



Towards structure learning

Causality

Christina Heinze-Deml

Spring 2019

Last week

- Counterfactuals
- Potential outcomes
- (Propensity score) matching

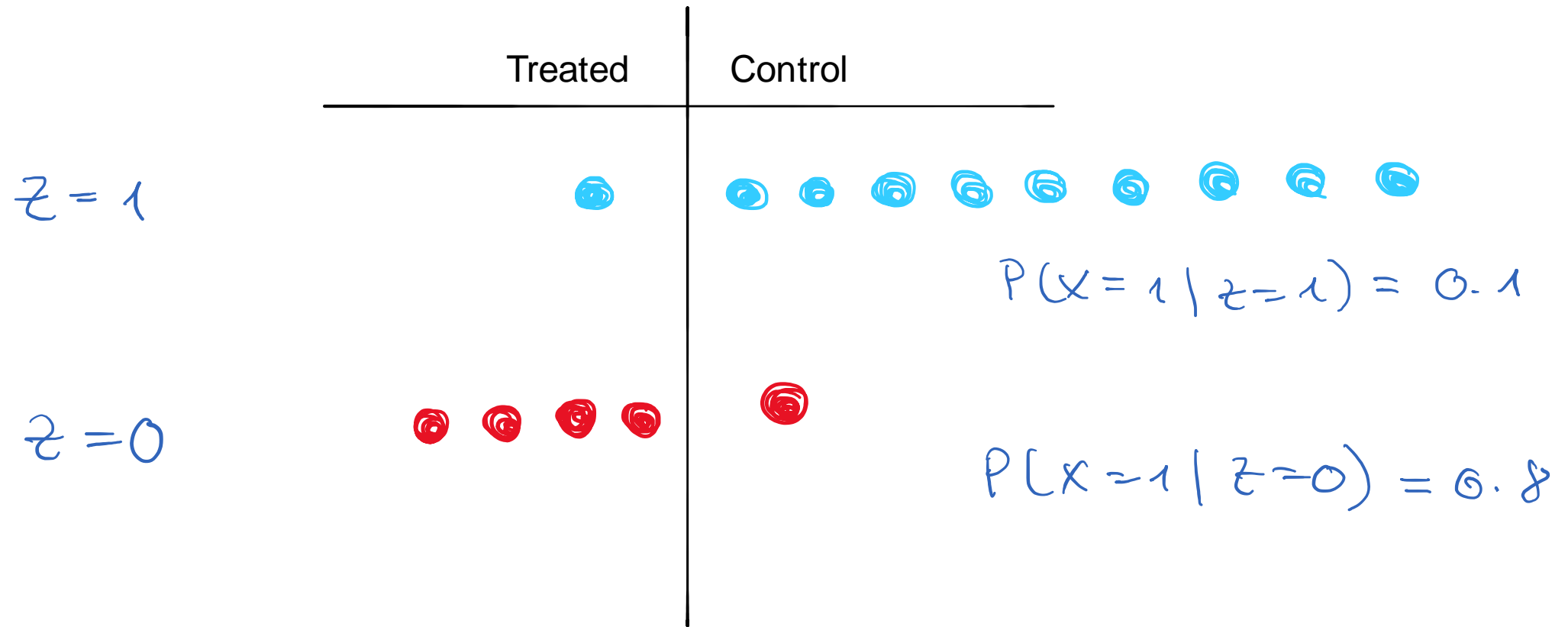
Today

- Inverse probability weighting
- “Towards structure learning”

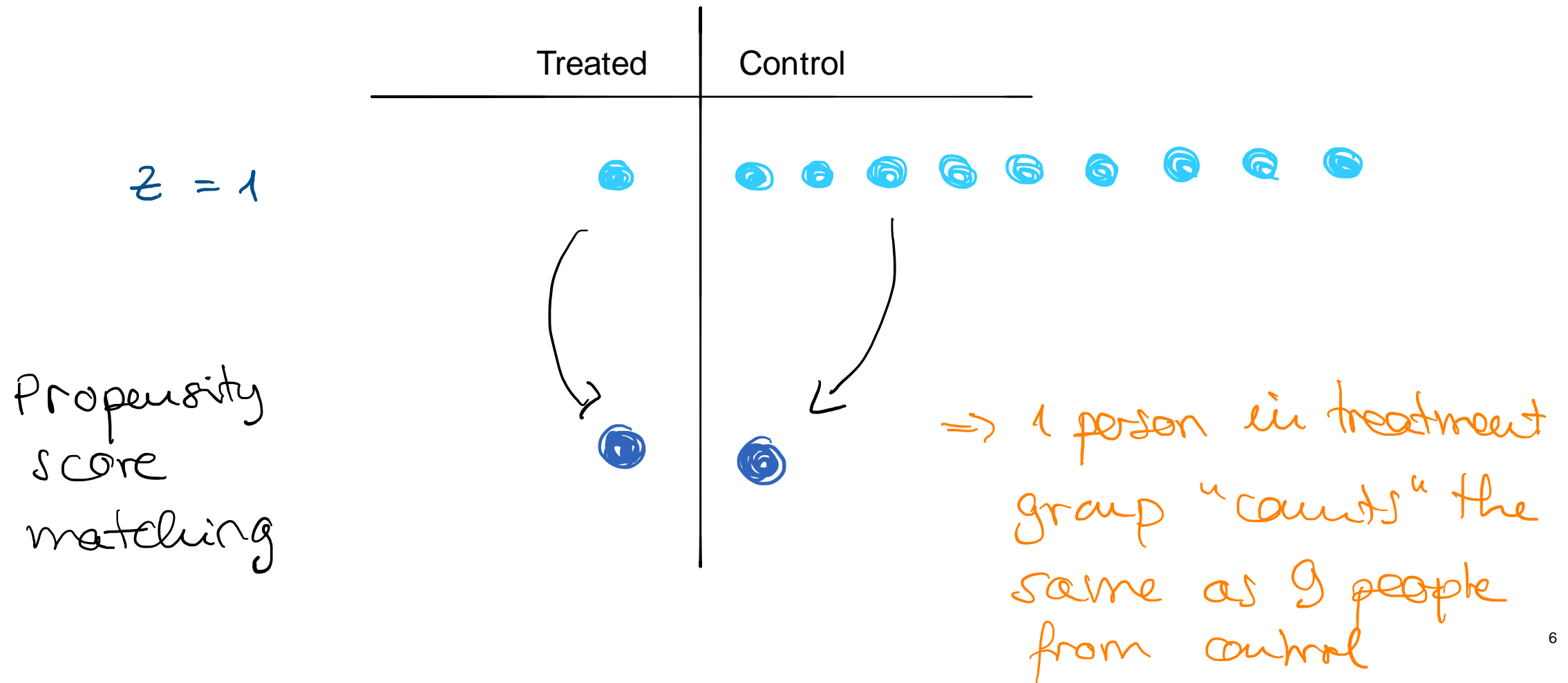
Inverse probability weighting

- Example
 - Single binary confounder Z
 - Suppose propensity score $P(X = 1|Z = 1) = 0.1$
 - Among people with $Z = 1$, only 10% receive the treatment
 - Suppose propensity score $P(X = 1|Z = 0) = 0.8$
 - Among people with $Z = 0$, 80% receive the treatment

Inverse probability weighting

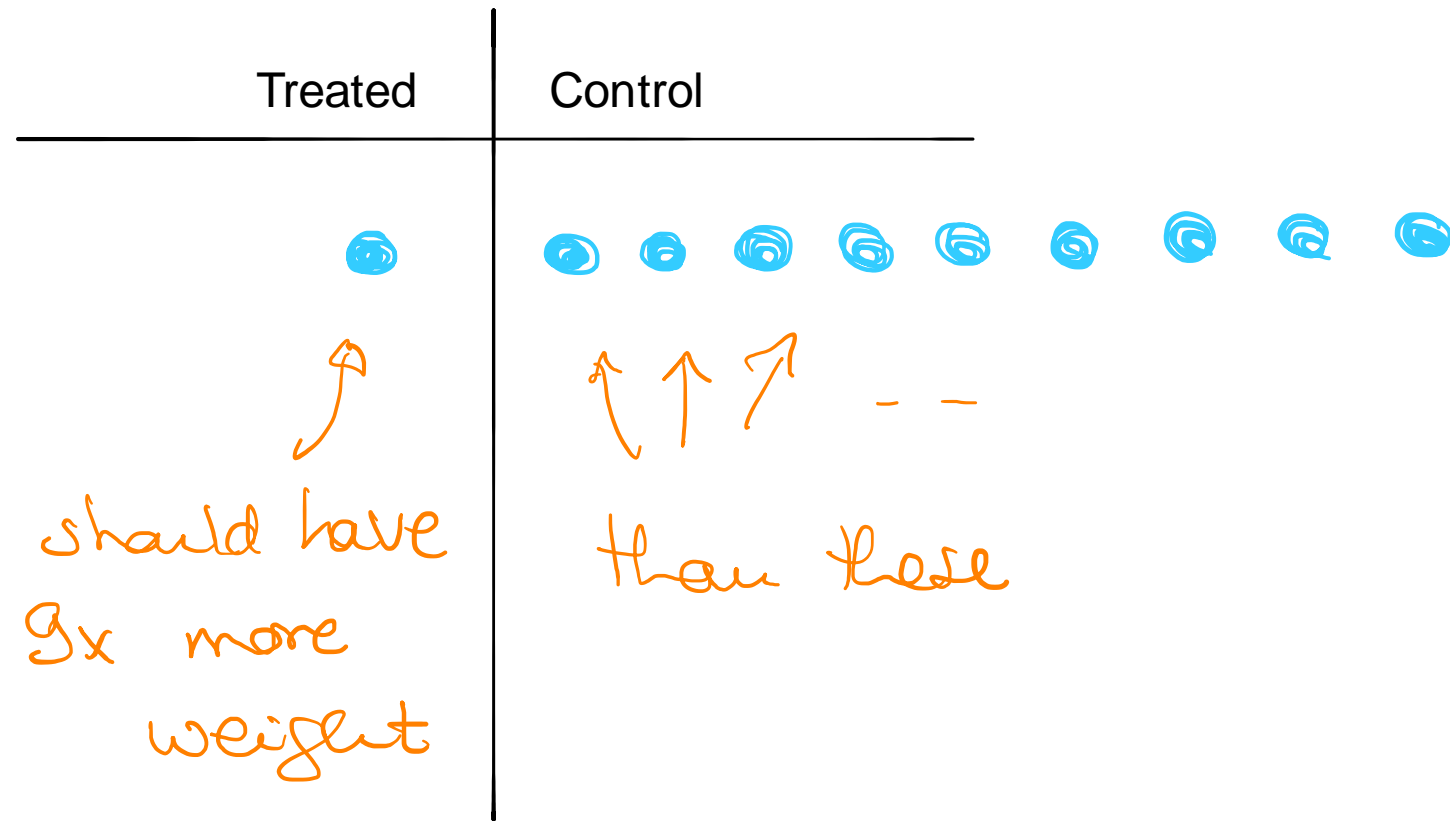


Inverse probability weighting



Inverse probability weighting

$$e = 1$$



Inverse probability weighting

- **Idea:** rather than match, use all data but downweigh/upweigh observations
 - Weighting by the inverse of the probability of treatment **received**
 - For treated: weigh by the inverse of the propensity score $\pi = P(X = 1|Z)$
 - For control: weigh by the inverse of $1 - \pi = P(X = 0|Z)$
- Known as **inverse probability of treatment weighting (IPTW)**

Example

Inverse probability weighting

- Estimator

$$\hat{E}(Y|do(X = 1)) = \frac{1}{n} \sum_i Y_i 1\{X_i = 1\} w_i$$

where $w_i = \frac{1}{\hat{\pi}_i} = \frac{1}{\hat{P}(X=1|Z_i)}$

- Equivalently for $\hat{E}(Y|do(X = 0))$
- If propensity score $\pi_i = P(X = 1|Z_i)$ is very small, weight will be very large
- Small estimation errors in $\hat{\pi}_i$ can lead to large estimation errors in $\hat{E}(Y|do(X = x))$

Inverse probability weighting

- More generally, consider
 - Observational distribution $p(x_V)$
 - Interventional distribution $p(x_V | do(X_k \leftarrow \tilde{N}_k))$ with imperfect intervention
 - Shorthand: $p(x_V | do(X_k \leftarrow \tilde{N}_k)) = \tilde{p}(x_V)$
- Factorizations agree except for the term of the intervened variable
 - Assume strictly positive densities

Inverse probability weighting

- Observational distribution $p(x_V)$
- Interventional distribution $p(x_V | do(X_k \leftarrow \tilde{N}_k)) = \tilde{p}(x_V)$
- Interested in certain aspect $l(x_V)$, then:

$$\begin{aligned} & E \left(l(x_V) | do(X_k \leftarrow \tilde{N}_k) \right) \\ &= \int l(x_V) \tilde{p}(x_V) dx_V \\ &= \int l(x_V) \frac{\tilde{p}(x_V)}{p(x_V)} p(x_V) dx_V \\ &= \int l(x_V) \frac{\tilde{p}(x_k | x_{\text{pa}(k)})}{p(x_k | x_{\text{pa}(k)})} p(x_V) dx_V \end{aligned}$$

Inverse probability weighting

- Given observations x_V^1, \dots, x_V^n drawn from **observational distribution** $p(x_V)$, can construct estimator for expectation under **interventional distribution**:

$$\hat{E} \left(l(x_V) | do(X_k \leftarrow \tilde{N}_k) \right) = \frac{1}{n} \sum_i l(x_V^i) w_i$$

where $w_i = \frac{\tilde{p}(x_k^i | x_{\text{pa}(k)}^i)}{p(x_k^i | x_{\text{pa}(k)}^i)}$

- [See Series 4.]
- Related to survey sampling, importance sampling, reinforcement learning
 - See Elements of Causal Inference, Chapter 8.2.

Towards structure learning

- Markov properties
- Causal minimality
- Faithfulness
- Markov equivalence

Markov properties

- Given a DAG $G = (V, E)$, a distribution P with density p on X_V is said to satisfy:
 - The **global Markov property** wrt G if for all pairwise disjoint subsets A, B and S of V :

$$A \text{ and } B \text{ are d-separated by } S \text{ in } G \Rightarrow X_A \perp\!\!\!\perp X_B | X_S \text{ in } P$$

- The **local Markov property** with respect to G if for all $j \in V$:

$$X_j \perp\!\!\!\perp X_{\text{nondesc}(j) \setminus \text{pa}(j)} | X_{\text{pa}(j)}$$

- The **Markov factorization property** with respect to G if

$$p(x_V) = \prod_{j \in V} p(x_j | x_{\text{pa}(j)})$$

Markov properties

- If P has a density (with respect to a product measure), then all three Markov properties are equivalent. We then simply say that:
 - P is **Markov** with respect to G
 - Equivalently: G is an **independence map (I-map)** of P

Why is the Markov property important?

- The Markov property connects a distribution and a DAG
- In a Bayesian network (G, p) , the Markov property holds by definition
- If X_V is generated from a SEM with DAG $G = (V, E)$, then the distribution of X_V is Markov with respect to G
 - Regardless of the choice of the structural equations

Minimal I-map

- Every distribution is Markov with respect to a full DAG
 - Equivalently: a full DAG is an I-map of any distribution
- We are not interested in such DAGs, but in DAGs that are in some sense sparse. This leads to the following definition:
- **Definition:** A DAG $G = (V, E)$ is a **minimal I-map** of a distribution P if:
 - G is an I-map of P , and
 - $G' = (V, E')$ with $E' \subset E$ is not an I-map of P
 - P is then said to satisfy **causal minimality** with respect to G

Minimal I-map

- Minimal I-maps can be easily constructed (see week 2):
 - Take any ordering of the variables
 - Write out the corresponding full factorization
 - Simplify the terms as much as possible and draw the corresponding DAG

Example

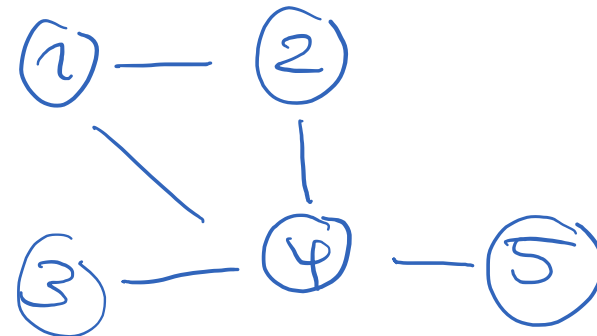
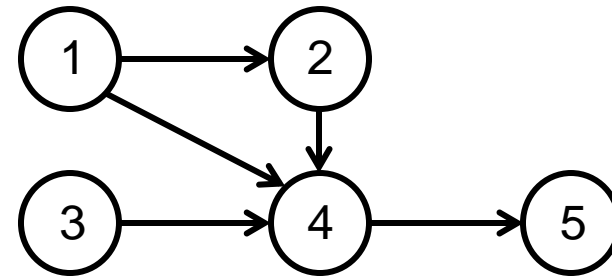
- Consider (X_1, X_2, X_3) and suppose that $X_1 \perp\!\!\!\perp X_3 | X_2$ is the only (conditional) independence:

$$f(x_1|x_2, x_3) = f(x_1|x_2) \text{ and } f(x_3|x_1, x_2) = f(x_3|x_2)$$

- Then $f(x_1, x_2, x_3) = f(x_1)f(x_2|x_1)f(x_3|x_1, x_2) = f(x_1)f(x_2|x_1)f(x_3|x_2)$
- DAG: $1 \rightarrow 2 \rightarrow 3$

Graph terminology

- The **skeleton** of G does not take the directions of the edges into account: it is the graph $G = (V, \tilde{E})$ with $(i, j) \in \tilde{E}$ if $i \rightarrow j$ or $j \rightarrow i$ in G



Clicker question – I-maps and minimal I-maps

- Which of the following statements are correct?
 - A distribution does not have a unique I-map.
 - A distribution does not have a unique minimal I-map.
 - Different minimal I-maps of a distribution can have different skeletons.
- We will discuss the answers at hand of the following example:
Consider (X_1, X_2, X_3) and suppose that $X_1 \perp\!\!\!\perp X_2 | X_3$ is the only (conditional) independence.

Faithfulness and perfect maps

- Given a DAG $G = (V, E)$, a distribution P on X_V is said to be **faithful** with respect to G if for all pairwise disjoint subsets A, B and S of V :

$$X_A \perp\!\!\!\perp X_B | X_S \text{ in } P \Rightarrow A \text{ and } B \text{ are d-separated by } S \text{ in } G$$

- If a distribution P is **Markov and faithful** with respect to a DAG G , then G is said to be a **perfect map** of P . In this case, we have that for all pairwise disjoint subsets A, B and S of V :

$$X_A \perp\!\!\!\perp X_B | X_S \text{ in } P \Leftrightarrow A \text{ and } B \text{ are d-separated by } S \text{ in } G$$

Faithfulness and perfect maps

- If a distribution P is **Markov and faithful** with respect to a DAG G , then G is said to be a **perfect map** of P . In this case, we have that for all pairwise disjoint subsets A, B and S of V :

$$X_A \perp\!\!\!\perp X_B | X_S \text{ in } P \iff A \text{ and } B \text{ are d-separated by } S \text{ in } G$$

- Combination of the Markov and faithfulness assumptions creates one-to-one link between d-separation in the DAG and conditional independence in P
- This will turn out to be very convenient for structure learning
- Not every distribution has a perfect map

Examples

Perfect maps and Markov equivalence

- **Definition:** Two DAGs G_1 and G_2 are **Markov equivalent** if they describe the same set of d-separation relationships, i.e., for all pairwise disjoint subsets A, B and S of V , we have:

A and B are d-separated by S in $G_1 \iff A$ and B are d-separated by S in G_2

- A perfect map (if it exists) is unique up to Markov equivalence

Recap

- Concepts to know:
 - Inverse probability weighting
 - Markov properties
 - Causal minimality
 - Faithfulness
 - Markov equivalence

References and acknowledgments

- Slides adapted from M. Maathuis
- Inverse probability weighting
 - Peters, Janzing and Schölkopf (2017). Elements of Causal Inference. Chapter 8.2.1
- Markov properties, faithfulness and causal minimality
 - Peters, Janzing and Schölkopf (2017). Elements of Causal Inference. Chapter 6.5
- Optional reading:
 - Pearl (2009). Causality: Models, Reasoning and Inference. Chapter 1.
 - Lauritzen (1996). Graphical Models. Chapter 3.