

A hand in a suit sleeve is lighting a fuse that is positioned next to a row of seven wooden blocks of varying heights on a wooden surface. The background is a dark, textured wood.

Causal Inference Theory and Applications in Enterprise Computing

Dr. Matthias Uflacker, Johannes Huegle, Christopher Schmidt

April 16, 2019

- **Recap Causal Inference in a Nutshell**
- **Introduction to Structural Causal Models**
 1. Preliminaries
 2. Structural Causal Models
 3. (Local) Markov Condition
 4. Factorization
 5. Global Markov Condition
 6. Functional Model and Markov Conditions
 7. Faithfulness
 8. Constraint-based Causal Inference
 9. Markov Equivalence Class
 10. Summary
 11. Structural Causal Models in Application
 12. Excursion: Maximal Ancestral Graphs

Causal Inference
Theory and Applications
in Enterprise Computing

Uflacker, Huegle,
Schmidt

Slide 2

A hand in a dark suit jacket and white shirt cuff is shown striking a matchstick. The matchstick is lit, with a bright yellow flame and a trail of white smoke. The hand is positioned over a row of seven wooden blocks of varying heights, which are standing on a wooden surface. The background is a blurred wooden texture. A semi-transparent red banner is overlaid at the bottom of the image, containing the title text.

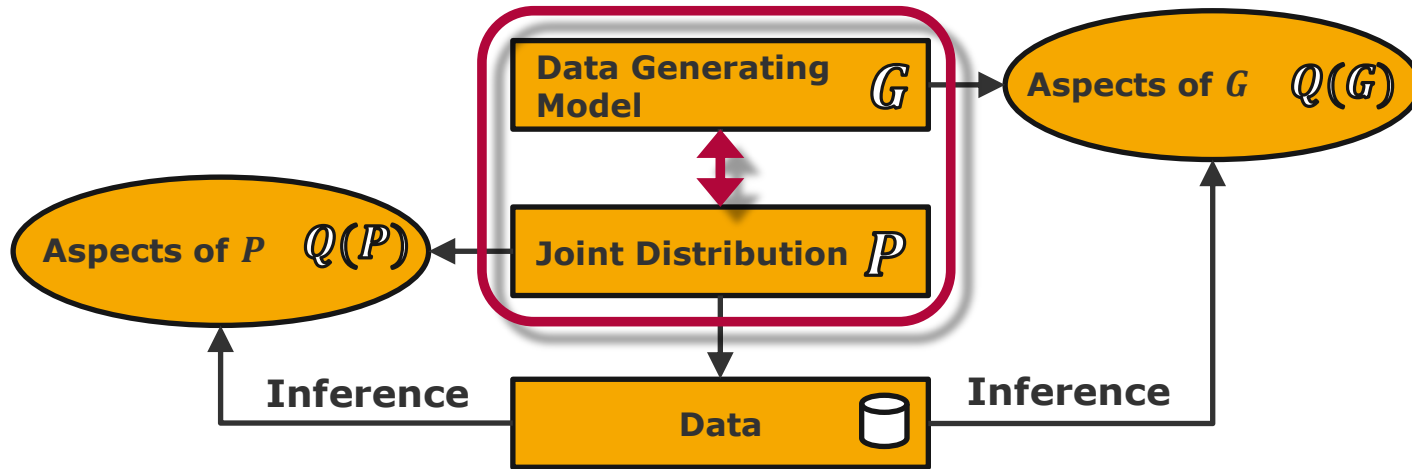
Causal Inference in a Nutshell

Causal Inference in a Nutshell

Recap: The Concept

Traditional Statistical Inference Paradigm

Paradigm of Structural Causal Models



Causal Inference
Theory and Applications
in Enterprise Computing

Uflacker, Huegle,
Schmidt

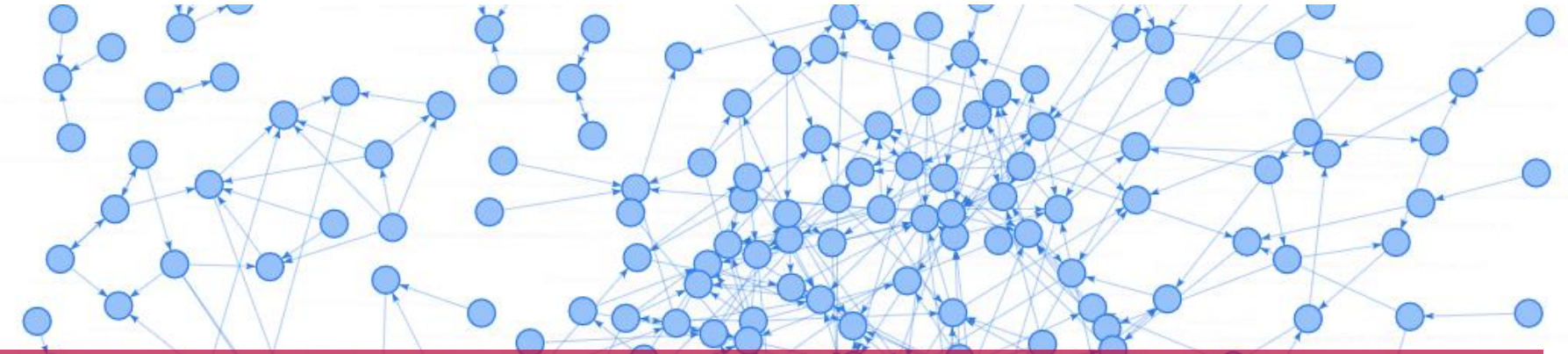
Slide 4

E.g., what is the sailors' probability of recovery when **we see** a treatment with lemons?

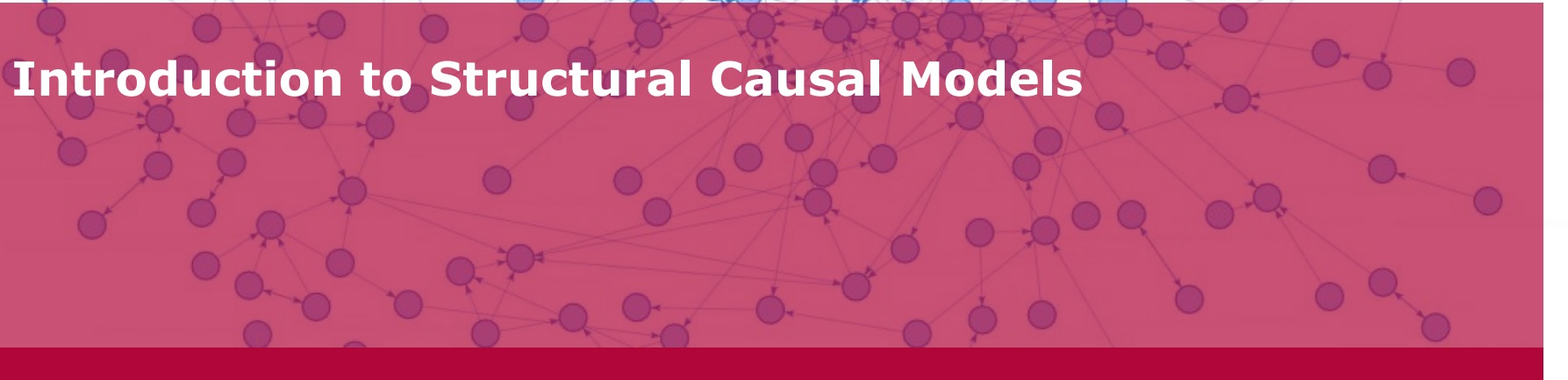
$$Q(P) = P(\text{recovery}|\text{lemons})$$

E.g., what is the sailors' probability of recovery if **we do** treat them with lemons?

$$Q(G) = P(\text{recovery}|\text{do}(\text{lemons}))$$



Introduction to Structural Causal Models



Introduction to Causal Graphical Models

Content

1. Preliminaries
2. Structural Causal Models
3. (Local) Markov Condition
4. Factorization
5. Global Markov Condition
6. Functional Model and Markov Conditions
7. Faithfulness
8. Constraint-based Causal Inference
9. Markov Equivalence Class
10. Summary
11. Structural Causal Models in Application
12. Excursion: Maximal Ancestral Graphs

Causal Inference
Theory and Applications
in Enterprise Computing

Uflacker, Huegle,
Schmidt

Slide 6

1. Preliminaries

Notation

- A, B events
- X, Y, Z random variables
- x value of random variable

- Pr probability measure
- P_X probability distribution of X
- p density
- $p(X)$ density of P_X
- $p(x)$ density of P_X evaluated at the point x

- $X \perp Y$ independence of X and Y
- $X \perp Y \mid Z$ conditional independence of X and Y given Z

1. Preliminaries

Independence of Events

- Two events A and B are called *independent* if

$$\Pr(A \cap B) = \Pr(A) \cdot \Pr(B),$$

or - rewritten in *conditional probabilities* - if

$$\Pr(A) = \frac{\Pr(A \cap B)}{\Pr(B)} = \Pr(A|B),$$

$$\Pr(B) = \frac{\Pr(A \cap B)}{\Pr(A)} = \Pr(B|A).$$

- A_1, \dots, A_n are called (*mutually*) *independent* if for every subset $S \subset \{1, \dots, n\}$ we have

$$\Pr\left(\bigcap_{i \in S} A_i\right) = \prod_{i \in S} \Pr(A_i).$$

- Note:**

for $n \geq 3$, pairwise independence $\Pr(A_i \cap A_j) = \Pr(A_i) \cdot \Pr(A_j)$ for all i, j does not imply (mutual) independence.

1. Preliminaries

Independence of Random Variables

- Two real-valued random variables X and Y are called *independent*,

$$X \perp Y,$$

if for every $x, y \in \mathbb{R}$, the events $\{X \leq x\}$ and $\{Y \leq y\}$ are independent,

Or, in terms of densities: for all x, y ,

$$p(x, y) = p(x)p(y).$$

- Note:**

If $X \perp Y$, then $E[XY] = E[X]E[Y]$, and $cov(X, Y) = E[XY] - E[X]E[Y] = 0$.

The converse is not true: ~~If $cov(X, Y) = 0$, then $X \perp Y$.~~

No correlation does not imply independence

However, we have, for large \mathcal{F} : $(\forall f, g \in \mathcal{F}: cov(f(X), g(Y)) = 0)$, then $X \perp Y$.

Causal Inference
Theory and Applications
in Enterprise Computing

Uflacker, Huegle,
Schmidt

1. Preliminaries

Conditional Independence of Random Variables

- Two real-valued random variables X and Y are called *conditionally independent* given Z ,

$$X \perp Y \mid Z \text{ or } (X \perp Y \mid Z)_P$$

if

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

For all x, y and for all z s.t. $p(z) > 0$.

- Note:**

It is possible to find X, Y which are conditionally independent given a variable Z but unconditionally dependent, and vice versa.

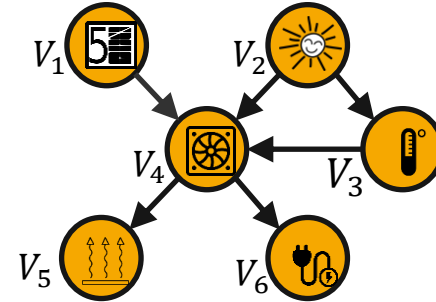
2. Structural Causal Models

Definition (Pearl)

- Directed Acyclic Graph (DAG) $G = (V, E)$
 - *Vertices* V_1, \dots, V_n
 - *Directed edges* $E = (V_i, V_j)$, i.e., $V_i \rightarrow V_j$,
 - *No cycles*
- Use kinship terminology, e.g., for path $V_i \rightarrow V_j \rightarrow V_k$
 - $V_i = Pa(V_j)$ *parent* of V_j
 - $\{V_i, V_j\} = Ang(V_k)$ *ancestors* of V_k
 - $\{V_j, V_k\} = Des(V_i)$ *descendants* of V_i
- Directed Edges encode *direct causes* via
 - $V_j = f_j(Pa(V_j), N_j)$ with independent noise N_1, \dots, N_n

➔ This forms the Causal Graphical Model

Cooling House Example:



- $V_1 = N(0,1)$
- $V_2 = N(0,1)$
- $V_3 = 3 V_2 + N(0,1)$
- $V_4 = 4 V_1 + 5 V_2 + 0.7 V_3 + N(0,1)$
- $V_5 = V_4 + N(0,1)$
- $V_6 = 1.2 V_4 + N(0,1)$

Causal Inference
Theory and Applications
in Enterprise Computing

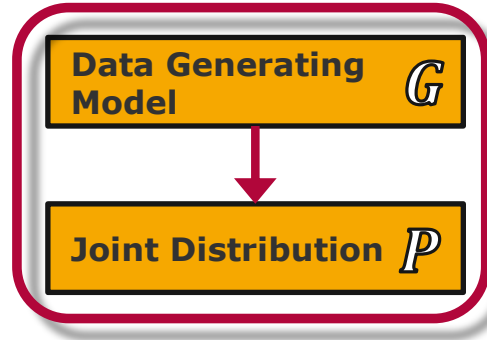
Uflacker, Huegle,
Schmidt

Slide 11

2. Structural Causal Models

Connecting G and P

- Basic Assumption: *Causal Sufficiency*
 - All relevant variables are included in the DAG G



$$(X \perp Y|Z)_G \Rightarrow (X \perp Y|Z)_P$$

- Key Postulate: *(Local) Markov Condition*
- Essential mathematical concept: *d-separation*
(describes the conditional independences required by a causal DAG)

Causal Inference
Theory and Applications
in Enterprise Computing

Uflacker, Huegle,
Schmidt

Slide **12**

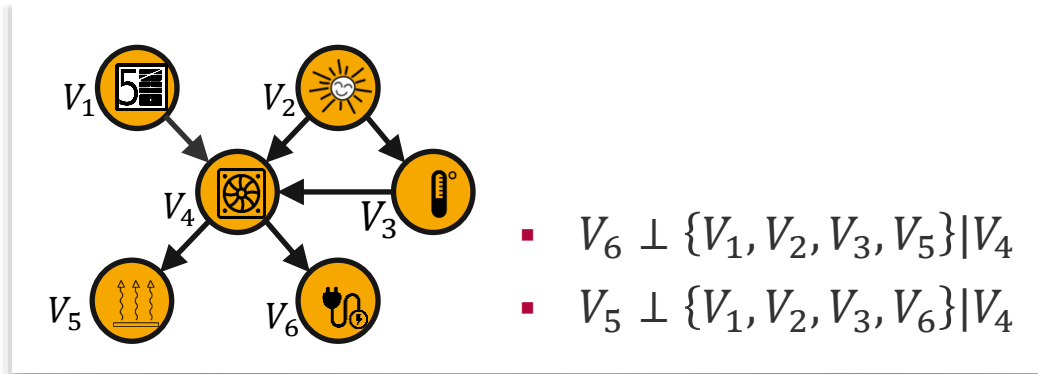
3. (Local) Markov Condition Theorem

(Local) Markov Condition:

V_j independent of nondescendants $ND(V_j)$, given parents $Pa(V_j)$, i.e.,

$$V_j \perp V_{V/(Des(V_j) \cup Pa(V_j))} | Pa(V_j).$$

- I.e., every information exchange with its nondescendants involves its parents
- Example:



3. (Local) Markov Condition

Supplement (Lauritzen 1996)

- Assume V_n has no descendants, then $ND(V_n) = \{V_1, \dots, V_{n-1}\}$.
- Thus the local Markov condition implies

$$V_n \perp \{V_1, \dots, V_{n-1}\} / Pa(V_n) \mid Pa(V_n).$$

- Hence, the general decomposition

$$p(v_1, \dots, v_n) = p(v_n | v_1, \dots, v_{n-1}) p(v_1, \dots, v_{n-1})$$

becomes

$$p(v_1, \dots, v_n) = p(v_n | Pa(v_n)) p(\{v_1, \dots, v_{n-1}\} / Pa(v_n)).$$

- Induction over n yields to

$$p(v_1, \dots, v_n) = \prod_{i=1}^n p(v_i | Pa(v_i)).$$

- I.e., the graph shows us how to factor the joint distribution P_V .

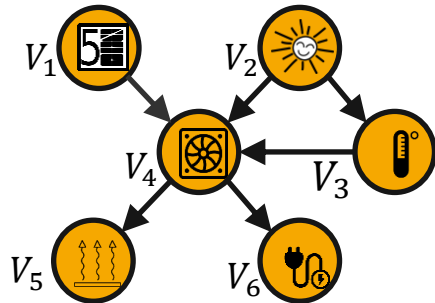
4. Factorization

Definition

Factorization:

$$p(v_1, \dots, v_n) = \prod_{i=1}^n p(v_i | Pa(v_i)).$$

- I.e., conditionals as causal mechanisms generating statistical dependence
- Example:



$$\begin{aligned} p(V) &= p(v_1, \dots, v_n) \\ &= p(v_1) \cdot p(v_2) \\ &\quad \cdot p(v_3 | v_2) \cdot p(v_4 | v_1, v_2, v_3) \\ &\quad \cdot p(v_5 | v_4) \cdot p(v_6 | v_4) \\ &= \prod_{i=1}^n p(v_i | Pa(v_i)) \end{aligned}$$

5. Global Markov Condition

D-Separation (Pearl 1988)

- *Path* = sequence of pairwise distinct vertices where consecutive ones are adjacent
- A path q is said to be *blocked* by a set S if
 - q contains a *chain* $V_i \rightarrow V_j \rightarrow V_k$ or a *fork* $V_i \leftarrow V_j \rightarrow V_k$ such that the middle node is in S , or
 - q contains a *collider* $V_i \rightarrow V_j \leftarrow V_k$ such that the middle node is not in S and such that no descendant of V_j is in S .

D-separation:

S is said to **d-separate** X and Y in the DAG G , i.e.,

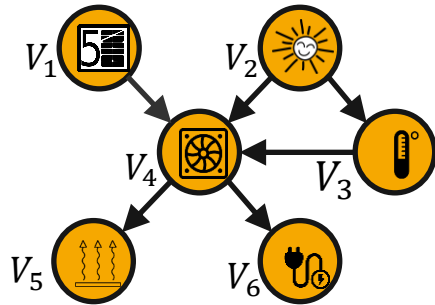
$$(X \perp Y | S)_G,$$

if S blocks every path from a vertex in X to a vertex in Y .

5. Global Markov Condition

Examples of d-Separation

■ Example:



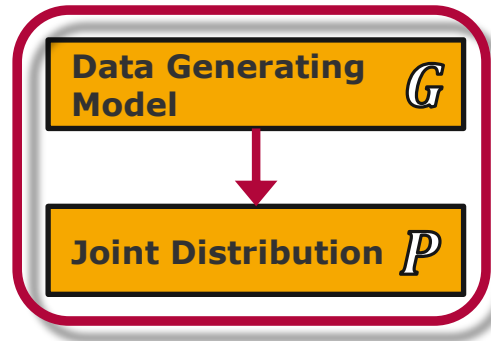
- The path from V_1 to V_6 is blocked by V_4 .
- V_1 and V_6 are d-separated by V_4 .
- The path $V_2 \rightarrow V_3 \rightarrow V_4 \rightarrow V_6$ is blocked by V_3 or V_4 or both.
- But: V_2 and V_6 are d-separated only by V_4 or $\{V_3, V_4\}$.
- V_1 and V_2 are not blocked by V_4 .
- V_4 is a fork in $V_5 \leftarrow V_4 \rightarrow V_6$.
- V_5 and V_6 are d-separated by V_4 .

5. Global Markov Condition Theorem

Global Markov Condition:

For all disjoint subsets of vertices X, Y and Z we have that
 X, Y d-separated by $Z \Rightarrow (X \perp Y | Z)_P$.

- I.e., we have $(X \perp Y | Z)_G \Rightarrow (X \perp Y | Z)_P$



Causal Inference
Theory and Applications
in Enterprise Computing

Uflacker, Huegle,
Schmidt

Slide **18**

6. Functional Model and Markov Conditions

Theorem (Lauritzen 1996, Pearl 2000)

Theorem:

The following are equivalent:

- Existence of a *functional causal model* G ;
- *Local Causal Markov condition*: V_j statistically independent of nondescendants, given parents
(i.e.: every information exchange with its nondescendants involves its parents)
- *Global Causal Markov condition*: d-separation
(characterizes the set of independences implied by local Markov condition)
- *Factorization*: $p(v_1, \dots, v_n) = \prod_{i=1}^n p(v_i | Pa(v_i))$.

(subject to technical conditions)

$$\text{I.e., } (X \perp Y | Z)_G \Rightarrow (X \perp Y | Z)_P$$

7. Causal Faithfulness

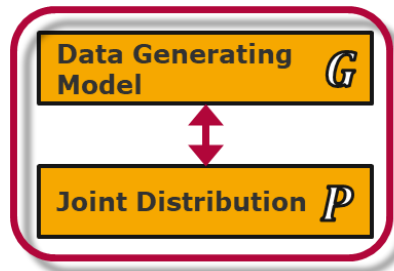
The key-postulate

Causal Faithfulness:

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

$$(X \perp Y | Z)_G \Leftarrow (X \perp Y | Z)_P$$

- I.e., we assume that any population P produced by this causal graph G has the independence relations obtained by applying d-separation to it
- Seems like a hefty assumption, but it really isn't: It assumes that whatever independencies occur in it arise not from incredible coincidence but rather from structure, i.e., data generating model G .
- Hence:



Causal Inference
Theory and Applications
in Enterprise Computing

Uflacker, Huegle,
Schmidt

Slide 20

8. Constraint-based Causal Inference

Concept (Spirtes, Glymor, Scheines and Pearl)

■ Assumptions:

- Causal Sufficiency
- Global Markov Condition
- Causal Faithfulness

■ Causal Structure Learning:

- Accept only those DAG's G as causal hypothesis for which
$$(X \perp Y | Z)_G \Leftrightarrow (X \perp Y | Z)_P.$$
- Defines the basis of constraint-based causal structure learning
- Identifies causal DAG up to Markov equivalence class (DAGs that imply the same conditional independencies)

Causal Inference
Theory and Applications
in Enterprise Computing

Uflacker, Huegle,
Schmidt

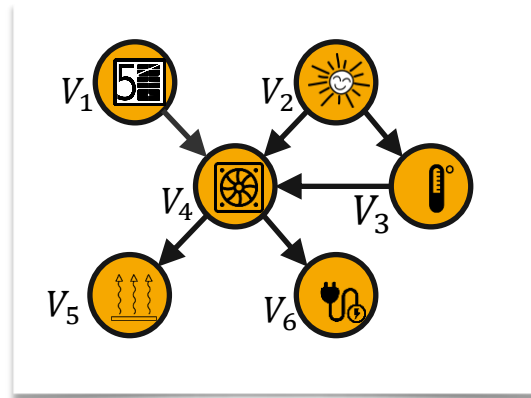
9. Markov Equivalence Class

Theorem (Verma and Pearl)

Theorem:

Two DAGs are Markov equivalent if and only if they have the same skeleton and the same v -structures

- *Skeleton:*
corresponding undirected graph
- *v -structure:*
substructure $X \rightarrow Y \leftarrow Z$ with no edges between X and Z .



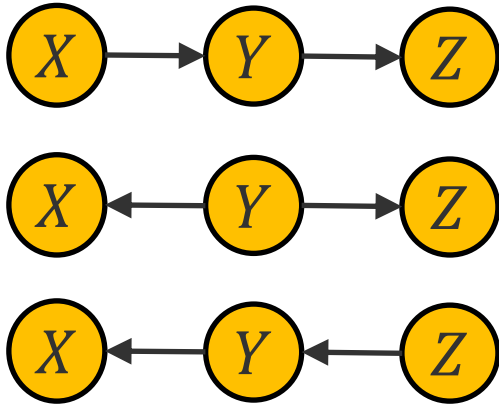
Causal Inference
Theory and Applications
in Enterprise Computing

Uflacker, Huegle,
Schmidt

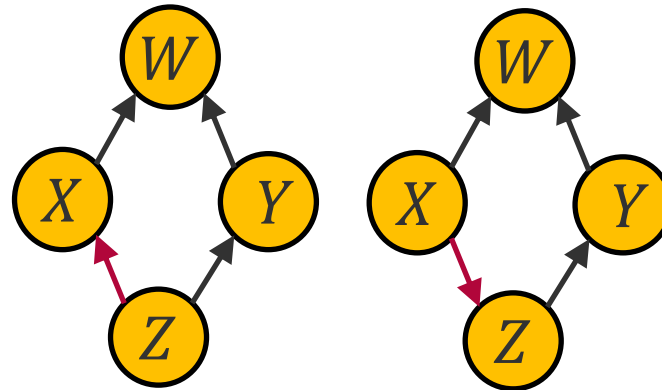
9. Markov Equivalence Class

Examples

- Same skeleton, no v -structure
- Same skeleton, same v -structure at W



$$X \perp Z \mid Y$$



10. Summary

Causal Structural Models

- Causal Structures formalized by DAG (directed acyclic graph) G with random variables V_1, \dots, V_n as vertices.

- Causal Sufficiency, Causal Faithfulness and Markov Condition imply

$$(X \perp Y | Z)_G \Leftrightarrow (X \perp Y | Z)_P.$$

- Local Markov Condition states that the density $p(v_1, \dots, v_n)$ then factorizes into

$$p(v_1, \dots, v_n) = \prod_{i=1}^n p(v_i | Pa(v_i)).$$

- Causal conditional $p(v_j | Pa(v_j))$ represent causal mechanisms.

11. Structural Causal Models in Application

Cooling House Example

```
Causal_Inference_in_Applica X
+ - 🔍 📄 ▶ ⌂ Markdown - R
In [ ]: n <- 10000
coolingData <- rmvDAG(n,coolingDAG)
#head(coolingData)
plot(density(coolingData[,1]), main="Density Plot")
```

This introduces the functional mechanisms in our system, which are described by the following equations

- $V_1 = N(0, 1)$
- $V_2 = N(0, 1)$
- $V_3 = 3 \cdot V_2 + N(0, 1)$
- $V_4 = 4 \cdot V_1 + 5 \cdot V_2 + 0.7 \cdot V_3 + N(0, 1)$
- $V_5 = V_4 + N(0, 1)$
- $V_6 = 1.2 \cdot V_4 + N(0, 1)$

In the following, we assume that these functional mechanisms are not known such that the goal remains to derive the causal relationships and the causal effects.

When looking at the correlationmatrix as a first examination step, we see that all variables are highly correlated:

```
In [ ]: round(cor(coolingData), 2)
```

3. Causal Graphical Models

In the framework causal graphical models, a directed edge $V_i \rightarrow V_j$ in our DAG represents a direct causal relationship of V_i to V_j .

3.A. D-Separation

Causal Sufficiency, Causal Faithfulness and Markov Condition imply that $(X \perp Y|Z)_G \Leftrightarrow (X \perp Y|Z)_P$. The essential mathematical concept is to find the d-separating sets S , e.g.

```
In [ ]: # Are V2 and V6 are d-separated by an empty set
dsep("V2","V6",NULL,coolingDAG)

# Are V2 and V6 are d-separated by V3 and V4?
dsep("V2","V6",c("V3","V4"),coolingDAG)
```

3.B. Conditional Independence

Then causal faithfulness and the Markov condition imply that two vertices V_i, V_j are conditionally independent given a set $S(V_i, V_j)$ if and only if the vertices V_i and V_j are d-separated by the set $S(V_i, V_j)$, e.g.:

```
In [ ]: # Are V2 and V6 are independent?
x <- 2
y <- 6
S <- c()
condIndFisherZ(x,y,S,cor(coolingData),n,qnorm(1- 0.05/2))

# Are V2 and V6 are independent given V3 and V4?
x <- 2
y <- 6
S <- c("V3","V4")
condIndFisherZ(x,y,S,cor(coolingData),n,qnorm(1- 0.05/2))
```

12. Excursion: Maximal Ancestral Graphs

Motivating Example

- Suppose, we are given the following list of conditional independencies among X, Y, Z and W :

- $X \perp\!\!\!\perp Z$,
- $Y \perp\!\!\!\perp W$,
- $X \perp\!\!\!\perp W$.
- $X \not\perp\!\!\!\perp Y$,
- $Y \not\perp\!\!\!\perp Z$,
- $Z \not\perp\!\!\!\perp W$.

- Which DAG could have generated these, and only these, independencies and dependencies?
- The pattern of dependencies must be:

$$X \text{ --- } Y \text{ --- } Z \text{ --- } W$$

- And there must be the following colliders:

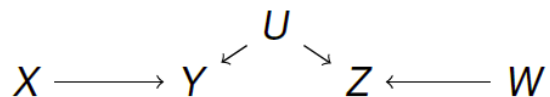
$$X \longrightarrow Y \longleftarrow Z$$

$$Y \longrightarrow Z \longleftarrow W$$

- There is no orientation of Y - Z that is consistent with the independencies.

12. Excursion: Maximal Ancestral Graphs DAG Models and Marginalization

- Let's include an additional variable U :



- This DAG model generates a probability distribution $P_{\{X,Y,Z,W,U\}}$ in which:
 - $X \perp\!\!\!\perp Z$,
 - $Y \perp\!\!\!\perp W$,
 - $X \perp\!\!\!\perp W$.
 - $X \not\perp\!\!\!\perp Y$,
 - $Y \not\perp\!\!\!\perp Z$,
 - $Z \not\perp\!\!\!\perp W$.
- The marginal distribution $P_{\{X,Y,Z,W\}} = P_{\{X,Y,Z,W,U\}} du$ must adhere the same independencies. But: this marginal distribution cannot be faithfully generated by any DAG.

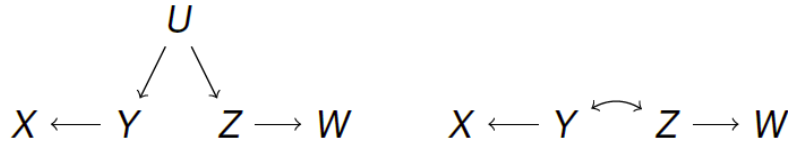
➔ **DAG models are not closed under marginalization!**

12. Excursion: Maximal Ancestral Graphs

Ancestral Graphs (informally)

- Ancestral Graph (AG)**

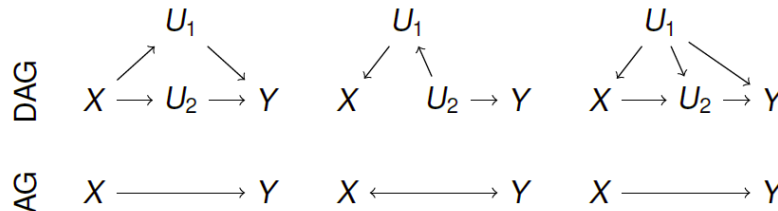
is a graph containing both directed and bi-directed edges, where the bi-directed edges stand for *latent variables*, e.g.,



- m-Separation**

If S m-separates X and Y in an ancestral graph M , then $X \perp Y \mid S$ in every density p that factorizes according to any DAG G that is represented by the AG M .

- Example**



12. Excursion: Maximal Ancestral Graphs

DAGs vs. AGs

■ Advantages of AGs

- AGs can faithfully represent more probability distributions than DAGs.
- AG models are closed under marginalization.
- AGs can (implicitly) represent unobserved variables, which exist in many (possibly almost all) applications.

■ Disadvantages of AGs

- Parameterization is difficult in the general case.
- Markov equivalence is difficult.

Literature

- Pearl, J. (2009). *Causal inference in statistics: An overview*. Statistics Surveys, 3:96-146.
- Pearl, J. (2009). *Causality: Models, Reasoning, and Inference*. Cambridge University Press.
- Spirtes, P., Glymour, C., and Scheines, R. (2000). Causation, Prediction, and Search. The MIT Press.

Causal Inference
Theory and Applications
in Enterprise Computing

Uflacker, Huegle,
Schmidt

Slide **30**

Thank you
for your attention!