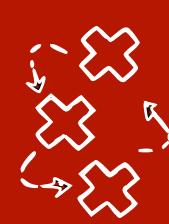


# Causal Data Science

## Lecture 8.2: Introduction to causal discovery

Lecturer: Sara Magliacane

UvA - Spring 2024



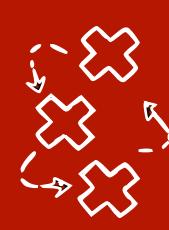
# Overview on where we are in the course

1	Introduction
2	Probability recap
3	Causal graphs
4	Interventions
5	Covariate adjustment
6	Frontdoor criterion, Instrumental variables
7	Counterfactuals and potential outcomes
8	Estimating causal effects, <a href="#">Missing data</a>
9	Constraint based structure learning
10	Score based structure learning
11	<a href="#">Advanced structure learning and transportability</a>
12	<a href="#">Causality-inspired ML</a>

Background on  
causal graphs

We know the causal  
graph, how do we  
estimate causal effects?

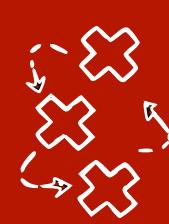
What happens if the  
graph is unknown?



# Common assumptions

- If  $P$  is **Markov** and **faithful** to  $G$ , then for any disjoint  $\mathbf{A}, \mathbf{B}, \mathbf{C} \subseteq \mathbf{V}$ :

$$\mathbf{A} \perp_G \mathbf{B} | \mathbf{C} \iff X_{\mathbf{A}} \perp\!\!\!\perp_P X_{\mathbf{B}} | X_{\mathbf{C}}$$

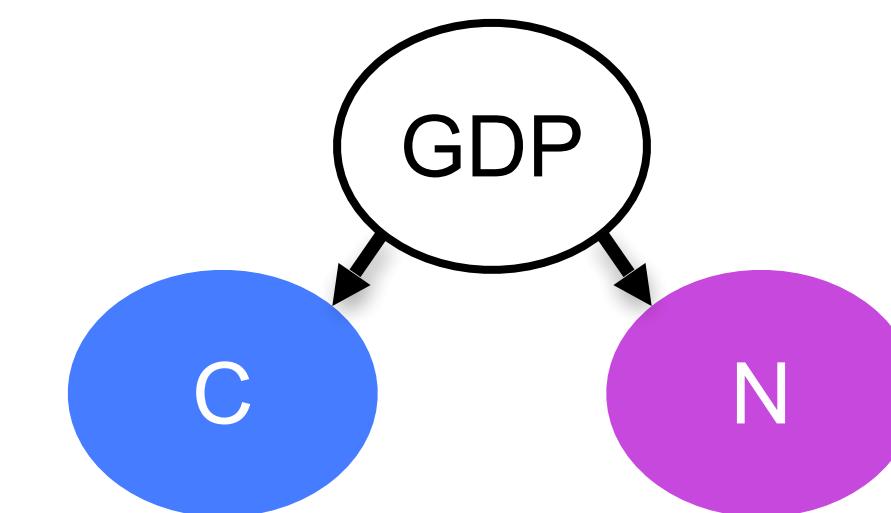
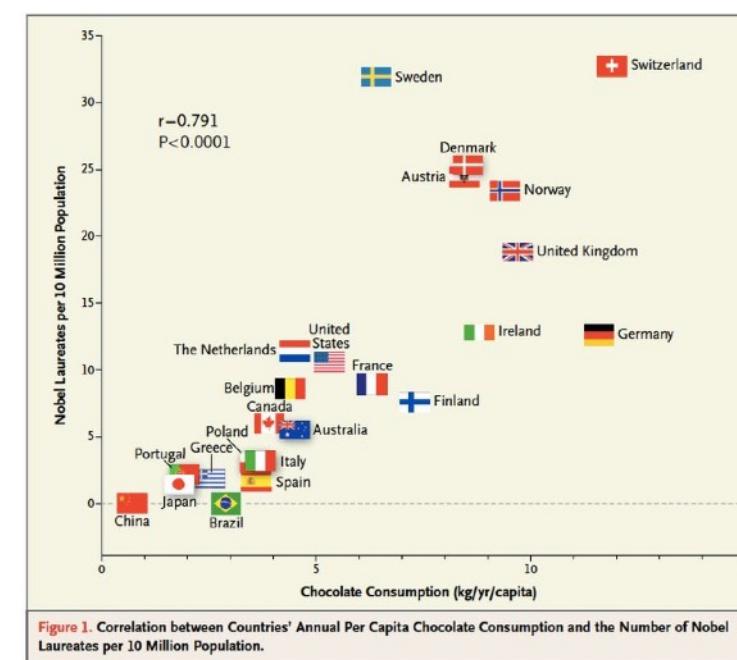


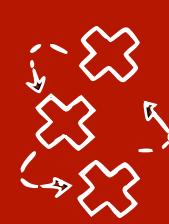
# Common assumptions

- If  $P$  is **Markov** and **faithful** to  $G$ , then for any disjoint  $A, B, C \subseteq V$ :

$$A \perp_G B | C \iff X_A \perp_P X_B | X_C$$

- **Causal sufficiency** - no **latent** confounders (common causes), no selection bias



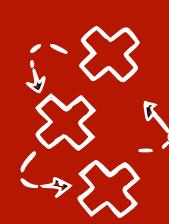


# Common assumptions

- If  $P$  is **Markov** and **faithful** to  $G$ , then for any disjoint  $\mathbf{A}, \mathbf{B}, \mathbf{C} \subseteq \mathbf{V}$ :

$$\mathbf{A} \perp_G \mathbf{B} | \mathbf{C} \iff X_{\mathbf{A}} \perp_P X_{\mathbf{B}} | X_{\mathbf{C}}$$

- **Causal sufficiency** - no latent confounders (common causes), no selection bias
- **Acyclicity** - the underlying graph is acyclic
- Cycles + causal insufficiency: sigma separation, Joint Causal Inference



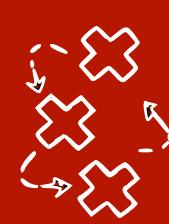
# Recap: Factorizing joint distributions

- Given any **ordering** of the variables  $(X_1, \dots, X_p)$  we can write:

$$P(X_1, \dots, X_p) = P(X_1)P(X_2 | X_1)\dots P(X_p | X_1, \dots, X_{p-1})$$

- For example  $P(X, Y, Z)$  can be equivalently factorized as:

- $P(X, Y, Z) = P(X)P(Y | X)P(Z | X, Y)$
- $P(X, Z, Y) = P(X)P(Z | X)P(Y | X, Z)$
- $P(Z, Y, X) = P(Z)P(Y | Z)P(X | Y, Z) \dots$



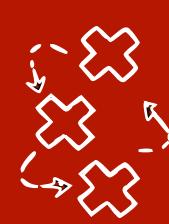
# Recap: Exploiting conditional independences

- We can **simplify** the factorisation by using **conditional independences**:

$$X_i \perp\!\!\!\perp X_j | X_Z \implies P(X_i | X_j, X_Z) = P(X_i | X_Z)$$

- For example:  $X \perp\!\!\!\perp Y | Z$ :

- $P(X, Y, Z) = P(X)P(Y | X)P(Z | X, Y)$
- $P(X, Z, Y) = P(X)P(Z | X)P(Y | X, Z)$
- $P(Z, Y, X) = P(Z)P(Y | Z)P(X | Y, Z) \dots$



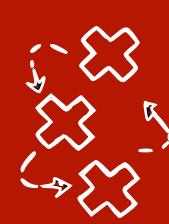
# Recap: Exploiting conditional independences

- We can **simplify** the factorisation by using **conditional independences**:

$$X_i \perp\!\!\!\perp X_j | X_Z \implies P(X_i | X_j, X_Z) = P(X_i | X_Z)$$

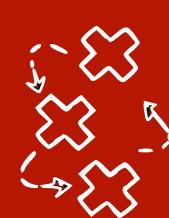
- For example:  $X \perp\!\!\!\perp Y | Z$ :

- $P(X, Y, Z) = P(X)P(Y | X)P(Z | X, Y)$
- $P(X, Z, Y) = P(X)P(Z | X)P(Y | \cancel{X}, Z) = P(X)P(Z | X)P(Y | Z)$
- $P(Z, Y, X) = P(Z)P(Y | Z)P(X | \cancel{X}, Z) = P(Z)P(Y | Z)P(X | Z)$



# Recap: Bayesian networks

- We have a set of random variables  $X_1, \dots, X_p$  with joint  $P(X_1, \dots, X_p)$
- We have a DAG  $G$ , s.t. **each random variable  $X_i$  is represented by node  $i$**
- We then say  $P(X_1, \dots, X_p)$  **factorizes over  $G$**  if
$$P(X_1, \dots, X_p) = \prod_{i \in V} P(X_i | \mathbf{X}_{\text{Pa}_G(i)})$$
- A **Bayesian network** (BN) is the tuple  $(G, P)$  s.t.  **$P$  factorizes over  $G$**



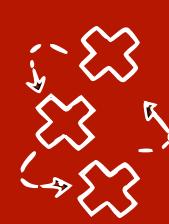
# Global Markov Property

- If  $(G, P)$  is a Bayesian network with a DAG  $G = (V, E)$ , i.e.  **$P$  factorizes according to  $G$** , then for any disjoint  $A, B, C \subseteq V$ :

$$A \perp_G B | C \implies X_A \perp_P X_B | X_C$$

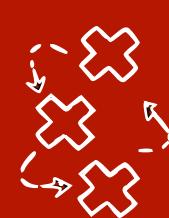
If  $P$  has a density (e.g.  
no deterministic  
relations)  
[Lauritzen 1996]

- **d-separations** that can be read purely from a graph imply **conditional independences** in the random variables and data generated by the graph
- In general some of the independences in the data are not represented in the graph, i.e. we have more edges that create “extra d-connections”



# Terminology - I-maps

- $(G, P)$  is a Bayesian network with a DAG  $G = (V, E)$ . We say that  $G$  is **an I-map of  $P$  (independence map)**, for any disjoint  $A, B, C \subseteq V$ :
  - $A \perp_G B | C \implies X_A \perp\!\!\!\perp_P X_B | X_C$  (*Global Markov Property*)



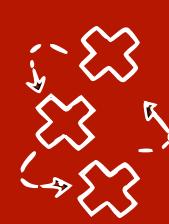
# Terminology - I-maps

- $(G, P)$  is a Bayesian network with a DAG  $G = (V, E)$ . We say that  $G$  is **an I-map of  $P$  (independence map)**, for any disjoint  $A, B, C \subseteq V$ :
  - $A \perp_G B | C \implies X_A \perp\!\!\!\perp_P X_B | X_C$  (*Global Markov Property*)

Many distributions factorize according to many graphs, e.g. any  $P$  factorizes according to all fully connected graphs

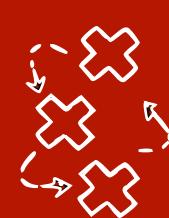
We want to narrow down the candidate graphs to a smaller set containing the true causal graph

What about sparsity/minimality?



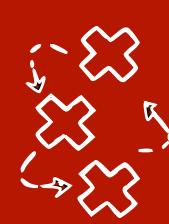
# Terminology - I-maps

- $(G, P)$  is a Bayesian network with a DAG  $G = (V, E)$ . We say that  **$G$  is an I-map of  $P$  (independence map)**, for any disjoint  $A, B, C \subseteq V$ :
  - $A \perp_G B | C \implies X_A \perp\!\!\!\perp_P X_B | X_C$  (*Global Markov Property*)
- **Causal minimality/Minimal I-map:** an I-map in which we cannot remove edges, otherwise it stops being an I-map



# Terminology - I-maps

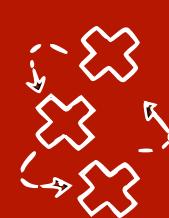
- $(G, P)$  is a Bayesian network with a DAG  $G = (V, E)$ . We say that  **$G$  is an I-map of  $P$**  (independence map), for any disjoint  $A, B, C \subseteq V$ :
  - $A \perp_G B | C \implies X_A \perp\!\!\!\perp_P X_B | X_C$  (*Global Markov Property*)
- **Causal minimality/Minimal I-map:** an I-map in which we cannot remove edges, otherwise it stops being an I-map
  - $G = (V, E)$  is an I-map of  $P$
  - $G' = (V, E')$  such that  $E' \subset E$  is **not** an I-map of  $P$   
 $\exists A, B, C \text{ s.t. } A \perp_{G'} B | C \text{ and } X_A \perp\!\!\!\perp_P X_B | X_C$



# Constructing a Minimal I-Map

1. Choose any **ordering** of the variables  $(X_1, \dots, X_p)$
2. Write the factorisation via chain rule, **simplify** factorisation as much as possible with **conditional independences**:

$$X_i \perp\!\!\!\perp X_j | X_{\mathbf{Z}} \implies P(X_i | X_j, X_{\mathbf{Z}}) = P(X_i | X_{\mathbf{Z}})$$

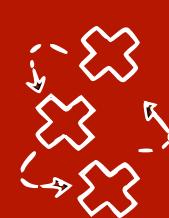


# Constructing a Minimal I-Map

1. Choose any **ordering** of the variables  $(X_1, \dots, X_p)$
2. Write the factorisation via chain rule, **simplify** factorisation as much as possible with **conditional independences**:

$$X_i \perp\!\!\!\perp X_j | X_{\mathbf{Z}} \implies P(X_i | X_j, X_{\mathbf{Z}}) = P(X_i | X_{\mathbf{Z}})$$

3. **Create a BN** in which each  $X_i$  is caused by the variables in  $P(X_i | X_{\text{Pa}(i)})$



# Constructing a Minimal I-Map

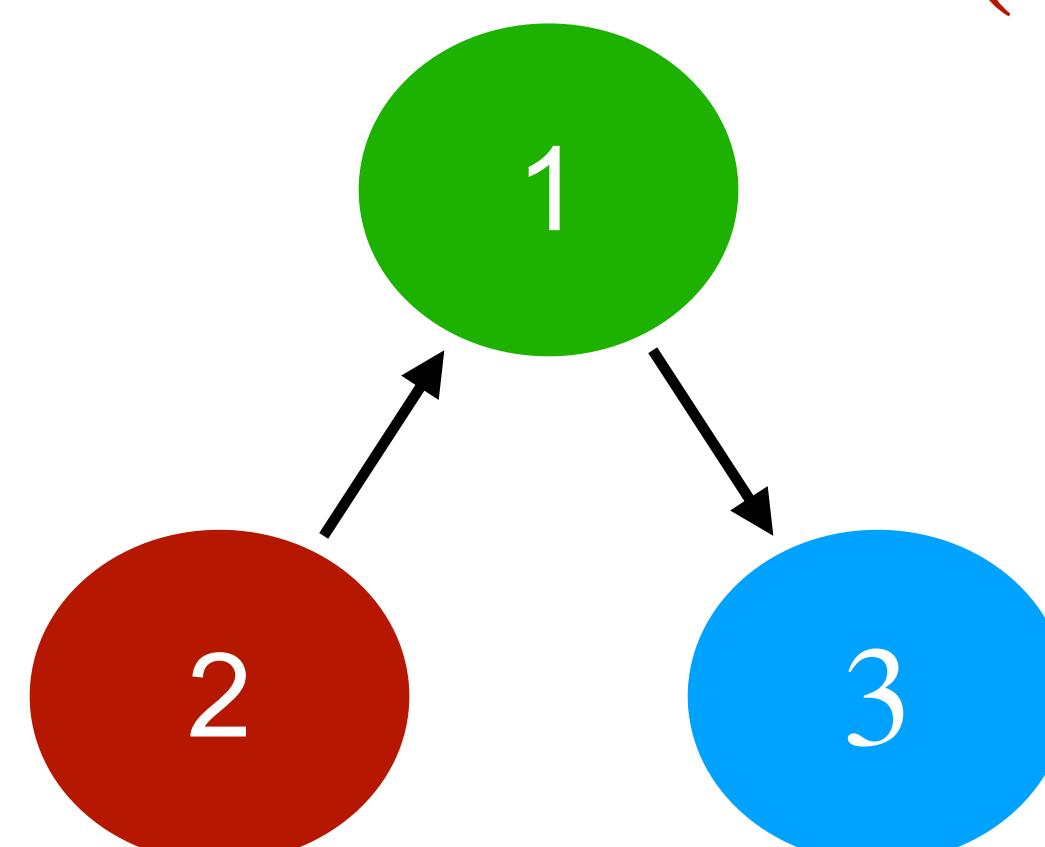
1. Choose any **ordering** of the variables  $(X_1, \dots, X_p)$
2. Write the factorisation via chain rule, **simplify** factorisation as much as possible with **conditional independences**:

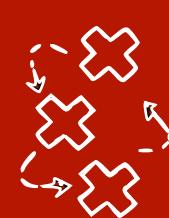
$$X_i \perp\!\!\!\perp X_j | X_{\mathbf{Z}} \implies P(X_i | X_j, X_{\mathbf{Z}}) = P(X_i | X_{\mathbf{Z}})$$

3. Create a BN in which each  $X_i$  is caused by the variables in  $P(X_i | X_{\text{Pa}(i)})$

If  $X_3 \perp\!\!\!\perp X_2 | X_1$

$$P(X_2, X_1, X_3) = P(X_2)P(X_1 | X_2)P(X_3 | X_1, \cancel{X_2})$$





# Constructing a Minimal I-Map (not unique!)

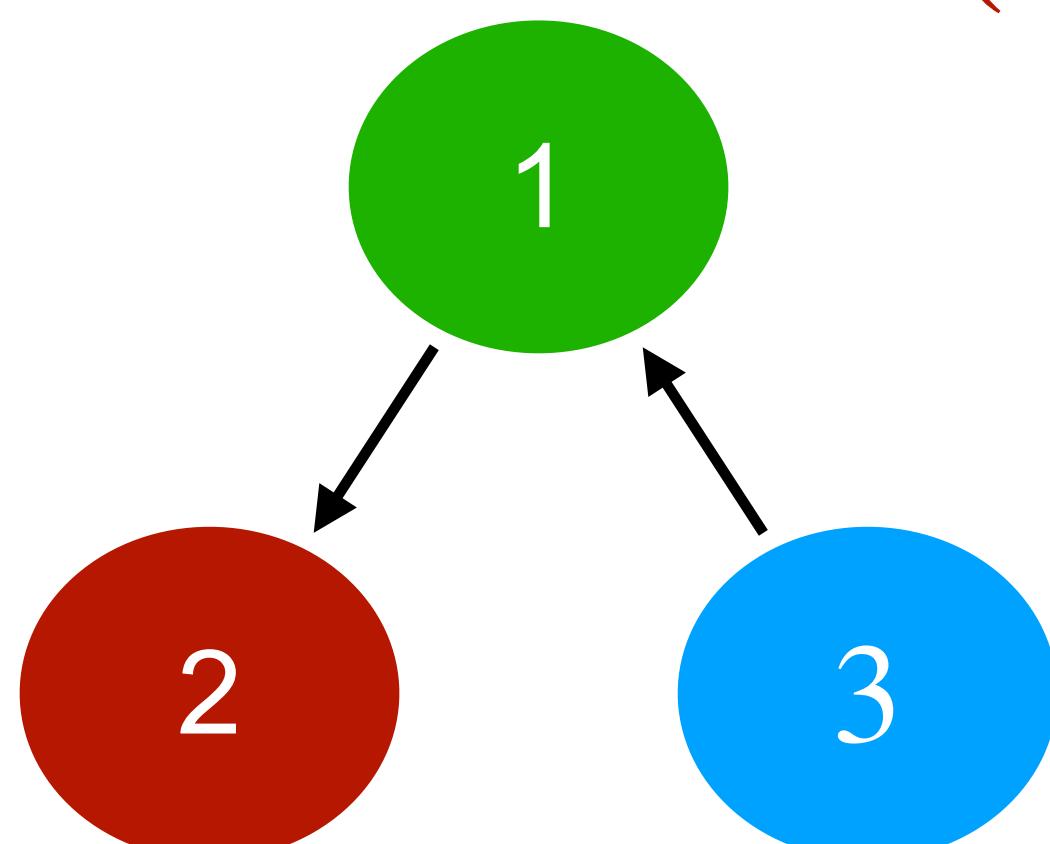
1. Choose any **ordering** of the variables  $(X_1, \dots, X_p)$
2. Write the factorisation via chain rule, **simplify** factorisation as much as possible with **conditional independences**:

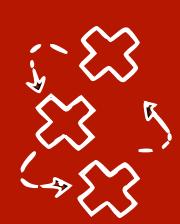
$$X_i \perp\!\!\!\perp X_j | X_{\mathbf{Z}} \implies P(X_i | X_j, X_{\mathbf{Z}}) = P(X_i | X_{\mathbf{Z}})$$

3. Create a BN in which each  $X_i$  is caused by the variables in  $P(X_i | X_{\text{Pa}(i)})$

If  $X_3 \perp\!\!\!\perp X_2 | X_1$

$$P(X_3, X_1, X_2) = P(X_3)P(X_1 | X_3)P(X_2 | X_1, \cancel{X_3})$$

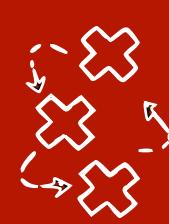




# Minimal I-Map example

- Modified classic example: Discipline (D), Intelligence (I), SAT score (S), Grade (G)

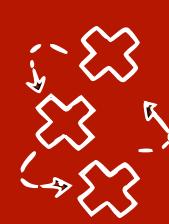
$$\begin{array}{lll} p(D, I, S, G) & D \perp\!\!\!\perp I & \\ & D \perp\!\!\!\perp G & D \perp\!\!\!\perp I | G \quad D \perp\!\!\!\perp G | I, S \\ & G \perp\!\!\!\perp S | I & D \perp\!\!\!\perp G | I \quad S \perp\!\!\!\perp G | I, D \end{array}$$



# Minimal I-Map example

- Modified classic example: Discipline (D), Intelligence (I), SAT score (S), Grade (G)

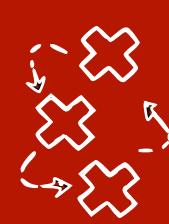
$$\begin{array}{lll} p(D, I, S, G) & D \perp\!\!\!\perp I & p(D) \cdot p(I|D) \cdot p(S|I,D) \cdot p(G|S,I,D) \\ & D \perp\!\!\!\perp G & D \perp\!\!\!\perp I | G \quad D \perp\!\!\!\perp G | I, S \\ & G \perp\!\!\!\perp S | I & D \perp\!\!\!\perp G | I \quad S \perp\!\!\!\perp G | I, D \end{array}$$



# Minimal I-Map example

- Modified classic example: Discipline (D), Intelligence (I), SAT score (S), Grade (G)

$$\begin{array}{lll} p(D, I, S, G) & D \perp\!\!\!\perp I & p(D) \cdot p(I|D) \cdot p(S|I,D) \cdot p(G|S,I,D) \\ & D \perp\!\!\!\perp G & D \perp\!\!\!\perp I | G \quad D \perp\!\!\!\perp G | I, S \\ & G \perp\!\!\!\perp S | I & D \perp\!\!\!\perp G | I \quad S \perp\!\!\!\perp G | I, D \\ & & p(D) \cdot p(I) \cdot p(S|I,D) \cdot p(G|I) \end{array}$$



# Minimal I-Map example

- Modified classic example: Discipline (D), Intelligence (I), SAT score (S), Grade (G)

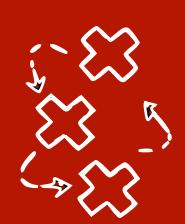
$$p(D, I, S, G) \quad D \perp\!\!\!\perp I \quad p(D) \cdot p(I|D) \cdot p(S|I,D) \cdot p(G|S,I,D)$$

$$D \perp\!\!\!\perp G \quad D \perp\!\!\!\perp I|G \quad D \perp\!\!\!\perp G|I,S$$

$$G \perp\!\!\!\perp S|I \quad D \perp\!\!\!\perp G|I \quad S \perp\!\!\!\perp G|I,D$$

$$p(D) \cdot p(I) \cdot p(S|I,D) \cdot p(G|I)$$

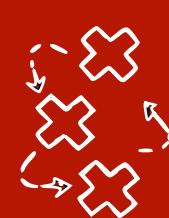




# Minimal I-Map example

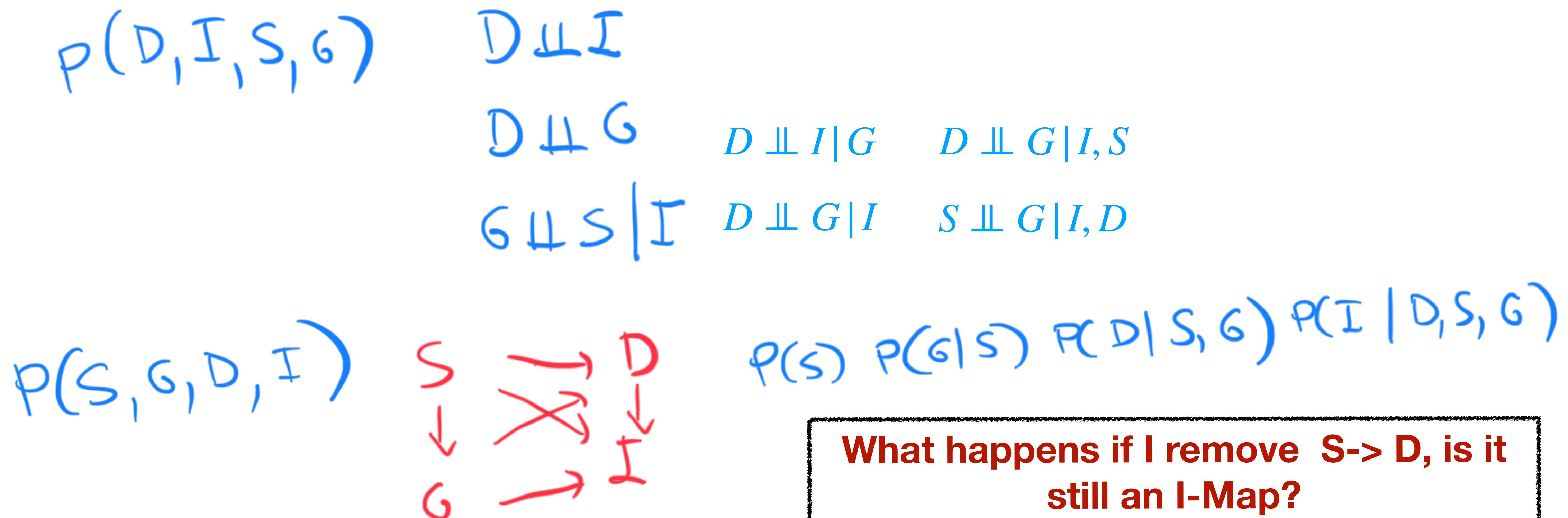
- Modified classic example: Discipline (D), Intelligence (I), SAT score (S), Grade (G)

$$\begin{array}{lll} p(D, I, S, G) & D \perp\!\!\!\perp I & \\ & D \perp\!\!\!\perp G & D \perp\!\!\!\perp I | G \quad D \perp\!\!\!\perp G | I, S \\ & G \perp\!\!\!\perp S | I & D \perp\!\!\!\perp G | I \quad S \perp\!\!\!\perp G | I, D \end{array}$$
$$p(S, G, D, I) \quad \begin{array}{c} S \rightarrow D \\ \downarrow \quad \times \quad \downarrow \\ G \rightarrow I \end{array} \quad \begin{array}{l} p(S) \quad p(G|S) \quad p(D|S, G) \quad p(I | D, S, G) \end{array}$$

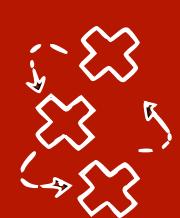


# Minimal I-Map example

- Modified classic example: Discipline (D), Intelligence (I), SAT score (S), Grade (G)



Now we have an extra d-separation  $S \perp\!\!\!\perp D | G$  not matching the conditional independences

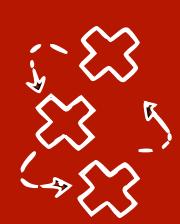


# Minimal I-Map example

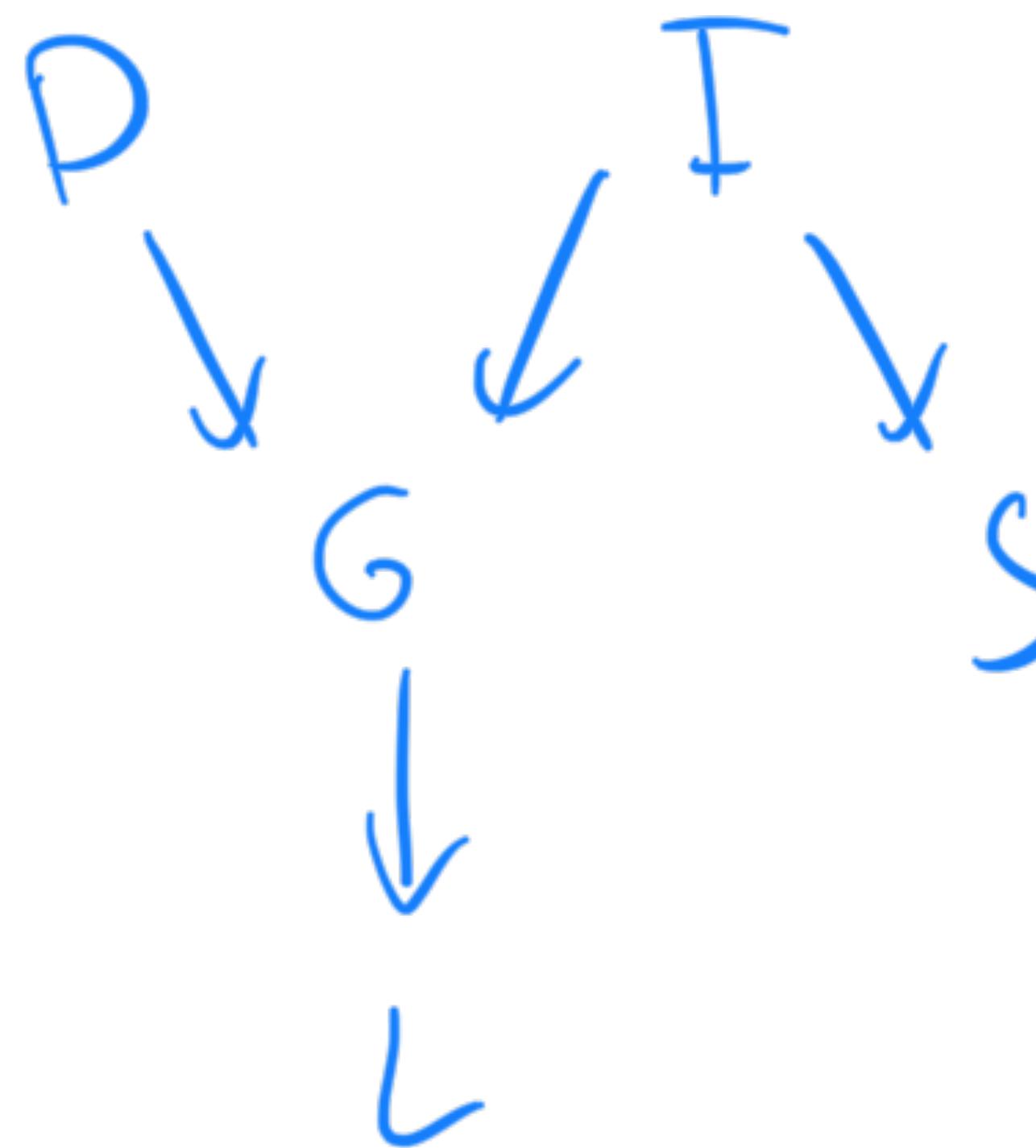
- Modified classic example: Discipline (D), Intelligence (I), SAT score (S), Grade (G)

$$\begin{array}{lll} p(D, I, S, G) & D \perp\!\!\!\perp I & \\ & D \perp\!\!\!\perp G & D \perp\!\!\!\perp I | G \quad D \perp\!\!\!\perp G | I, S \\ & G \perp\!\!\!\perp S | I & D \perp\!\!\!\perp G | I \quad S \perp\!\!\!\perp G | I, D \end{array}$$
$$p(S, G, D, I) \quad \begin{array}{c} S \rightarrow D \\ \downarrow \quad \diagup \\ G \rightarrow I \end{array} \quad \begin{array}{l} p(S) \quad p(G|S) \quad p(D|S, G) \quad p(I | D, S, G) \end{array}$$

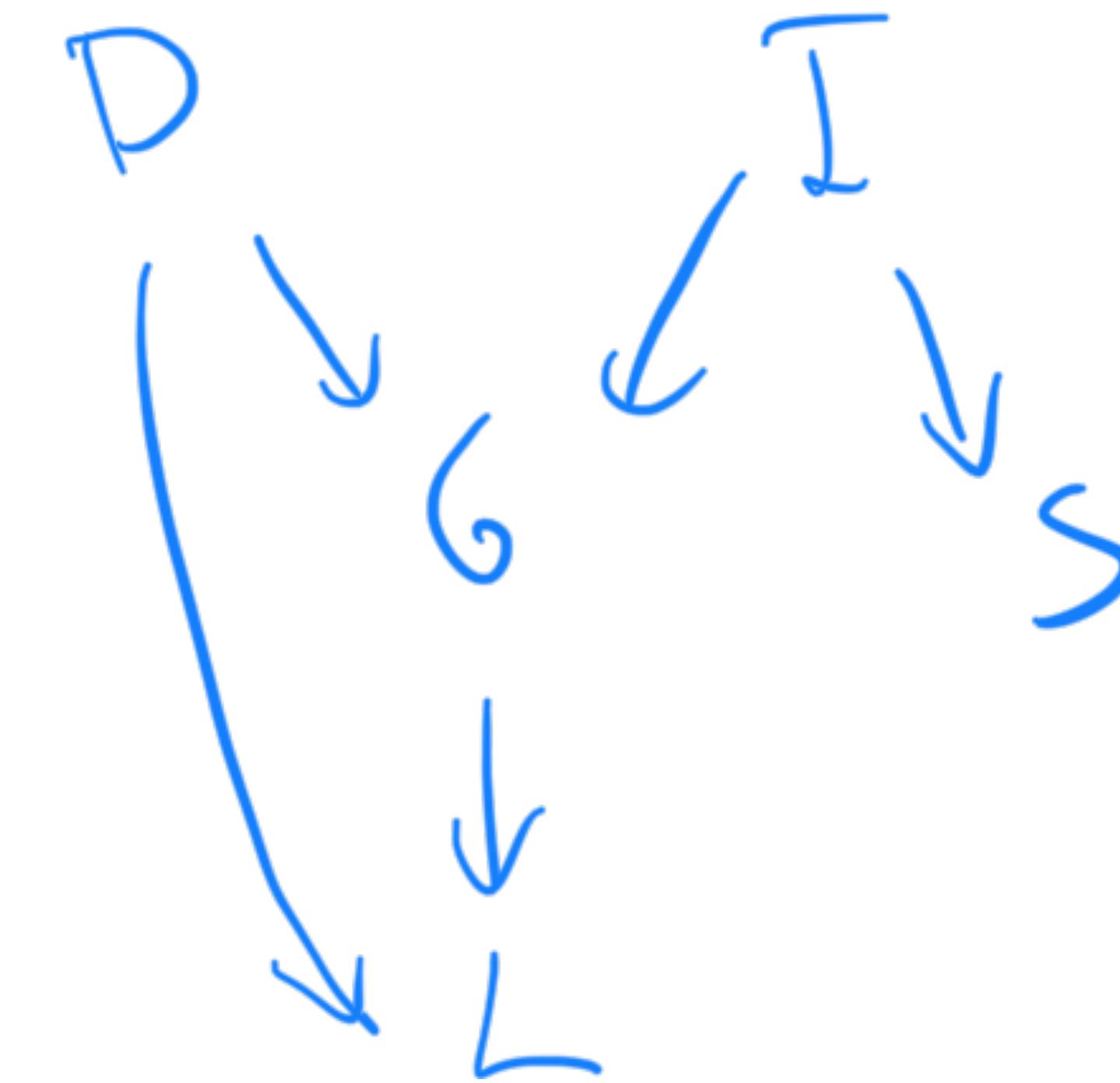
We can show that this is also a minimal I-Map  
(we cannot remove any edge without creating  
an extra d-separation not matching  
independences in the distribution)



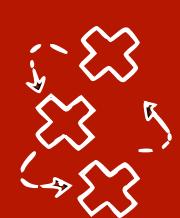
# D-separation reprise and minimal I-maps: Canvas quiz



Graph 1



Graph 2



# Minimal I-Maps vs independences:

$$P(D, I, S, G)$$

$$D \perp\!\!\!\perp I$$

$$D \perp\!\!\!\perp G$$

$$G \perp\!\!\!\perp S \mid I$$

$$D \perp\!\!\!\perp I \mid G$$

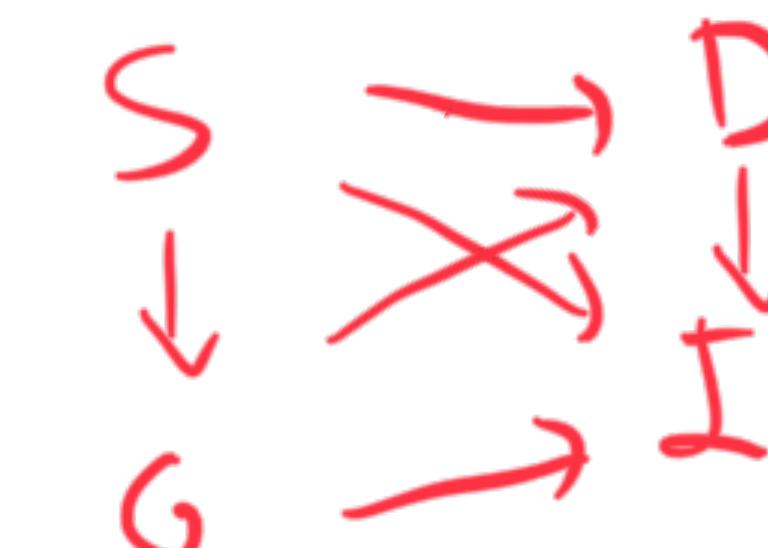
$$D \perp\!\!\!\perp G \mid I$$

$$D \perp\!\!\!\perp G \mid I, S$$

$$S \perp\!\!\!\perp G \mid I, D$$

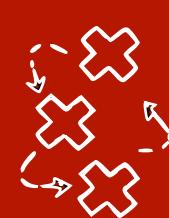


$$\begin{aligned} D \perp\!\!\!\perp_d I &\checkmark \\ D \perp\!\!\!\perp_d G &\checkmark \\ G \perp\!\!\!\perp_d S \mid I &\checkmark \end{aligned}$$



$$\begin{aligned} D \not\perp\!\!\!\perp_d I &\times \\ D \not\perp\!\!\!\perp_d G &\times \\ G \not\perp\!\!\!\perp_d S \mid I &\times \end{aligned}$$

Can we reduce the set of minimal I-maps even further, possibly making them fit “perfectly” the conditional independences?



# Global Markov Property and Faithfulness

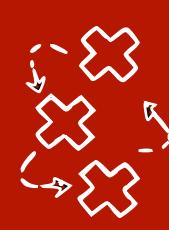
- If  $(G, P)$  is a Bayesian network with a DAG  $G = (V, E)$ , i.e.  $P$  factorizes according to  $G$ , then for any disjoint  $A, B, C \subseteq V$ :

$$A \perp B | C \implies X_A \perp\!\!\!\perp X_B | X_C$$

- The reverse implication is not true in general, but if it is  $P$  is faithful to  $G$

$$X_A \perp\!\!\!\perp X_B | X_C \implies A \perp B | C$$

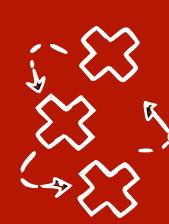
- If  $P$  is Markov and faithful to  $G$ , we say that  $G$  is a perfect map of  $P$ 
  - In general minimal maps can also not be perfect maps



# Perfect maps

- If  $P$  is Markov and faithful to  $G$ , we say that  $G$  is a **perfect map of  $P$** .  
Then, for any disjoint  $A, B, C \subseteq V$ :

$$A \perp_G B | C \iff X_A \perp\!\!\!\perp_P X_B | X_C$$

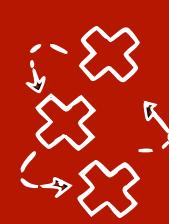


# Perfect maps - existence

- Not every distribution  $P$  has a perfect map

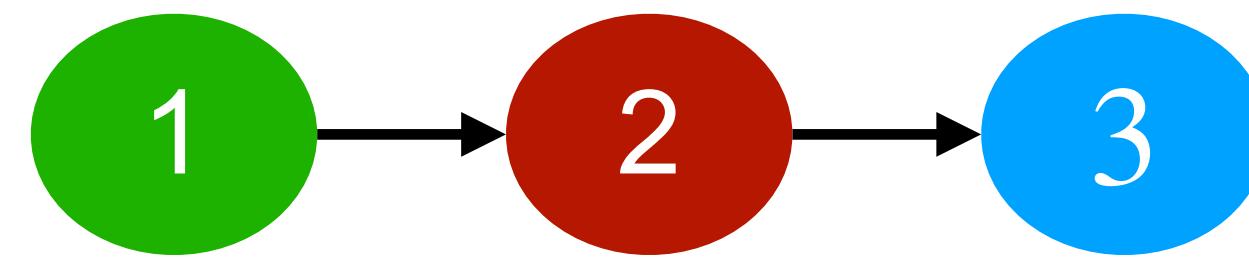


$$\begin{cases} X_1 = \epsilon_1 \\ X_2 = 3X_1 \\ X_3 = 4X_2 + \epsilon_3 \\ \epsilon_1, \epsilon_3 \sim N(0,1) \end{cases}$$

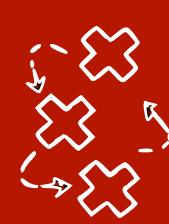


# Perfect maps - existence

- Not every distribution  $P$  has a perfect map

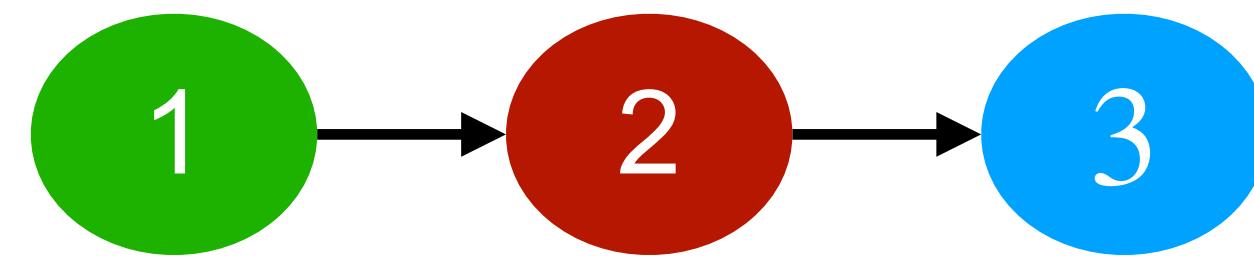


$$\left\{ \begin{array}{l} X_1 = \epsilon_1 \\ X_2 = 3X_1 \\ X_3 = 4X_2 + \epsilon_3 \\ \epsilon_1, \epsilon_3 \sim N(0,1) \end{array} \right. \quad \begin{aligned} X_2 \perp\!\!\!\perp X_3 | X_1 &\iff \forall x_1, x_2, x_3 : p(X_2 = x_2 | X_3 = x_3, X_1 = x_1) = p(X_2 = x_2 | X_1 = x_1) \\ &\iff p(X_2 = 3x_1 | X_1 = x_1, X_3 = x_3) = p(X_2 = 3x_1 | X_1 = x_1) \end{aligned}$$



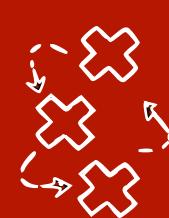
# Perfect maps - existence

- Not every distribution  $P$  has a perfect map



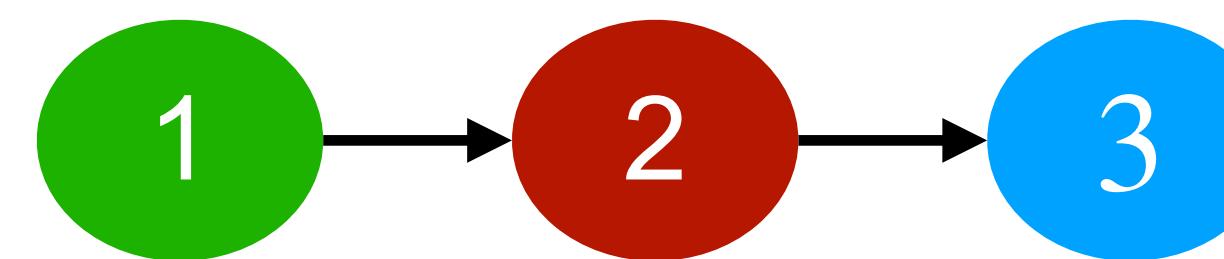
$$\begin{cases} X_1 = \epsilon_1 \\ X_2 = 3X_1 \\ X_3 = 4X_2 + \epsilon_3 \\ \epsilon_1, \epsilon_3 \sim N(0,1) \end{cases} \quad \begin{aligned} X_2 \perp\!\!\!\perp X_3 | X_1 &\iff \forall x_1, x_2, x_3 : p(X_2 = x_2 | X_3 = x_3, X_1 = x_1) = p(X_2 = x_2 | X_1 = x_1) \\ &\iff p(X_2 = 3x_1 | X_1 = x_1, X_3 = x_3) = p(X_2 = 3x_1 | X_1 = x_1) \end{aligned}$$

$X_2 \perp\!\!\!\perp X_3 | X_1$



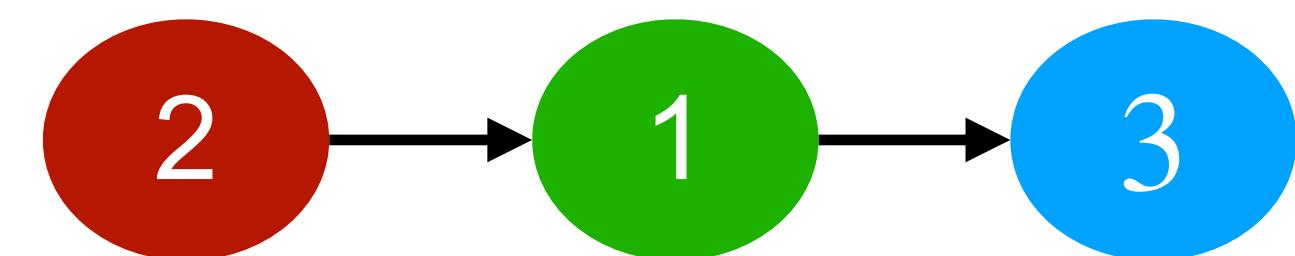
# Perfect maps - existence

- Not every distribution  $P$  has a perfect map



$$\begin{cases} X_1 = \epsilon_1 \\ X_2 = 3X_1 \\ X_3 = 4X_2 + \epsilon_3 \\ \epsilon_1, \epsilon_3 \sim N(0,1) \end{cases} \quad \begin{aligned} X_2 \perp\!\!\!\perp X_3 | X_1 &\iff \forall x_1, x_2, x_3 : p(X_2 = x_2 | X_3 = x_3, X_1 = x_1) = p(X_2 = x_2 | X_1 = x_1) \\ &\iff p(X_2 = 3x_1 | X_1 = x_1, X_3 = x_3) = p(X_2 = 3x_1 | X_1 = x_1) \\ X_2 \perp\!\!\!\perp X_3 | X_1 \end{aligned}$$

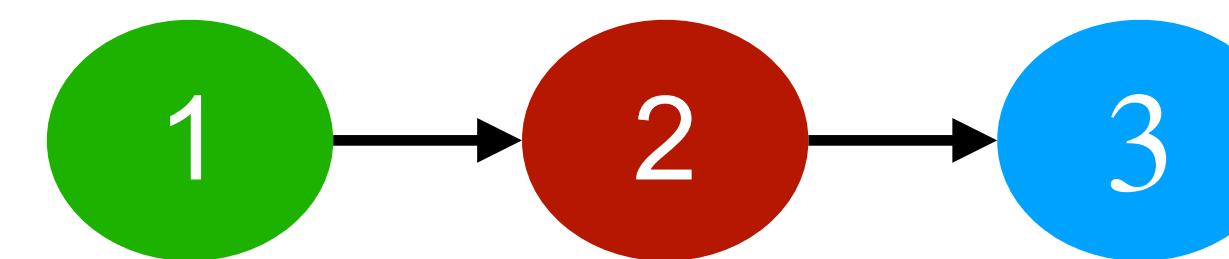
$$X_2 \perp\!\!\!\perp X_3 \quad X_1 \perp\!\!\!\perp X_3 \quad X_1 \perp\!\!\!\perp X_3 | X_2$$





# Perfect maps - existence

- Not every distribution  $P$  has a perfect map



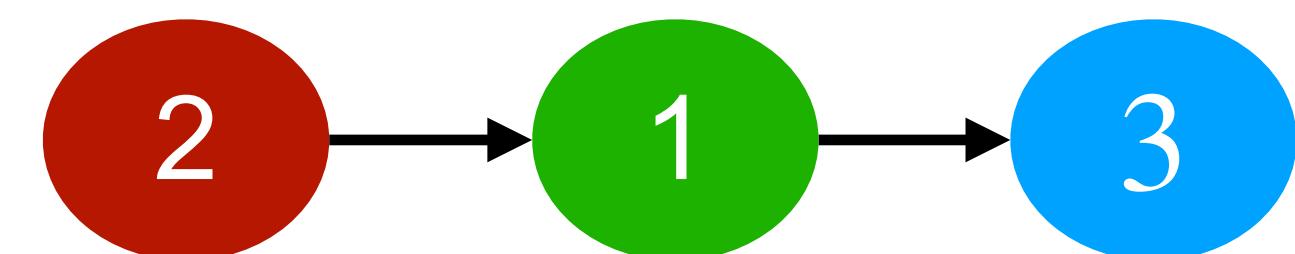
**There doesn't exist a DAG that is a perfect map for  $P$**

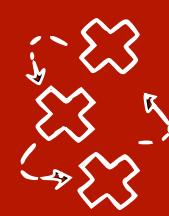
$$\begin{cases} X_1 = \epsilon_1 \\ X_2 = 3X_1 \\ X_3 = 4X_2 + \epsilon_3 \\ \epsilon_1, \epsilon_3 \sim N(0,1) \end{cases}$$

$$\begin{aligned} X_2 \perp\!\!\!\perp X_3 | X_1 &\iff \forall x_1, x_2, x_3 : p(X_2 = x_2 | X_3 = x_3, X_1 = x_1) = p(X_2 = x_2 | X_1 = x_1) \\ &\iff p(X_2 = 3x_1 | X_1 = x_1, X_3 = x_3) = p(X_2 = 3x_1 | X_1 = x_1) \end{aligned}$$

$$X_2 \not\perp\!\!\!\perp X_3 | X_1$$

$$X_2 \not\perp\!\!\!\perp X_3 \quad X_1 \not\perp\!\!\!\perp X_3 \quad X_1 \perp\!\!\!\perp X_3 | X_2$$



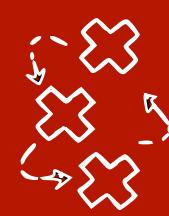


# Perfect maps - Markov equivalence

- If  $P$  is Markov and faithful to  $G$ , we say that  $G$  is a **perfect map of  $P$** .  
Then, for any disjoint  $A, B, C \subseteq V$ :

$$A \perp B | C \iff X_A \perp\!\!\!\perp X_B | X_C$$

- In general there are multiple DAGs that can describe the same d-separations (and independences), i.e. there are multiple perfect maps
- We call these DAGs **Markov equivalent** and we **cannot distinguish them from observational data alone** (or without further assumptions)



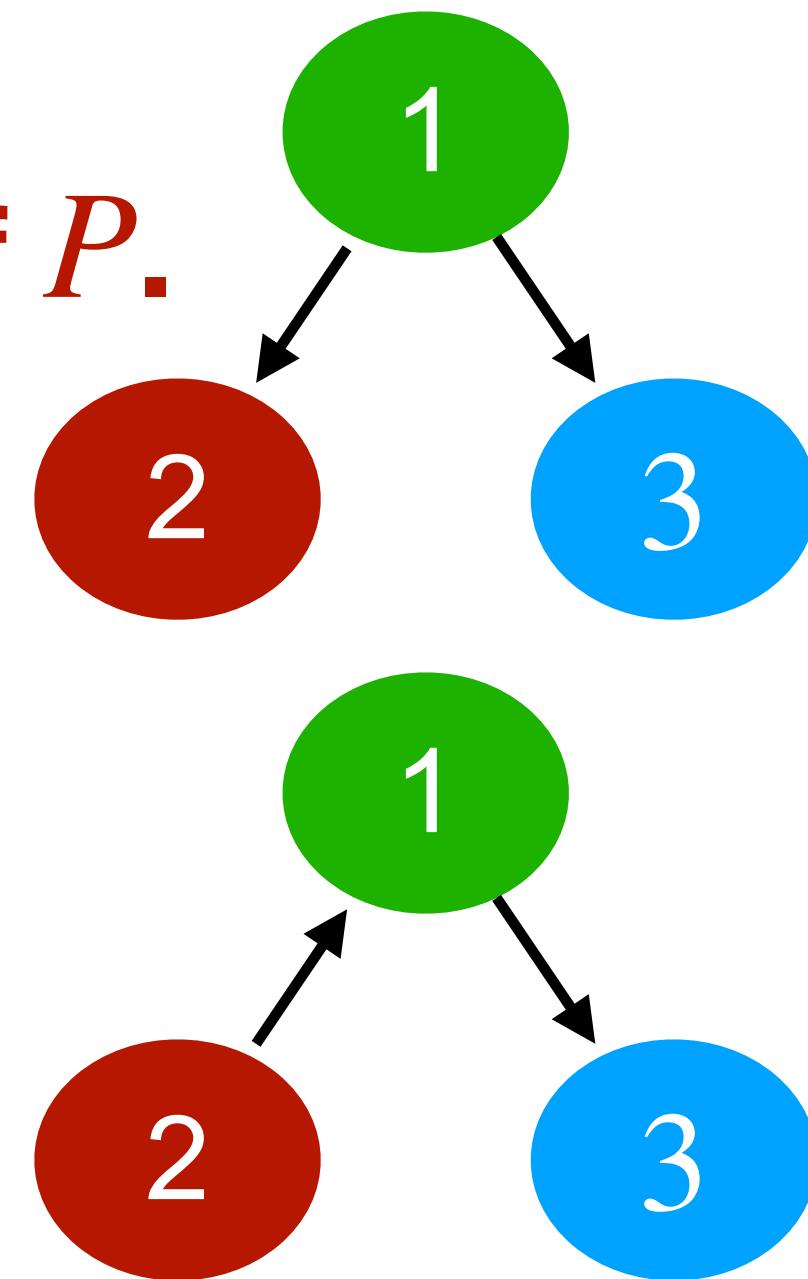
# Perfect maps - Markov equivalence

- If  $P$  is Markov and faithful to  $G$ , we say that  $G$  is a **perfect map of  $P$** .

Then, for any disjoint  $A, B, C \subseteq V$ :

$$A \perp B | C \iff X_A \perp\!\!\!\perp X_B | X_C$$

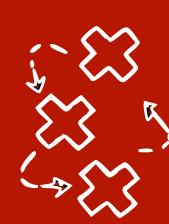
- In general there are multiple DAGs that can describe the same d-separations (and independences), i.e. there are multiple perfect maps
- We call these DAGs **Markov equivalent** and we **cannot distinguish them from observational data alone** (or without further assumptions)





# Next class - constraint-based causal discovery

- **In a nutshell:** we perform a set of conditional independence tests on the data and use them to constrain the possible graphs using d-separation
- In general, we can narrow down the possible graphs only up to their **Markov equivalence class (MEC)**
- The output of the algorithms we will see (e.g. SGS, PC) is a **CPDAG**, a mixed graph in which directed edges represent causal relations on which all DAGs in the MEC agree - these relations are **identifiable**



# Questions?

