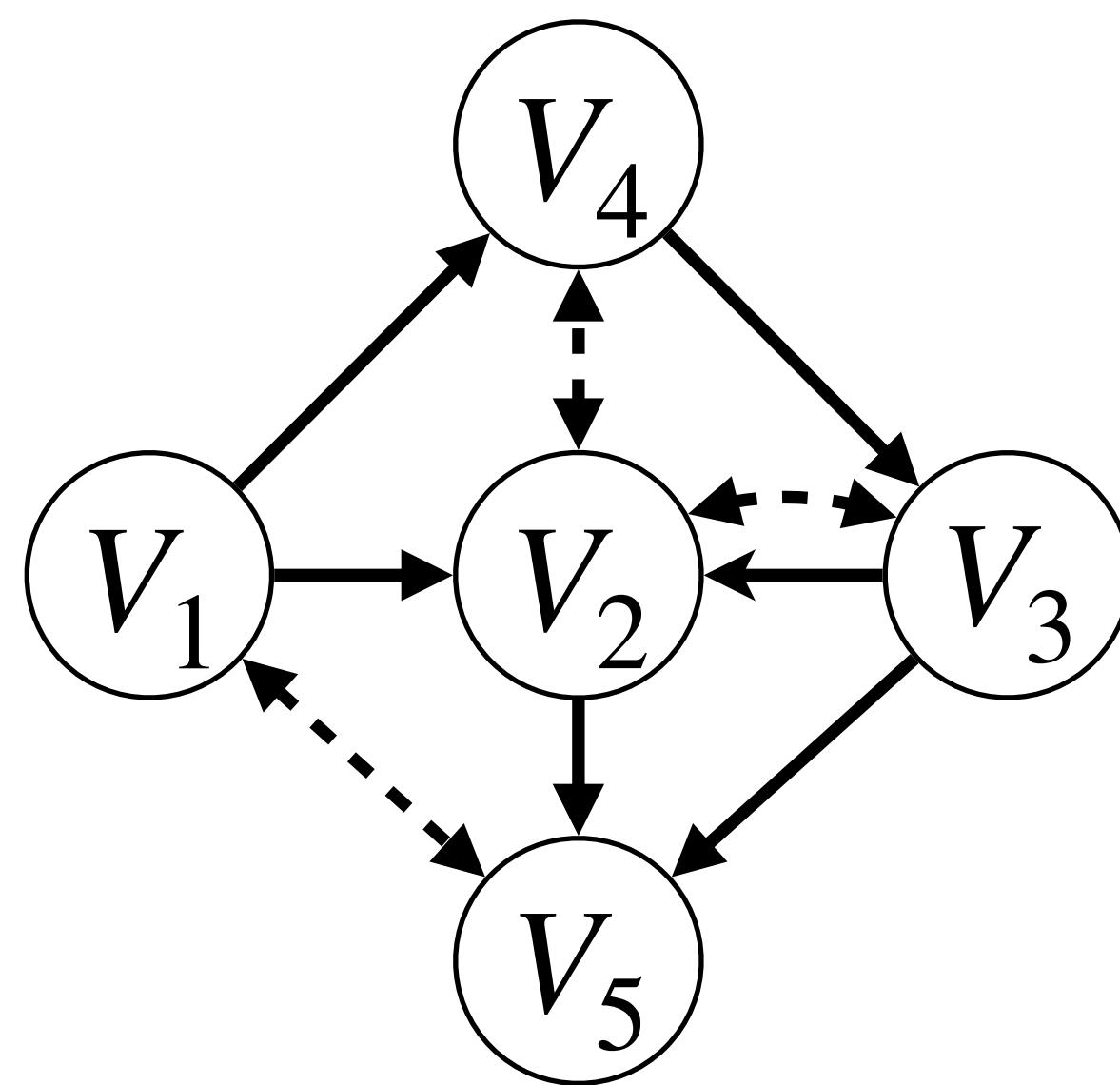


FCI - Orienting the Colliders

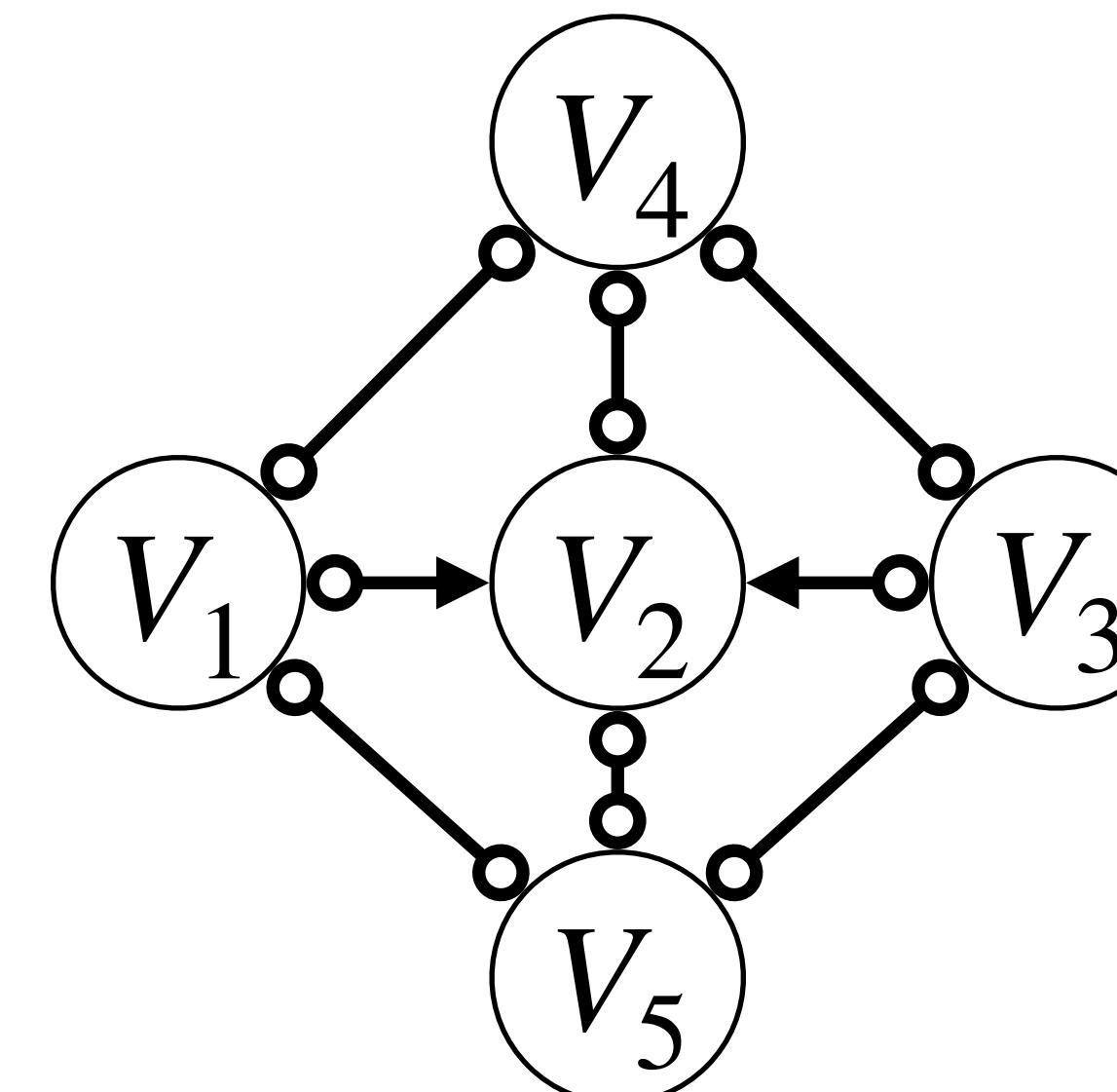
R0: If $\langle V_1, V_2, V_3 \rangle$ is unshielded and $V_2 \notin \text{Sepset}(V_1, V_3)$, then



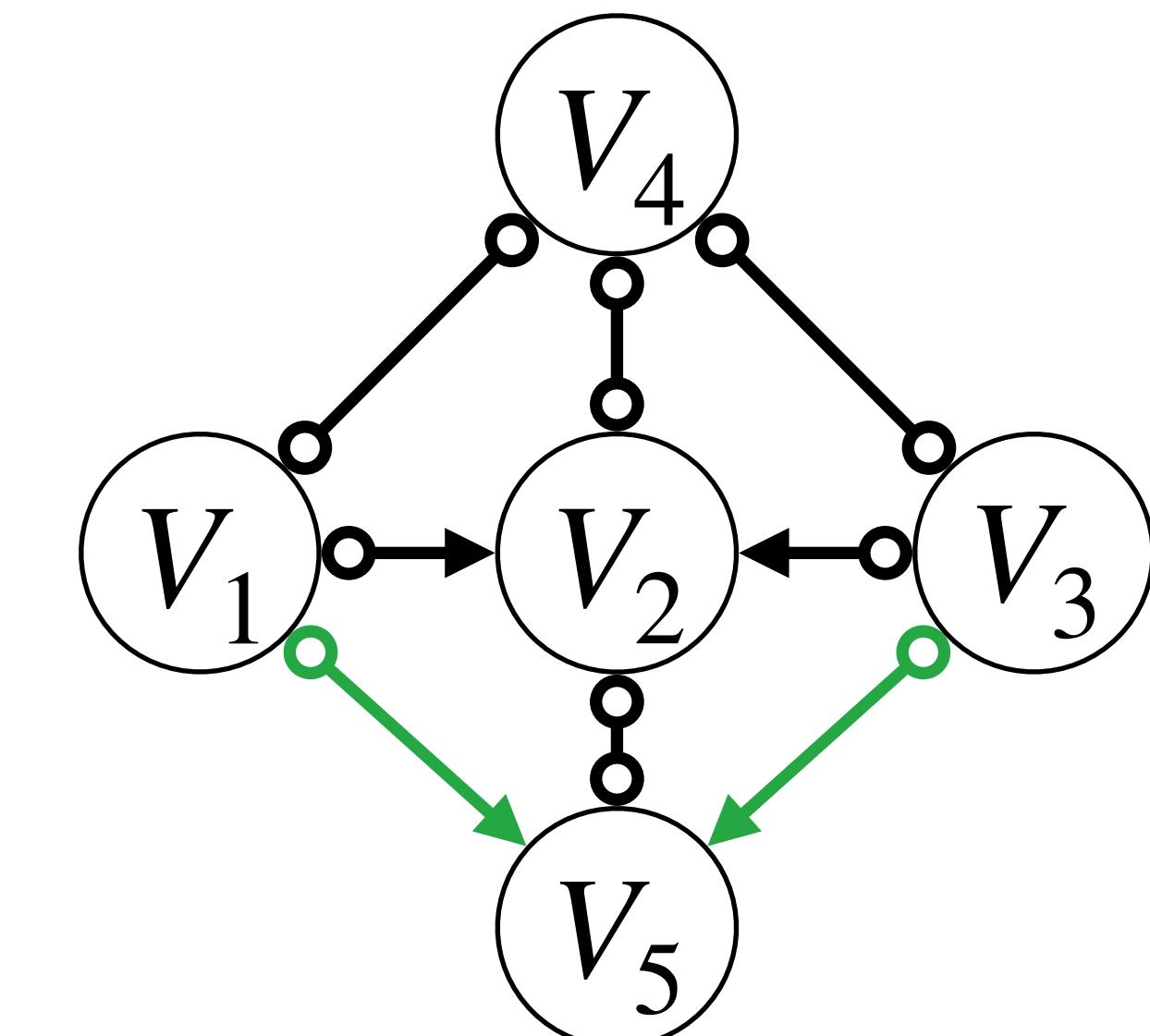
We apply R0 until no more collider can be oriented!



True, unknown ADMG



$V_1 \perp\!\!\!\perp V_3 | V_4$ and $V_1 \not\perp\!\!\!\perp V_3 | V_4, V_5$



Rules R1 to R10

Next, we apply Rules 1 to 10, in any order, whenever any of them applies.

Theorem by Ali et al, 2005 [Arrowhead Completeness]:

- R1 to R4, together with R0, are complete for arrowhead (non-ancestrality) orientation.

Theorem by Zhang J., 2008 [Tail Completeness]:

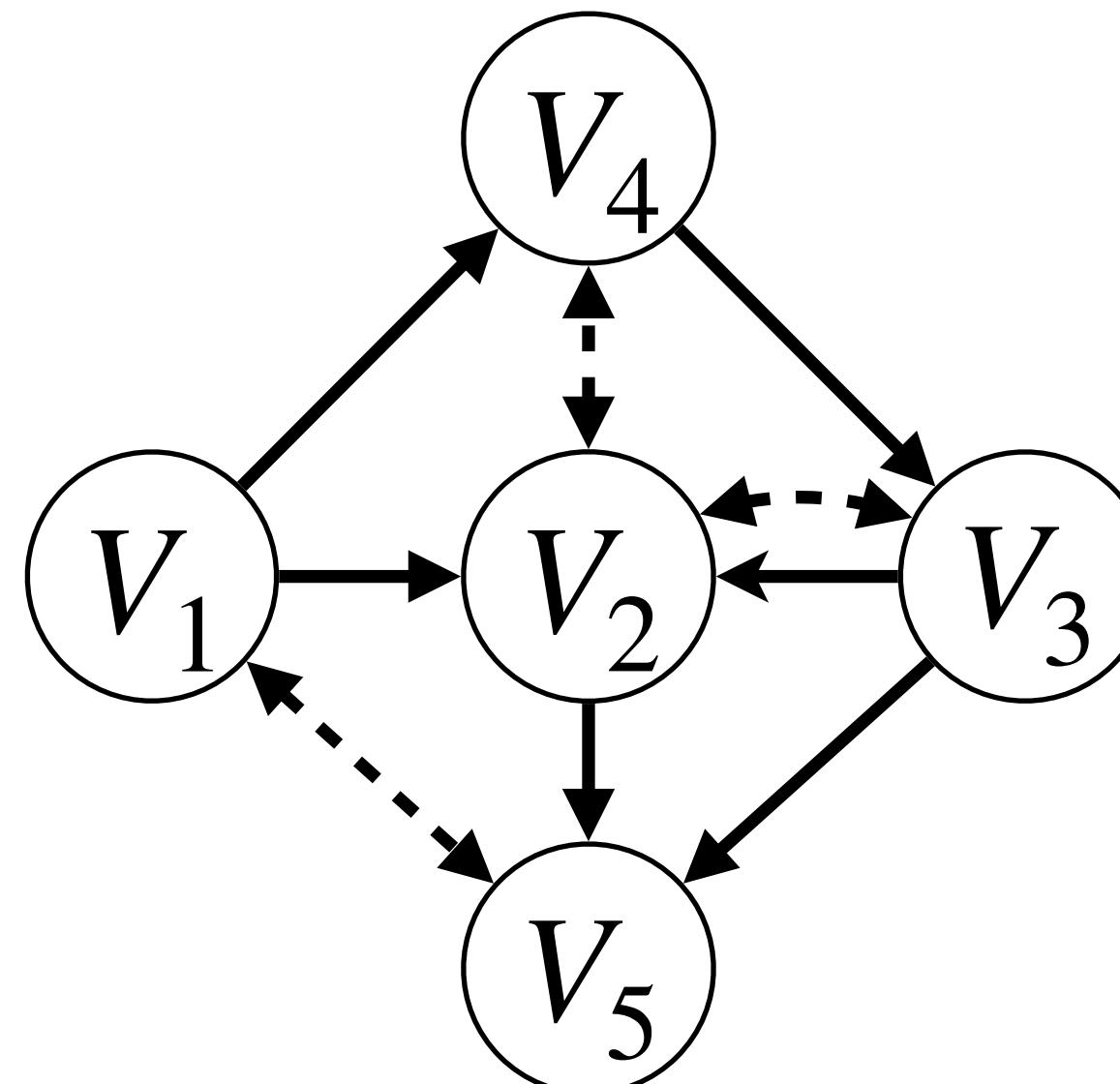
- R5 to R10, together with R0 to R4, are complete for arrowhead (non-ancestrality) and tail (ancestrality) orientation.
- If no selection variables are present, then rules R5 to R7 are not necessary.

Spirites, P., Glymour, C. N., Scheines, R., & Heckerman, D. (2000). *Causation, prediction, and search*. MIT press.

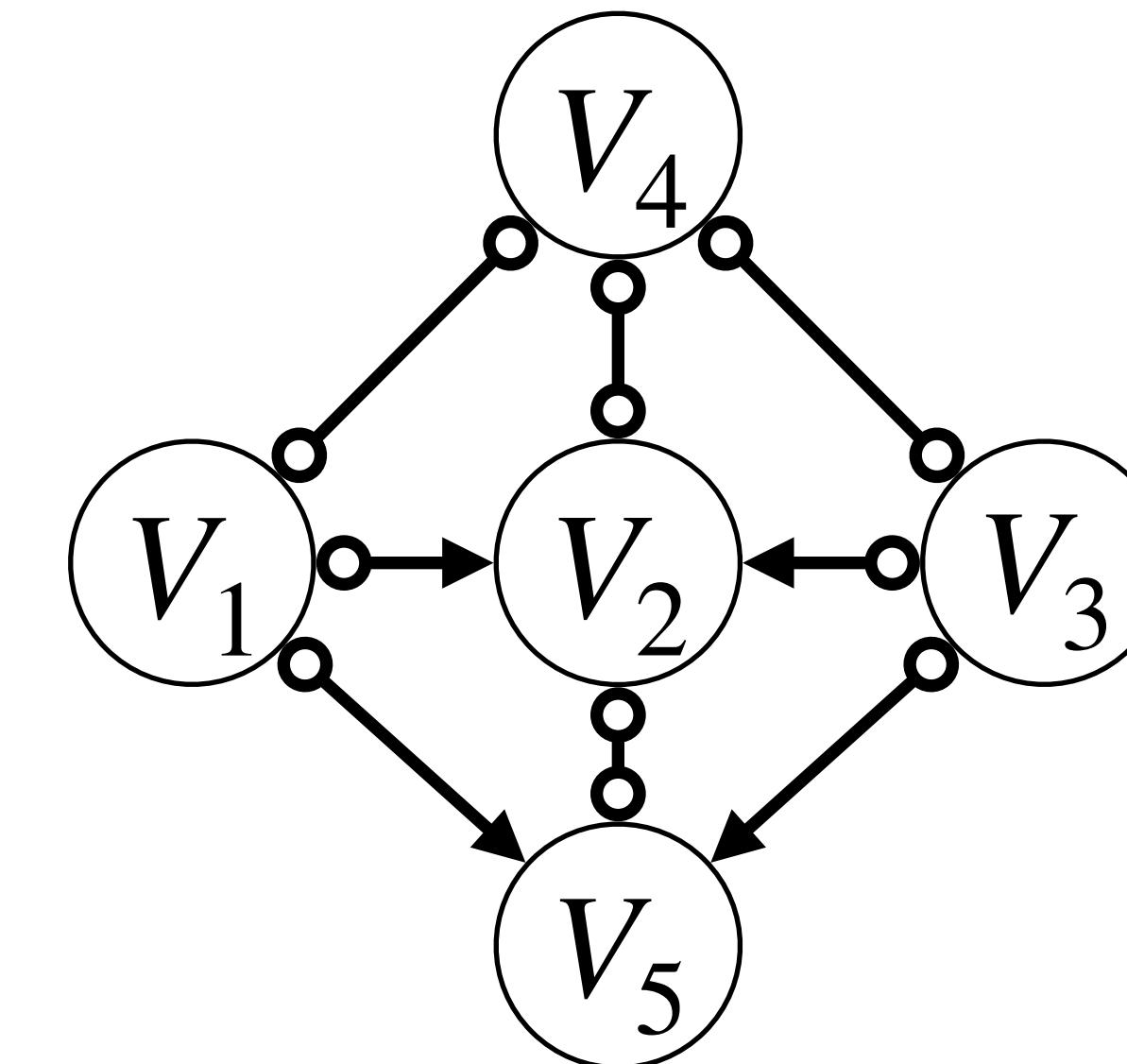
R.A. Ali, T. Richardson, P. Spirtes, J. Zhang, Towards characterizing Markov equivalence classes for directed acyclic graphs with latent variables. In Proceedings of the 21th Conference on Uncertainty in Artificial Intelligence, AUAI Press, 2005, pp. 10–17. ([Link](#))

Zhang, J., 2008. On the completeness of orientation rules for causal discovery in the presence of latent confounders and selection bias. *Artificial Intelligence*, 172(16-17), pp.1873-1896.

Going back to the example...



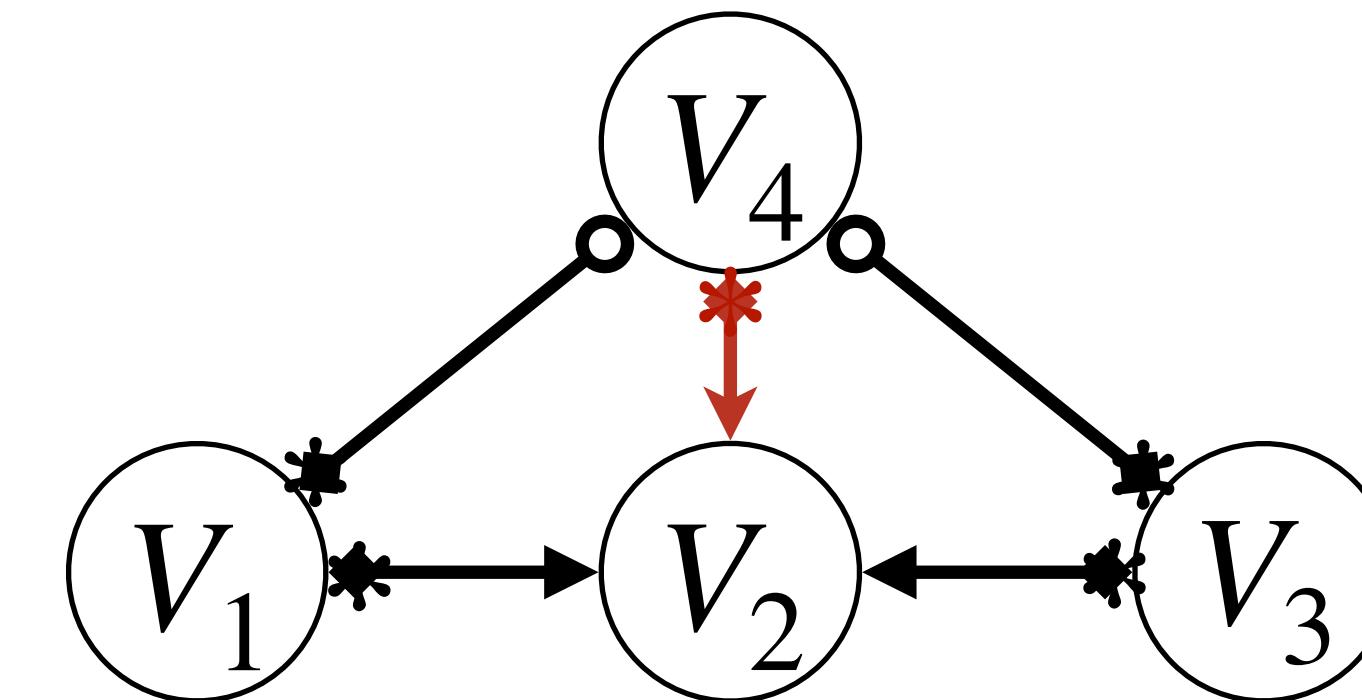
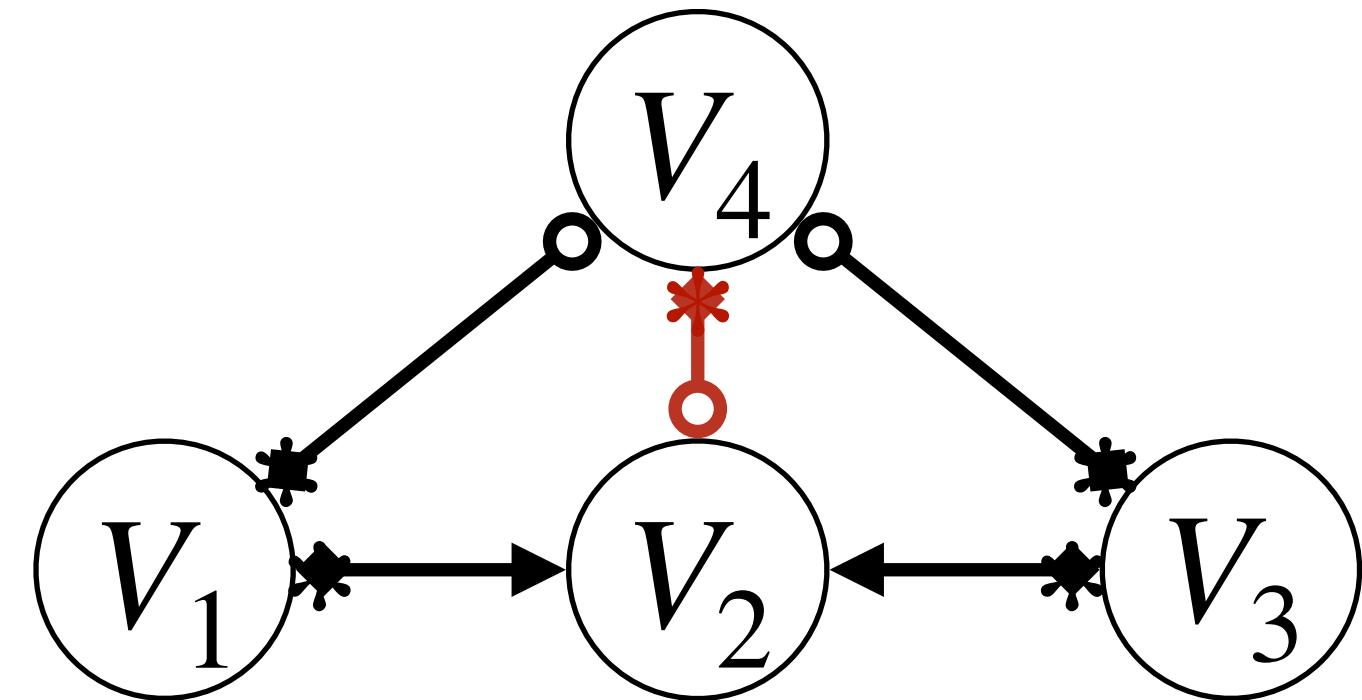
True, unknown ADMG



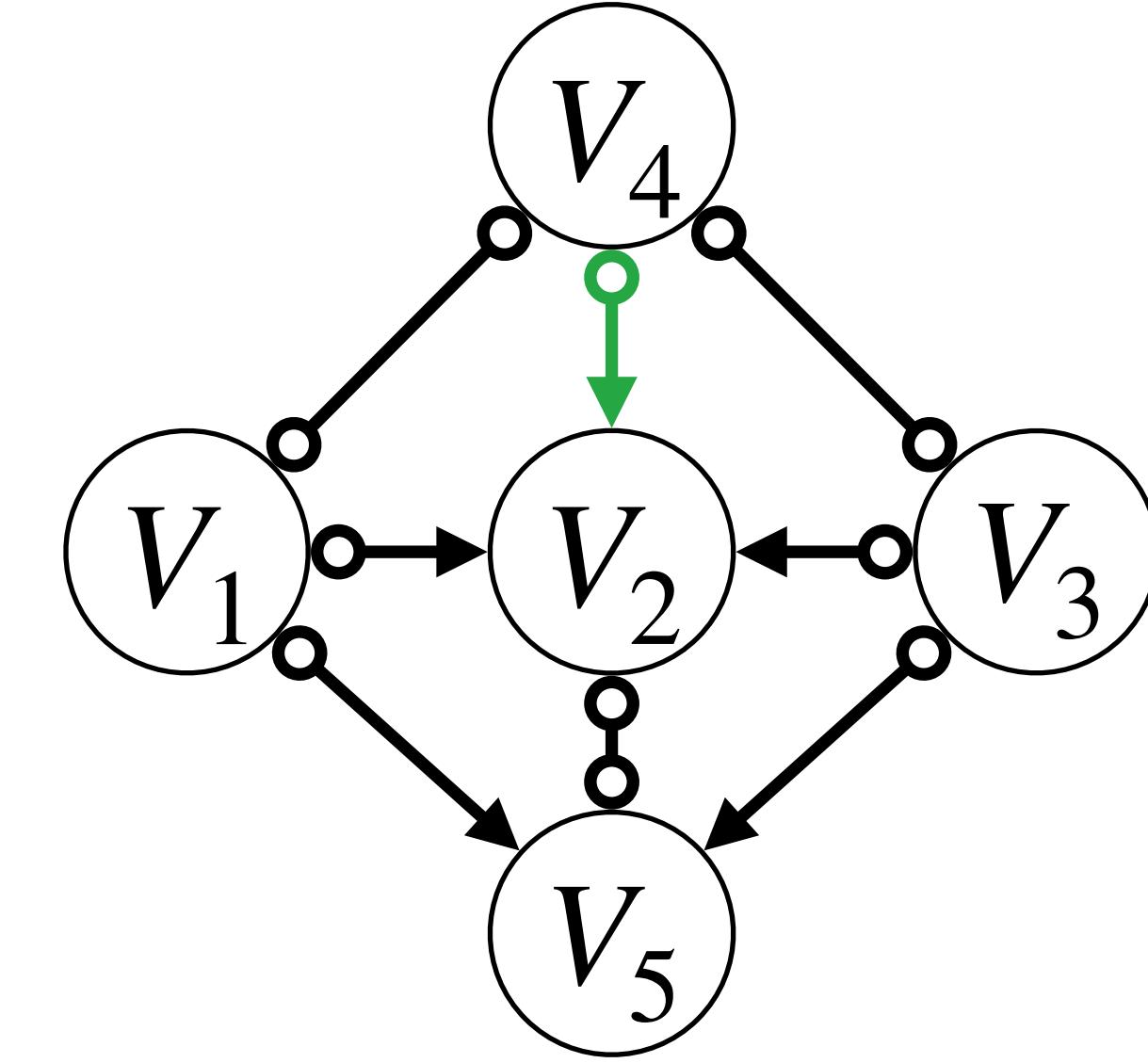
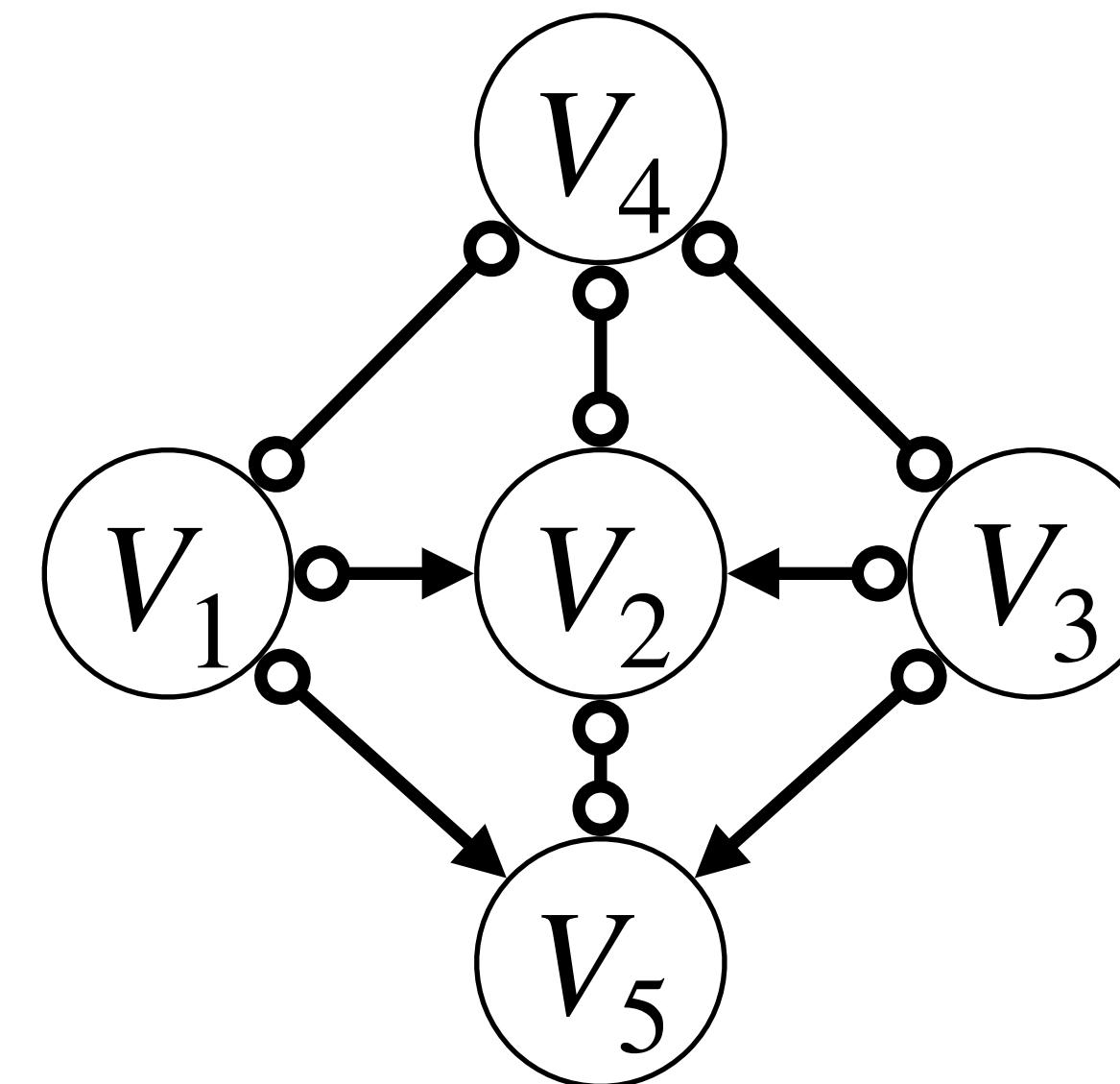
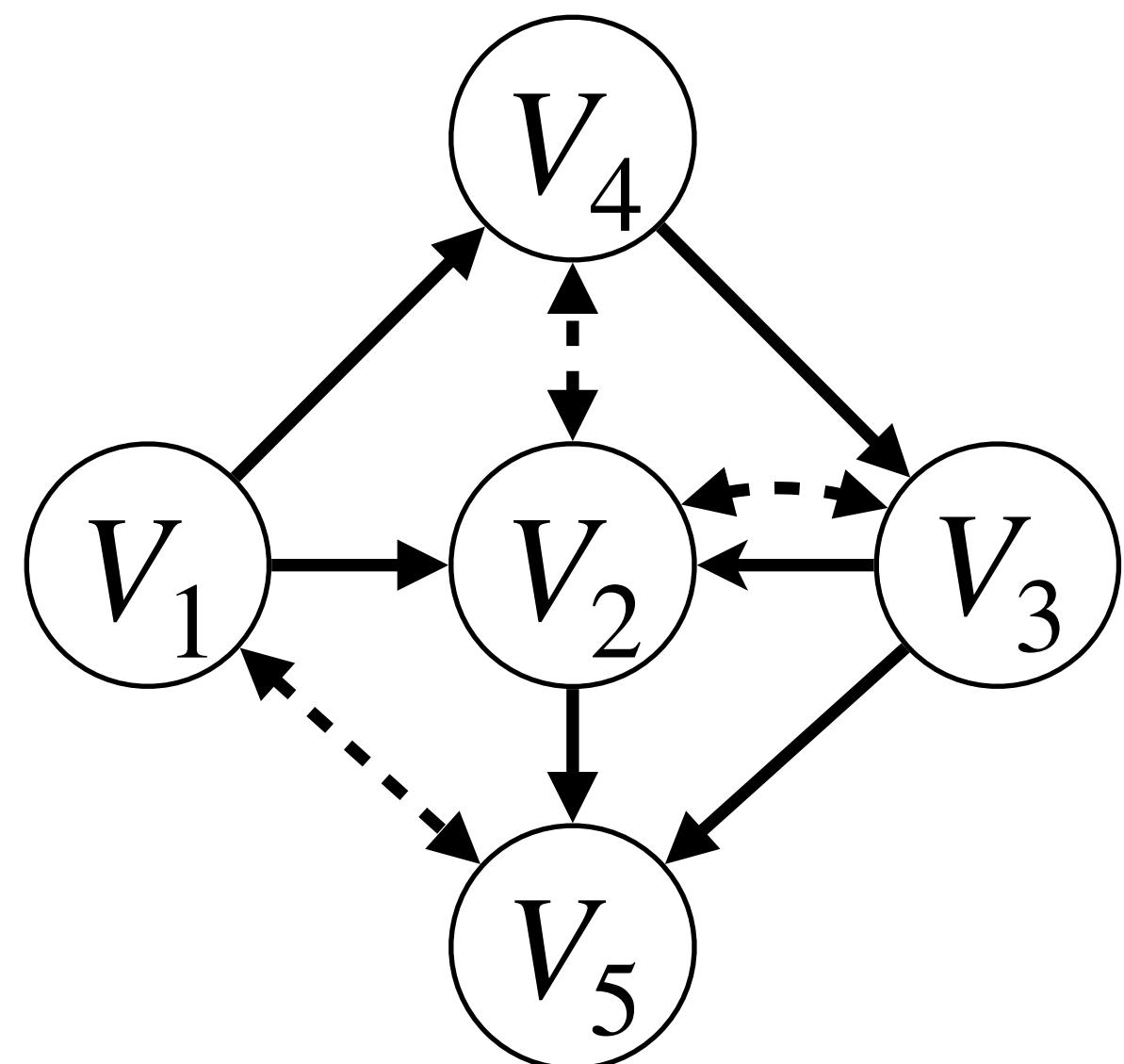
After Skeleton + R0

Applying Mark Inference Rules

R3:



where V_1 and V_3 are not adjacent



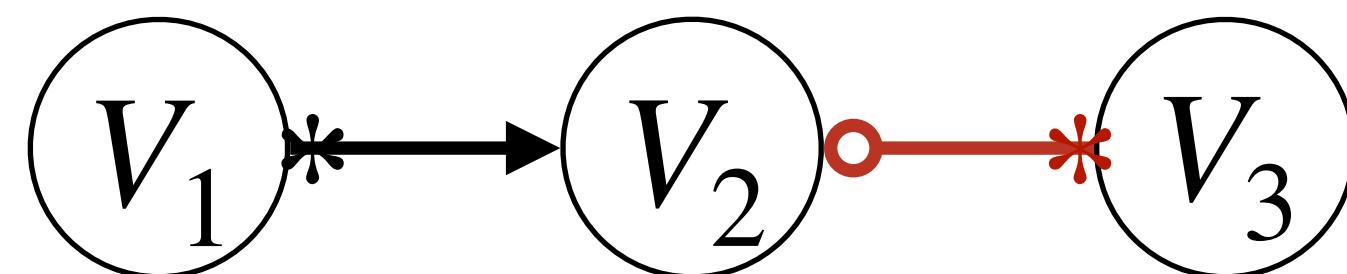
True, unknown ADMG

After Skeleton + R0

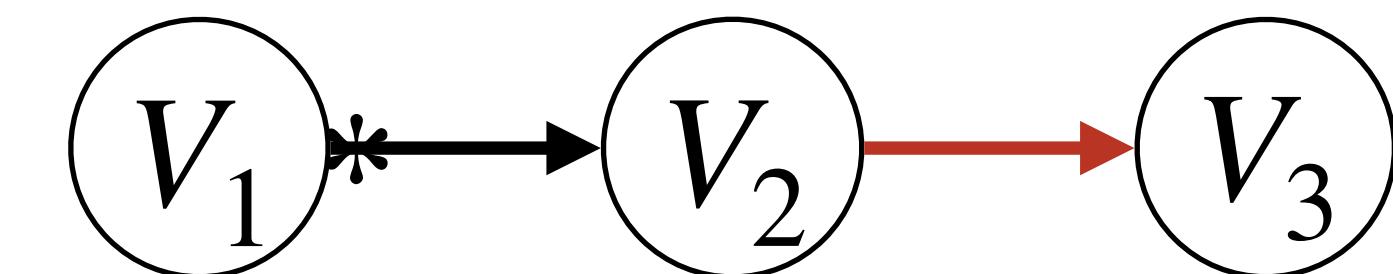
Applying R3

Applying Mark Inference Rules

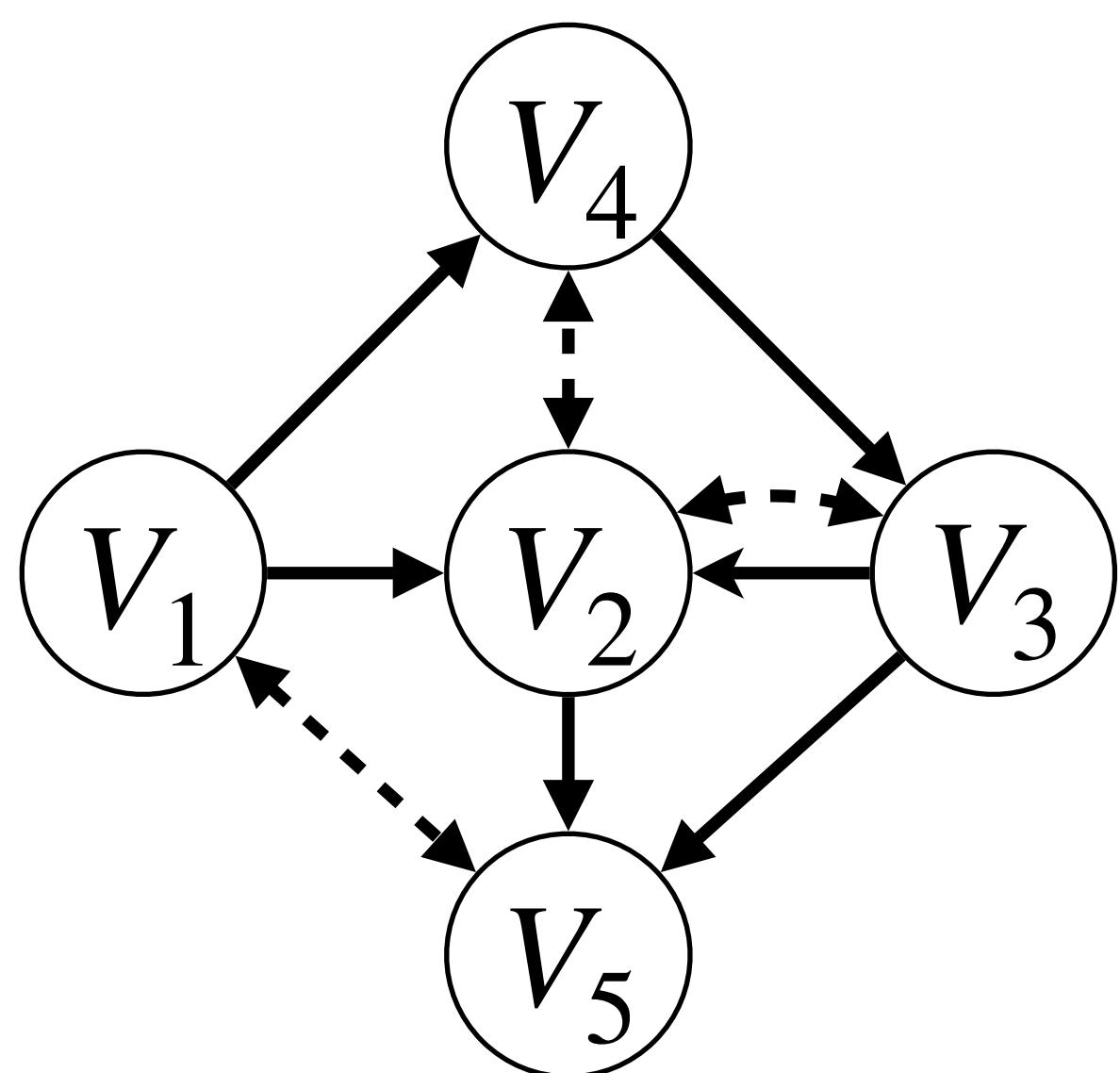
R1:



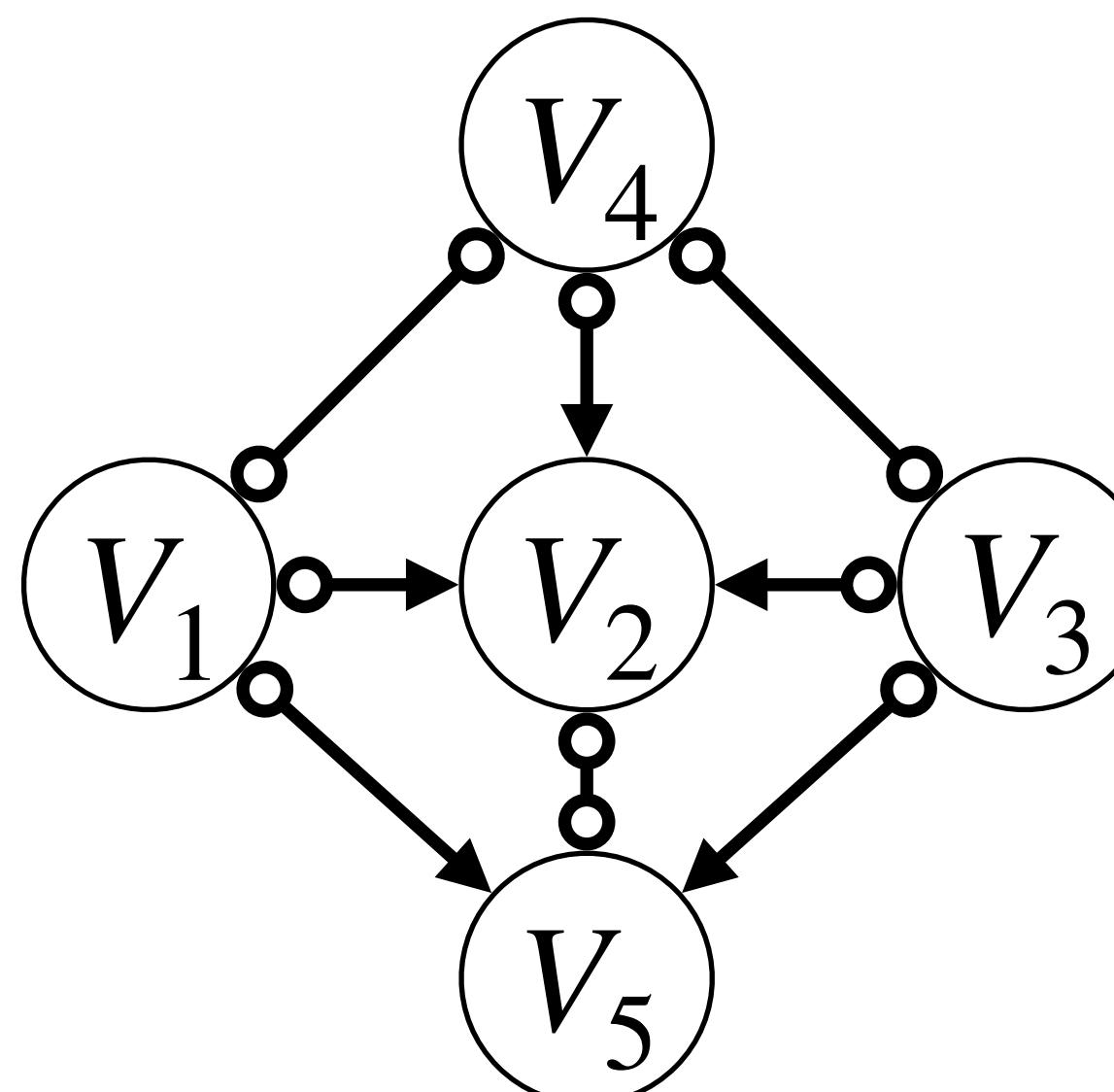
⇒



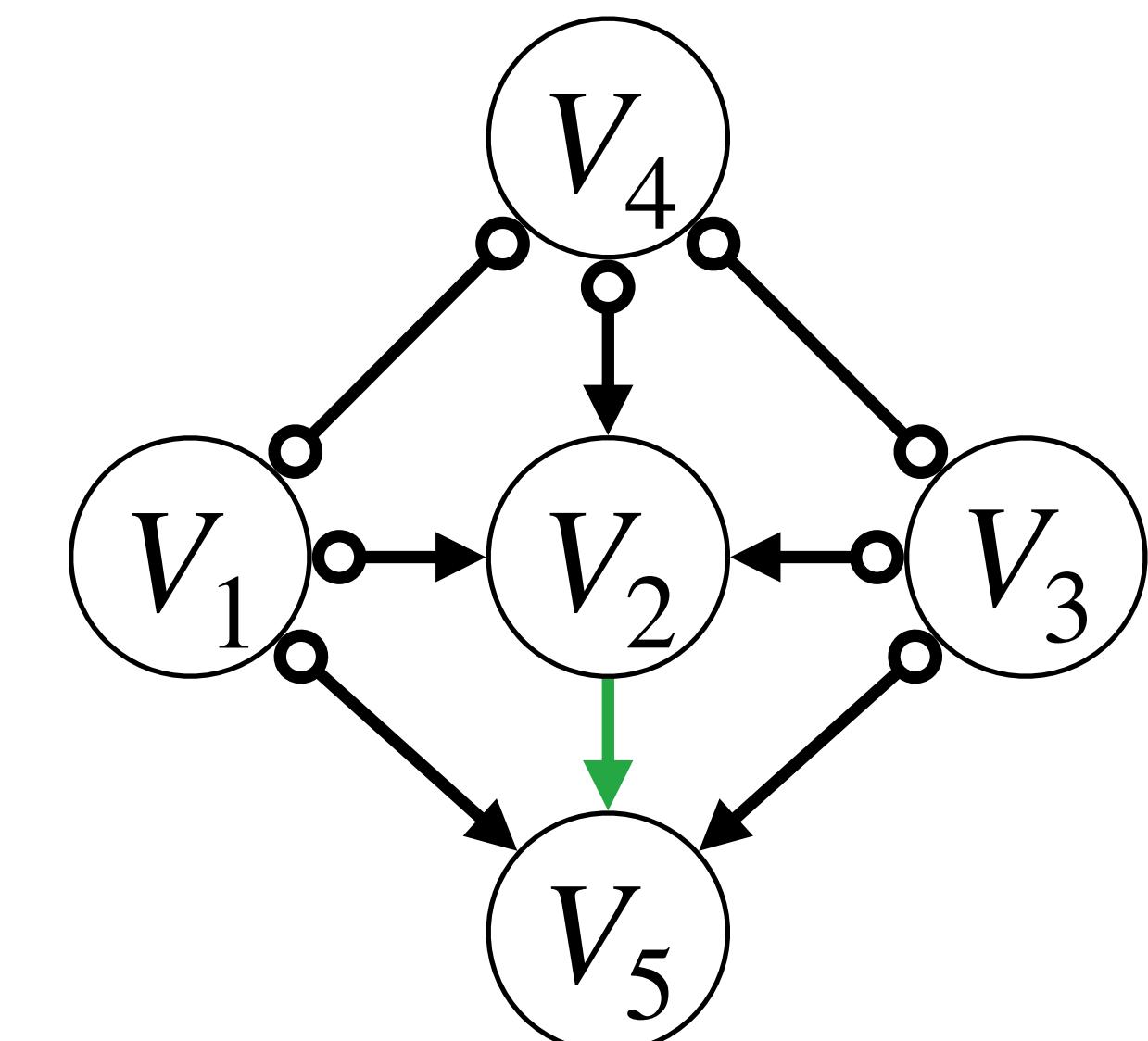
where V_1 and V_3 are not adjacent



True, unknown ADMG

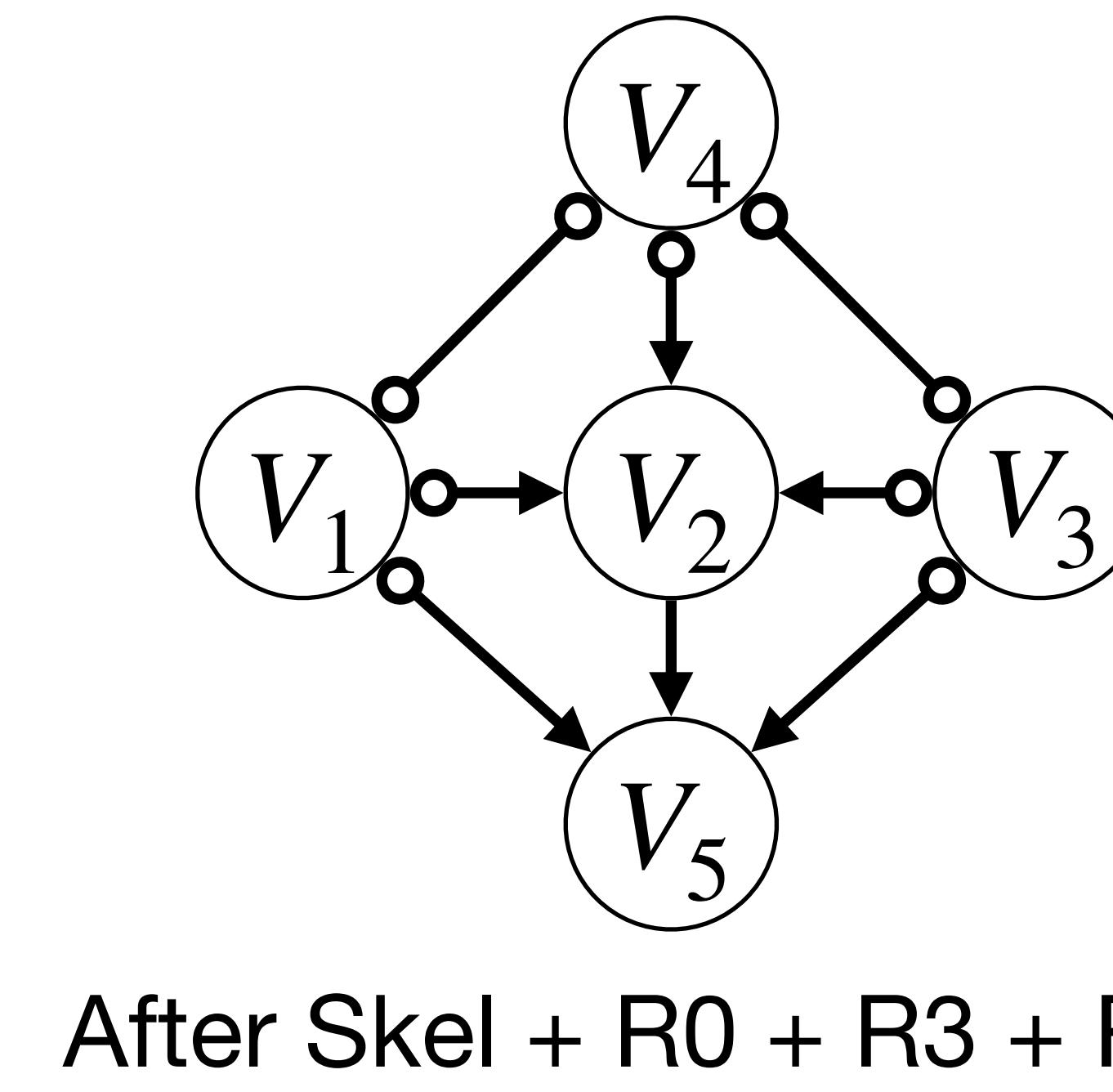
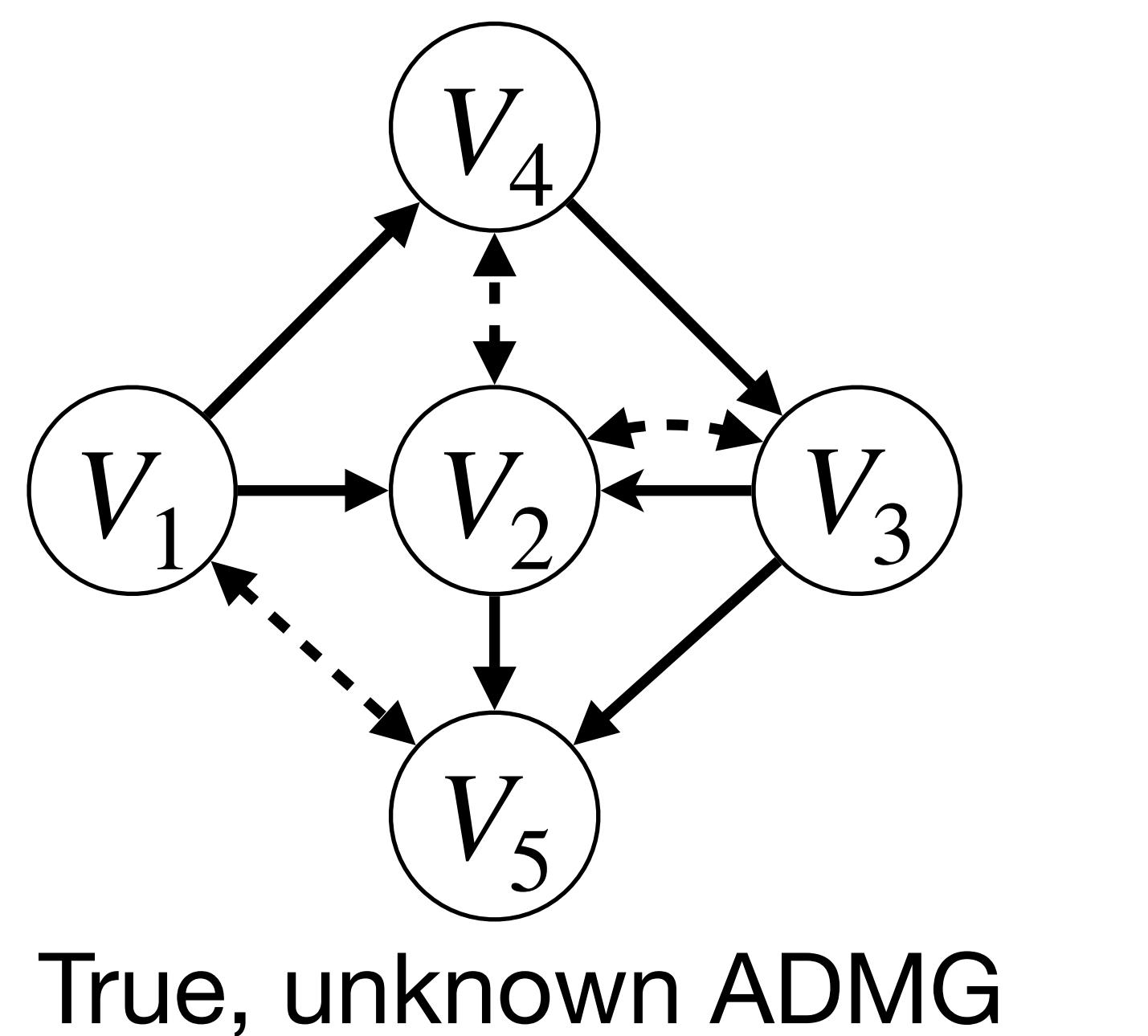
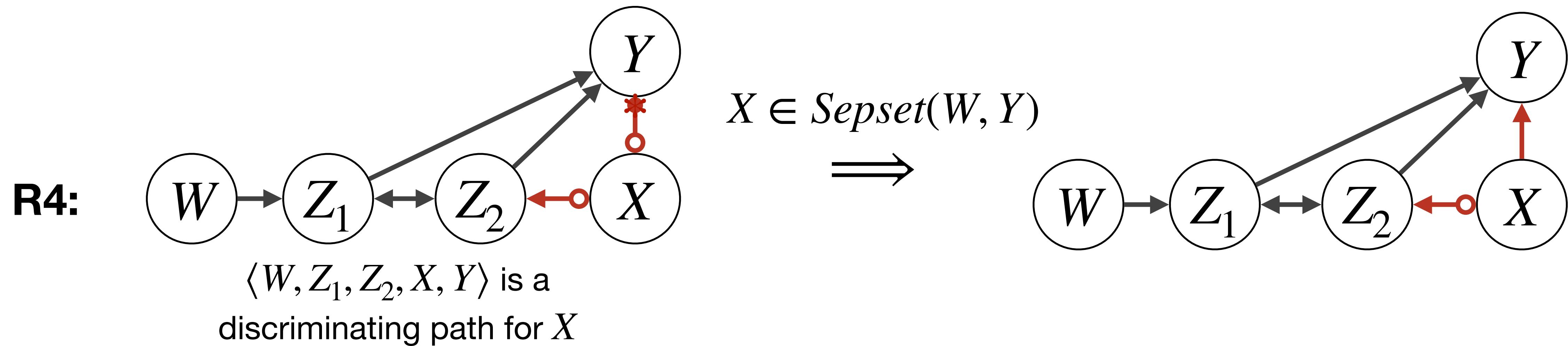


After Skel + R0 + R3

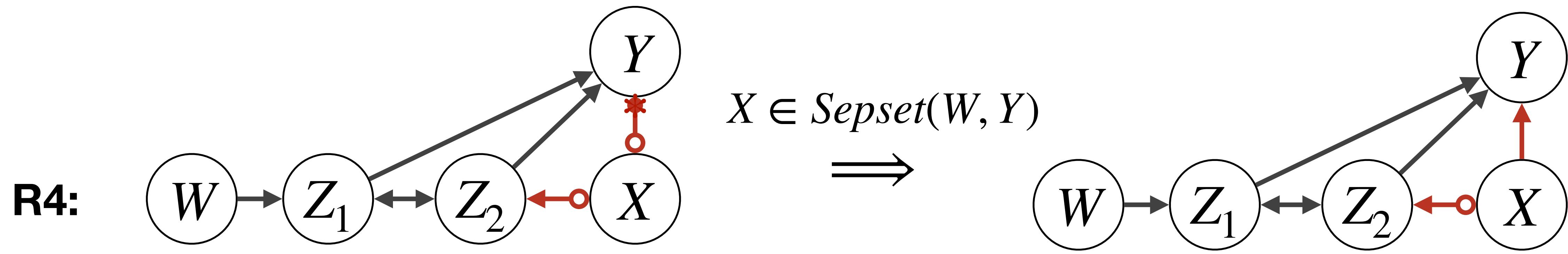


Applying R1

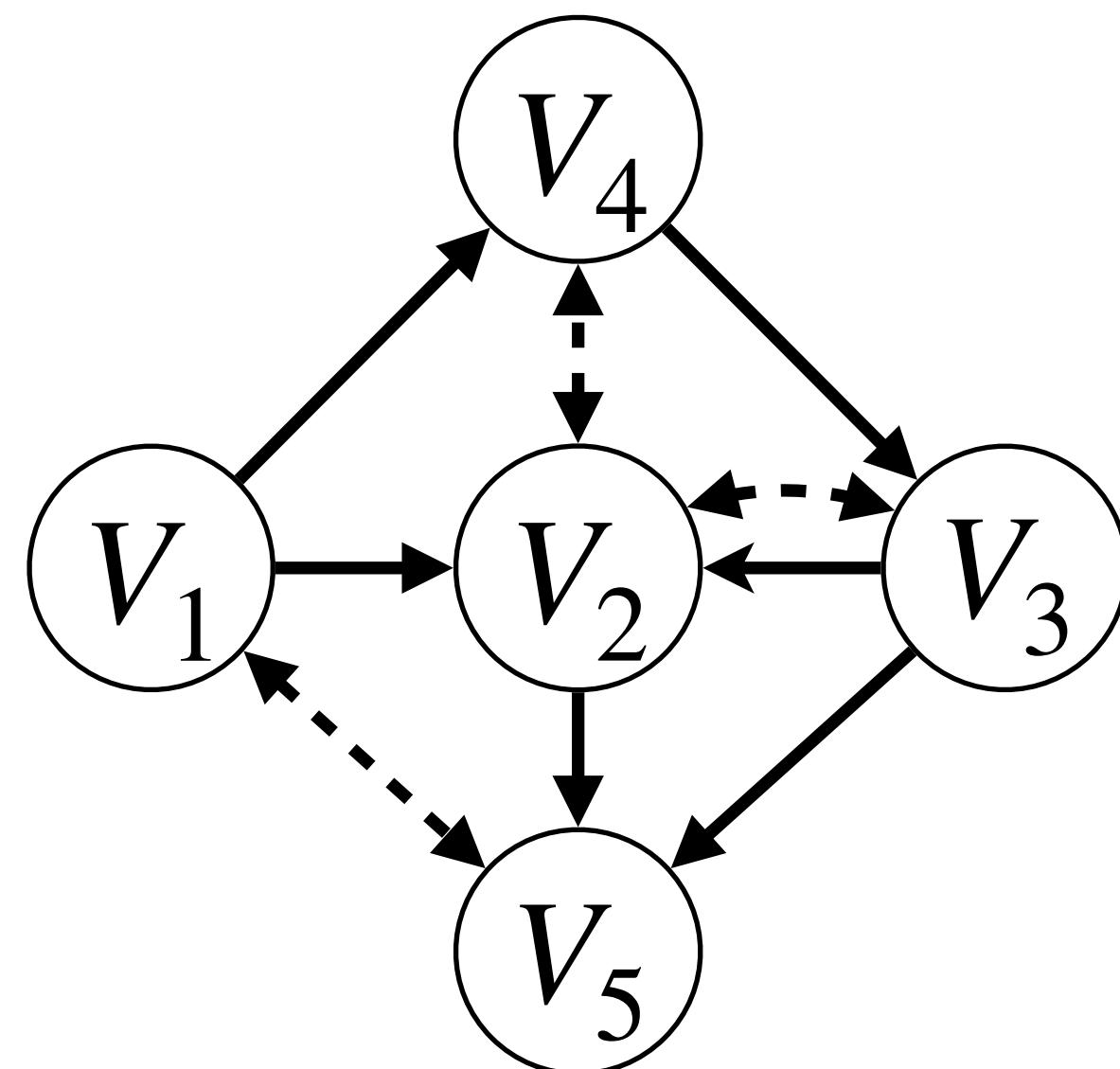
Applying Mark Inference Rules



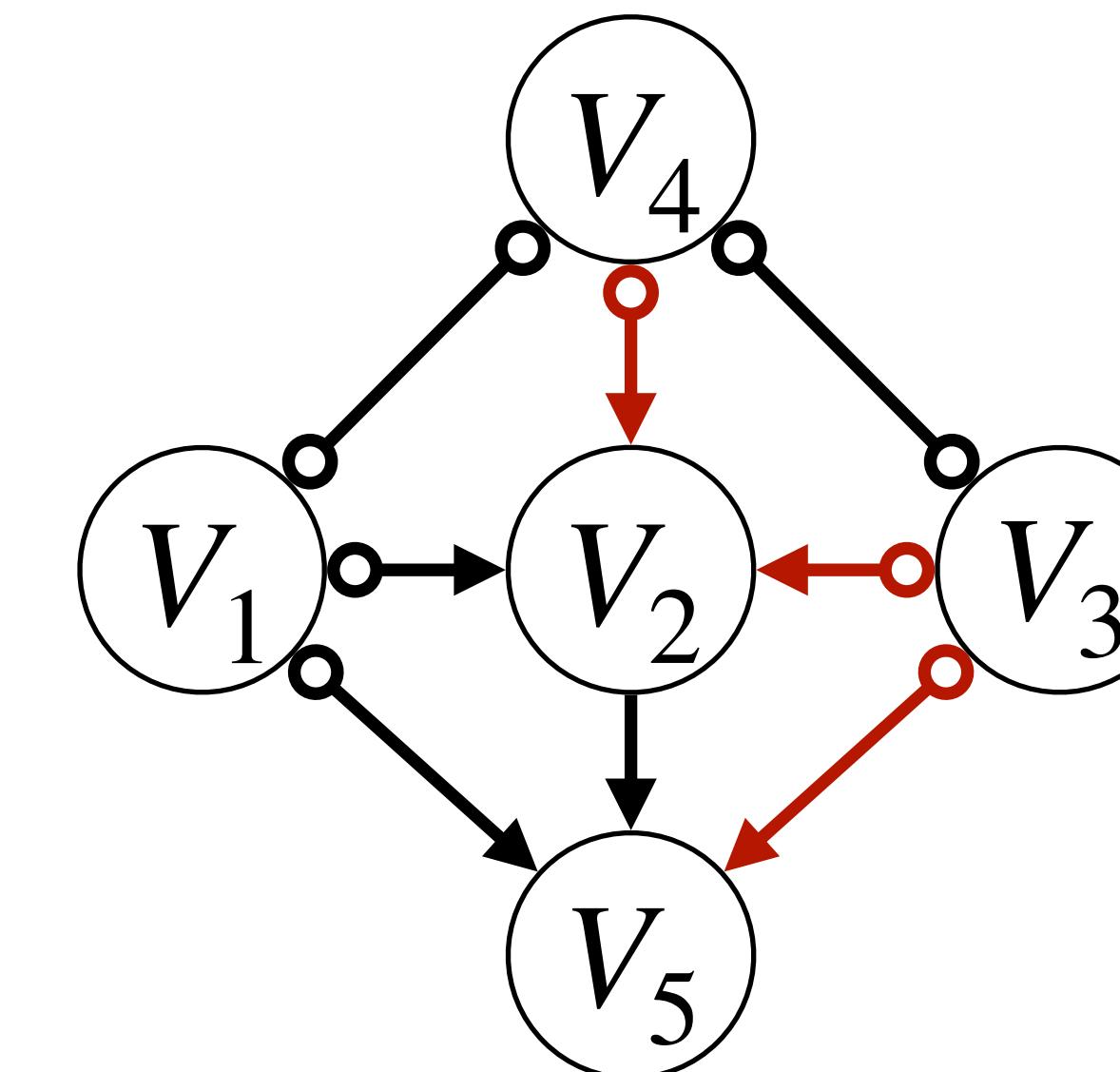
Applying Mark Inference Rules



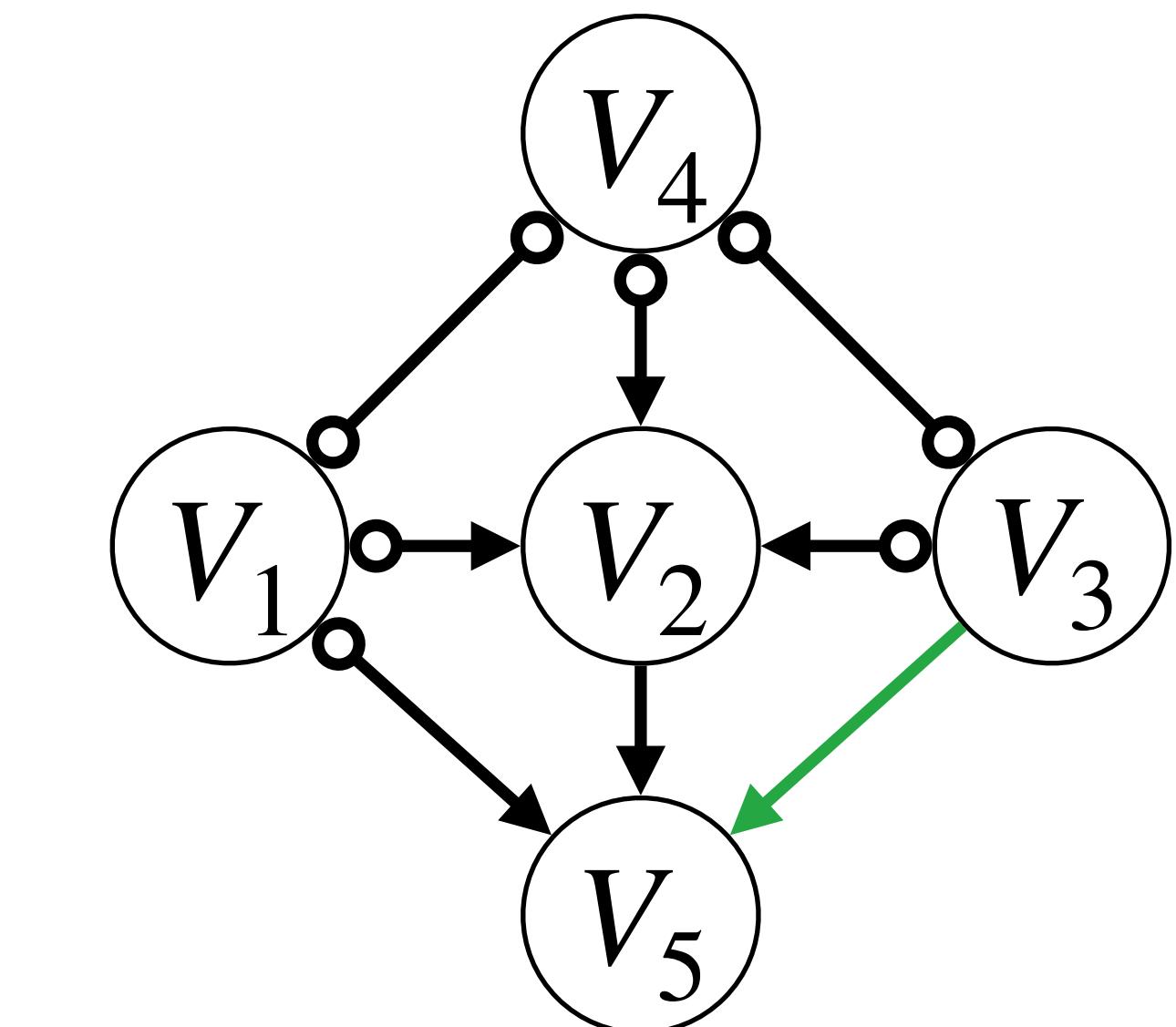
$\langle W, Z_1, Z_2, X, Y \rangle$ is a
discriminating path for X



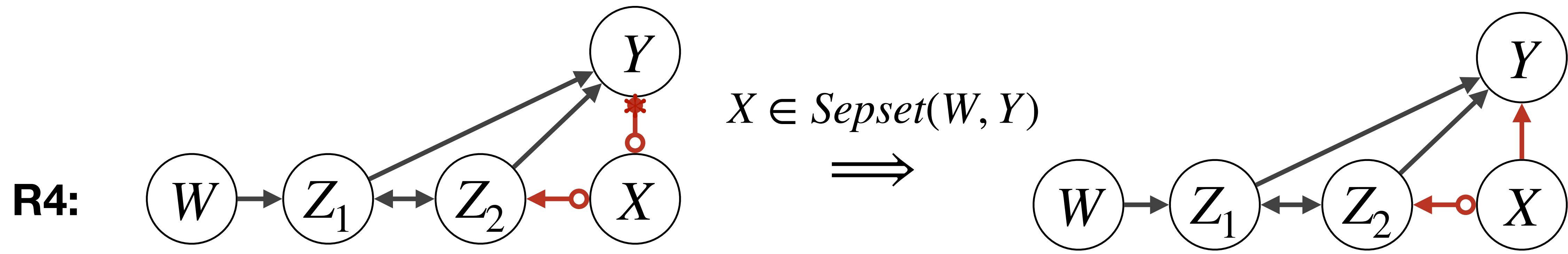
True, unknown ADMG



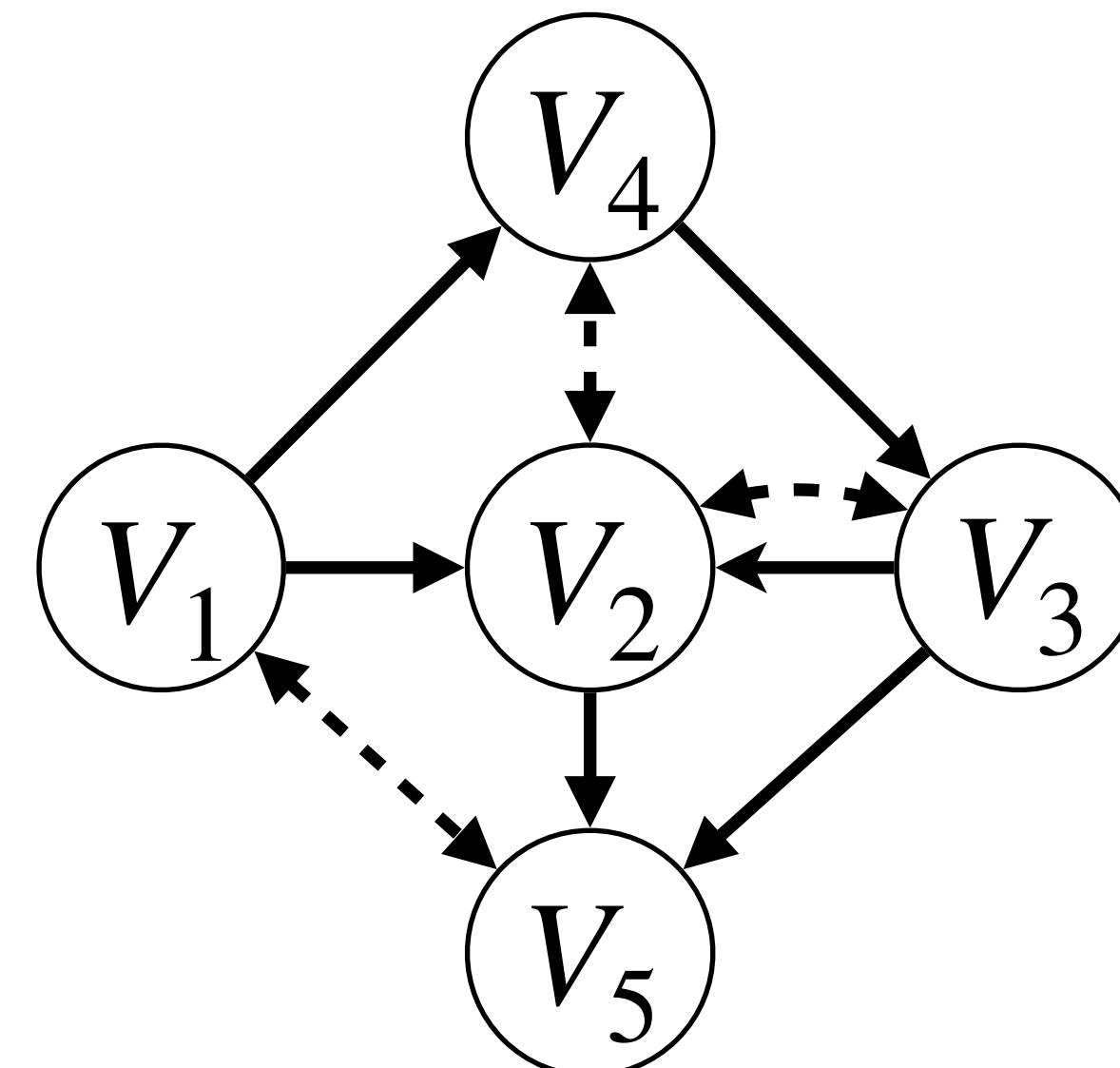
$\langle V_4, V_2, V_3, V_5 \rangle$ is a discriminating path for V_3 and
 $V_3 \in \text{Sepset}(V_4, V_5) - V_4 \perp\!\!\!\perp V_5 | V_1, V_2, V_3$



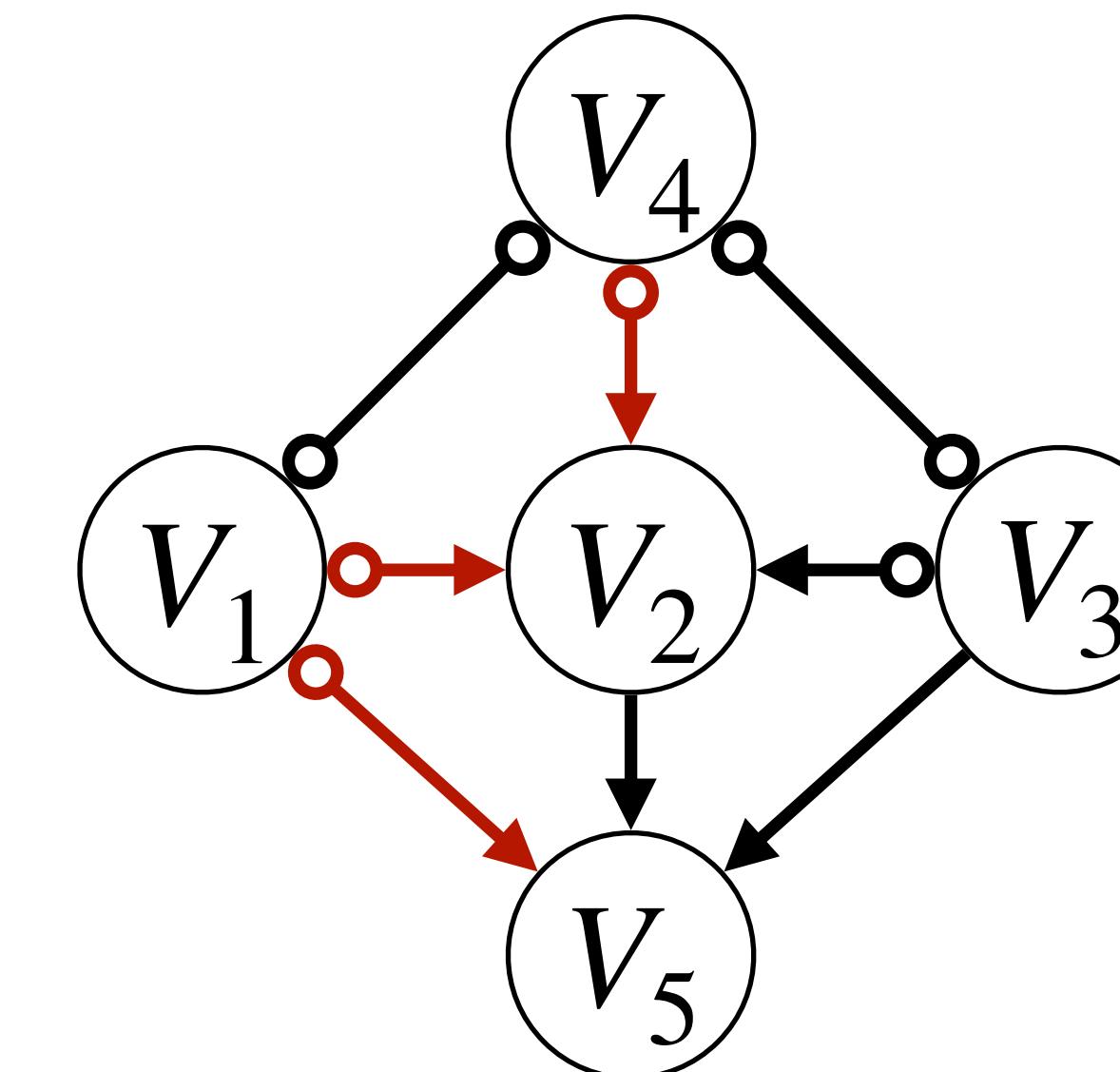
Applying Mark Inference Rules



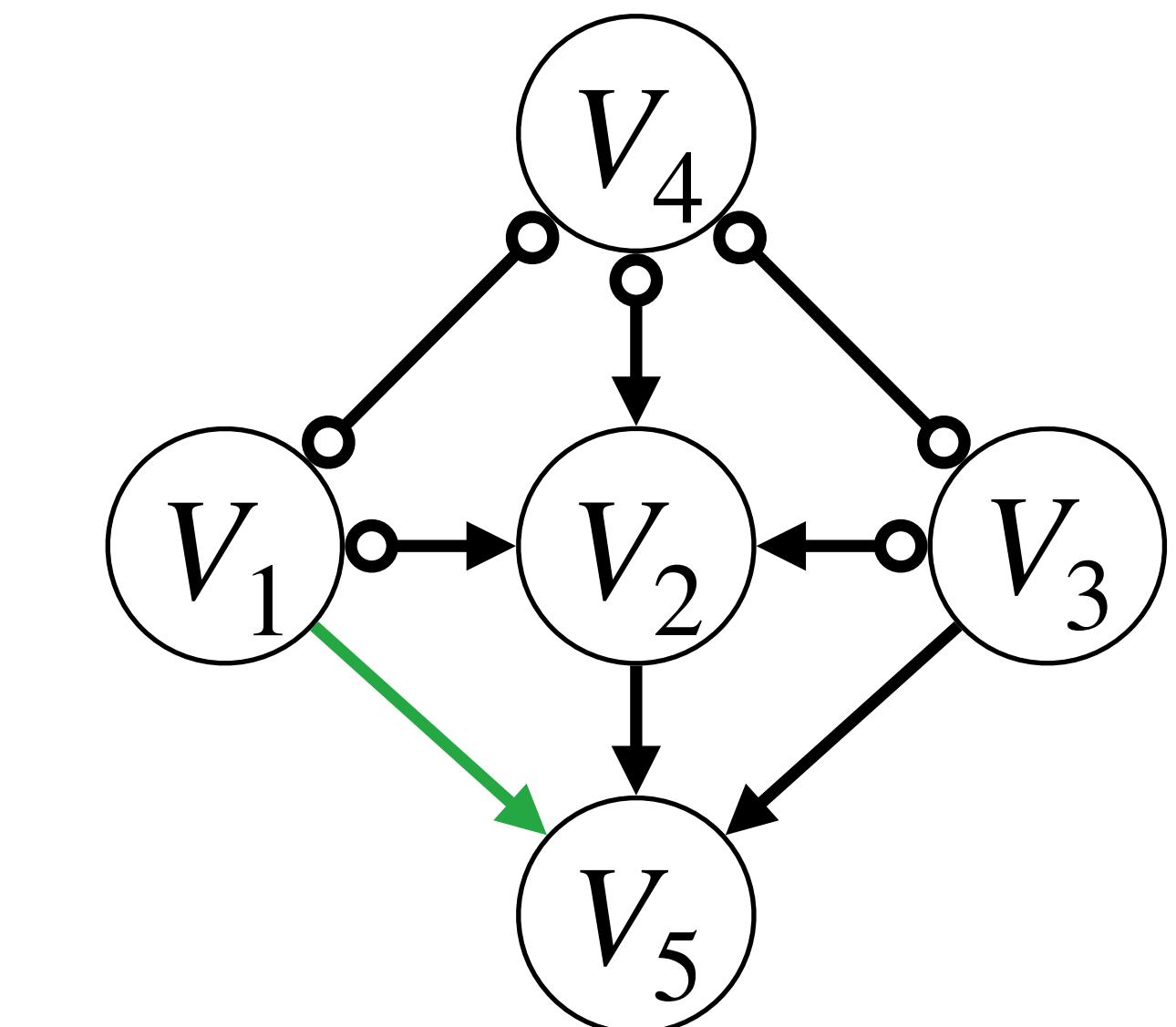
$\langle W, Z_1, Z_2, X, Y \rangle$ is a
discriminating path for X



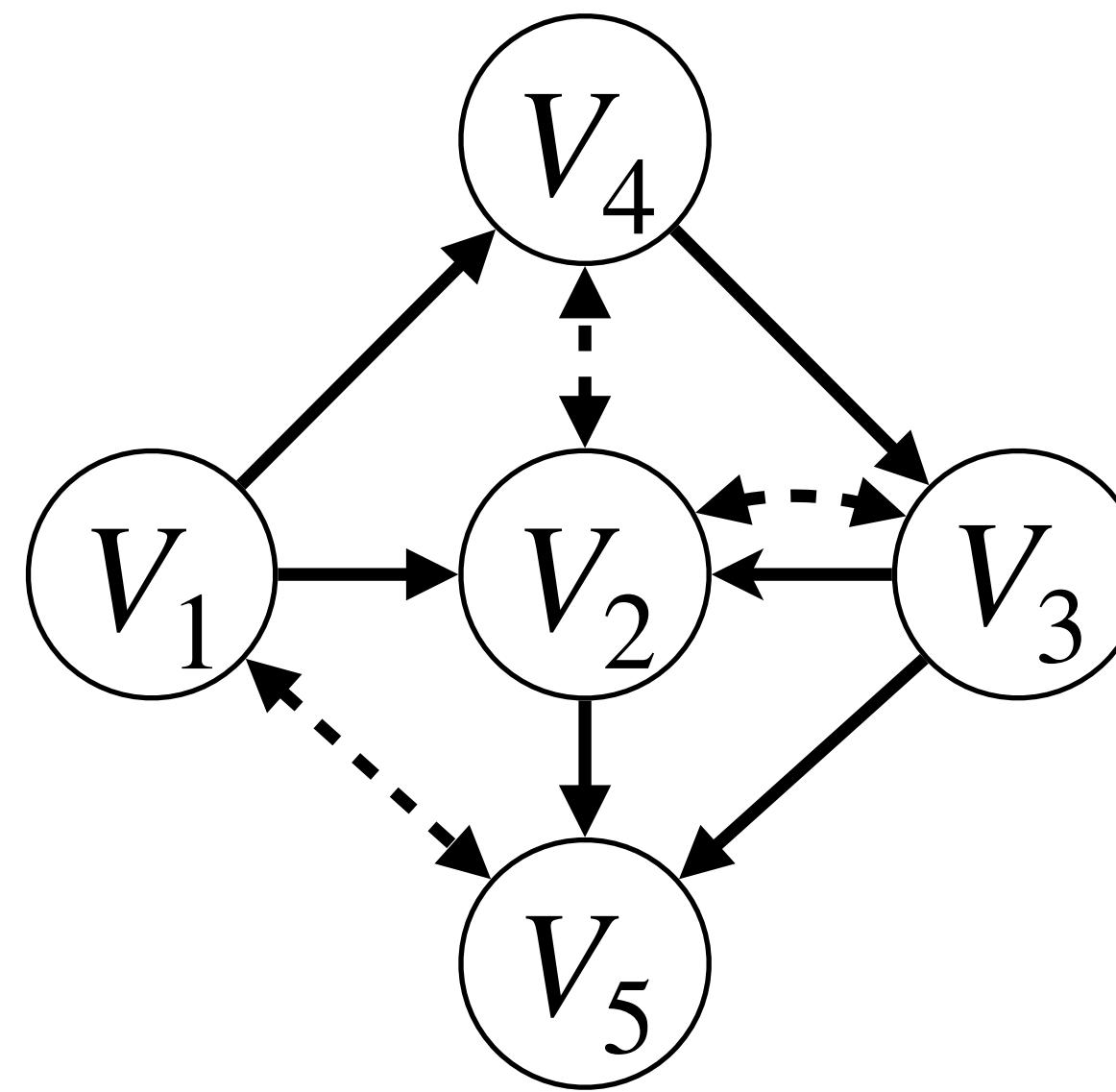
True, unknown ADMG



$\langle V_4, V_2, V_1, V_5 \rangle$ is a discriminating path for V_1 and
 $V_1 \in Sepset(V_4, V_5) - V_4 \perp\!\!\!\perp V_5 | V_1, V_2, V_3$



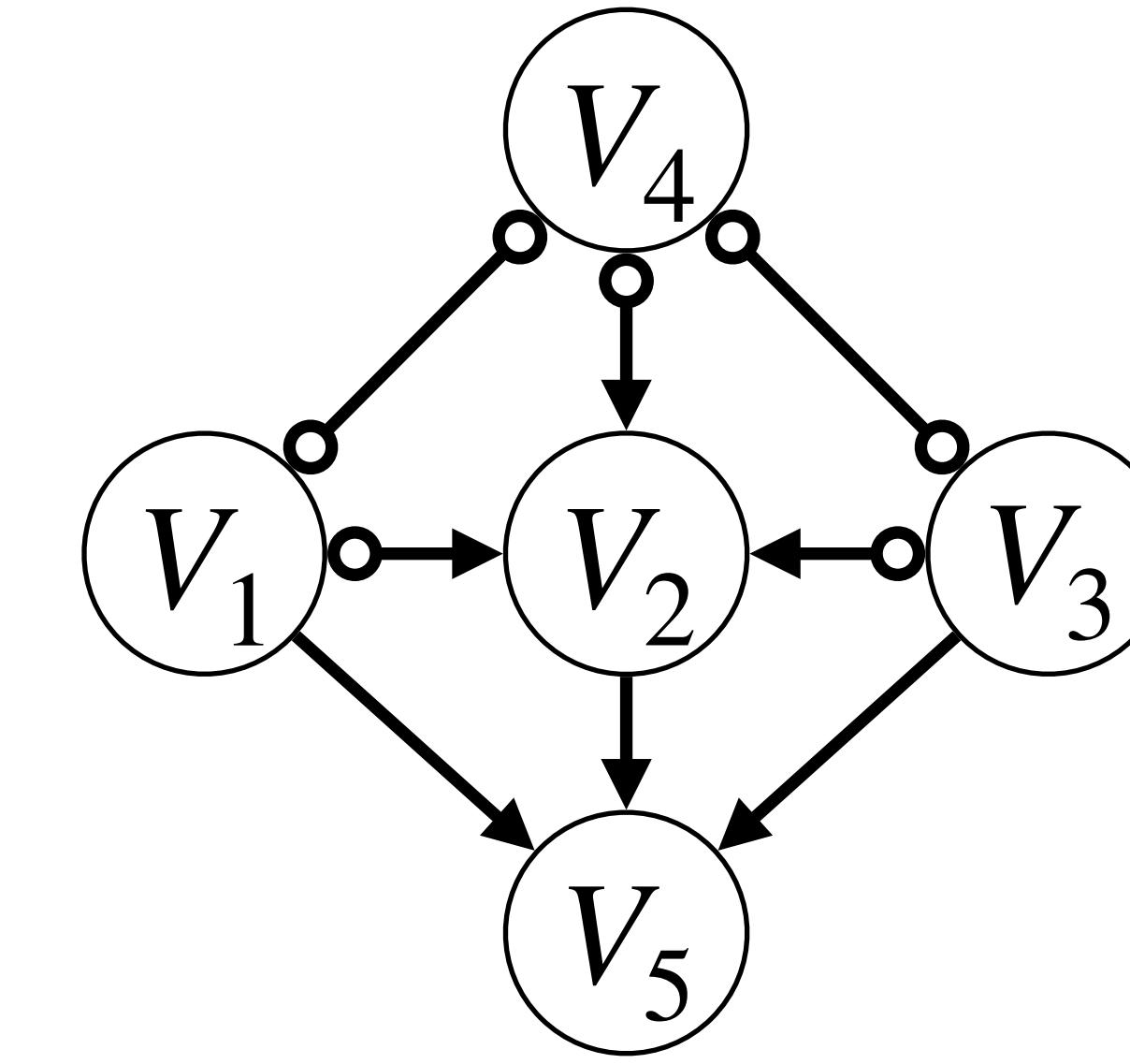
Final PAG



True, unknown ADMG

$$V_1 \perp\!\!\!\perp V_3 | V_4$$

$$V_4 \perp\!\!\!\perp V_5 | V_1, V_2, V_3$$



Final PAG

After Skel + R0 + R3 + R1 + R4 + R4

$$V_1 \perp\!\!\!\perp V_3 | V_4$$

$$V_4 \perp\!\!\!\perp V_5 | V_1, V_2, V_3$$

Causal Identification from PAGs



Can we identify causal effects from the equivalence class?

Effect Identification:

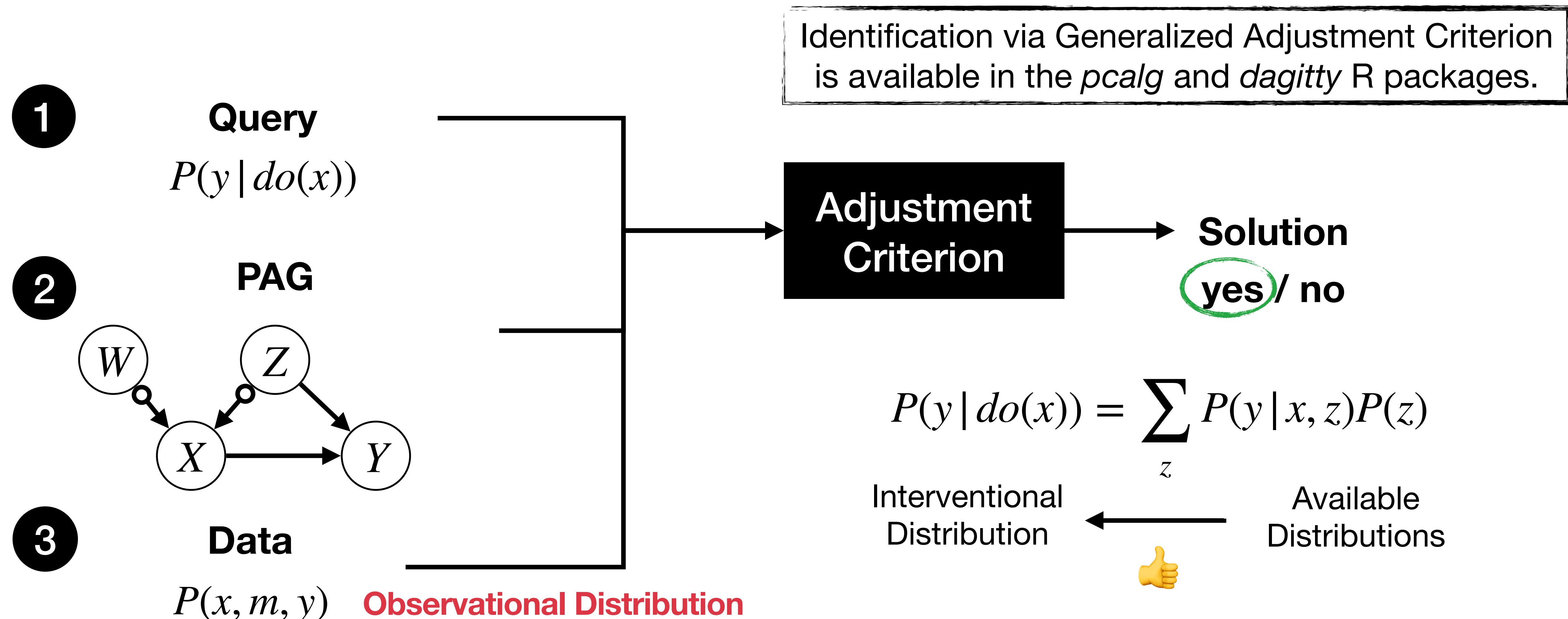
For Covariate Adjustment, we can use the Generalized Adjustment Criterion.

Recently, we proposed complete calculus and algorithms for the identification of marginal and conditional causal effect in PAGs!

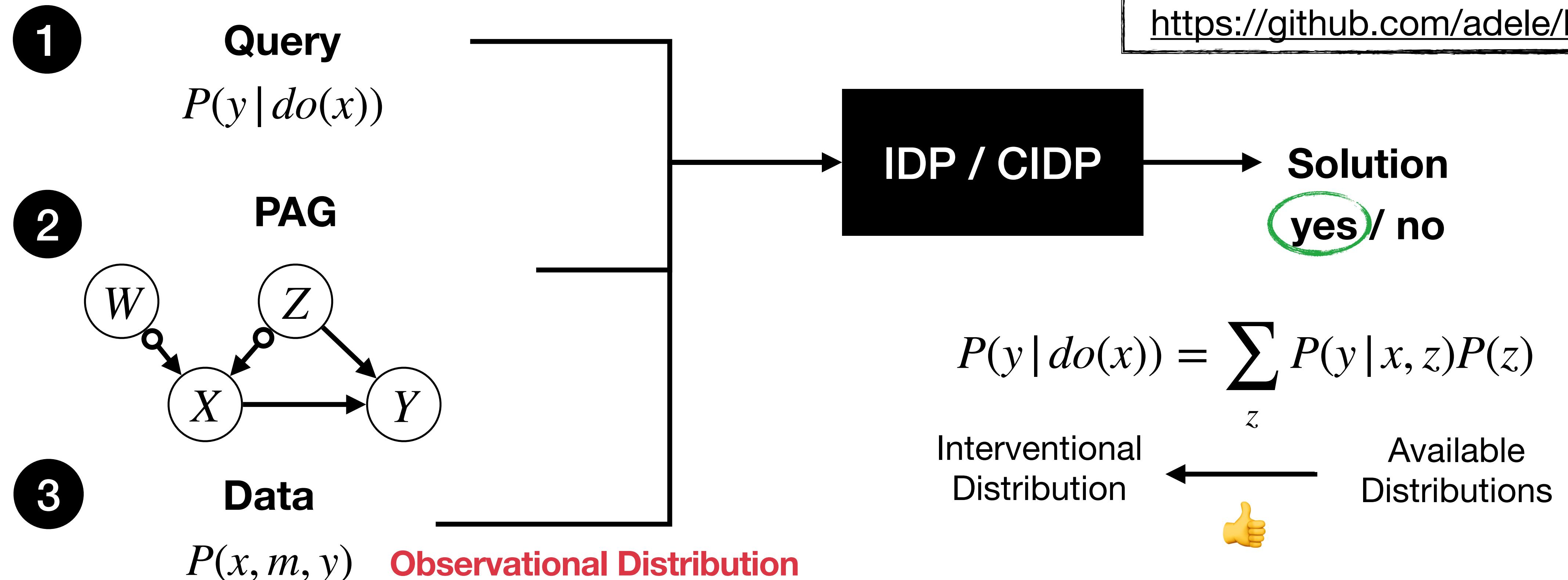
Perkovic, E., Textor, J. C., Kalisch, M., & Maathuis, M. H. (2018). Complete graphical characterization and construction of adjustment sets in Markov equivalence classes of ancestral graphs. Journal of Machine Learning Research 18 (2018) 1-62

Jaber A., **Ribeiro A. H.**, Zhang, J., Bareinboim, E. (2022) Causal Identification under Markov Equivalence - Calculus, Algorithm, and Completeness. In Proceedings of the 36th Annual Conference on Neural Information Processing Systems, NeurIPS. ([Link](#))

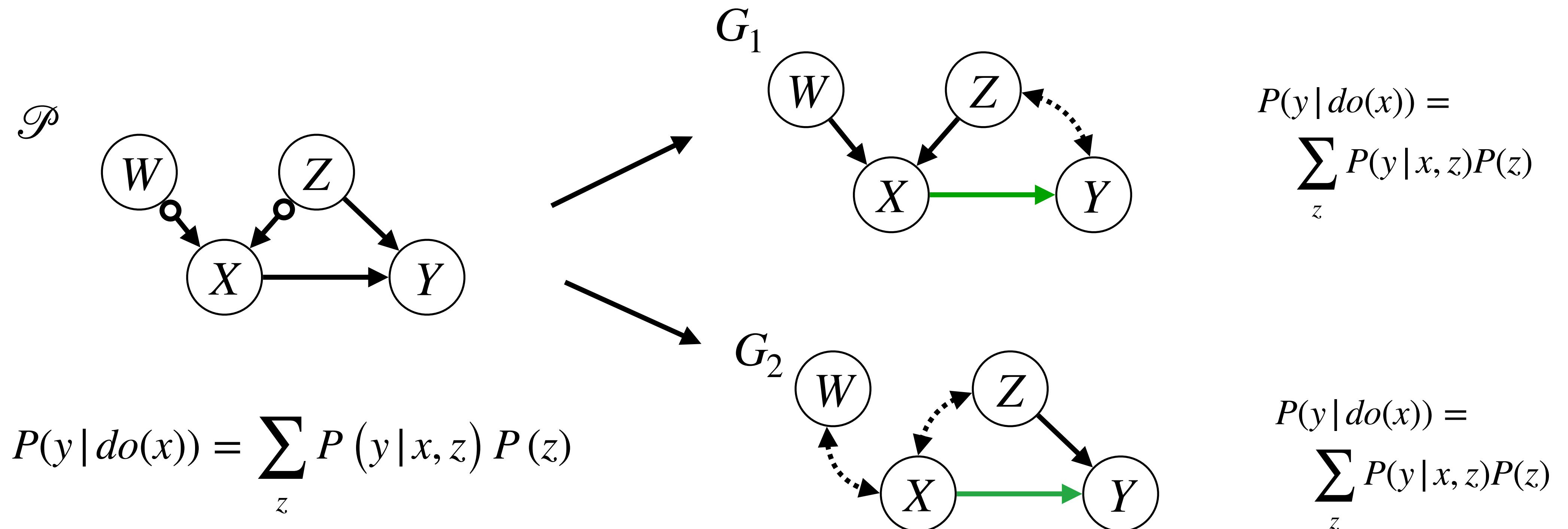
Adjustment in Markov Equivalence Classes



General Identification in Markov Equivalence Classes

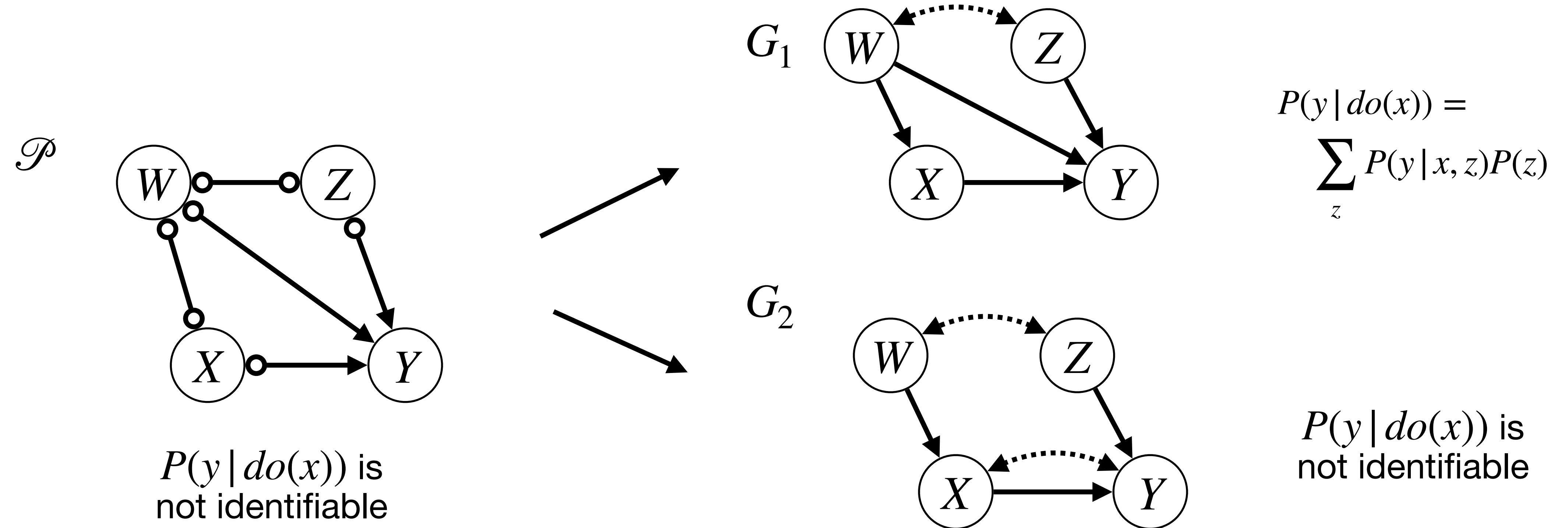


Effect Identifiability given a PAG



An effect identifiable in a PAG \mathcal{P} is identifiable in all causal diagrams G in the Markov Equivalence Class using the same identification formula!

Effect Non-Identifiability given a PAG



An effect not identifiable in a PAG \mathcal{P} is not identifiable in at least one causal diagrams G in the Markov Equivalence Class

Coding Exercises

Causality Tutorial:

- **Code:** <https://github.com/adele/Causality-Tutorial/blob/main/main.Rmd>
- **HTML Output:** <https://github.com/adele/Causality-Tutorial/blob/main/main.html>

Check Part II:

1. Causal Discovery
2. Causal Effect Identification from PAGs

Thank you! :)

Feel free to reach out to me if you have any questions:

adele.ribeiro@uni-marburg.de