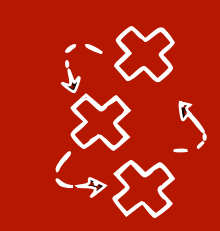


Constraint-based causal discovery

dr. Sara Magliacane (University of Amsterdam, MIT-IBM Watson AI Lab)

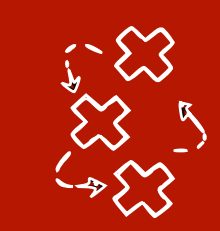




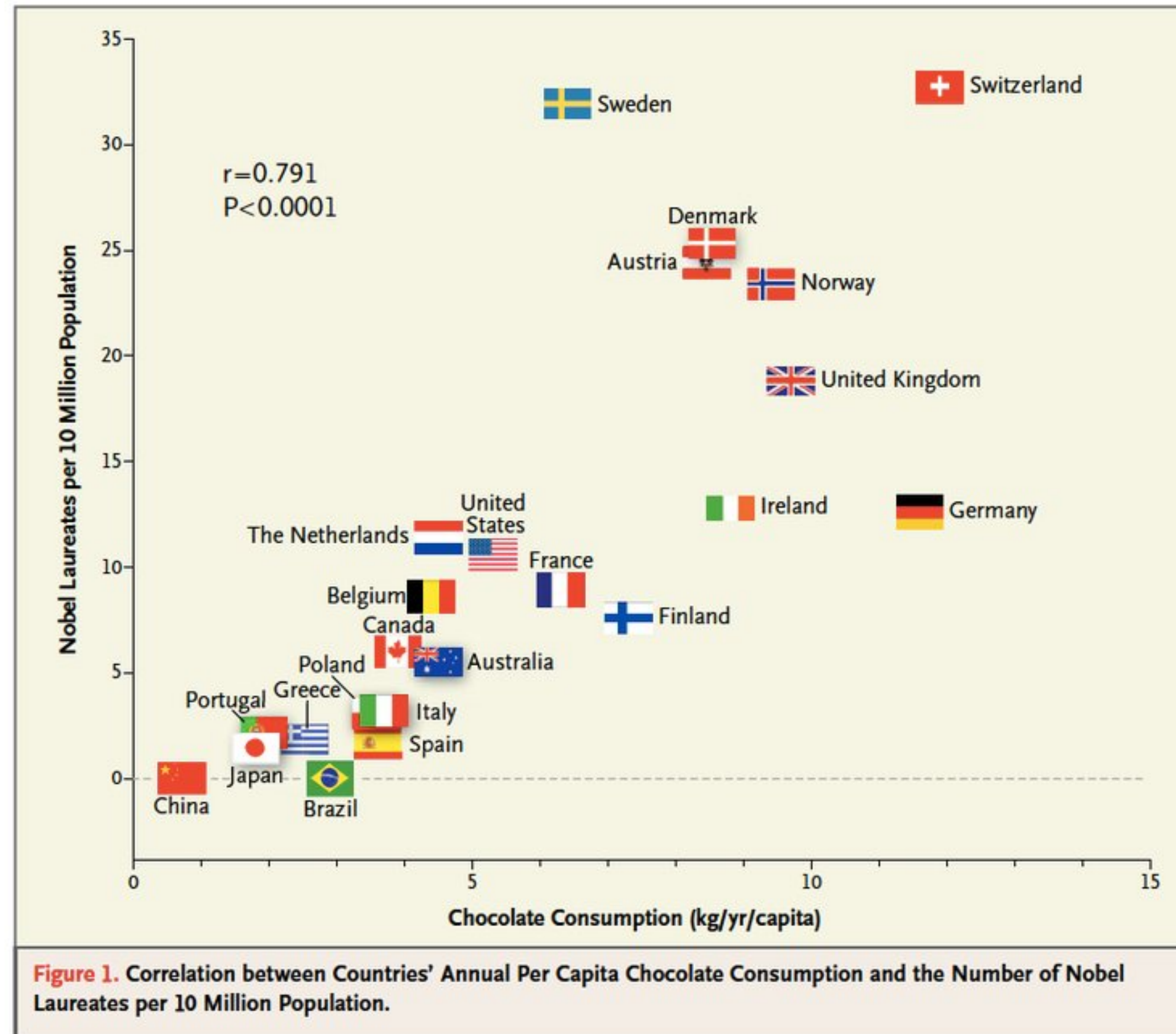
This class

- Introduction to causal discovery
 - Common assumptions: causal sufficiency, acyclicity, faithfulness
 - Constraint-based causal discovery on observational data (causal sufficiency)
 - SGS, PC
-
- Learning from multiple contexts or interventional data
 - Invariant Causal Prediction
 - Joint Causal Inference

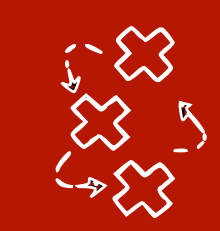
Inspired by <https://stat.ethz.ch/lectures/ss21/causality.php>



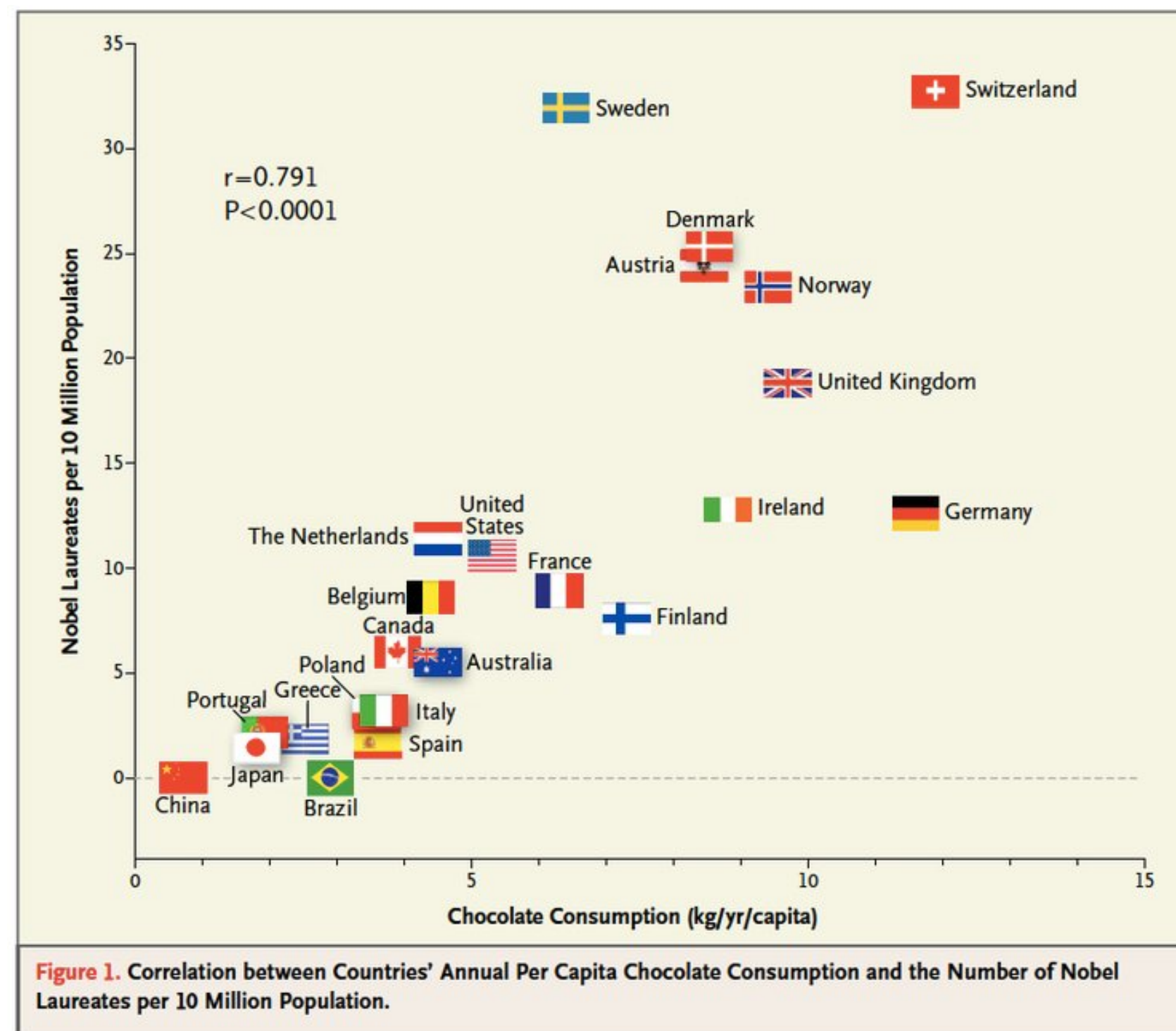
Does chocolate cause Nobel prizes?



[Messerli, 2012] <https://www.nejm.org/doi/full/10.1056/NEJMon1211064>



What is the causal graph here?



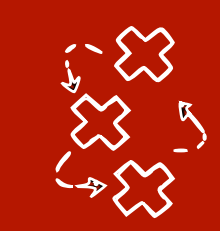
Reichenbach's principle of common cause:

A correlation between X and Y implies that

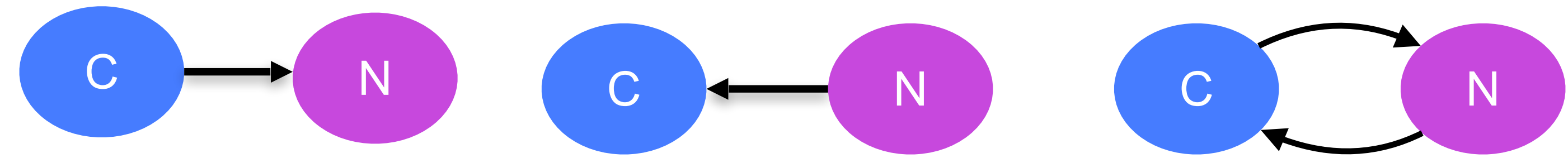
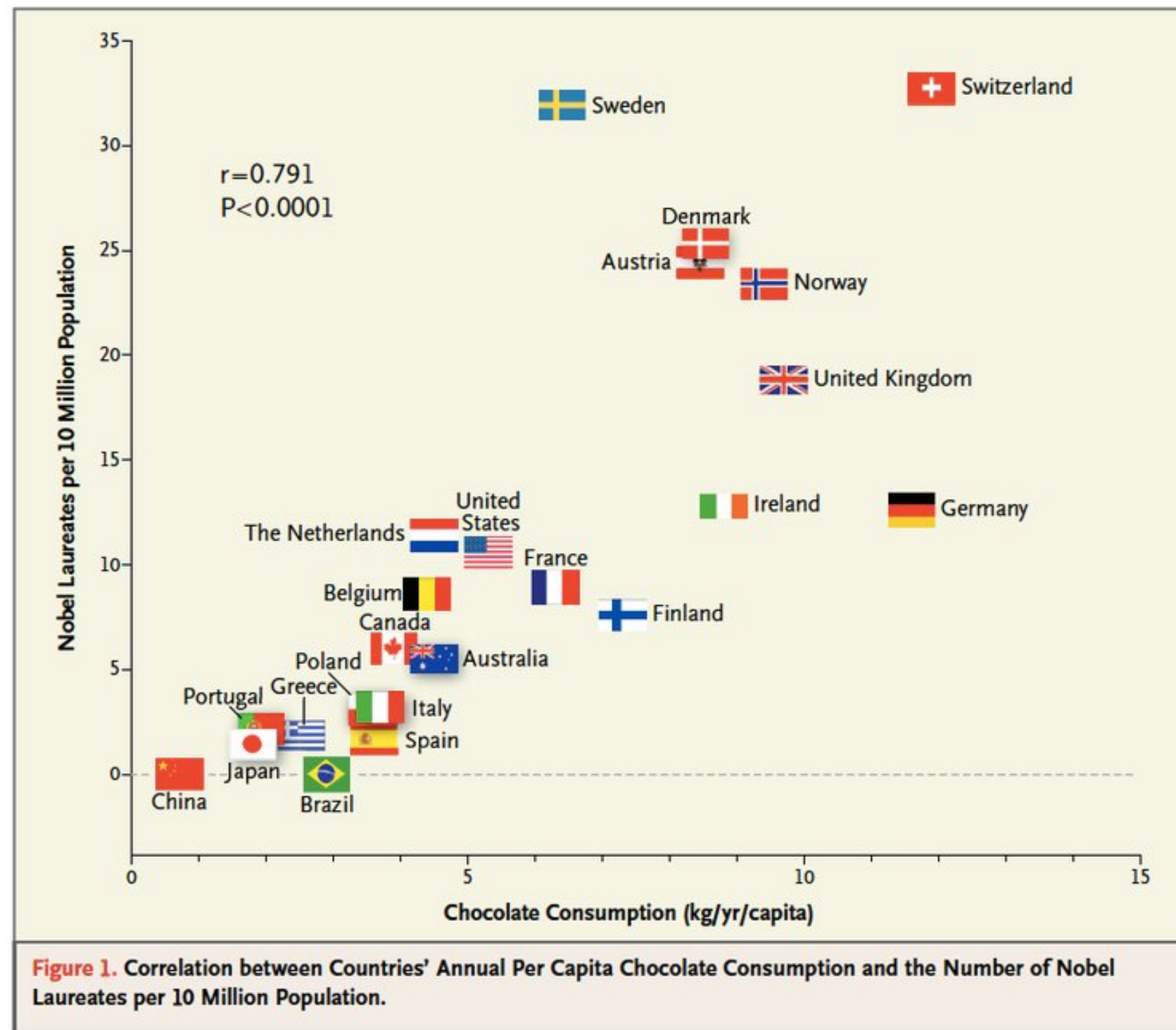
- X causes Y, or
- Y causes X, or
- There exists a common cause between X and Y (*a confounder*)

(or any combination of the above)

((we ignore selection bias*))

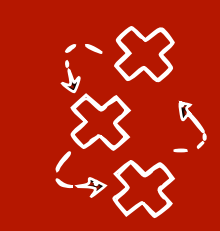


What is the causal graph here?

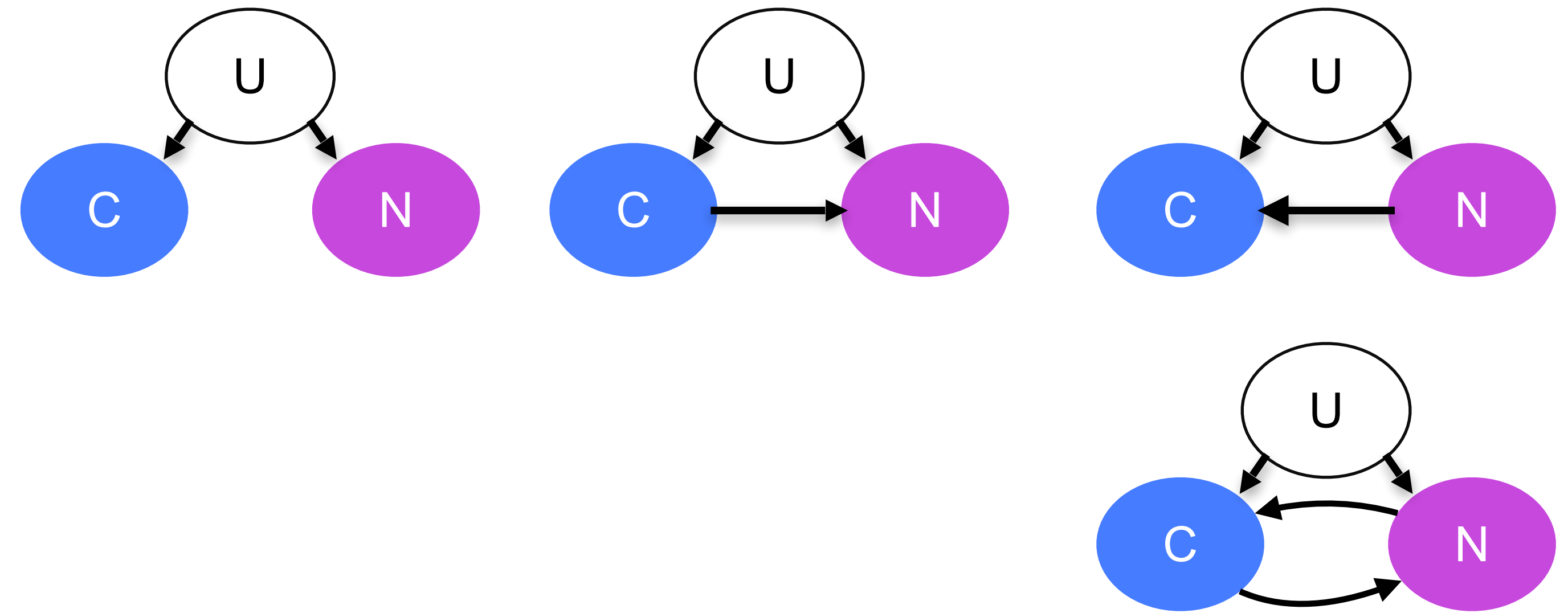
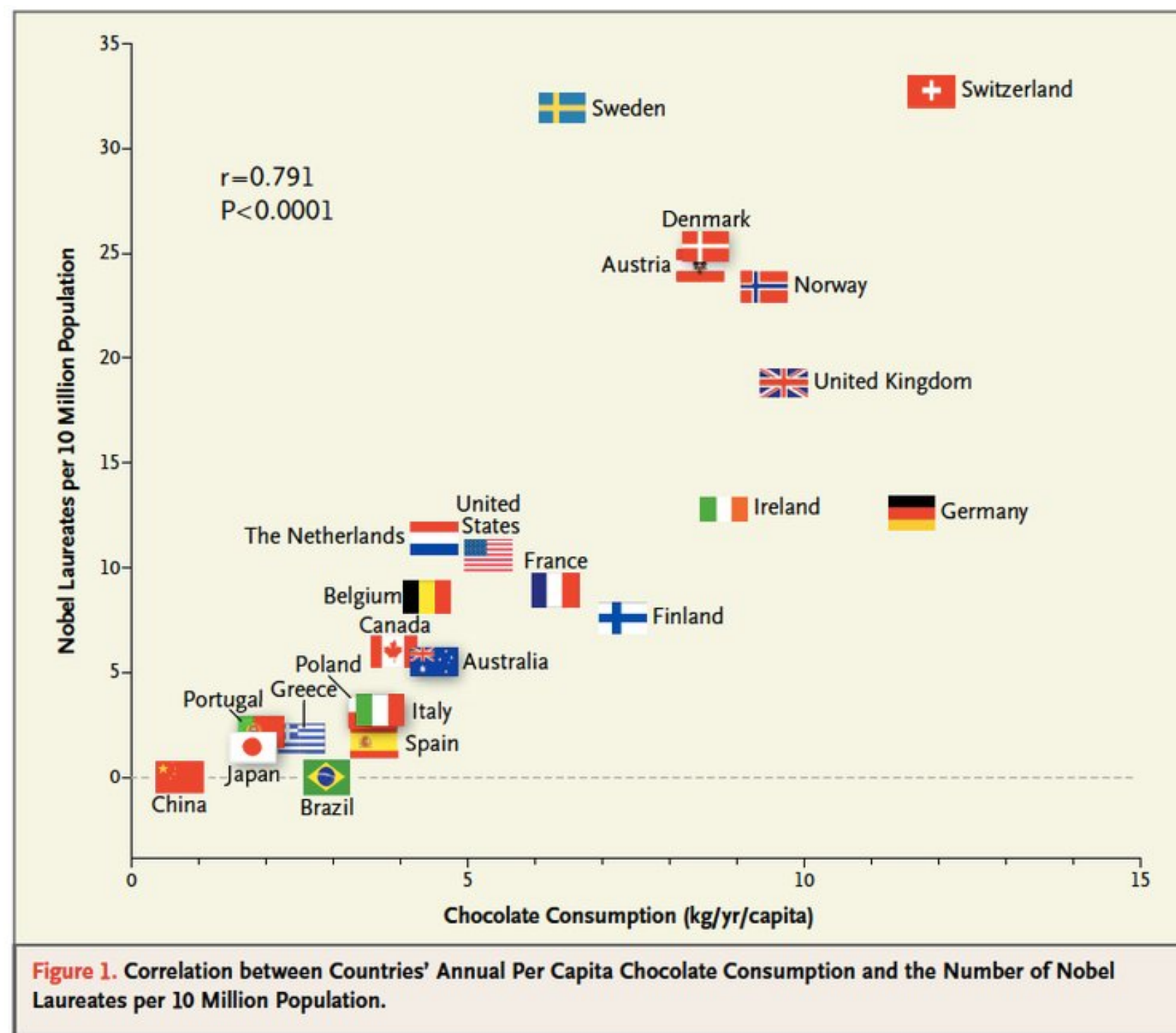


Reichenbach's principle of common cause:

A correlation between X and Y implies that X causes Y, Y causes X, or there exists a common cause between X and Y (or any combination)

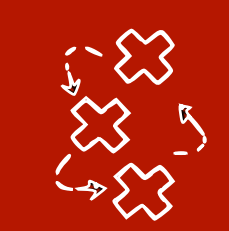


What is the causal graph here?

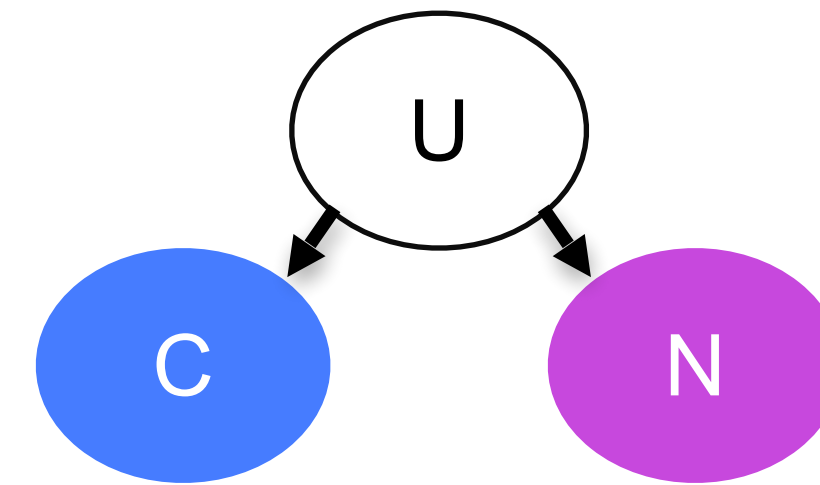
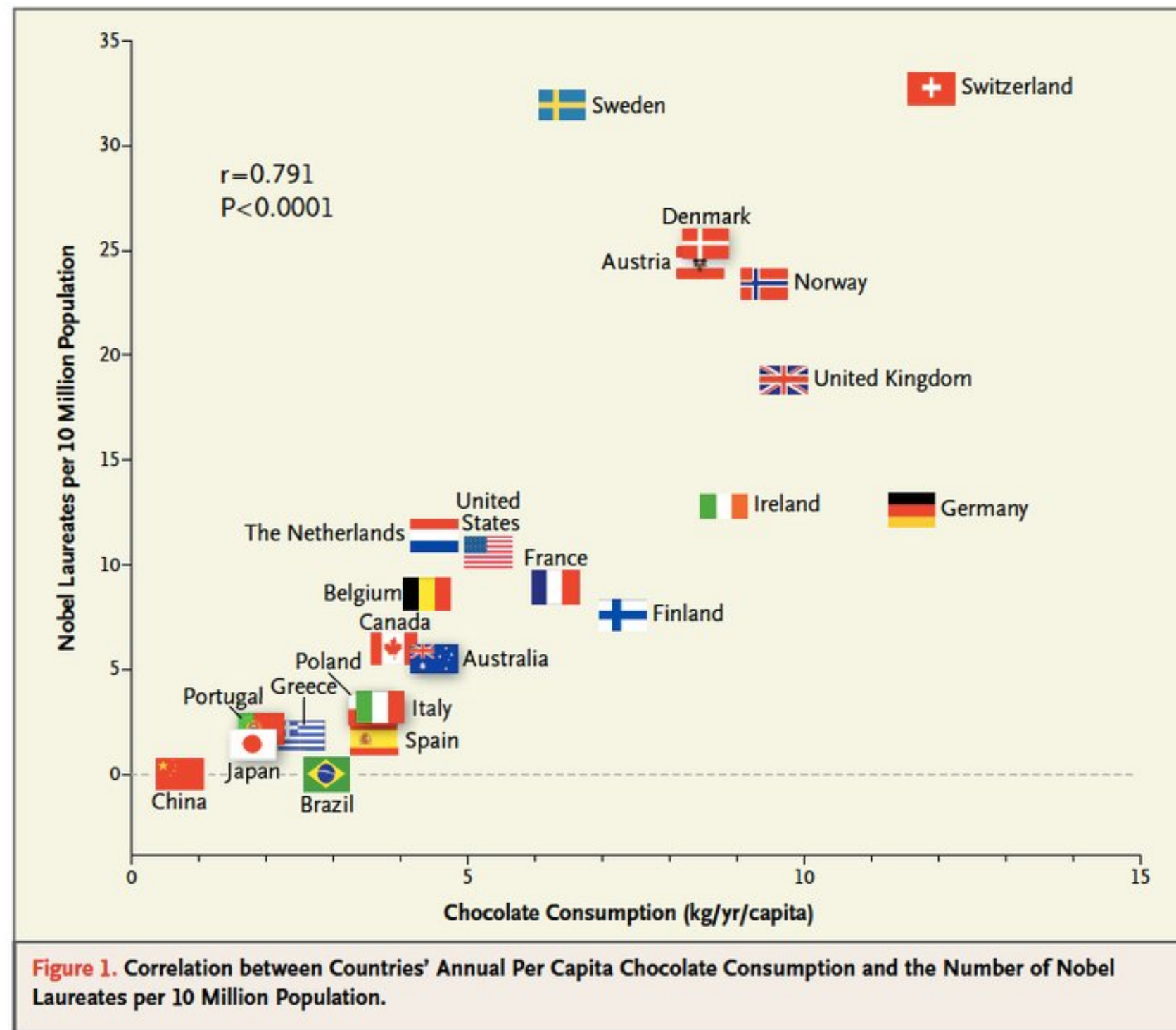


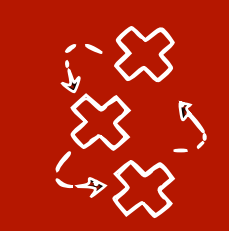
Reichenbach's principle of common cause:

A correlation between X and Y implies that X causes Y, Y causes X, or there exists a common cause between X and Y (or any combination)

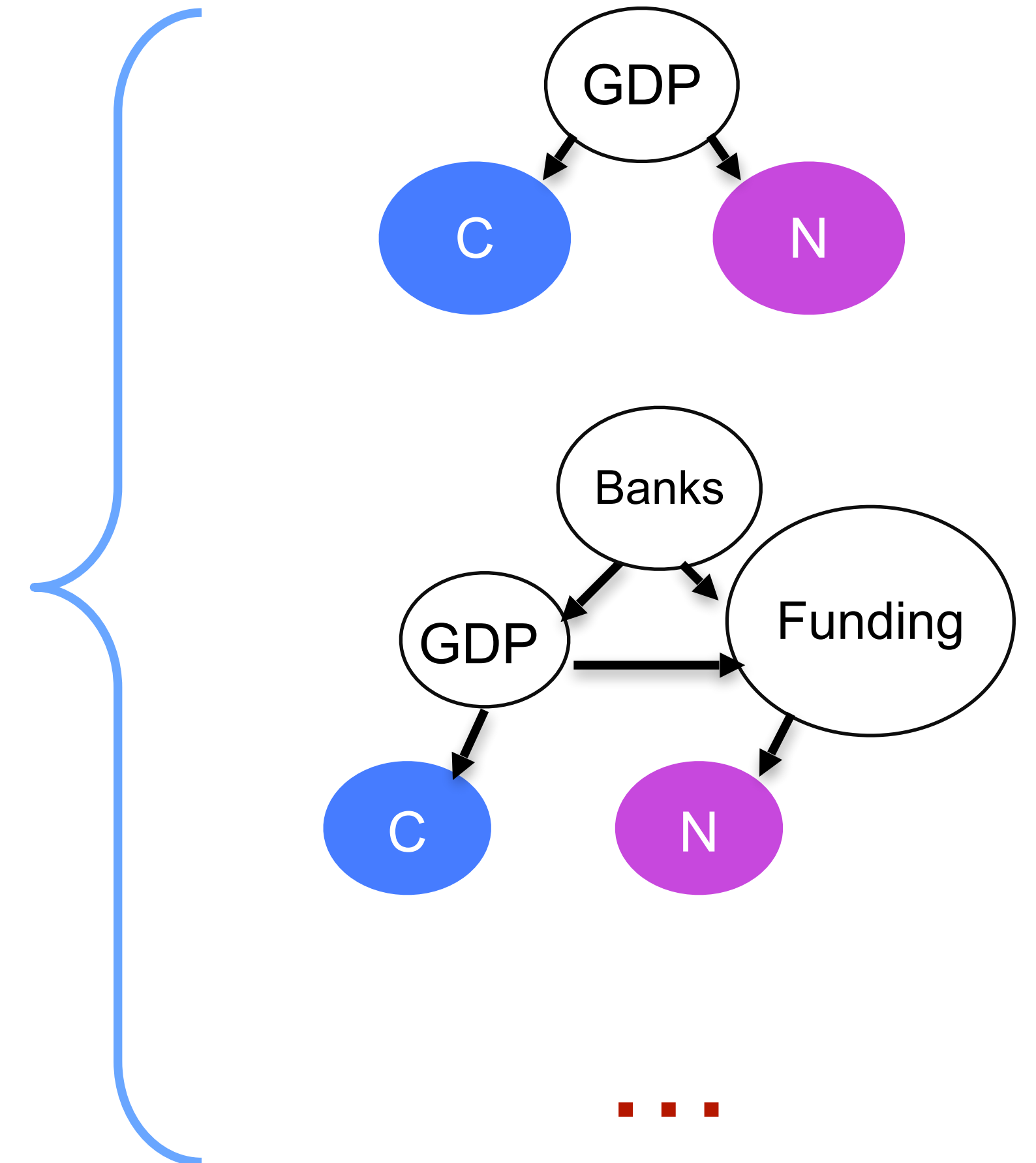
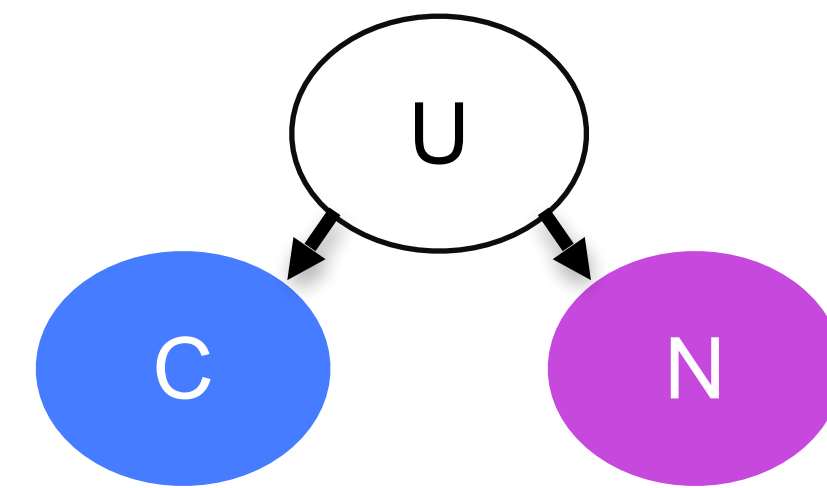
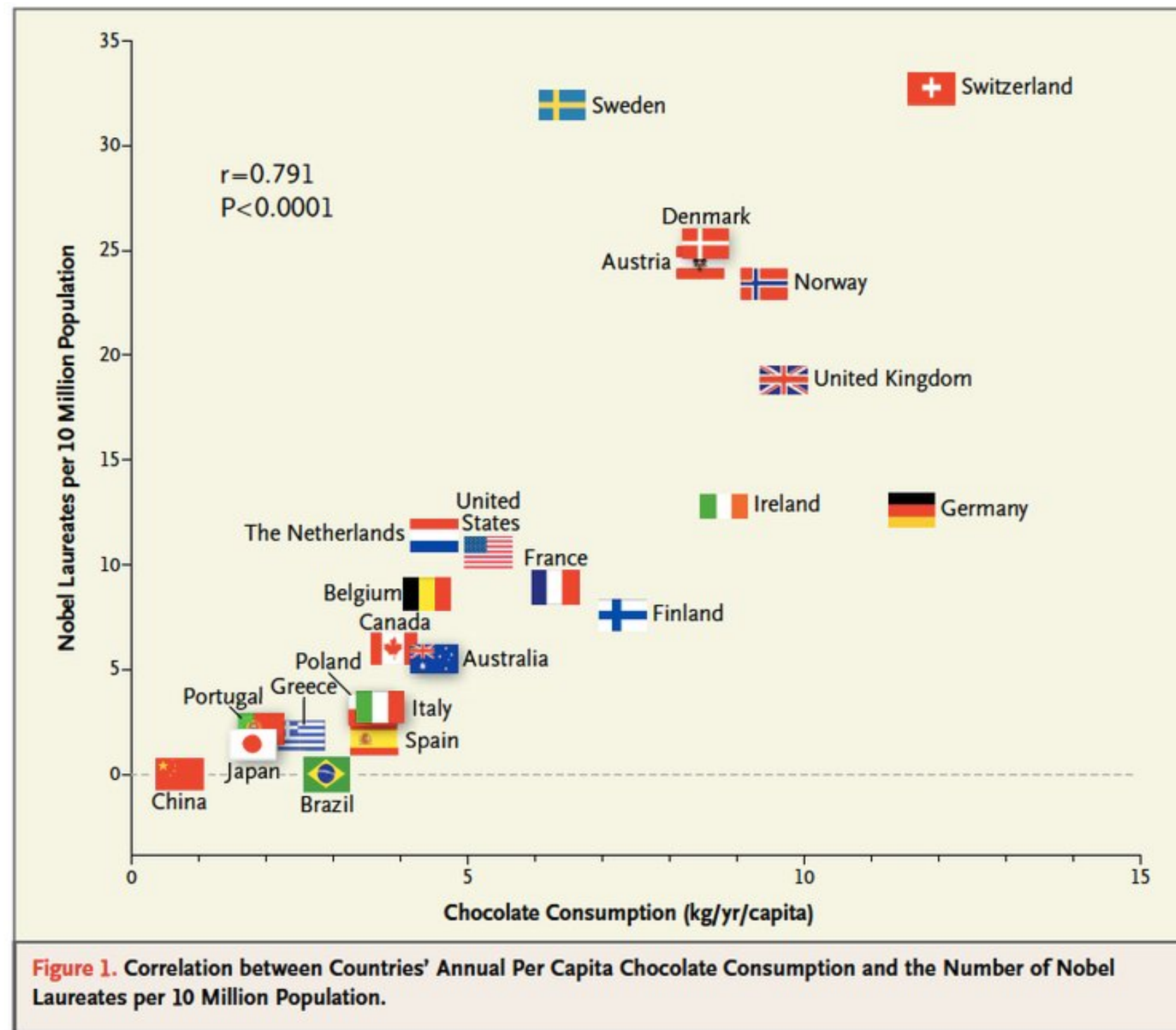


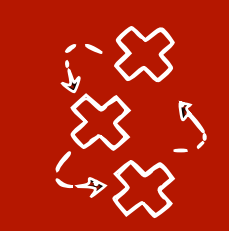
Latent confounding



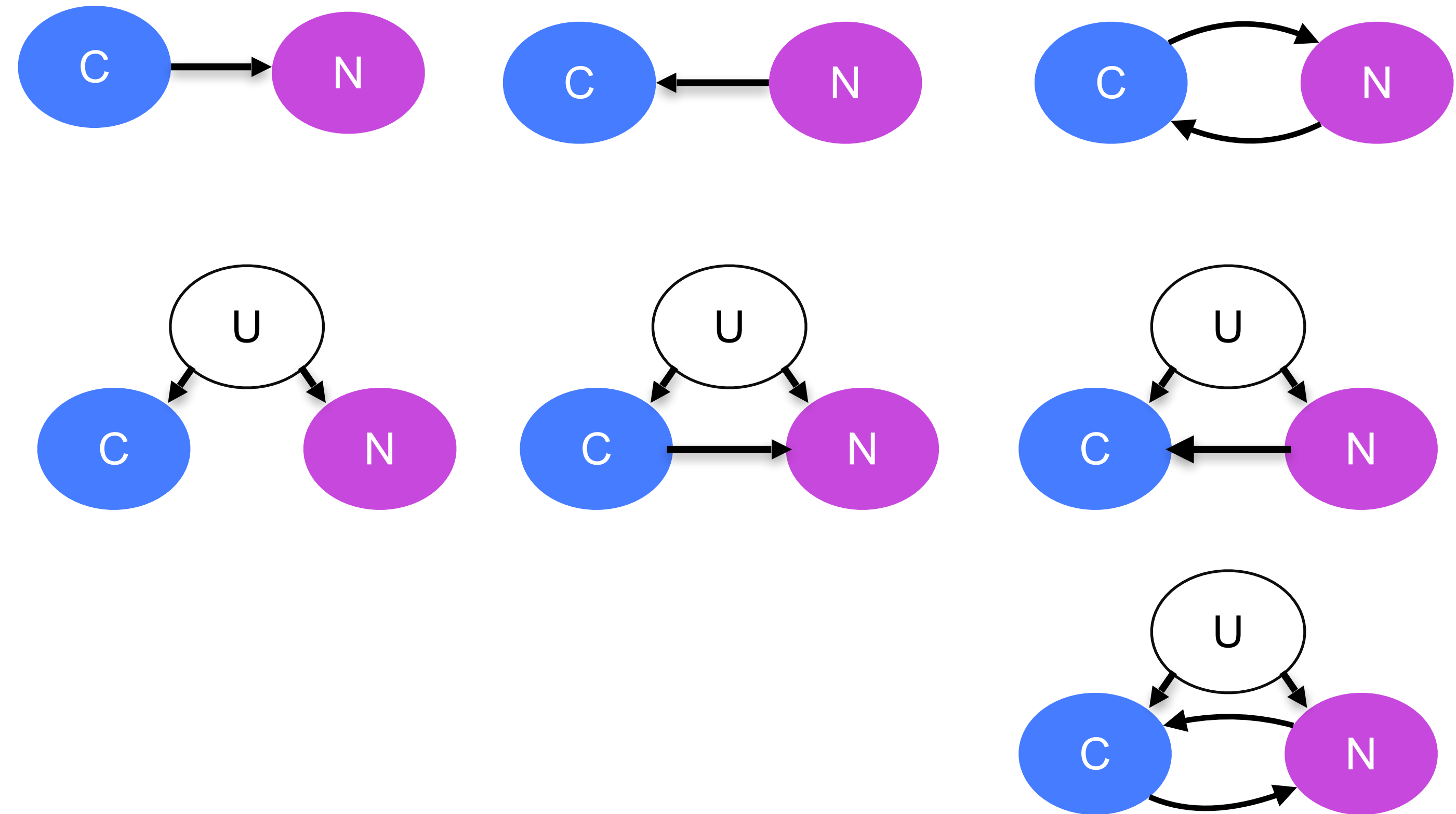
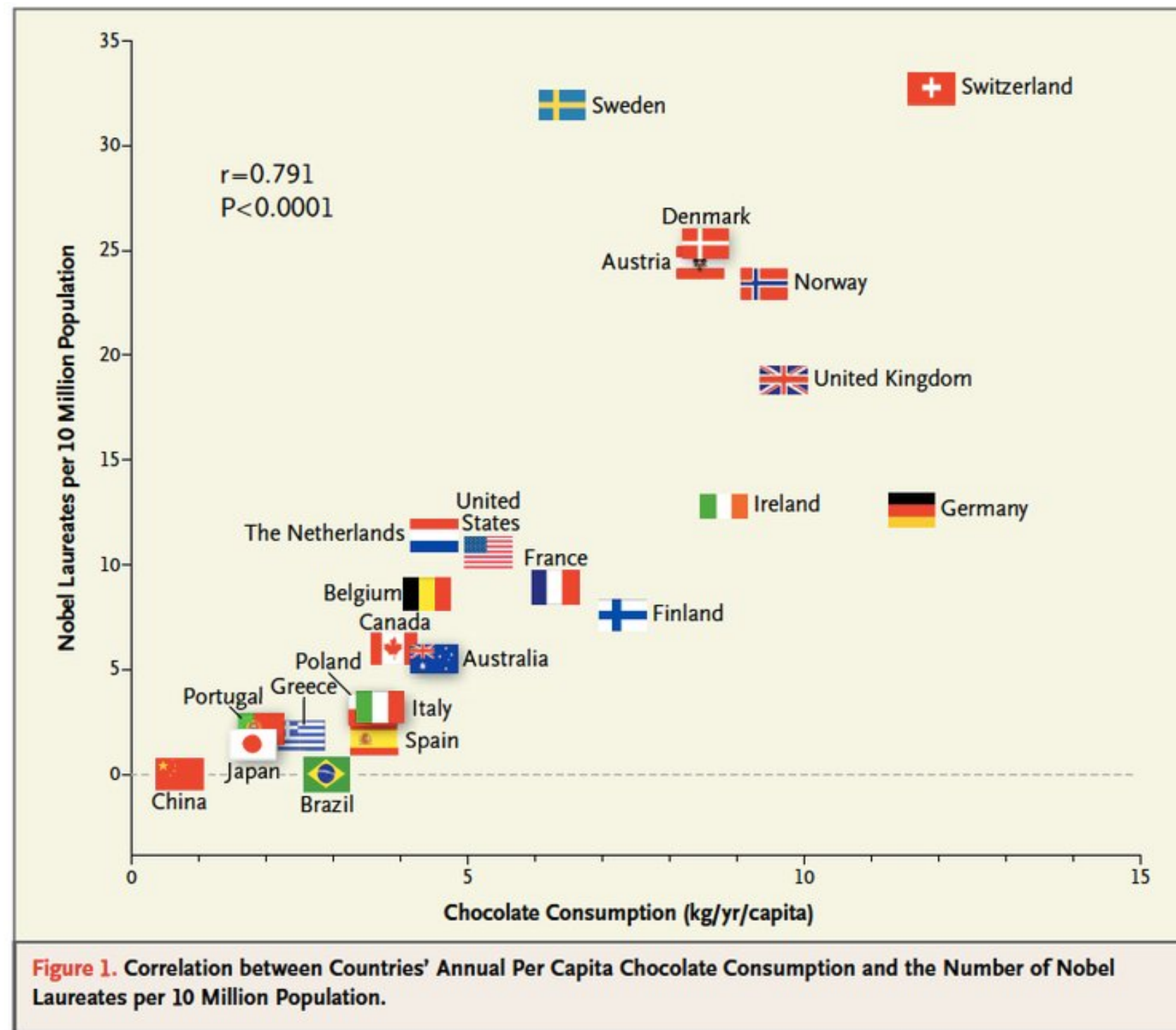


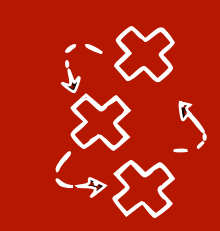
Latent confounding



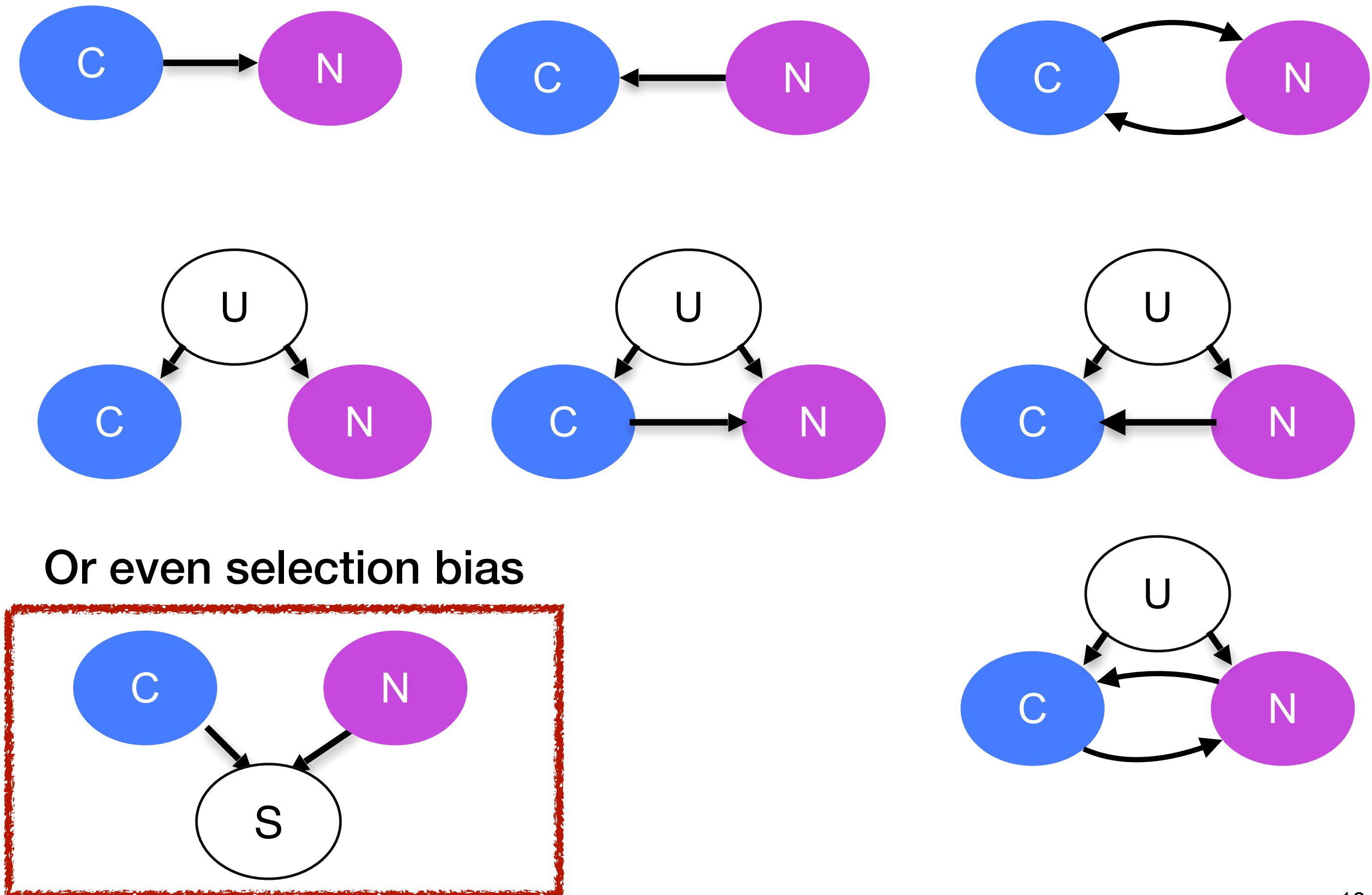
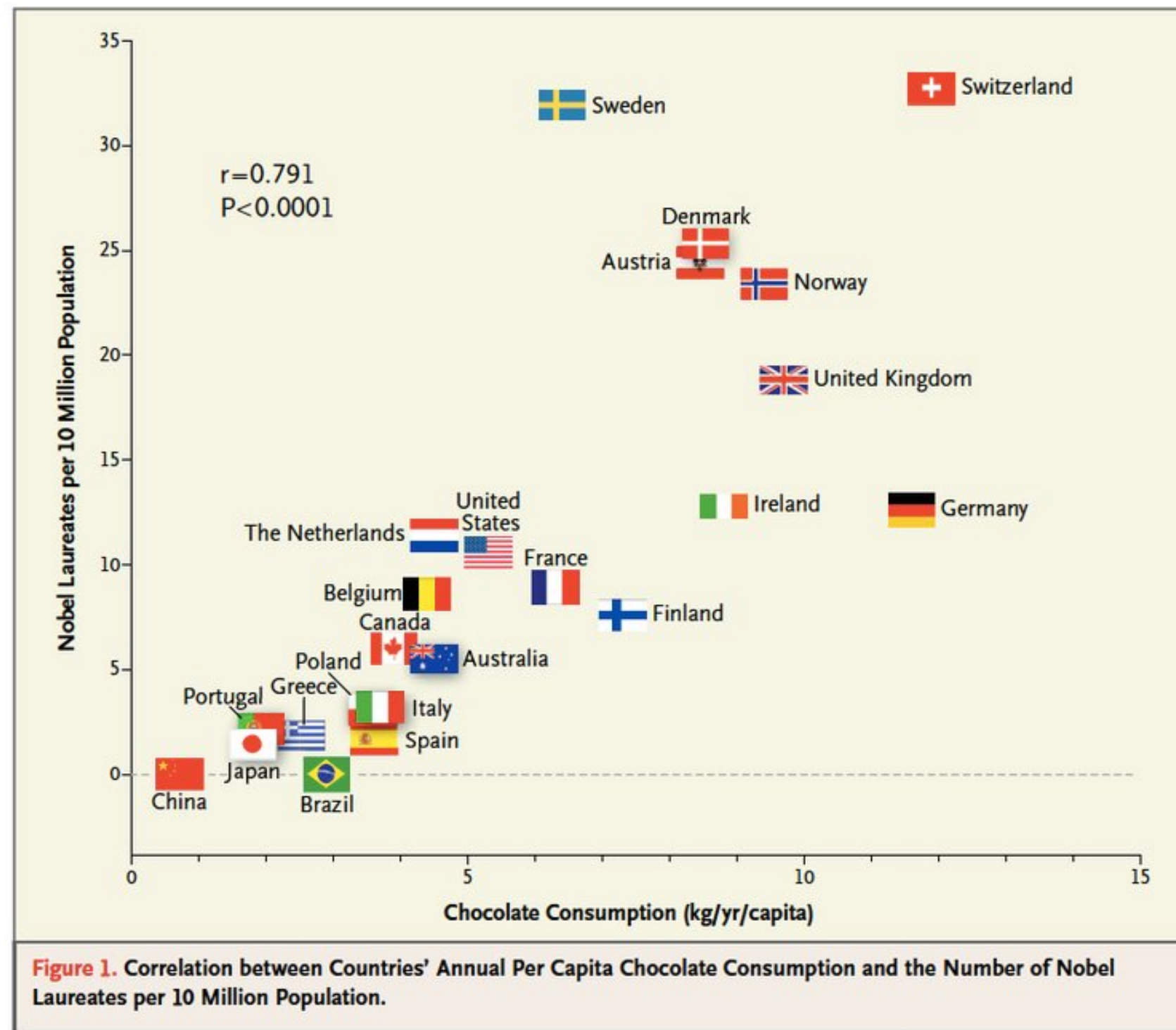


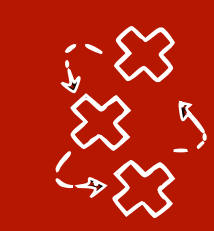
Statistically indistinguishable from only these data



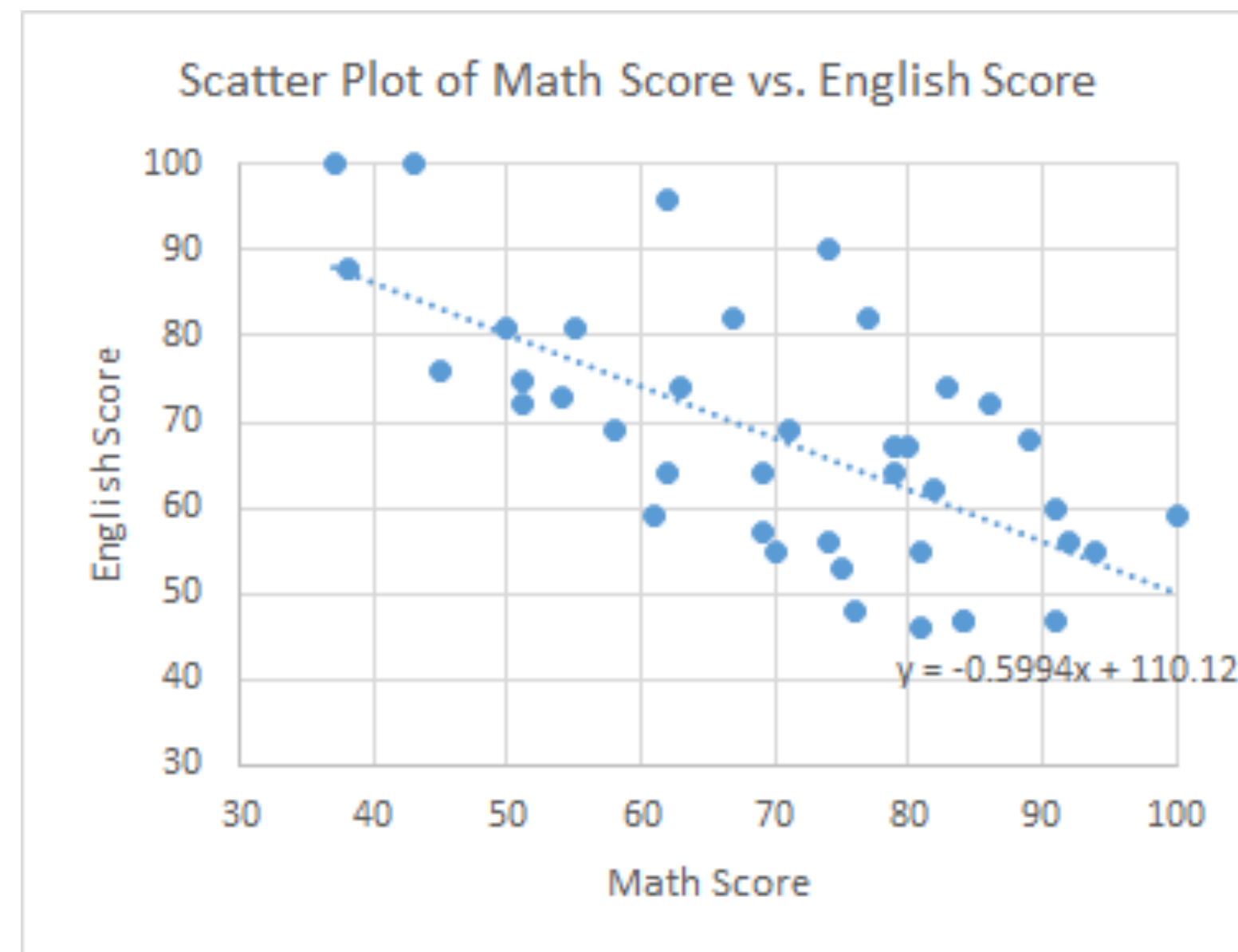


Statistically indistinguishable from only these data

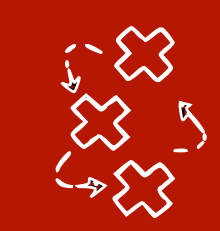




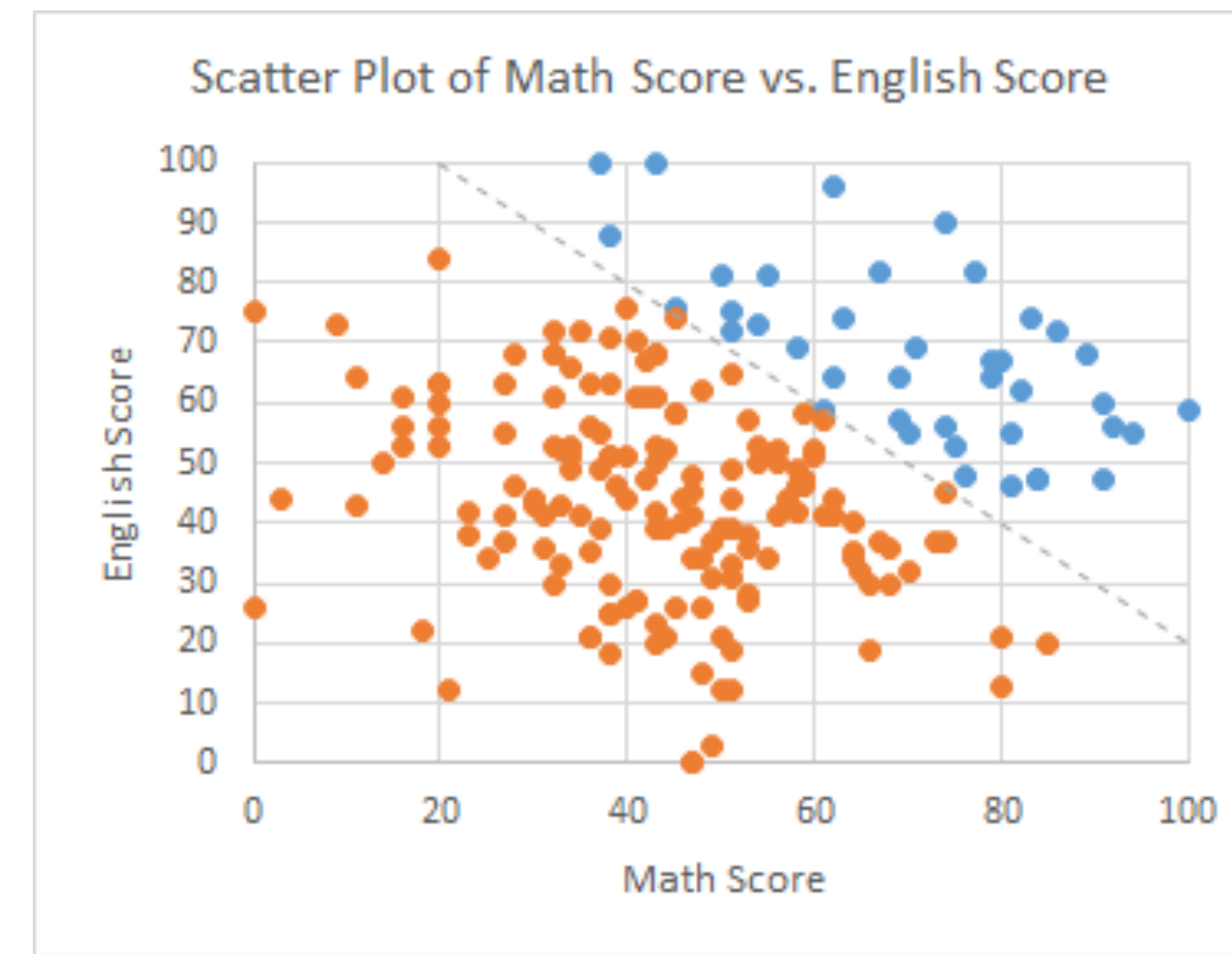
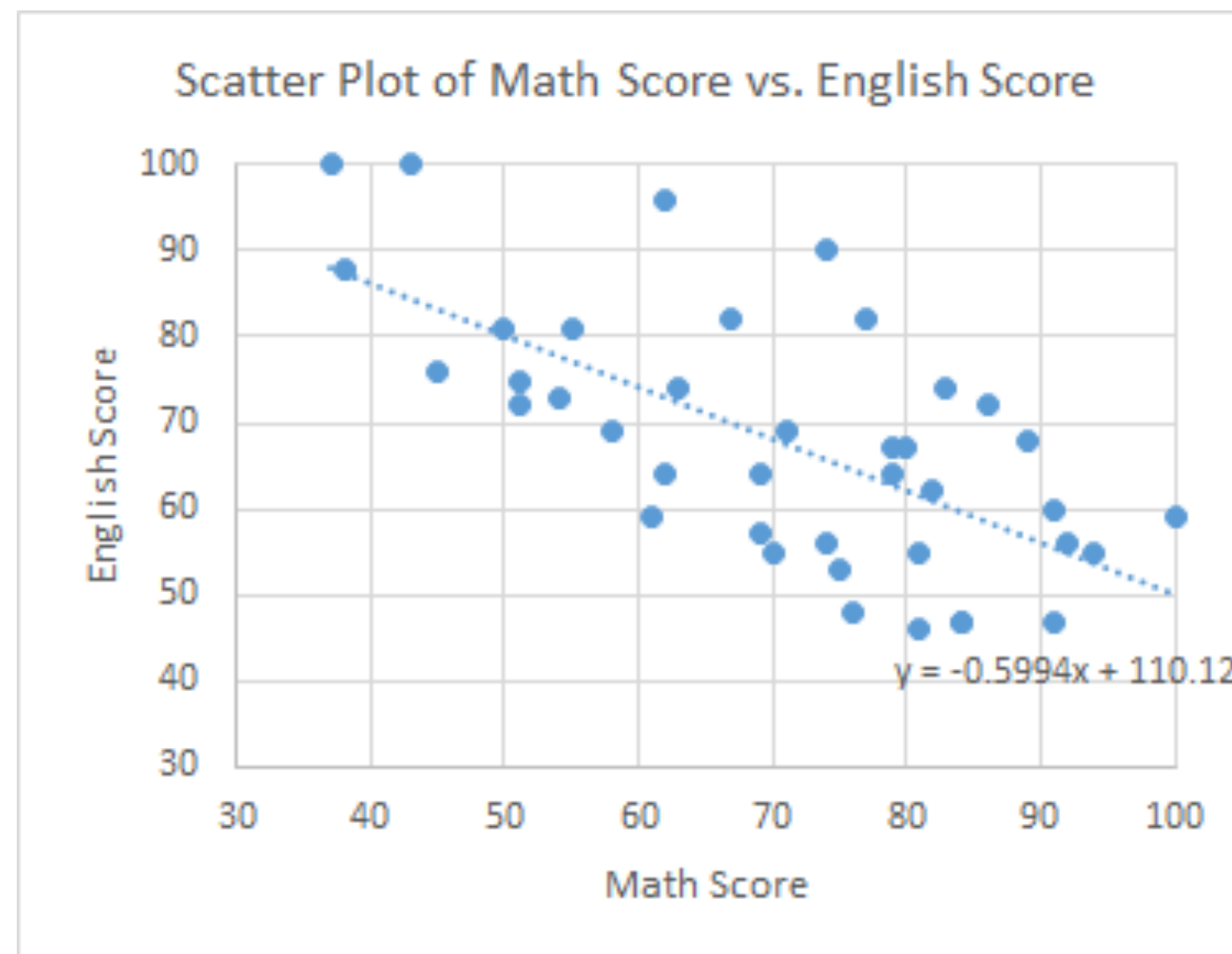
An example of selection bias



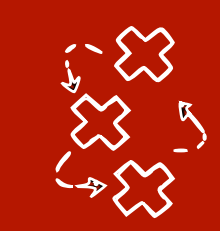
College admissions: English score is negatively correlated with Math Score



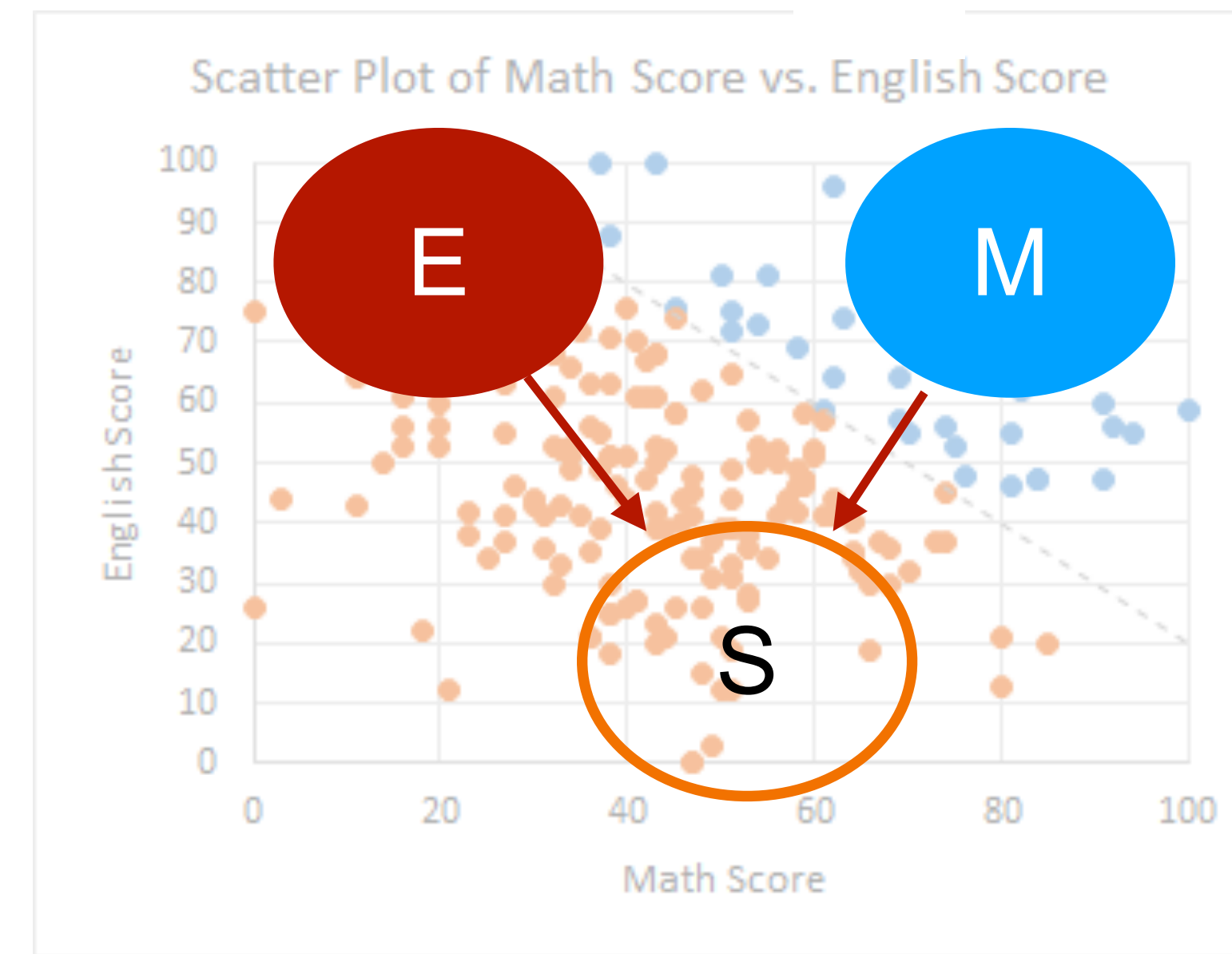
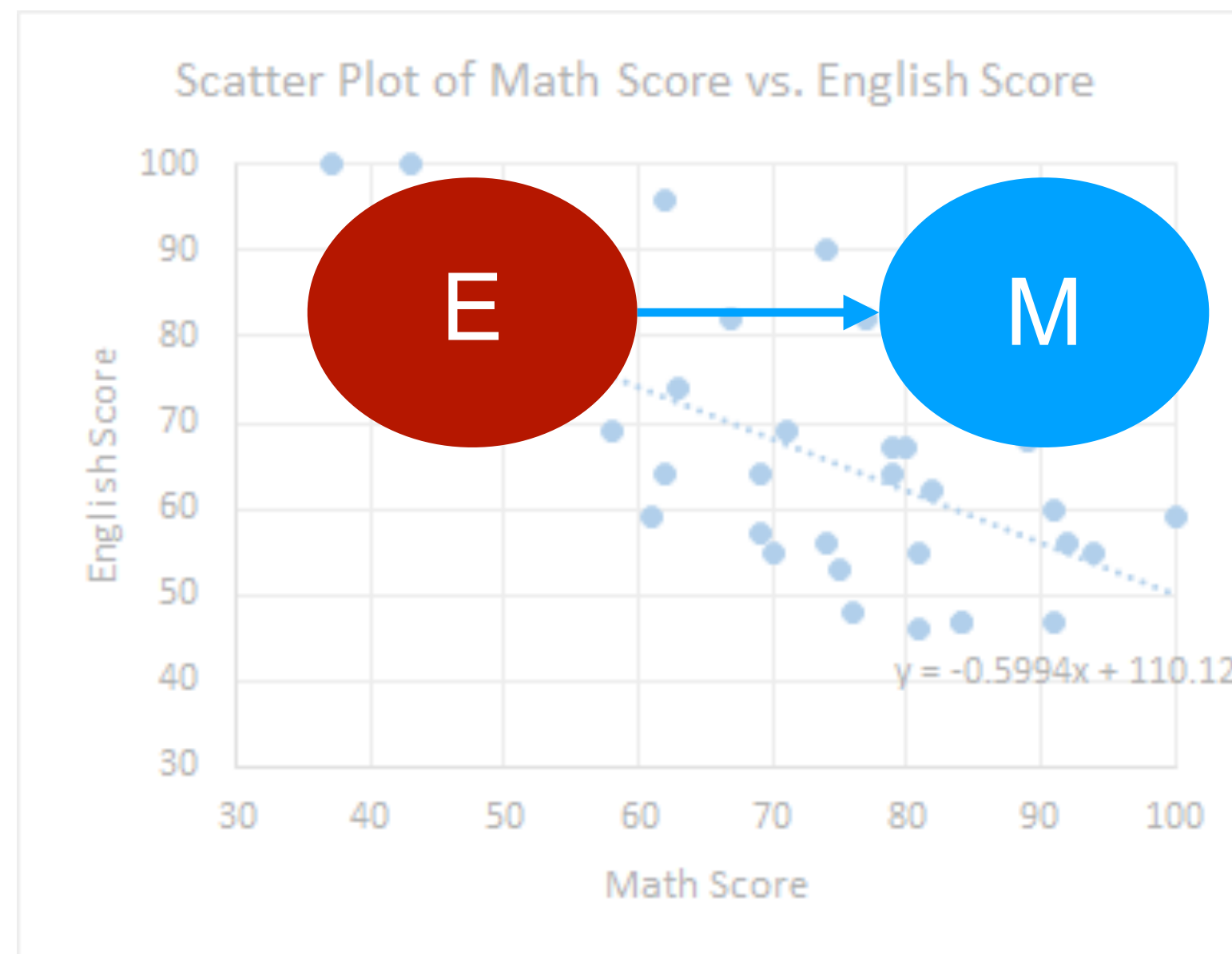
An example of selection bias



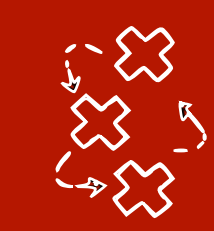
Selection rule: English score + Math Score > 120



An example of selection bias



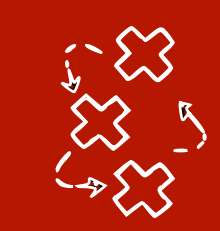
Selection rule: English score + Math Score > 120



Common assumptions

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

$$\mathbf{A} \perp_G \mathbf{B} \mid \mathbf{C} \iff X_{\mathbf{A}} \perp\!\!\!\perp_p X_{\mathbf{B}} \mid X_{\mathbf{C}}$$

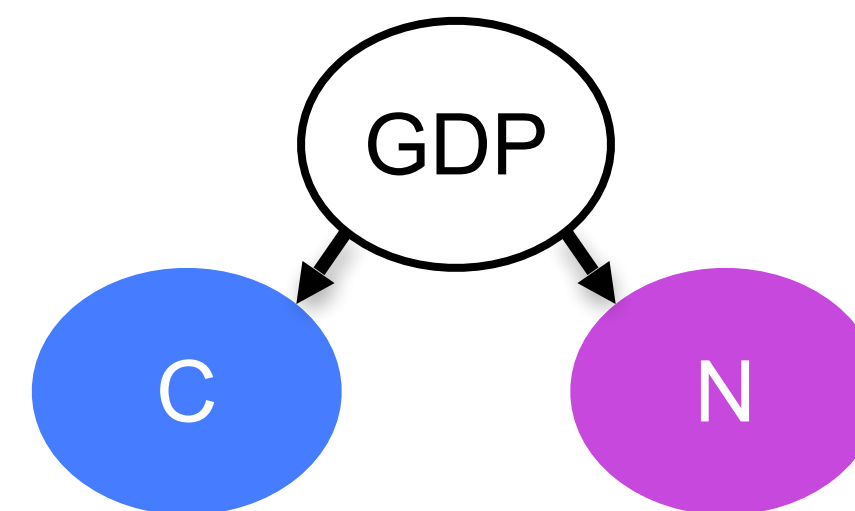
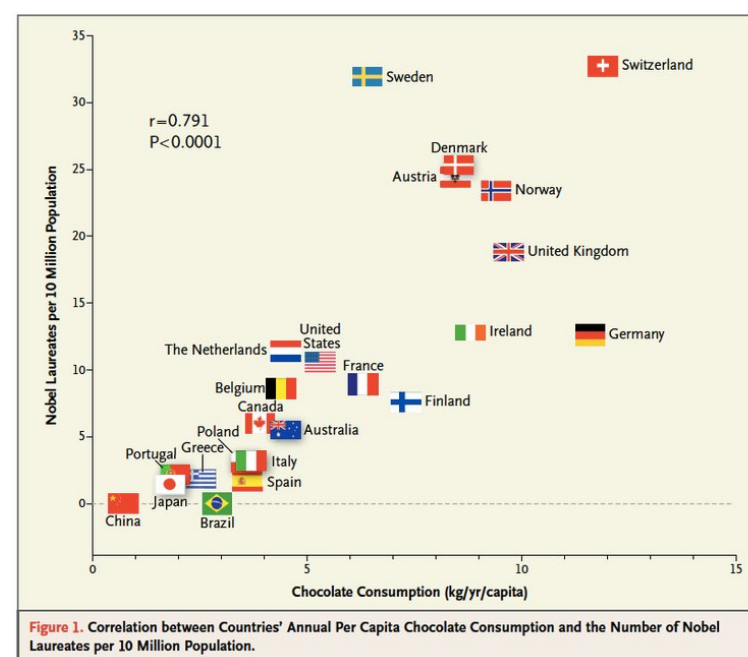


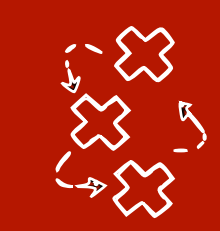
Common assumptions

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

$$\mathbf{A} \perp_G \mathbf{B} \mid \mathbf{C} \iff X_{\mathbf{A}} \perp\!\!\!\perp_p X_{\mathbf{B}} \mid X_{\mathbf{C}}$$

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





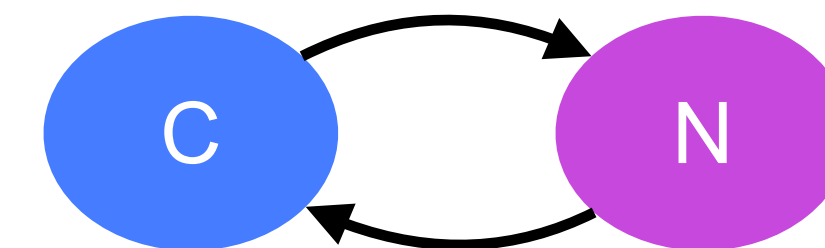
Common assumptions

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

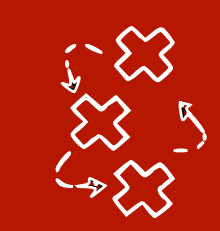
$$\mathbf{A} \perp_G \mathbf{B} \mid \mathbf{C} \iff X_{\mathbf{A}} \perp\!\!\!\perp_p X_{\mathbf{B}} \mid X_{\mathbf{C}}$$

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

- **Acyclicity** - the underlying graph is acyclic



- Cycles + causal insufficiency: sigma separation, Joint Causal Inference



Causal discovery (causal structure learning)

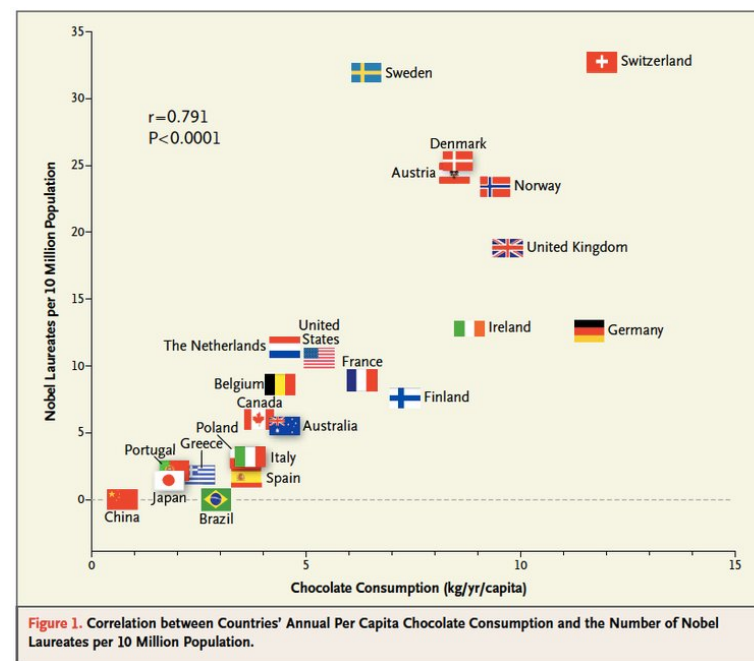


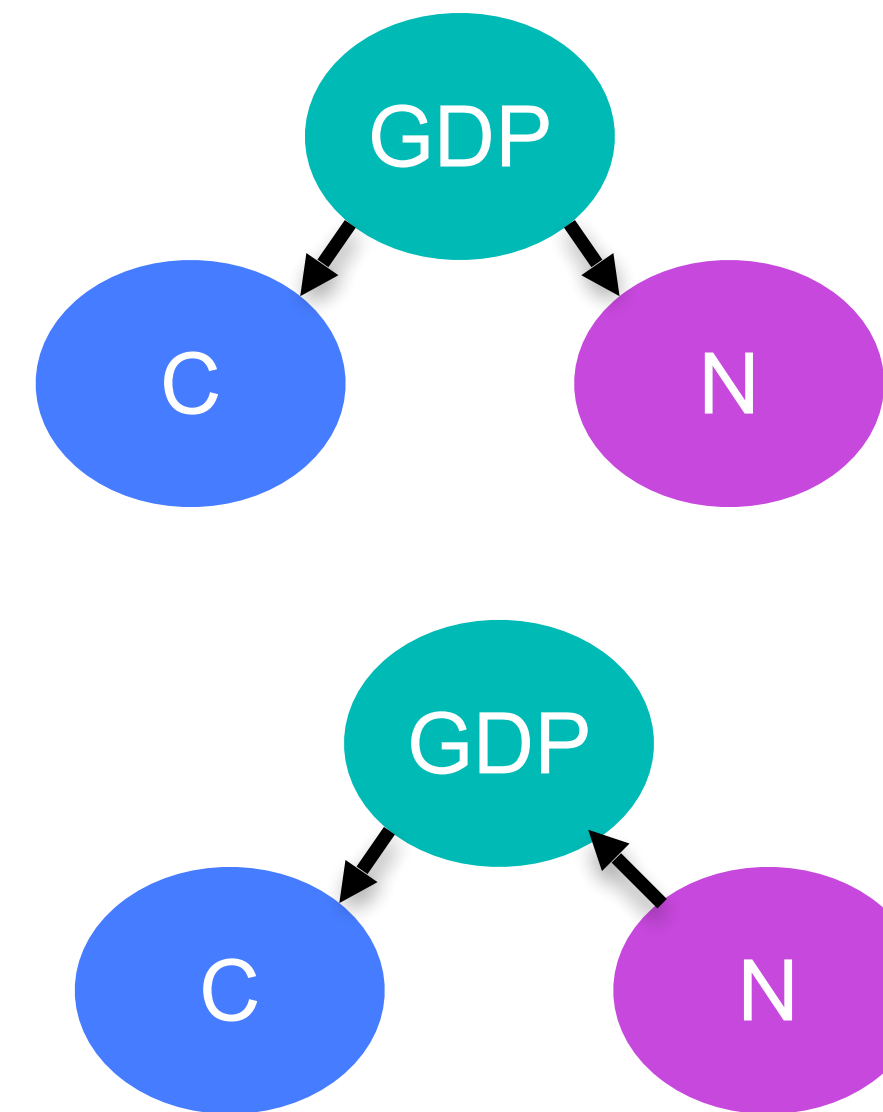
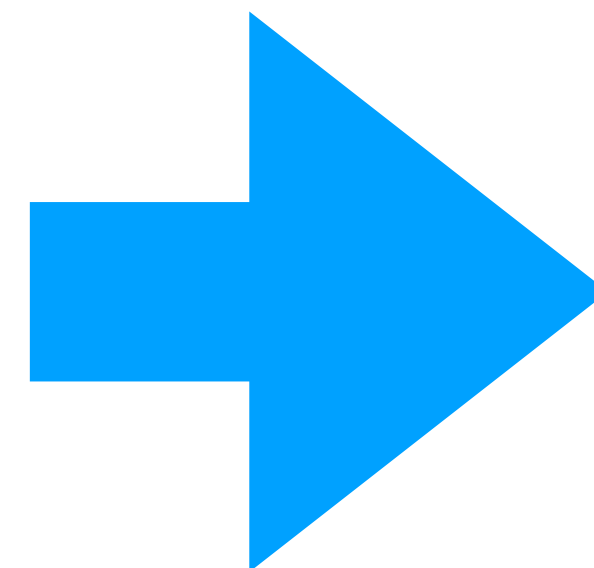
Figure 1. Correlation between Countries' Annual Per Capita Chocolate Consumption and the Number of Nobel Laureates per 10 Million Population.

C	N	GDP
4.5	5	33k
12	30	86k
10	20	46k
...

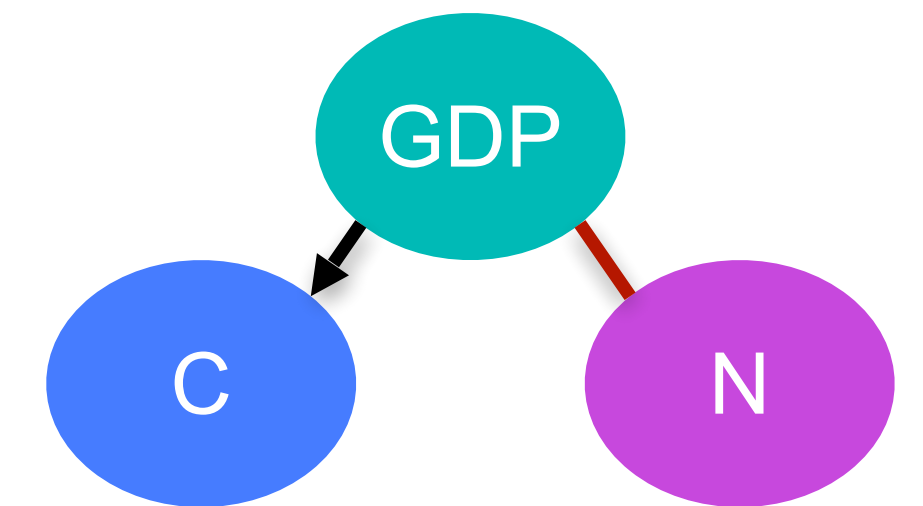
Observational data

$$C \nrightarrow GDP$$

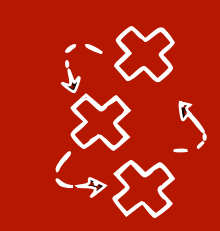
[Optional] Background knowledge



Sets of graphs that fit the data and background knowledge



Summary graph



Causal discovery simplified overview

Constraint-based causal discovery

- Conditional independence tests
- Observational data
- Output: MEC
- SGS, PC, FCI

Score-based causal discovery

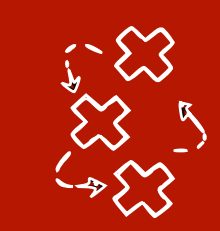
- Penalised likelihood
- Observational data
- Output: MEC
- GES, MMHC

Restricted models

- Nonlinear additive noise, Linear Non-Gaussianity
- Observational data
- Output: DAG
- RESIT, LINGAM

Interventional causal discovery / causal invariance

- Observational and Interventional data
- Output: parents of Y, I-MEC
- ICP, GIES, JCI



Causal discovery - this class

Constraint-based causal discovery

- Conditional independence tests
- Observational data
- Output: MEC
- SGS, PC

Score-based causal discovery

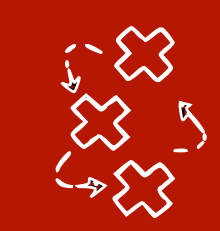
- Penalised likelihood
- Observational data
- Output: MEC
- GES, MMHC

Restricted models

- Nonlinear additive noise, Linear Non-Gaussianity
- Observational data
- Output: DAG
- RESIT, LINGAM

Interventional causal discovery / causal invariance

- Observational and Interventional data
- Output: parents of Y, I-MEC
- ICP, JCI



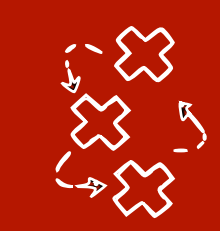
Recap: Global Markov Property & faithfulness

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

$$\mathbf{A} \perp_d \mathbf{B} \mid \mathbf{C} \implies X_{\mathbf{A}} \perp\!\!\!\perp X_{\mathbf{B}} \mid X_{\mathbf{C}}$$

**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



Recap: Global Markov Property & faithfulness

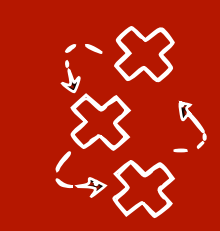
- If (G, p) is a Bayesian network with a DAG $G = (\mathbf{V}, \mathbf{E})$, i.e. **p factorizes according to G** , then for any disjoint $\mathbf{A}, \mathbf{B}, \mathbf{C} \subseteq \mathbf{V}$:

$$\mathbf{A} \perp_d \mathbf{B} \mid \mathbf{C} \implies X_{\mathbf{A}} \perp\!\!\!\perp X_{\mathbf{B}} \mid X_{\mathbf{C}}$$

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

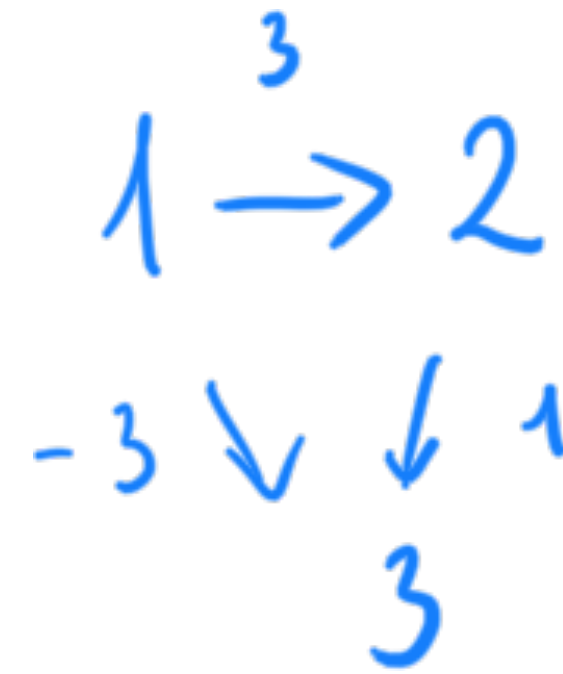
- **d-separations** that can be read purely from a graph imply **conditional independences** in the random variables and data generated by the graph
- The reverse implication is not true in general, but if it is, we say that **p is faithful to G** :

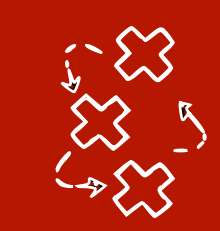
$$X_{\mathbf{A}} \perp\!\!\!\perp X_{\mathbf{B}} \mid X_{\mathbf{C}} \implies \mathbf{A} \perp_d \mathbf{B} \mid \mathbf{C}$$



Faithfulness violation example - cancelling paths

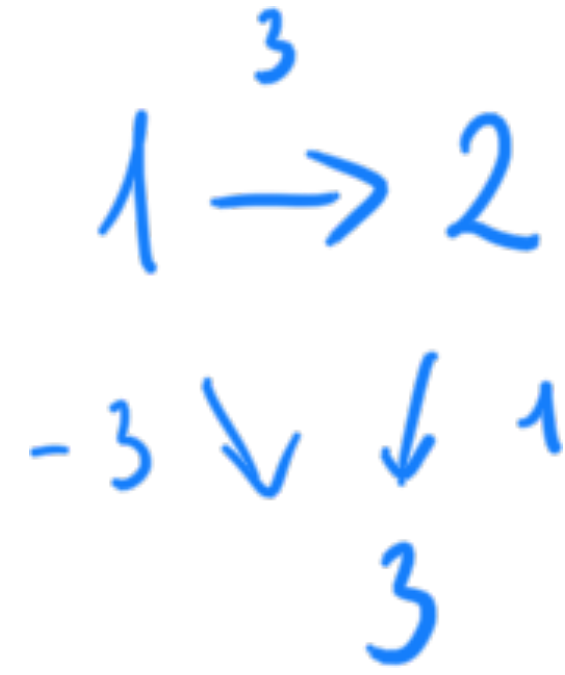
$$\left\{ \begin{array}{l} X_1 = \epsilon_1 \\ X_2 = 3X_1 + \epsilon_2 \\ X_3 = X_2 - 3X_1 + \epsilon_3 \\ \epsilon_1, \epsilon_2, \epsilon_3 \sim N(0,1) \\ \epsilon_1 \perp\!\!\!\perp \epsilon_2, \epsilon_1 \perp\!\!\!\perp \epsilon_3, \epsilon_2 \perp\!\!\!\perp \epsilon_3 \end{array} \right.$$



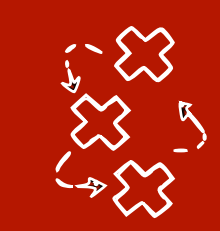


Faithfulness violation example - cancelling paths

$$\left\{ \begin{array}{l} X_1 = \epsilon_1 \\ X_2 = 3X_1 + \epsilon_2 \\ X_3 = X_2 - 3X_1 + \epsilon_3 \\ \epsilon_1, \epsilon_2, \epsilon_3 \sim N(0,1) \\ \epsilon_1 \perp\!\!\!\perp \epsilon_2, \epsilon_1 \perp\!\!\!\perp \epsilon_3, \epsilon_2 \perp\!\!\!\perp \epsilon_3 \end{array} \right.$$

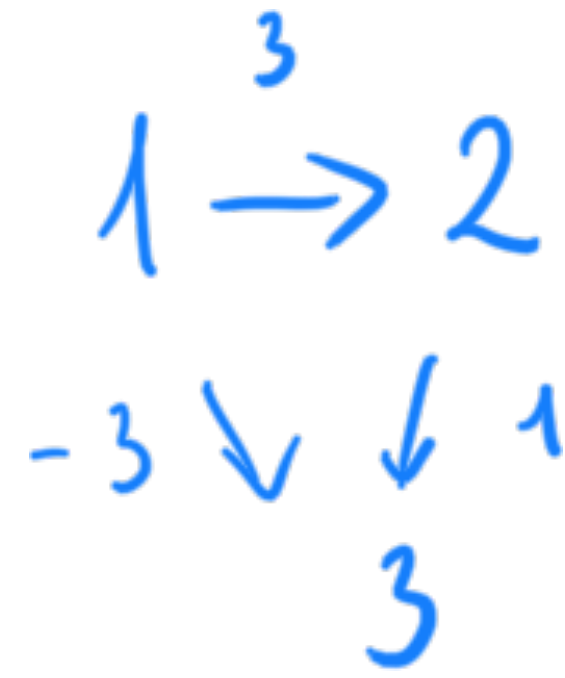


$$\left\{ \begin{array}{l} X_1 = \epsilon_1 \\ X_2 = 3\epsilon_1 + \epsilon_2 \\ X_3 = 3\epsilon_1 + \epsilon_2 - 3\epsilon_1 + \epsilon_3 \\ \epsilon_1, \epsilon_2, \epsilon_3 \sim N(0,1) \\ \epsilon_1 \perp\!\!\!\perp \epsilon_2, \epsilon_1 \perp\!\!\!\perp \epsilon_3, \epsilon_2 \perp\!\!\!\perp \epsilon_3 \end{array} \right.$$

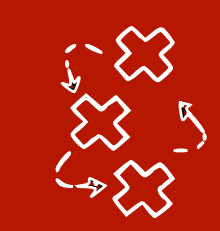


Faithfulness violation example - cancelling paths

$$\left\{ \begin{array}{l} X_1 = \epsilon_1 \\ X_2 = 3X_1 + \epsilon_2 \\ X_3 = X_2 - 3X_1 + \epsilon_3 \\ \epsilon_1, \epsilon_2, \epsilon_3 \sim N(0,1) \\ \epsilon_1 \perp\!\!\!\perp \epsilon_2, \epsilon_1 \perp\!\!\!\perp \epsilon_3, \epsilon_2 \perp\!\!\!\perp \epsilon_3 \end{array} \right.$$

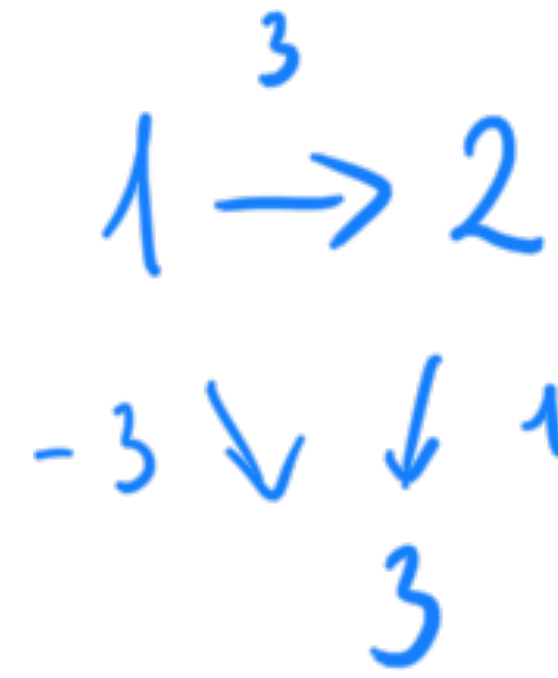


$$\left\{ \begin{array}{l} X_1 = \epsilon_1 \\ X_2 = 3\epsilon_1 + \epsilon_2 \\ X_3 = \cancel{3\epsilon_1} + \epsilon_2 - \cancel{3\epsilon_1} + \epsilon_3 \\ \epsilon_1, \epsilon_2, \epsilon_3 \sim N(0,1) \\ \epsilon_1 \perp\!\!\!\perp \epsilon_2, \epsilon_1 \perp\!\!\!\perp \epsilon_3, \epsilon_2 \perp\!\!\!\perp \epsilon_3 \end{array} \right.$$

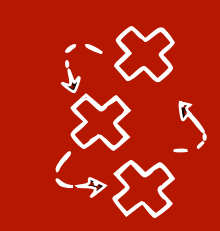


Faithfulness violation example - cancelling paths

$$\left\{ \begin{array}{l} X_1 = \epsilon_1 \\ X_2 = 3X_1 + \epsilon_2 \\ X_3 = X_2 - 3X_1 + \epsilon_3 \\ \epsilon_1, \epsilon_2, \epsilon_3 \sim N(0,1) \\ \epsilon_1 \perp\!\!\!\perp \epsilon_2, \epsilon_1 \perp\!\!\!\perp \epsilon_3, \epsilon_2 \perp\!\!\!\perp \epsilon_3 \end{array} \right.$$



$$\left\{ \begin{array}{l} X_1 = \epsilon_1 \\ X_2 = 3\epsilon_1 + \epsilon_2 \\ X_3 = \epsilon_2 + \epsilon_3 \\ \epsilon_1, \epsilon_2, \epsilon_3 \sim N(0,1) \\ \epsilon_1 \perp\!\!\!\perp \epsilon_2, \epsilon_1 \perp\!\!\!\perp \epsilon_3, \epsilon_2 \perp\!\!\!\perp \epsilon_3 \end{array} \right.$$



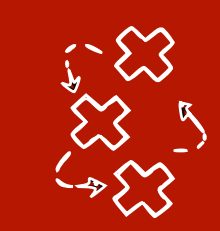
Faithfulness violation example - cancelling paths

$$\left\{ \begin{array}{l} X_1 = \epsilon_1 \\ X_2 = 3X_1 + \epsilon_2 \\ X_3 = X_2 - 3X_1 + \epsilon_3 \\ \epsilon_1, \epsilon_2, \epsilon_3 \sim N(0,1) \\ \epsilon_1 \perp\!\!\!\perp \epsilon_2, \epsilon_1 \perp\!\!\!\perp \epsilon_3, \epsilon_2 \perp\!\!\!\perp \epsilon_3 \end{array} \right.$$

$$\begin{array}{c} 3 \\ 1 \rightarrow 2 \\ -3 \downarrow \downarrow 1 \\ 3 \end{array}$$

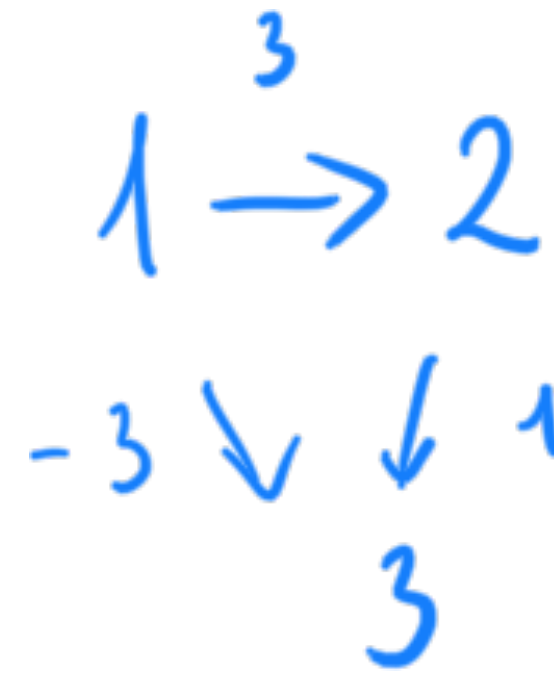
$$\left\{ \begin{array}{l} X_1 = \epsilon_1 \\ X_2 = 3\epsilon_1 + \epsilon_2 \\ X_3 = \epsilon_2 + \epsilon_3 \\ \epsilon_1, \epsilon_2, \epsilon_3 \sim N(0,1) \\ \epsilon_1 \perp\!\!\!\perp \epsilon_2, \epsilon_1 \perp\!\!\!\perp \epsilon_3, \epsilon_2 \perp\!\!\!\perp \epsilon_3 \end{array} \right.$$

$$X_1 \perp\!\!\!\perp X_2$$



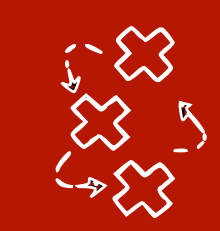
Faithfulness violation example - cancelling paths

$$\left\{ \begin{array}{l} X_1 = \epsilon_1 \\ X_2 = 3X_1 + \epsilon_2 \\ X_3 = X_2 - 3X_1 + \epsilon_3 \\ \epsilon_1, \epsilon_2, \epsilon_3 \sim N(0,1) \\ \epsilon_1 \perp\!\!\!\perp \epsilon_2, \epsilon_1 \perp\!\!\!\perp \epsilon_3, \epsilon_2 \perp\!\!\!\perp \epsilon_3 \end{array} \right.$$



$$\left\{ \begin{array}{l} X_1 = \epsilon_1 \\ X_2 = 3\epsilon_1 + \epsilon_2 \\ X_3 = \epsilon_2 + \epsilon_3 \\ \epsilon_1, \epsilon_2, \epsilon_3 \sim N(0,1) \\ \epsilon_1 \perp\!\!\!\perp \epsilon_2, \epsilon_1 \perp\!\!\!\perp \epsilon_3, \epsilon_2 \perp\!\!\!\perp \epsilon_3 \end{array} \right.$$

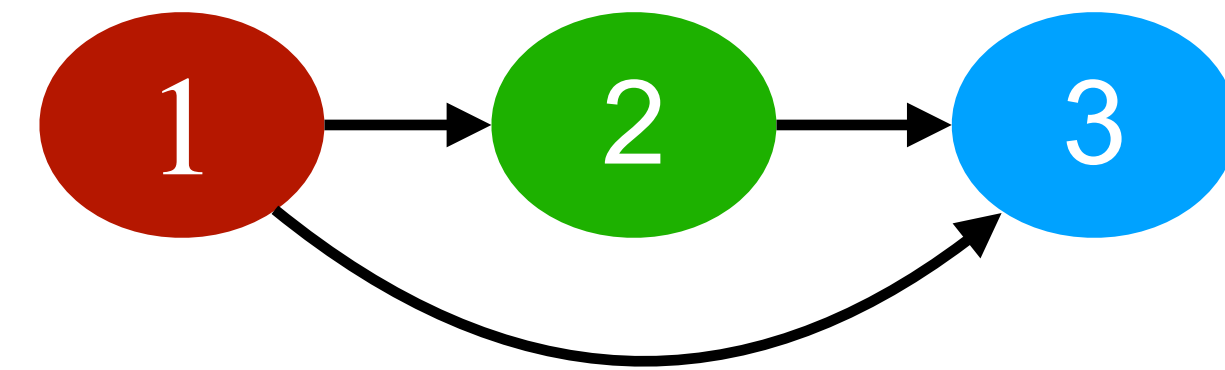
$$X_1 \perp\!\!\!\perp X_2 \quad X_2 \perp\!\!\!\perp X_3 \quad X_1 \perp\!\!\!\perp X_3$$



Faithfulness violation example - cancelling paths

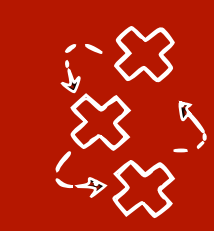
$$\left\{ \begin{array}{l} X_1 = \epsilon_1 \\ X_2 = 3X_1 + \epsilon_2 \\ X_3 = X_2 - 3X_1 + \epsilon_3 \\ \epsilon_1, \epsilon_2, \epsilon_3 \sim N(0,1) \\ \epsilon_1 \perp\!\!\!\perp \epsilon_2, \epsilon_1 \perp\!\!\!\perp \epsilon_3, \epsilon_2 \perp\!\!\!\perp \epsilon_3 \end{array} \right.$$

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



$$X_1 \not\perp_d X_2 \quad X_2 \not\perp_d X_3 \quad X_1 \not\perp_d X_3$$

**Faithfulness violation:
conditional independence does
not imply d-separation**



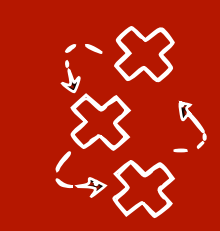
Perfect maps

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

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

$$\mathbf{A} \perp_d \mathbf{B} \mid \mathbf{C} \iff X_{\mathbf{A}} \perp\!\!\!\perp X_{\mathbf{B}} \mid X_{\mathbf{C}}$$

- This correspondence is the basis of learning causal graphs from data
- **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

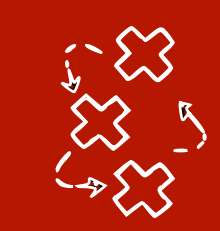


Perfect maps - existence

- Not every distribution p has a perfect map

$$\left\{ \begin{array}{l} X_1 = \epsilon_1 \\ X_2 = X_1 \quad \boxed{\text{Deterministic function}} \\ X_3 = X_2 + \epsilon_3 \\ \epsilon_1, \epsilon_3 \sim N(0,1) \\ \epsilon_1 \perp\!\!\!\perp \epsilon_3 \end{array} \right.$$





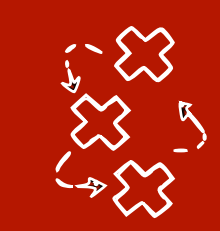
Perfect maps - existence

- Not every distribution p has a perfect map



$$\left\{ \begin{array}{l} X_1 = \epsilon_1 \\ X_2 = X_1 \\ X_3 = X_2 + \epsilon_3 \\ \epsilon_1, \epsilon_3 \sim N(0,1) \\ \epsilon_1 \perp\!\!\!\perp \epsilon_3 \end{array} \right. \quad \boxed{\text{Deterministic function}}$$

$$\left\{ \begin{array}{l} X_1 = \epsilon_1 \\ X_2 = \epsilon_1 \\ X_3 = \epsilon_1 + \epsilon_3 \\ \epsilon_1, \epsilon_3 \sim N(0,1) \\ \epsilon_1 \perp\!\!\!\perp \epsilon_3 \end{array} \right. \quad \begin{array}{l} X_1 \perp\!\!\!\perp X_2 \\ X_2 \perp\!\!\!\perp X_3 \\ X_1 \perp\!\!\!\perp X_3 \\ X_1 \perp\!\!\!\perp X_3 \mid X_2 \end{array}$$



Perfect maps - existence

- Not every distribution p has a perfect map



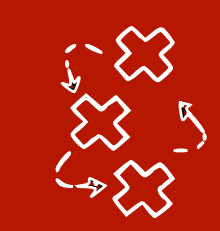
$$\left\{ \begin{array}{l} X_1 = \epsilon_1 \\ X_2 = X_1 \\ X_3 = X_2 + \epsilon_3 \\ \epsilon_1, \epsilon_3 \sim N(0,1) \\ \epsilon_1 \perp\!\!\!\perp \epsilon_3 \end{array} \right.$$

Deterministic function

$$\left\{ \begin{array}{l} X_1 = \epsilon_1 \\ X_2 = \epsilon_1 \\ X_3 = \epsilon_1 + \epsilon_3 \\ \epsilon_1, \epsilon_3 \sim N(0,1) \\ \epsilon_1 \perp\!\!\!\perp \epsilon_3 \end{array} \right. \quad \begin{array}{l} X_1 \perp\!\!\!\perp X_2 \\ X_2 \perp\!\!\!\perp X_3 \\ X_1 \perp\!\!\!\perp X_3 \\ X_1 \perp\!\!\!\perp X_3 \mid X_2 \end{array}$$

$$p(\epsilon_1 \mid \epsilon_1, \epsilon_1 + \epsilon_3) = p(\epsilon_1 \mid \epsilon_1)$$

$$X_2 \perp\!\!\!\perp X_3 \mid X_1$$



Perfect maps - existence

- Not every distribution p has a perfect map

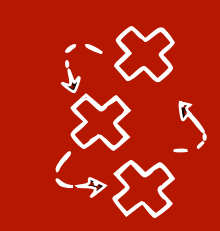


$$\left\{ \begin{array}{l} X_1 = \epsilon_1 \\ X_2 = X_1 \\ X_3 = X_2 + \epsilon_3 \\ \epsilon_1, \epsilon_3 \sim N(0,1) \\ \epsilon_1 \perp\!\!\!\perp \epsilon_3 \end{array} \right. \quad \boxed{\text{Deterministic function}}$$

$$\left\{ \begin{array}{l} X_1 = \epsilon_1 \\ X_2 = \epsilon_1 \\ X_3 = \epsilon_1 + \epsilon_3 \\ \epsilon_1, \epsilon_3 \sim N(0,1) \\ \epsilon_1 \perp\!\!\!\perp \epsilon_3 \end{array} \right. \quad \begin{array}{l} X_1 \perp\!\!\!\perp X_2 \\ X_2 \perp\!\!\!\perp X_3 \\ X_1 \perp\!\!\!\perp X_3 \\ X_1 \perp\!\!\!\perp X_3 \mid X_2 \end{array}$$

$$X_2 \perp_d X_3 \mid X_1 ?$$

$$X_2 \perp\!\!\!\perp X_3 \mid X_1$$



Perfect maps - existence

- Not every distribution p has a perfect map



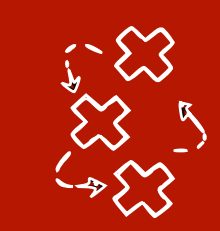
$$\left\{ \begin{array}{l} X_1 = \epsilon_1 \\ X_2 = X_1 \\ X_3 = X_2 + \epsilon_3 \\ \epsilon_1, \epsilon_3 \sim N(0,1) \\ \epsilon_1 \perp\!\!\!\perp \epsilon_3 \end{array} \right.$$

Deterministic function

$$\left\{ \begin{array}{l} X_1 = \epsilon_1 \\ X_2 = \epsilon_1 \\ X_3 = \epsilon_1 + \epsilon_3 \\ \epsilon_1, \epsilon_3 \sim N(0,1) \\ \epsilon_1 \perp\!\!\!\perp \epsilon_3 \end{array} \right. \quad \begin{array}{l} X_1 \perp\!\!\!\perp X_2 \\ X_2 \perp\!\!\!\perp X_3 \\ X_1 \perp\!\!\!\perp X_3 \\ X_1 \perp\!\!\!\perp X_3 \mid X_2 \\ X_2 \perp\!\!\!\perp X_3 \mid X_1 \end{array}$$

$$X_2 \perp\!\!\!\perp_d X_3 \mid X_1$$

$$X_2 \perp\!\!\!\perp X_3 \mid X_1$$



Perfect maps - existence

- Not every distribution p has a perfect map



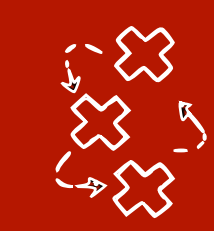
There exists no Bayesian network that can represent these conditional in/dependences as d-separations perfectly!

$$\left\{ \begin{array}{l} X_1 = \epsilon_1 \\ X_2 = X_1 \\ X_3 = X_2 + \epsilon_3 \\ \epsilon_1, \epsilon_3 \sim N(0,1) \\ \epsilon_1 \perp\!\!\!\perp \epsilon_3 \end{array} \right. \quad \boxed{\text{Deterministic function}}$$

$$\left\{ \begin{array}{l} X_1 = \epsilon_1 \\ X_2 = \epsilon_1 \\ X_3 = \epsilon_1 + \epsilon_3 \\ \epsilon_1, \epsilon_3 \sim N(0,1) \\ \epsilon_1 \perp\!\!\!\perp \epsilon_3 \end{array} \right. \quad \begin{array}{l} X_1 \perp\!\!\!\perp X_2 \\ X_2 \perp\!\!\!\perp X_3 \\ X_1 \perp\!\!\!\perp X_3 \\ X_1 \perp\!\!\!\perp X_3 \mid X_2 \\ X_2 \perp\!\!\!\perp X_3 \mid X_1 \end{array}$$

$$X_2 \not\perp\!\!\!\perp_d X_3 \mid X_1$$

$$X_2 \perp\!\!\!\perp X_3 \mid X_1$$



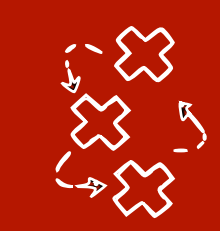
Markov equivalence class (MEC)

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

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

$$\mathbf{A} \perp_d \mathbf{B} \mid \mathbf{C} \iff X_{\mathbf{A}} \perp\!\!\!\perp X_{\mathbf{B}} \mid X_{\mathbf{C}}$$

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



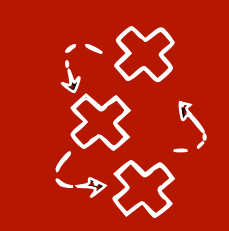
Question: Markov equivalence class

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



Which graphs on X, Y, Z are perfect maps of these conditional independences?

Hint: we can start by orienting this undirected graph



Question: Markov equivalence class

$X \not\perp_d Y$ $X \perp_d Y | Z$

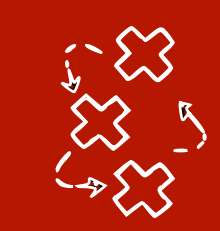
Which graphs on X, Y, Z are perfect maps of these conditional independences?

$X - Z - Y$

$X \rightarrow Z \rightarrow Y$

$X \not\perp_d Y \checkmark$

$X \perp_d Y | Z \checkmark$



Question: Markov equivalence class

$X \not\perp_d Y$ $X \perp_d Y | Z$

Which graphs on X, Y, Z are perfect maps of these conditional independences?

$X - Z - Y$

$X \rightarrow Z \rightarrow Y$

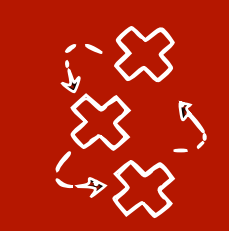
$X \leftarrow Z \leftarrow Y$

$X \not\perp_d Y \quad \checkmark$

$X \perp_d Y | Z \quad \checkmark$

$X \not\perp_d Y \quad \checkmark$

$X \perp_d Y | Z \quad \checkmark$



Question: Markov equivalence class

$X \not\perp_d Y$ $X \perp_d Y | Z$

Which graphs on X, Y, Z are perfect maps of these conditional independences?



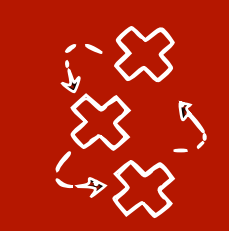
$X \not\perp_d Y \checkmark$
 $X \perp_d Y | Z \checkmark$



$X \not\perp_d Y \checkmark$
 $X \perp_d Y | Z \checkmark$

$X \not\perp_d Y \checkmark$

$X \perp_d Y | Z \checkmark$



Question: Markov equivalence class

$X \not\perp_d Y$ $X \perp_d Y | Z$

Which graphs on X, Y, Z are perfect maps of these conditional independences?



$X \not\perp_d Y \checkmark$
 $X \perp_d Y | Z \checkmark$



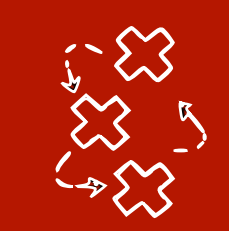
$X \not\perp_d Y \checkmark$
 $X \perp_d Y | Z \checkmark$



$X \not\perp_d Y \checkmark$
 $X \perp_d Y | Z \checkmark$



$X \perp_d Y \times$
 $X \not\perp_d Y | Z \times$



Question: Markov equivalence class

$X \not\perp_d Y$ $X \perp_d Y | Z$

Which graphs on X, Y, Z are perfect maps of these conditional independences?



$X \not\perp_d Y \checkmark$
 $X \perp_d Y | Z \checkmark$

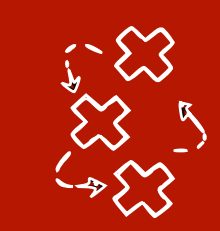


$X \not\perp_d Y \checkmark$
 $X \perp_d Y | Z \checkmark$



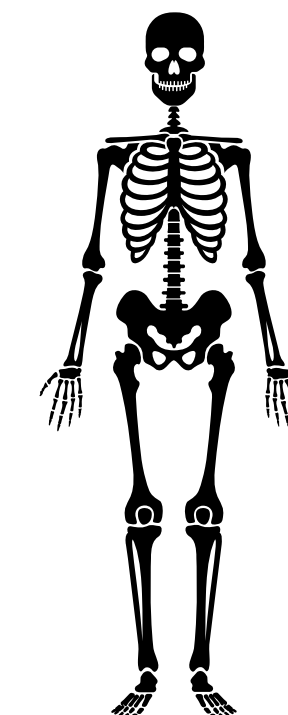
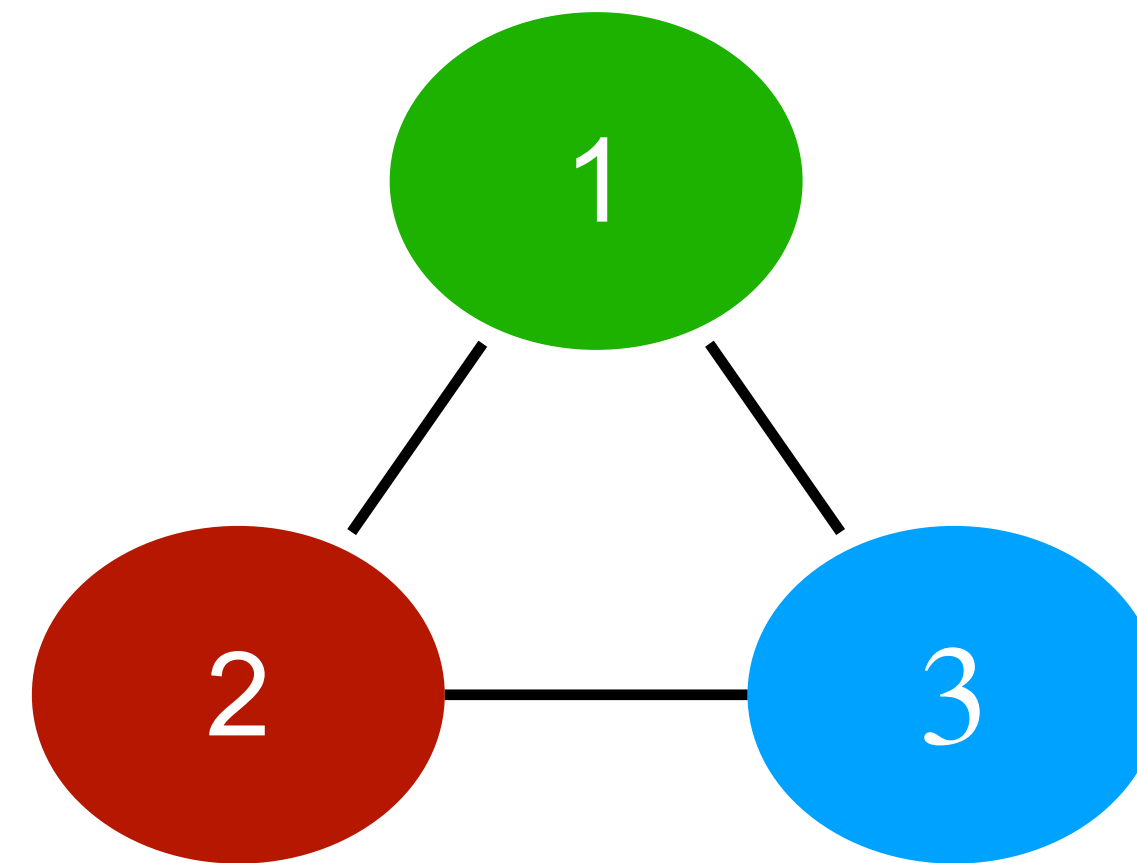
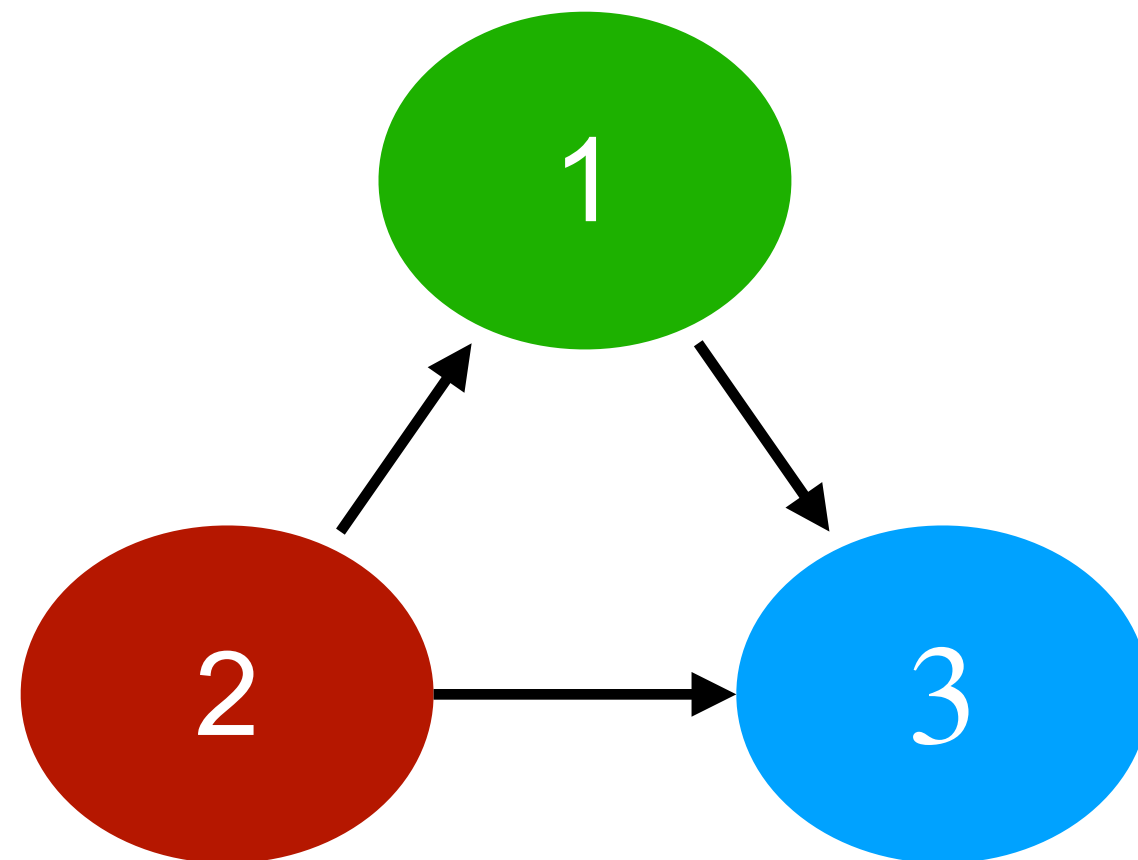
$X \not\perp_d Y \checkmark$
 $X \perp_d Y | Z \checkmark$

These three graphs represent a Markov Equivalence Class

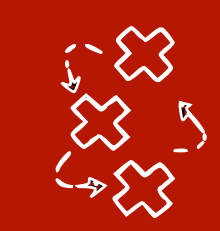


Graph terminology: skeletons

- The skeleton of a DAG $G = (\mathbf{V}, \mathbf{E})$ is the undirected graph $U = (\mathbf{V}, \mathbf{E}')$ that has **an undirected edge** $(i, j) \in \mathbf{E}'$ for every **directed edge** $i \rightarrow j \in \mathbf{E}$ (and no other edges)

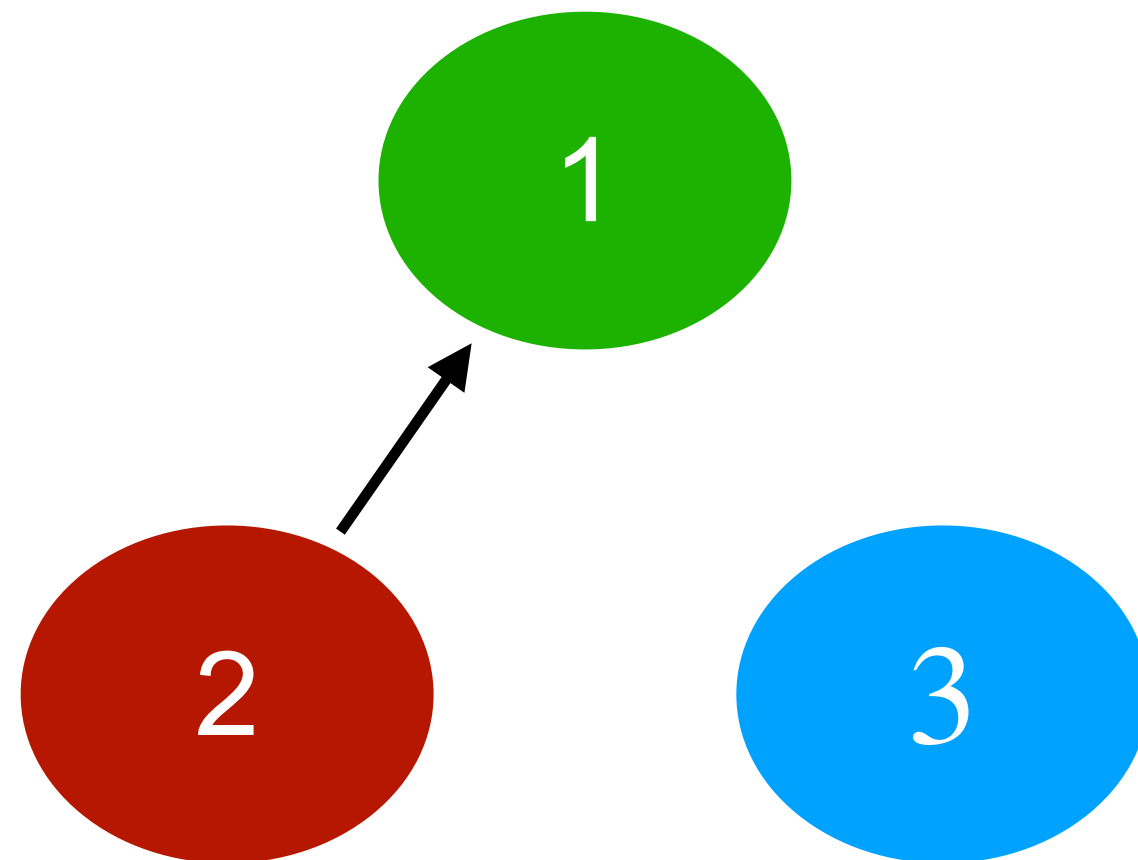


In other words a graph that is the same, but without the arrowheads

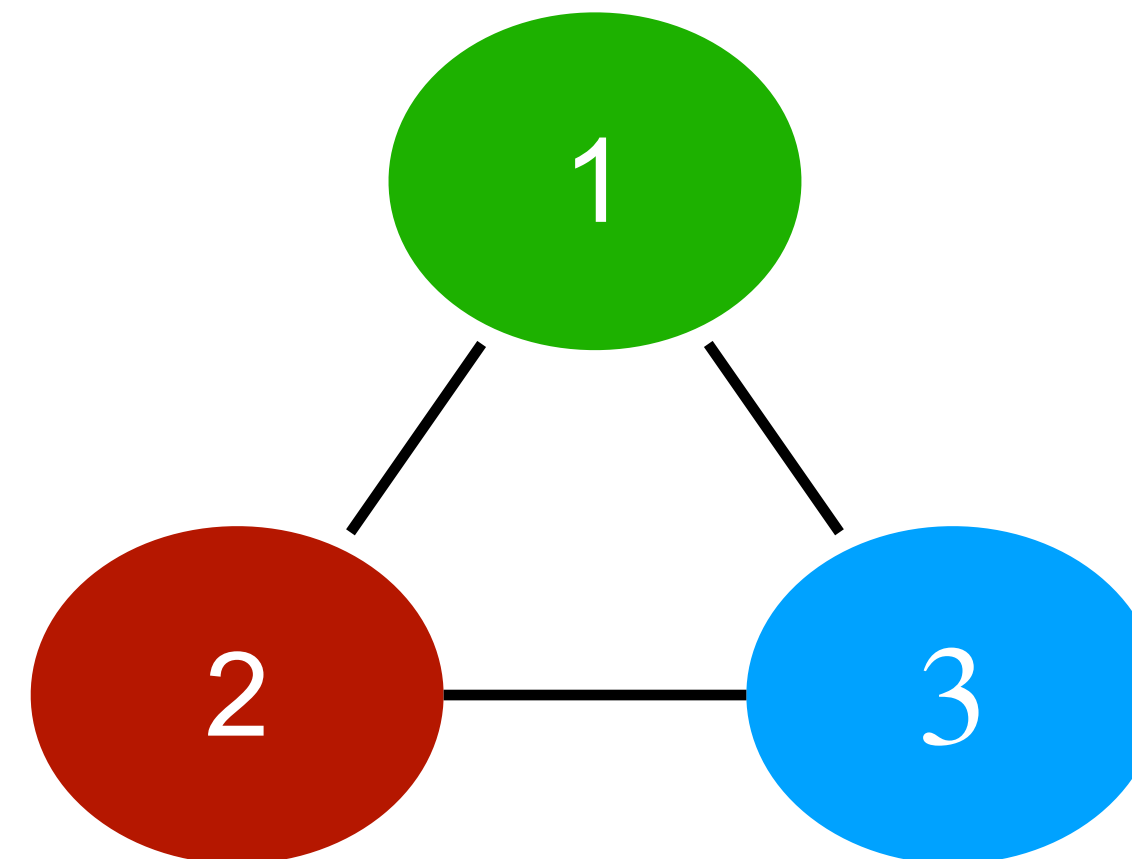


Skeleton exercise

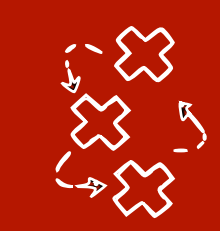
- The skeleton of a DAG $G = (\mathbf{V}, \mathbf{E})$ is the undirected graph $U = (\mathbf{V}, \mathbf{E}')$ that has **an undirected edge** $(i, j) \in \mathbf{E}'$ for every **directed edge** $i \rightarrow j \in \mathbf{E}$ and no other edges



G

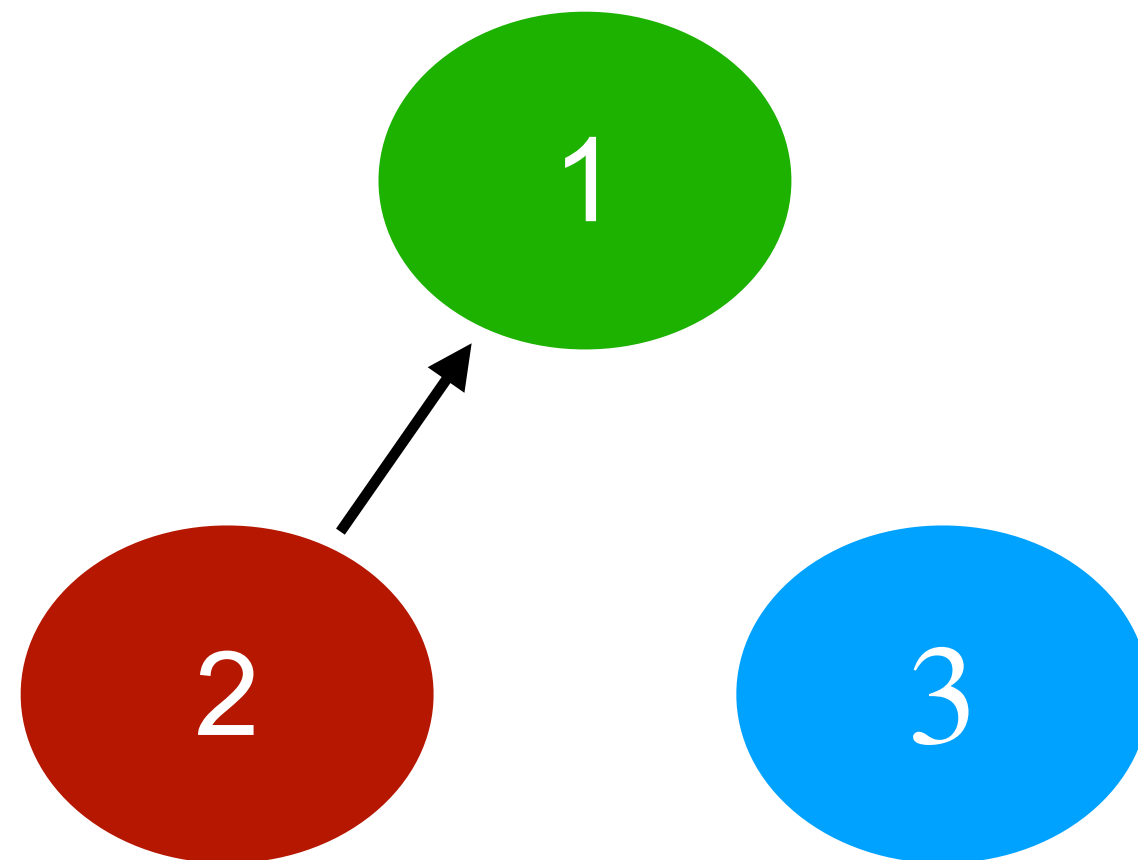


Is this the skeleton of G?

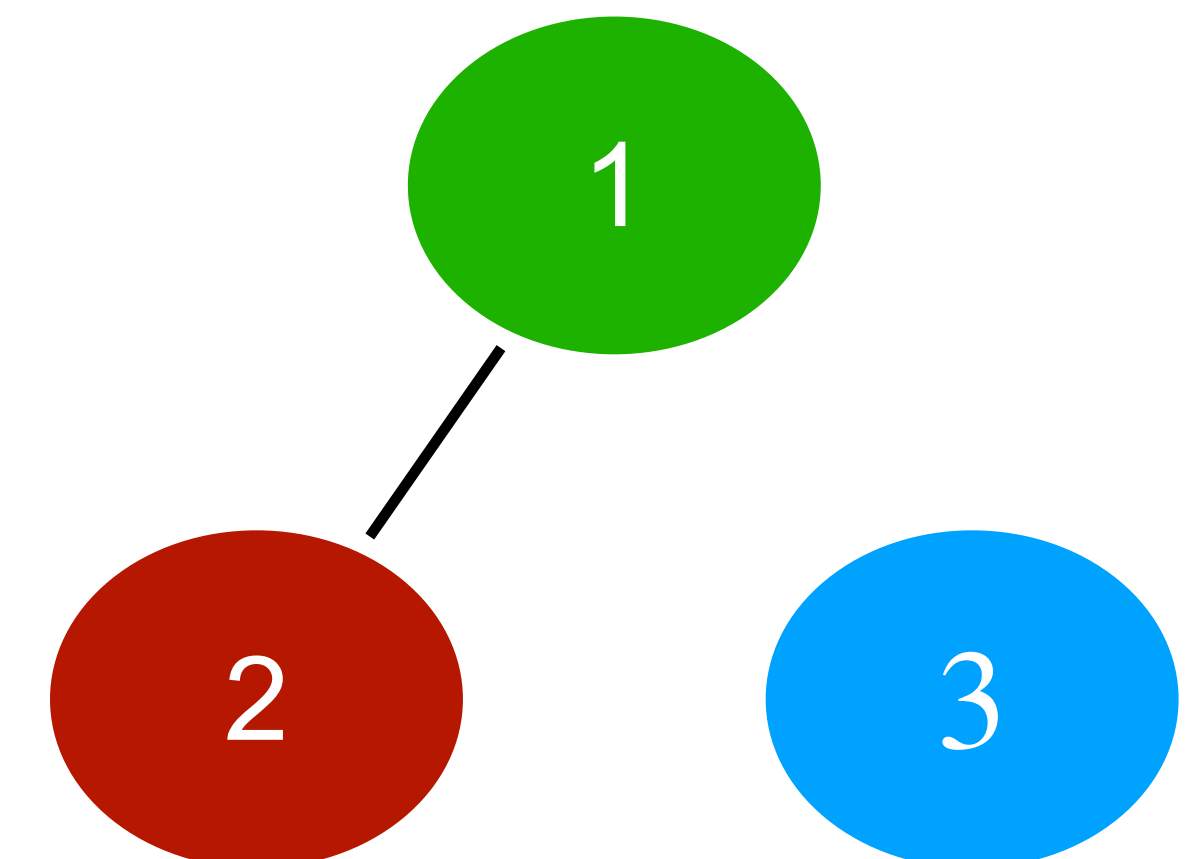
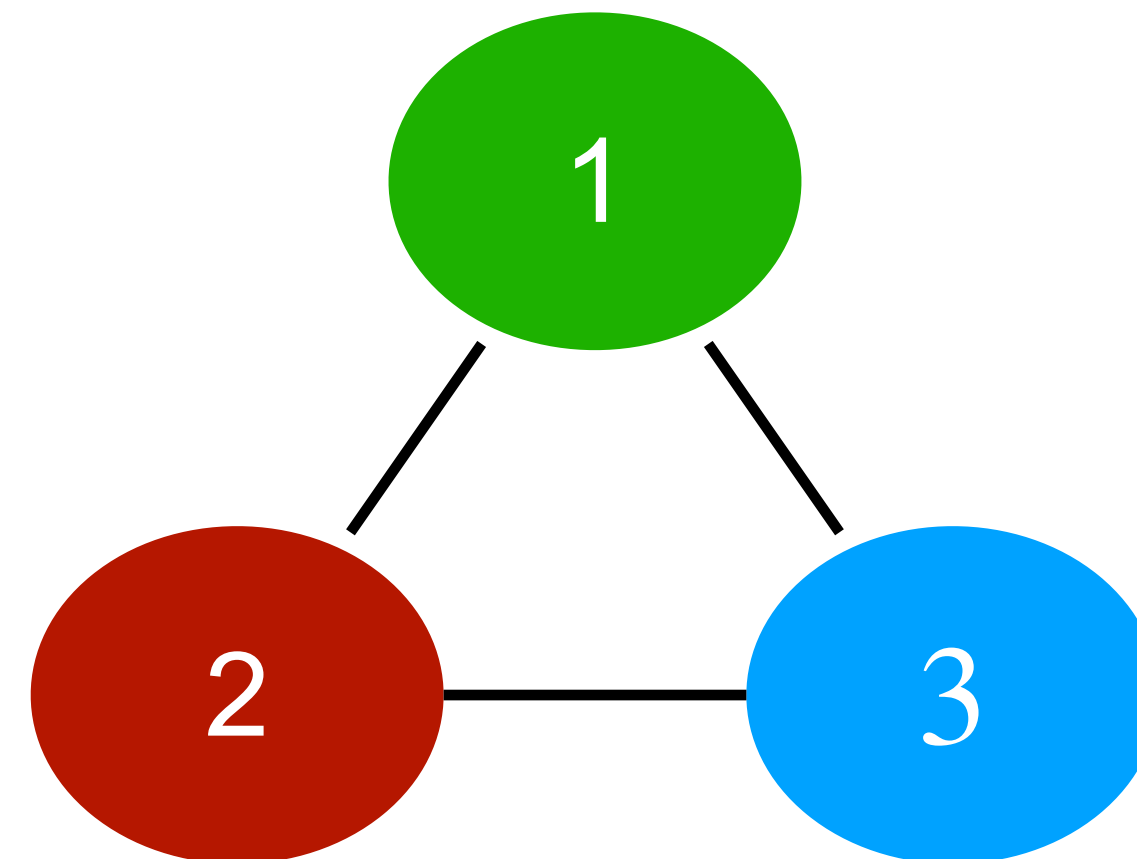


Skeleton exercise

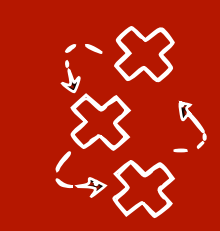
- The skeleton of a DAG $G = (\mathbf{V}, \mathbf{E})$ is the undirected graph $U = (\mathbf{V}, \mathbf{E}')$ that has **an undirected edge** $(i, j) \in \mathbf{E}'$ for every **directed edge** $i \rightarrow j \in \mathbf{E}$ and no other edges



G

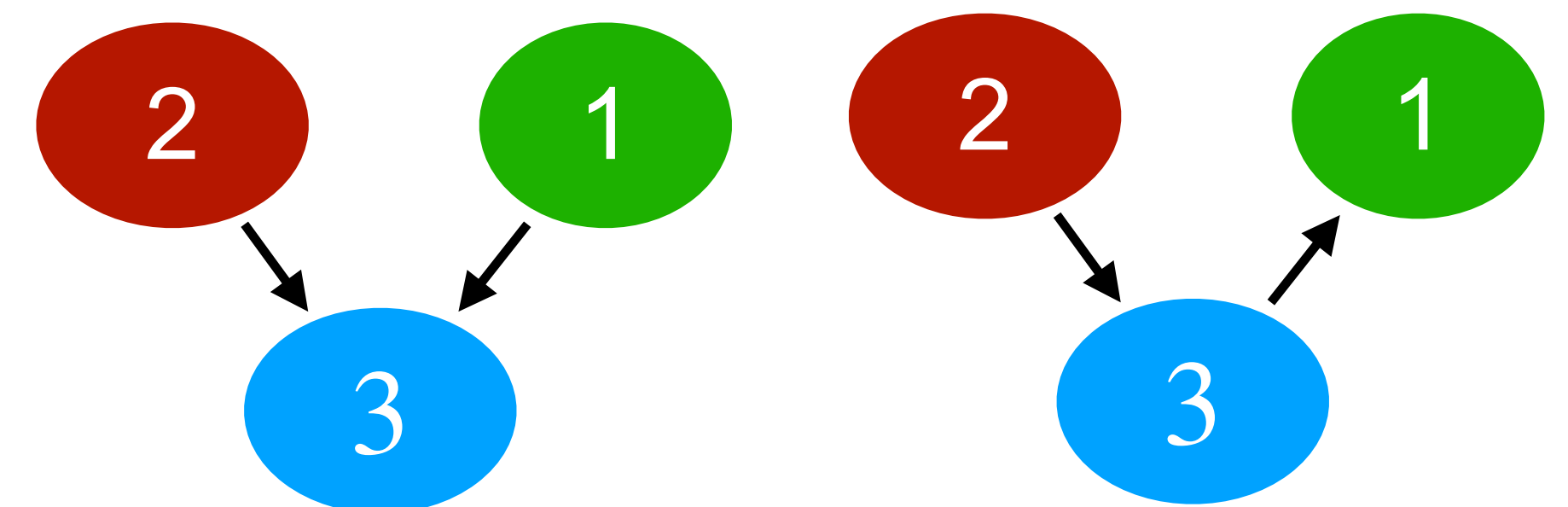


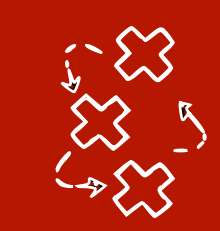
Is this the skeleton of G?



More graph terminology: adjacent

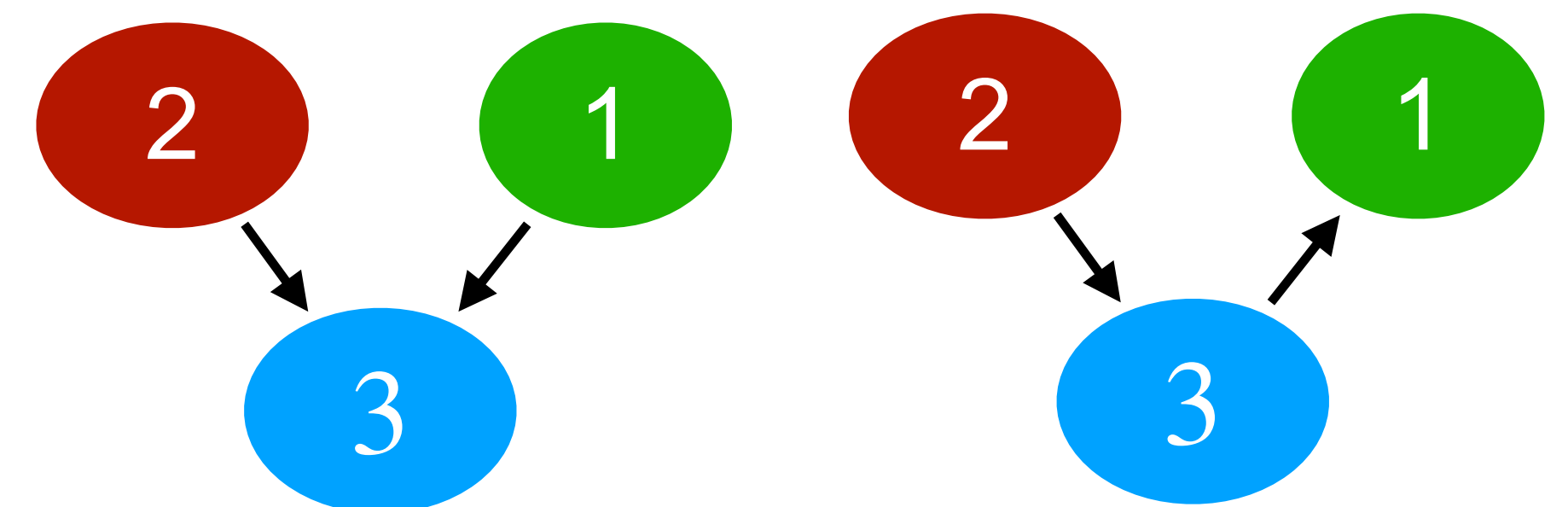
- Nodes i and j in a DAG G are **adjacent/neighbours** if $i \rightarrow j$ or $j \rightarrow i$ in G

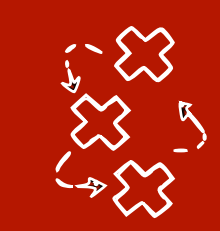




More graph terminology: adjacent

- Nodes i and j in a DAG G are **adjacent/neighbours** if $i \rightarrow j$ or $j \rightarrow i$ in G
 - i.e., they are connected by an undirected edge in the skeleton of G
 - We denote adjacency with $i - j$, while $i \not- j$ means non-adjacent

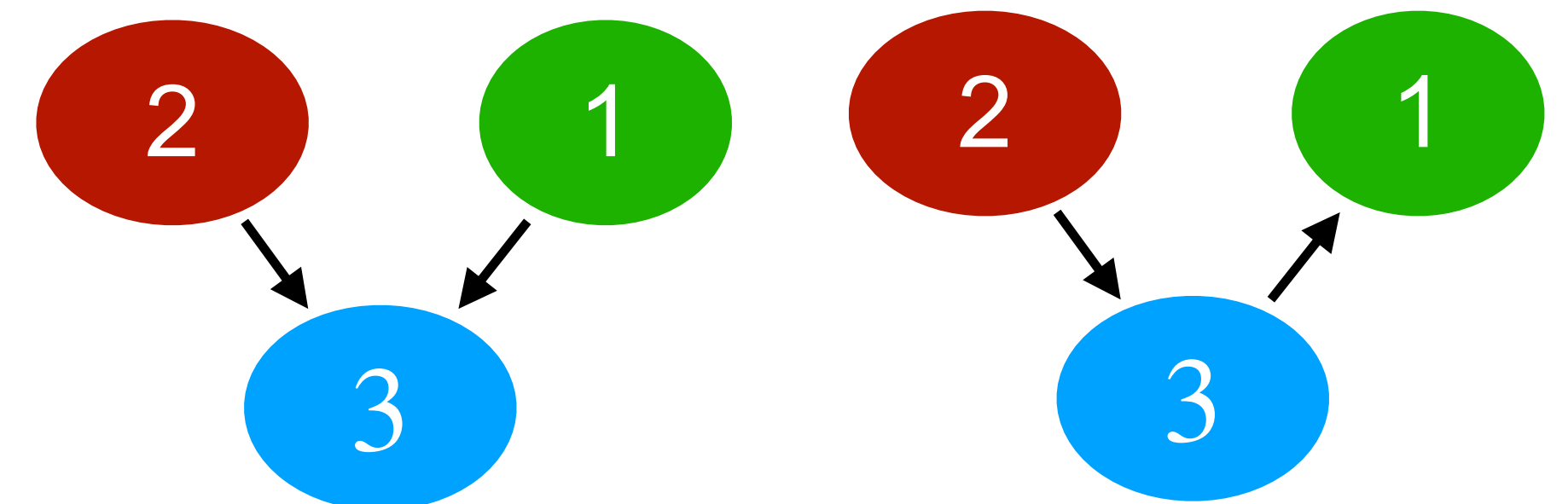


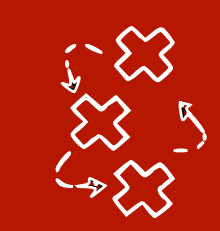


More graph terminology: adjacent

- Nodes i and j in a DAG G are **adjacent/neighbours** if $i \rightarrow j$ or $j \rightarrow i$ in G
 - i.e., they are connected by an undirected edge in the skeleton of G
 - We denote adjacency with $i - j$, while $i \not- j$ means non-adjacent

Is 2 adjacent to 3?



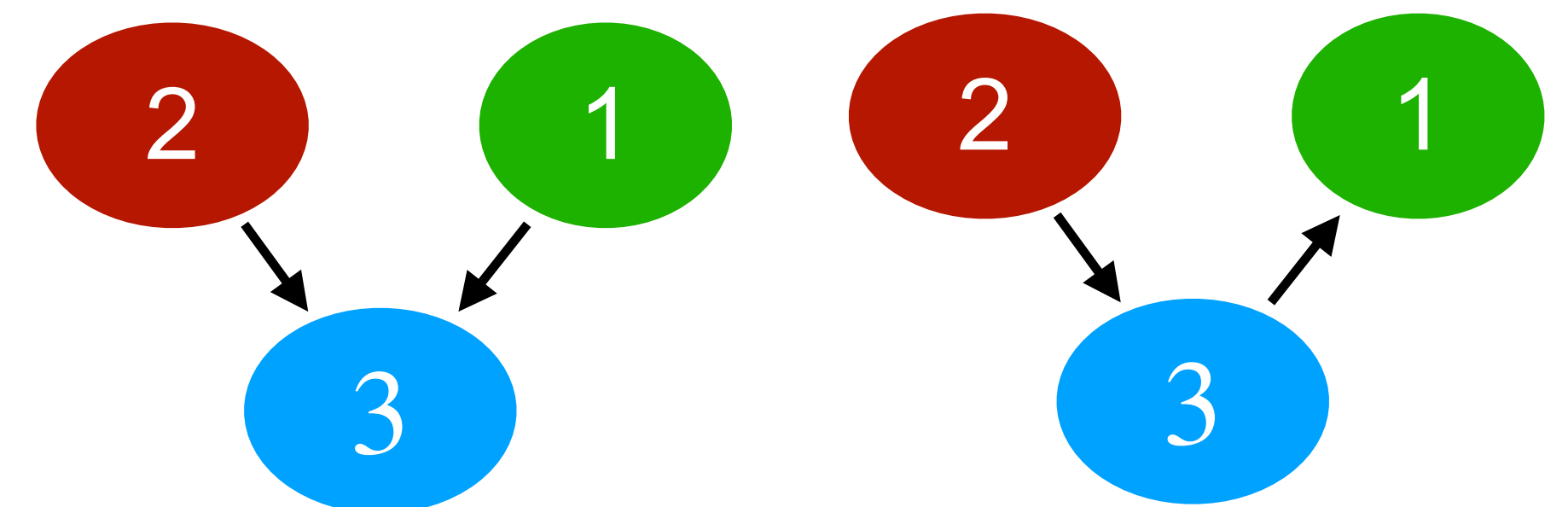


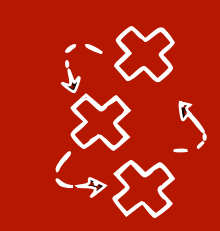
More graph terminology: adjacent

- Nodes i and j in a DAG G are **adjacent/neighbours** if $i \rightarrow j$ or $j \rightarrow i$ in G
 - I.e., they are connected by an undirected edge in the skeleton of G
 - We denote adjacency with $i - j$, while $i \not- j$ means non-adjacent

Is 2 adjacent to 3?

Is 2 adjacent to 1?



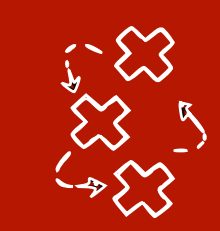


More graph terminology: v-structures/immoralities

- A triple of nodes (i, j, k) in a DAG G is a **v-structure (unshielded collider)** if $i \rightarrow j \leftarrow k$ in G and i is not adjacent to k

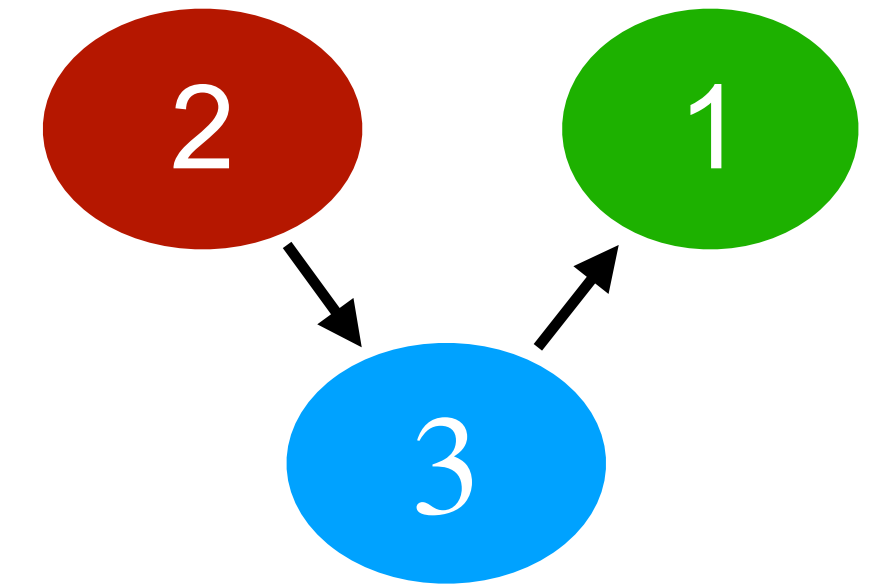
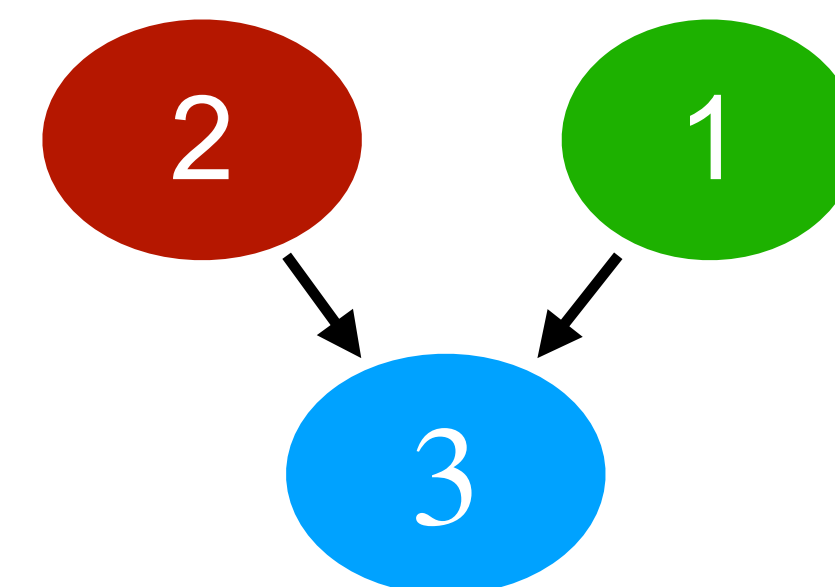
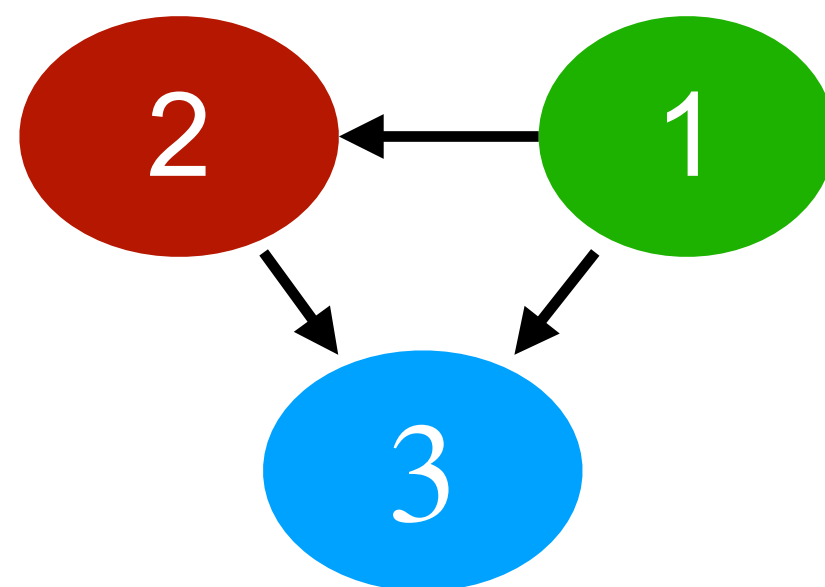
An edge from i to k would be called a shield

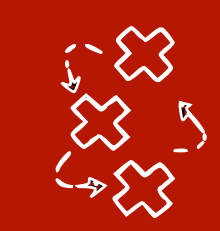
It's also called an immorality, because the two parents i and j who have a common child k are not “married” (adjacent)



More graph terminology: v-structures/immoralities

- A triple of nodes (i, j, k) in a DAG G is a **v-structure (unshielded collider)** if $i \rightarrow j \leftarrow k$ in G and i is not adjacent to k

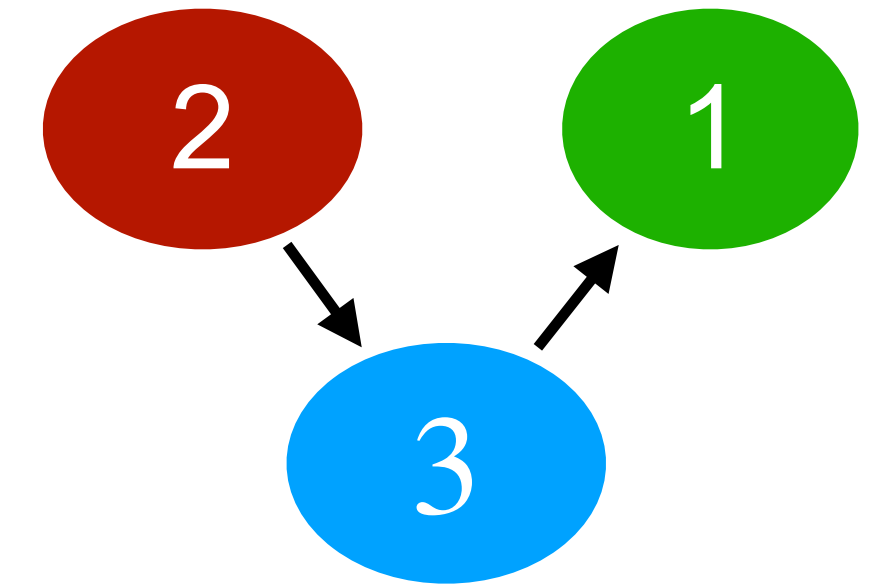
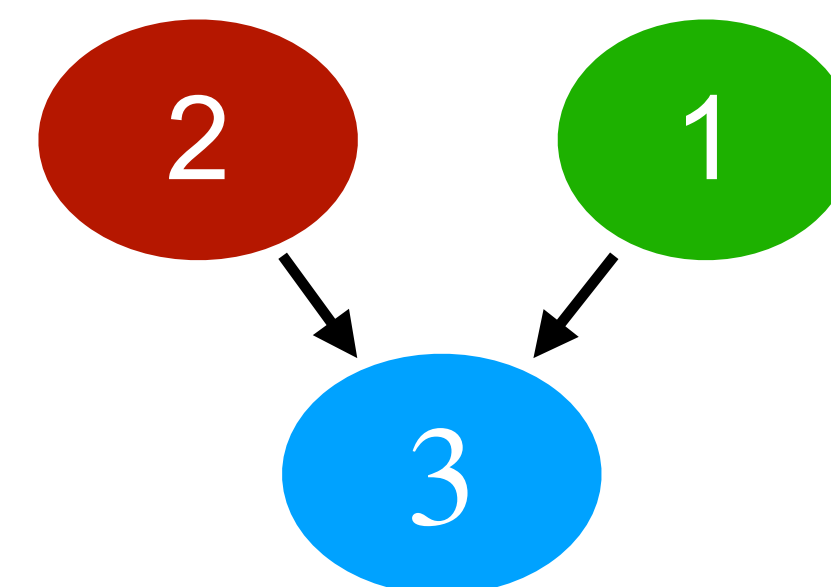
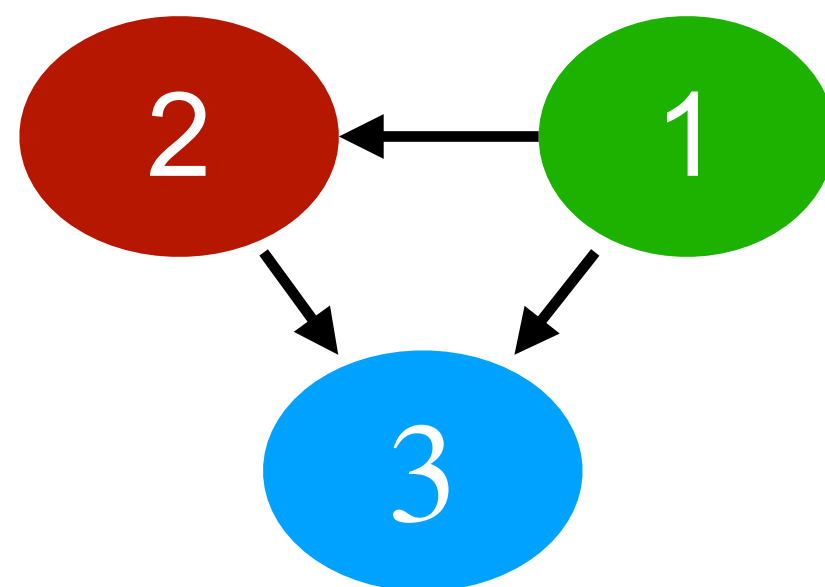


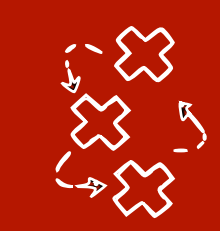


More graph terminology: v-structures/immoralities

- A triple of nodes (i, j, k) in a DAG G is a **v-structure (unshielded collider)** if $i \rightarrow j \leftarrow k$ in G and i is not adjacent to k

In which graphs is 3 a collider on the path between 2 and 1?

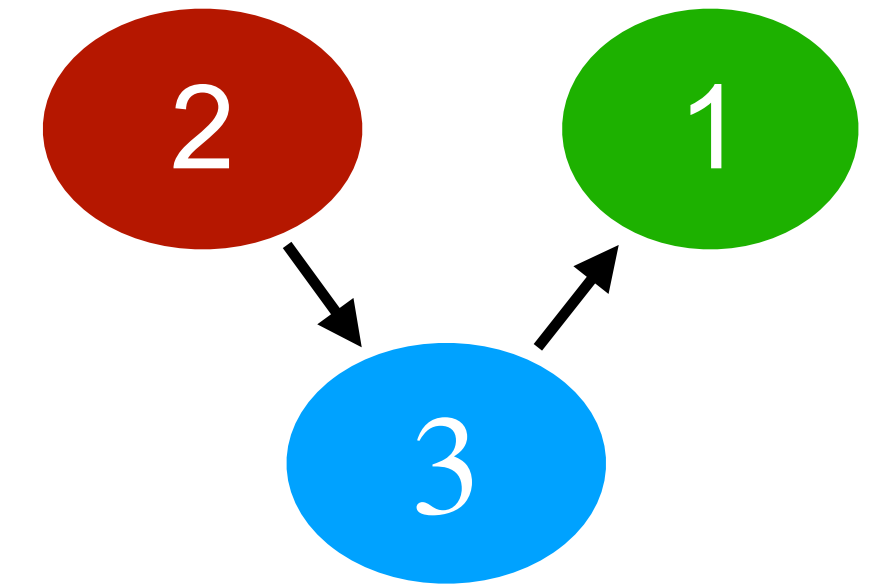
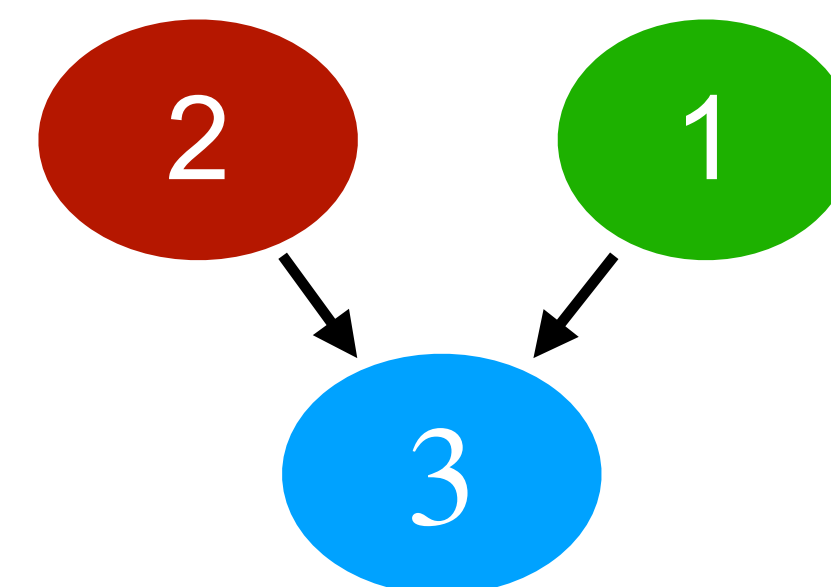
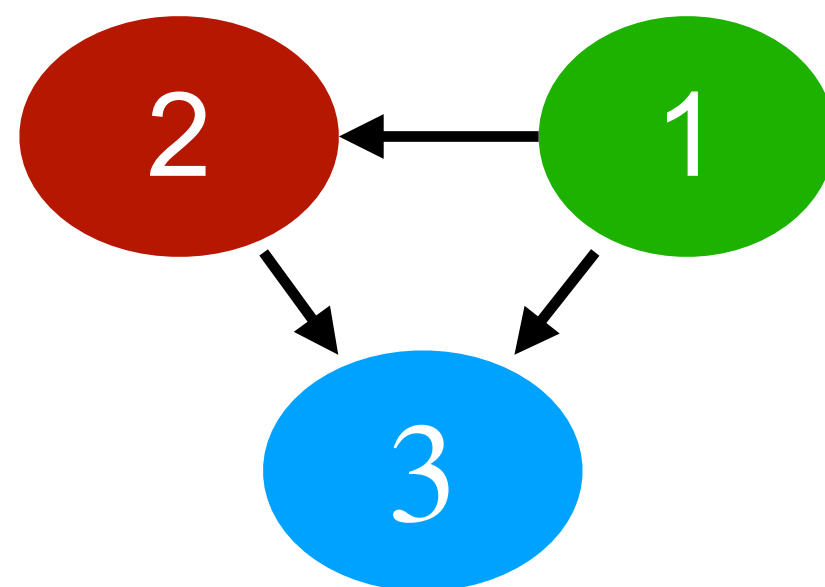


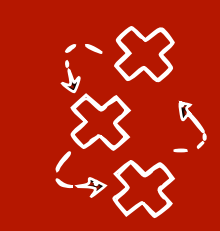


More graph terminology: v-structures/immoralities

- A triple of nodes (i, j, k) in a DAG G is a **v-structure (unshielded collider)** if $i \rightarrow j \leftarrow k$ in G and i is not adjacent to k

In which graphs is 3 part of a v-structure (2,3,1)?

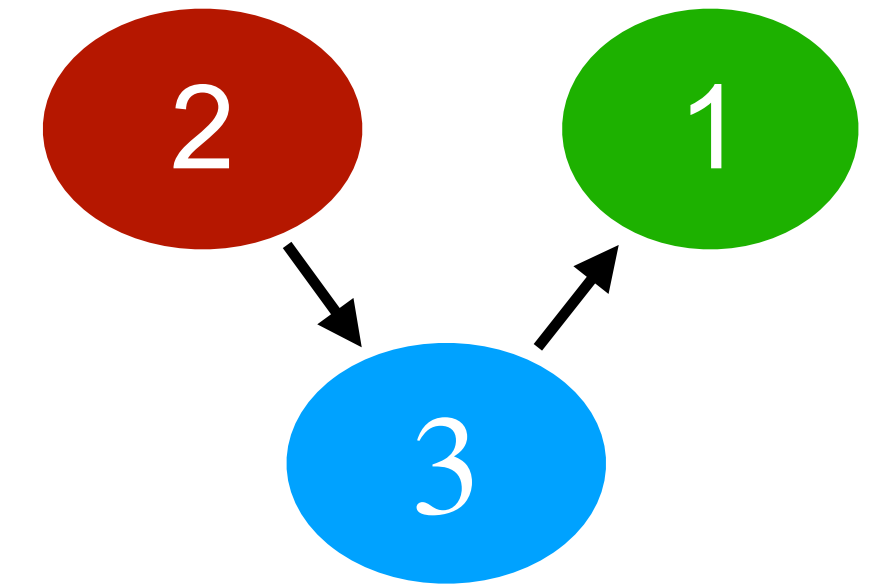
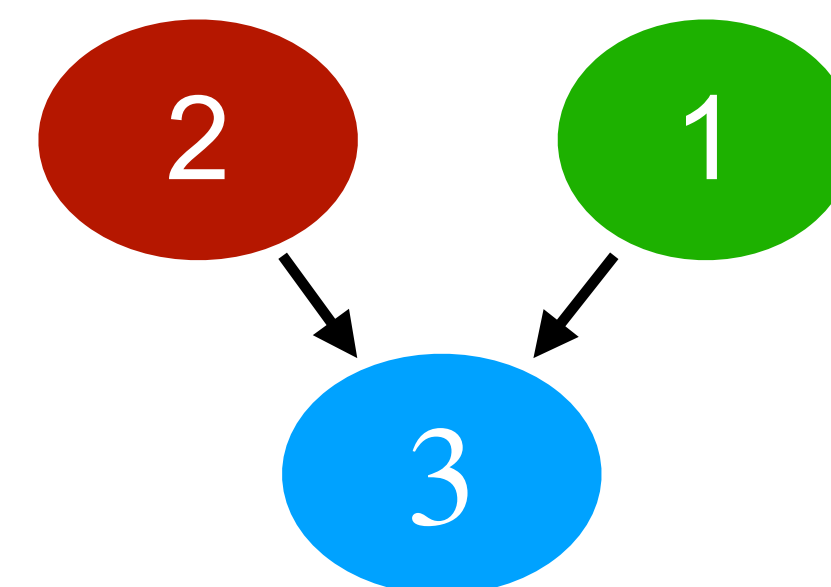
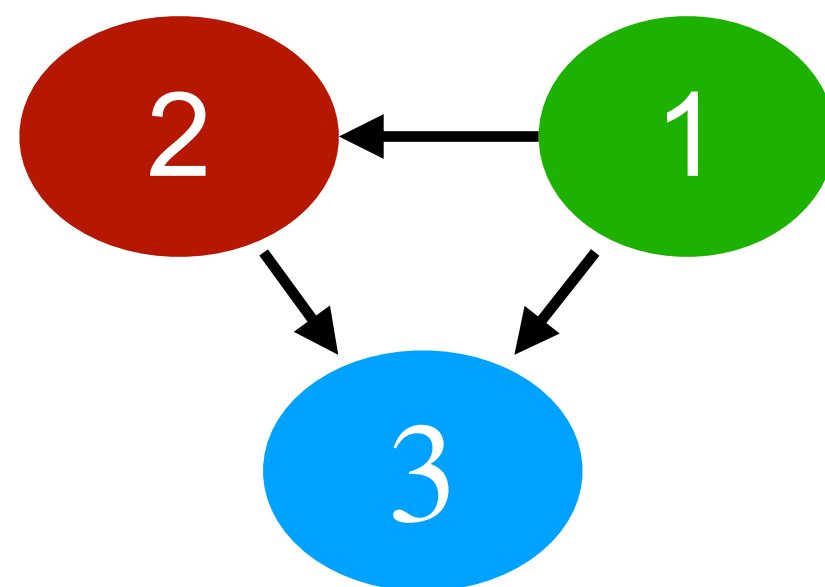


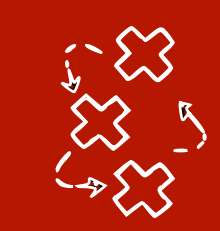


More graph terminology: unshielded triple

- A triple of nodes (i, j, k) in a DAG G is a **v-structure (unshielded collider)** if $i \rightarrow j \leftarrow k$ in G and i is not adjacent to k in G
- A triple of nodes (i, j, k) in a DAG G is a **an unshielded triple** if $i - j, j - k$ and i is not adjacent to k in G

In which graphs is $(2,3,1)$ an unshielded triple?

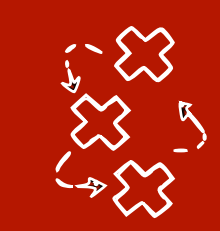




Markov equivalence class and CPDAGs

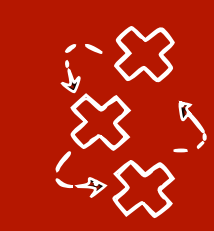
essential graphs *Summary graphs*

- (Verma and Pearl 1990) show that all DAGs in a Markov equivalence class have the **same skeleton** and the **same v-structures**



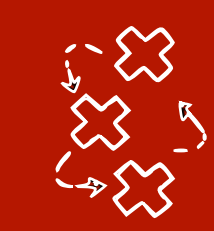
Markov equivalence class and CPDAGs

- (Verma and Pearl 1990) show that all DAGs in a Markov equivalence class have the **same skeleton** and the **same v-structures**
- We can represent the skeleton and the orientations (edge marks) all DAGs in a Markov equivalence class (MEC) have in common with a **Complete Partially Directed Acyclic Graph (CPDAG):**



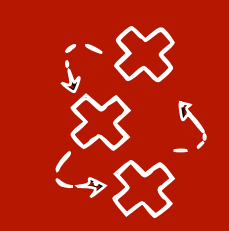
Markov equivalence class and CPDAGs

- (Verma and Pearl 1990) show that all DAGs in a Markov equivalence class have the **same skeleton** and the **same v-structures**
- We can represent the skeleton and the orientations (edge marks) all DAGs in a Markov equivalence class (MEC) have in common with a **Complete Partially Directed Acyclic Graph (CPDAG):**
 - We have a **directed** edge $i \rightarrow j$ if **all DAGs** in the MEC have $i \rightarrow j$



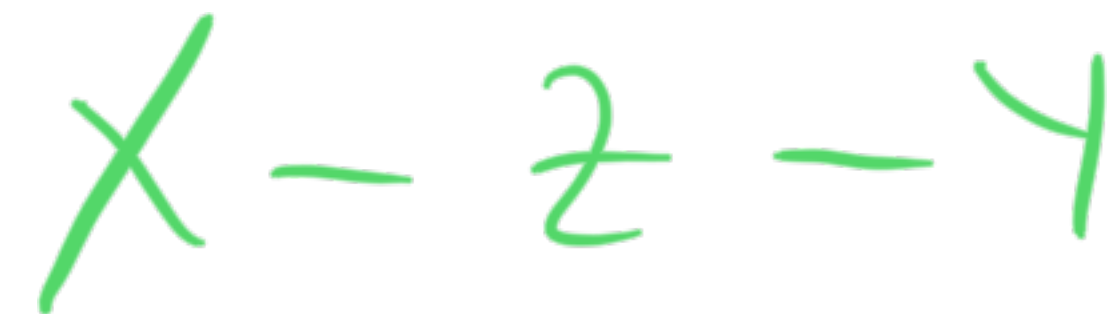
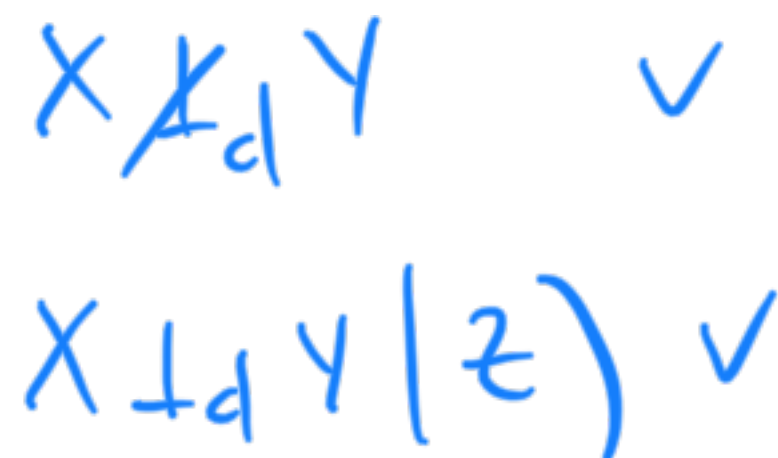
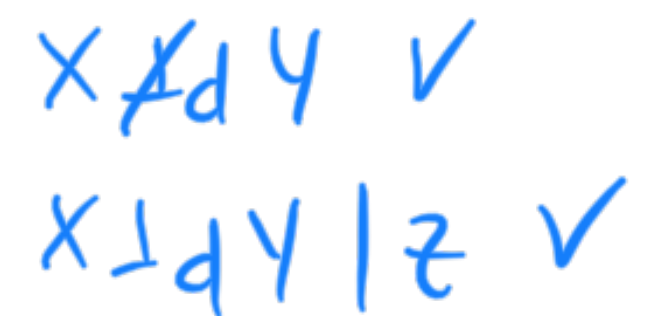
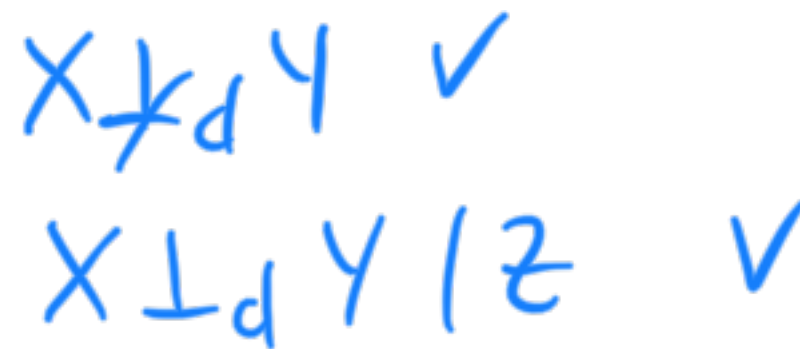
Markov equivalence class and CPDAGs

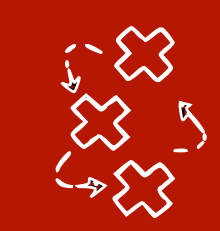
- (Verma and Pearl 1990) show that all DAGs in a Markov equivalence class have the **same skeleton** and the **same v-structures**
- We can represent the skeleton and the orientations (edge marks) all DAGs in a Markov equivalence class (MEC) have in common with a **Complete Partially Directed Acyclic Graph (CPDAG)**:
 - We have a **directed** edge $i \rightarrow j$ if **all DAGs** in the MEC have $i \rightarrow j$
 - We have an **undirected** edge $i - j$ if **some DAGs** in the MEC have $i \rightarrow j$ and others have $j \rightarrow i$



Markov equivalence class and CPDAGs

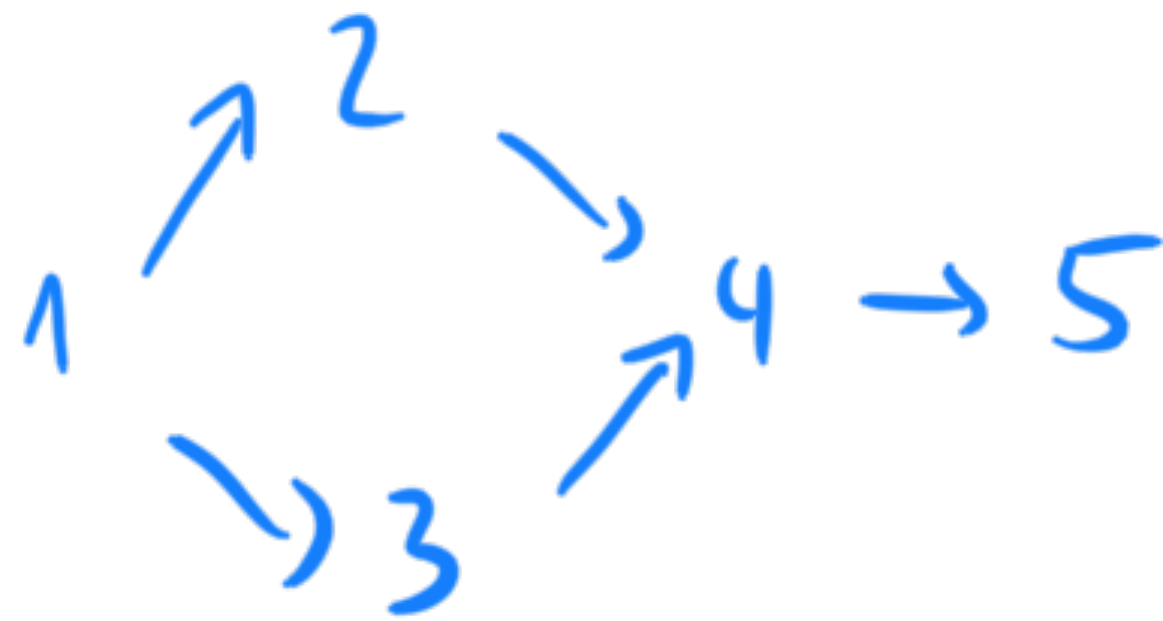
- **Complete Partially Directed Acyclic Graph (CPDAG):** *Summary graphs*
 - We have a directed edge $i \rightarrow j$ if all DAGs in the MEC have $i \rightarrow j$
 - We have an undirected edge $i - j$ if some DAGs have $i \rightarrow j$ and others $j \rightarrow i$



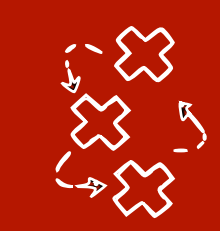


CPDAG question

- (Verma and Pearl 1990) show that all DAGs in a Markov equivalence class have the **same skeleton** and the **same v-structures**

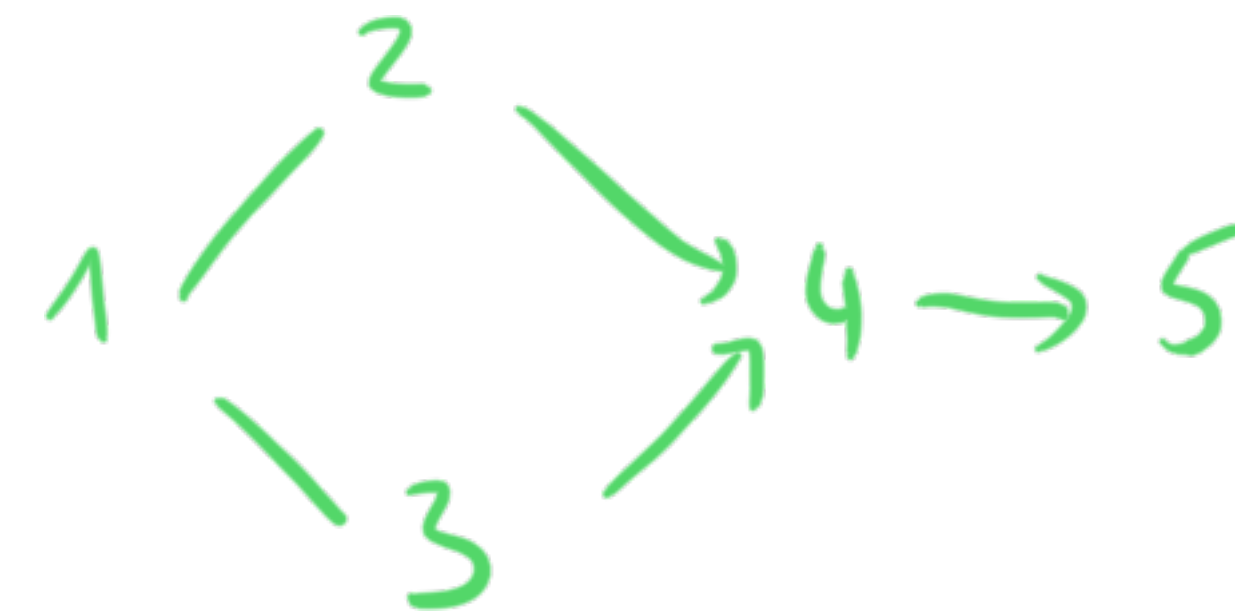
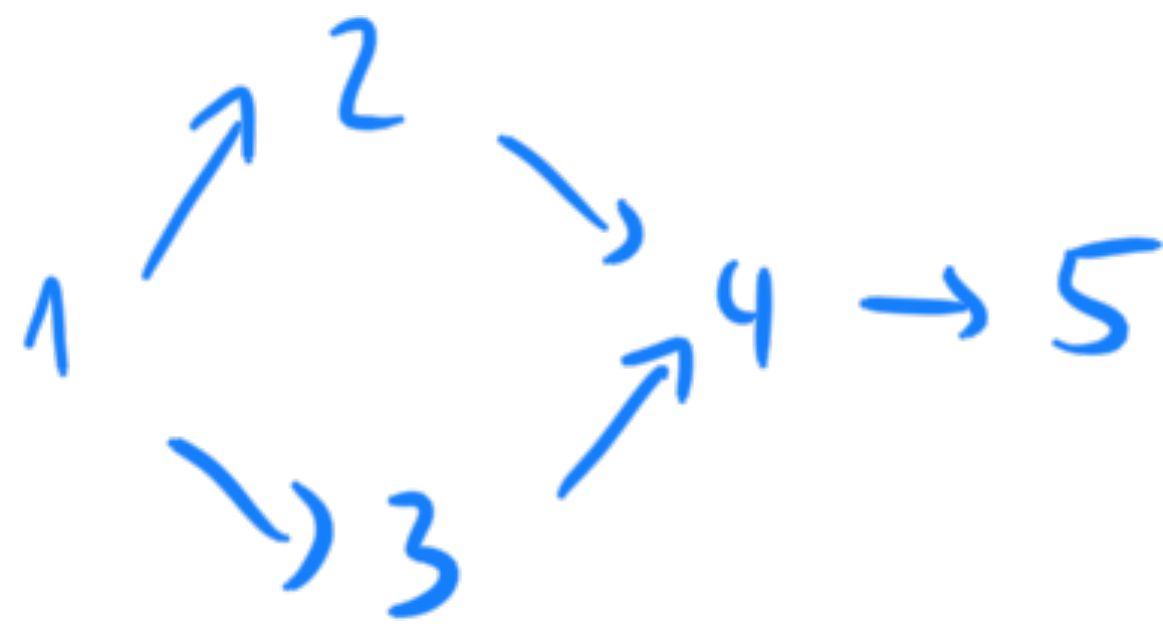


If the true graph is known we can compute the **CPDAG** representing the MEC



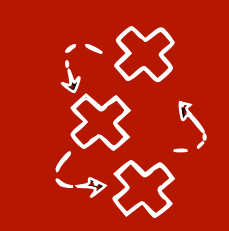
CPDAG question

- (Verma and Pearl 1990) show that all DAGs in a Markov equivalence class have the **same skeleton** and the **same v-structures**



If the true graph is known we can compute
the CPDAG representing the MEC

What if it isn't known?!

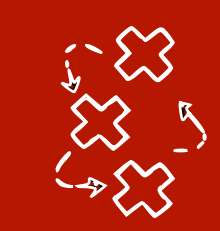


Constraint-based causal discovery

- **Idea:** we perform conditional independence tests on **observational** data and use them to constrain the possible graphs using d-separation

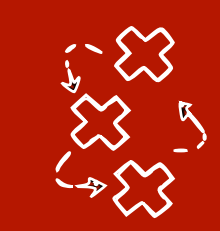
$X \not\perp\!\!\!\perp Y$ $X \perp\!\!\!\perp Y | Z$

$X - Z - Y$



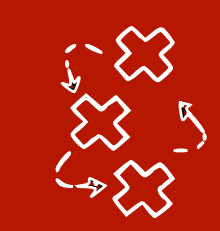
Constraint-based causal discovery

- **Idea:** we perform conditional independence tests on **observational** 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)**



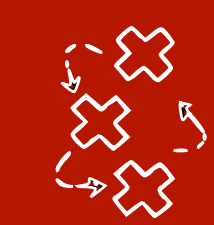
Constraint-based causal discovery

- **Idea:** we perform conditional independence tests on **observational** 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**



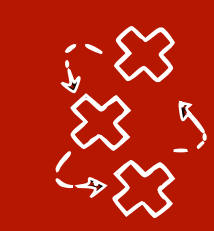
SGS algorithm (Spirtes, Glymour, Scheines)

- Assuming p is Markov and faithful to an unknown graph G
- We can estimate a CPDAG from samples of p in three steps:
 1. Determine the **skeleton**
 2. Determine the **v-structures**
 3. Direct as many remaining edges as possible



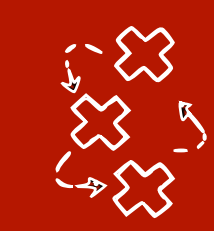
SGS algorithm (Spirtes, Glymour, Scheines)

- Assuming p is Markov and faithful to an unknown graph G
- We can estimate a CPDAG from samples of p in three steps:
 1. Determine the **skeleton**
 2. Determine the **v-structures** (*given the tests in the previous phase*)
 3. Direct as many remaining edges as possible
- **Note**: the directed parts of the CPDAG will agree with G , but some parts might stay undirected



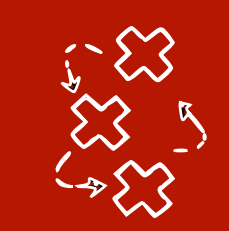
Step 1: Skeleton learning

- Given $G = (\mathbf{V}, \mathbf{E})$, nodes $i, j \in \mathbf{V}, i \neq j$ then:
 - If i is adjacent to j , they cannot be d-separated by any subset of remaining nodes (and vice-versa)



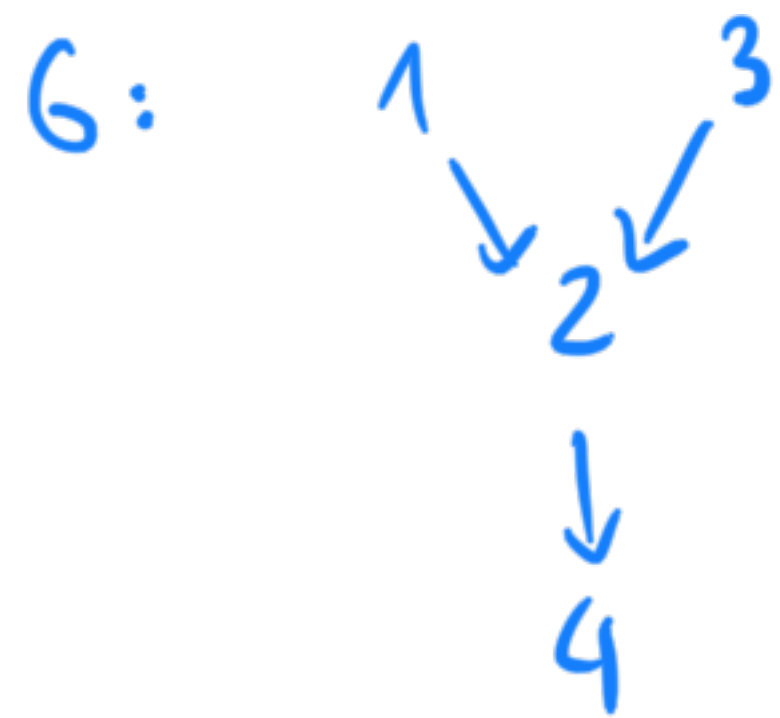
Step 1: Skeleton learning

- Given $G = (\mathbf{V}, \mathbf{E})$, nodes $i, j \in \mathbf{V}, i \neq j$ then:
 - If **i is adjacent to j** , they cannot be d-separated by any subset of remaining nodes (and vice-versa)
- 1. Start with **completely connected undirected graph U**
- 2. For each pair $i, j \in \mathbf{V}, i \neq j$, and for any subset $\mathbf{S} \subseteq \mathbf{V} \setminus \{i, j\}$
 - Check **if $X_i \perp\!\!\!\perp X_j \mid X_{\mathbf{S}}$** for any \mathbf{S} in data
 - If this is true, by faithfulness $i \perp_G j \mid \mathbf{S}$, so we can **remove $i - j$ in U**



Step 1: Skeleton learning - example

True causal graph

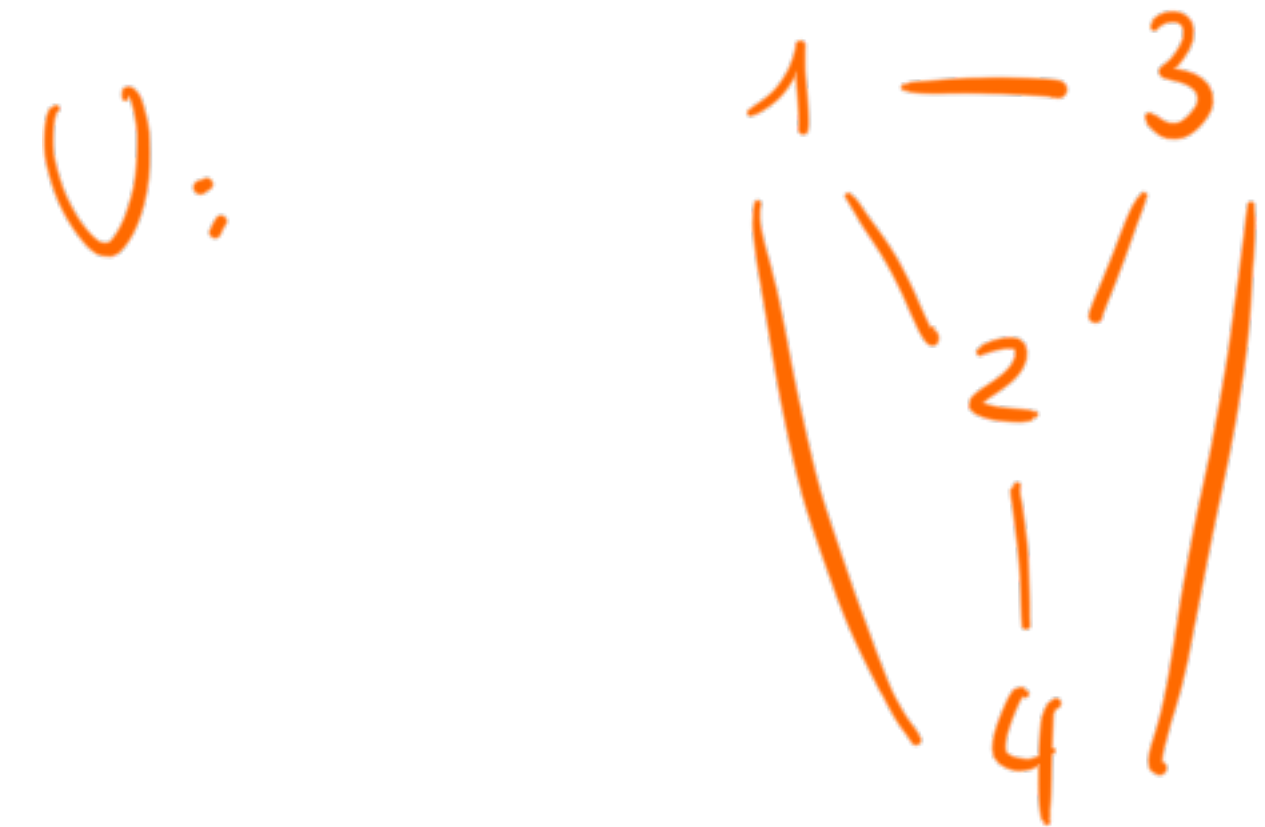


P: $P(x_1) \cdot P(x_3) \cdot P(x_2 | x_1, x_3) P(x_4 | x_2)$

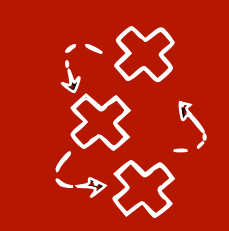
Data

x_1	x_2	x_3	x_4
1	1	1	0
1	0	0	1
0	1	0	1
1	0	1	0
1	0	0	0

Algorithm output

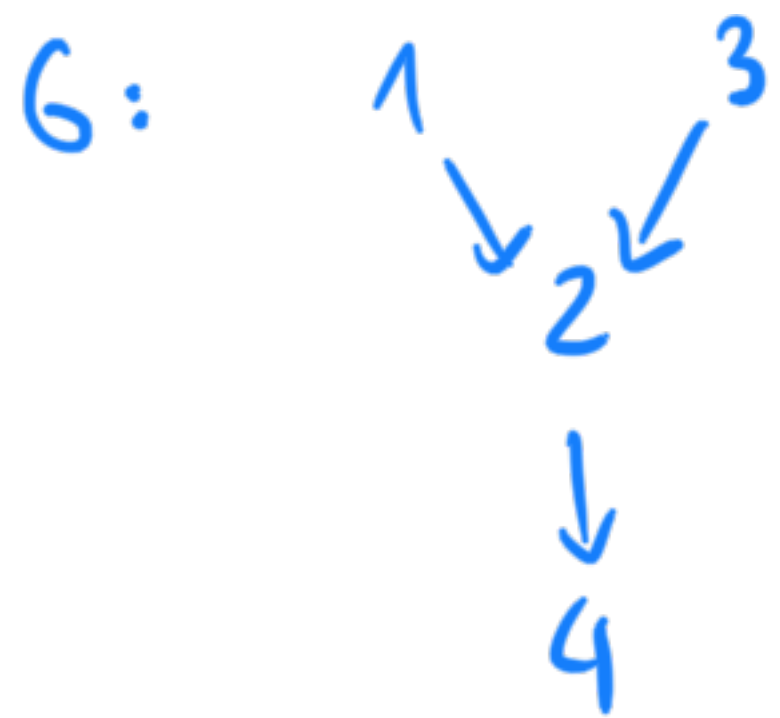


For this example, we will pretend we do not know the true graph, but only the CIs that we can derive from data



Step 1: Skeleton learning - example

True causal graph



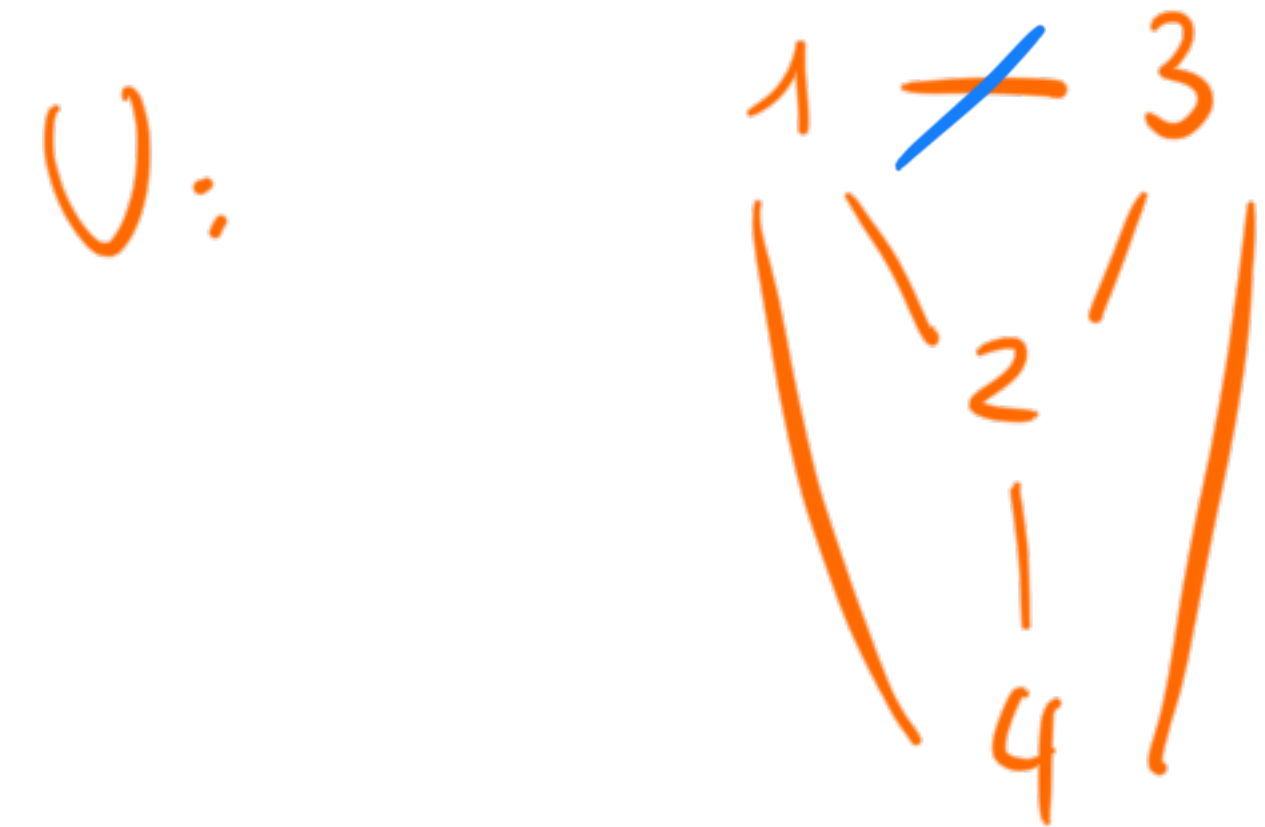
P: $P(x_1) \cdot P(x_3) \cdot P(x_2 | x_1, x_3) P(x_4 | x_2)$

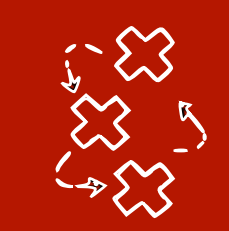
Data

x_1	x_2	x_3	x_4
1	1	1	0
1	0	0	1
0	1	0	1
1	0	1	0
1	0	0	0

$x_1 \perp\!\!\!\perp x_3$

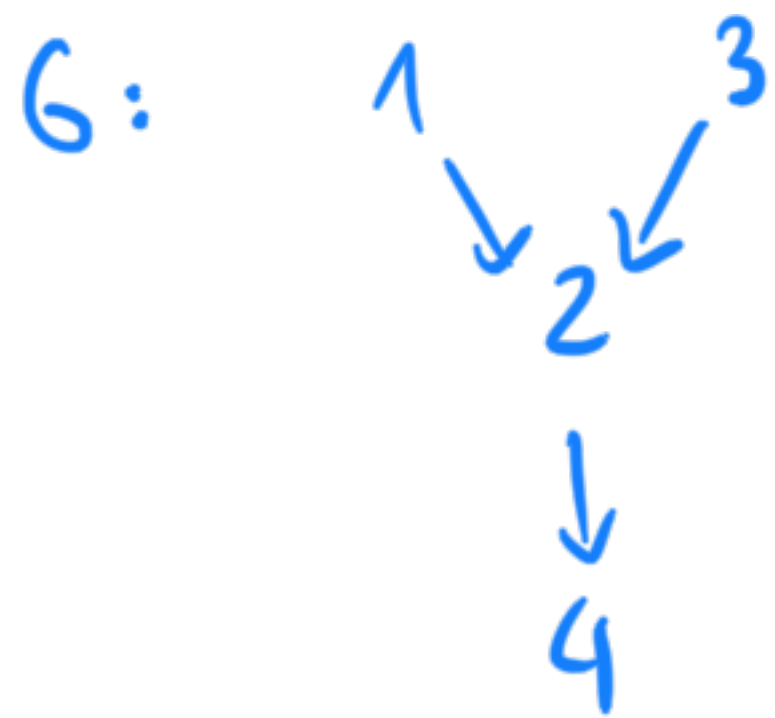
Algorithm output





Step 1: Skeleton learning - example

True causal graph



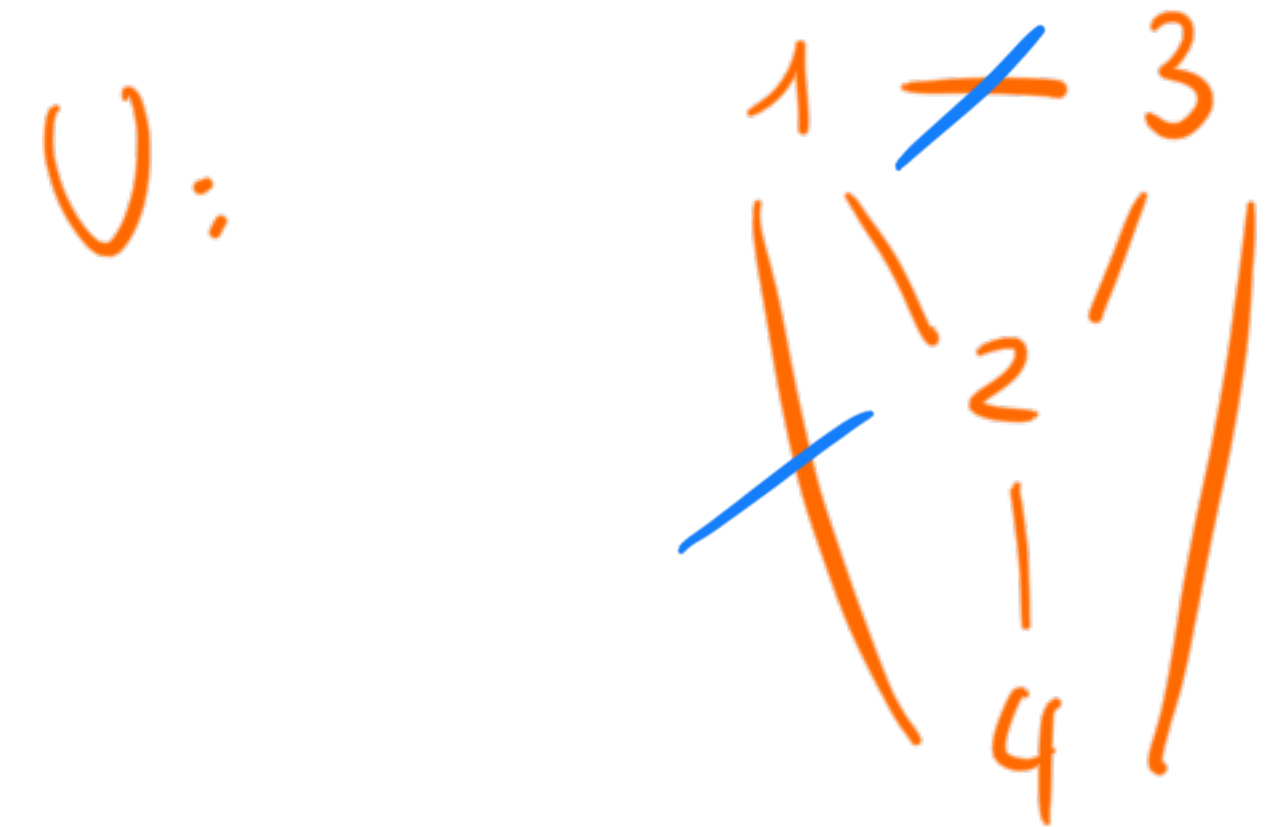
P: $P(x_1) \cdot P(x_3) \cdot P(x_2 | x_1, x_3) P(x_4 | x_2)$

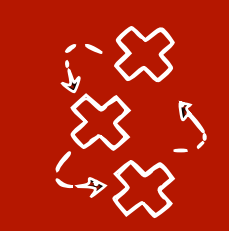
Data

x_1	x_2	x_3	x_4
1	1	1	0
1	0	0	1
0	1	0	1
1	0	1	0
1	0	0	0

$x_1 \perp\!\!\!\perp x_4 | x_2$

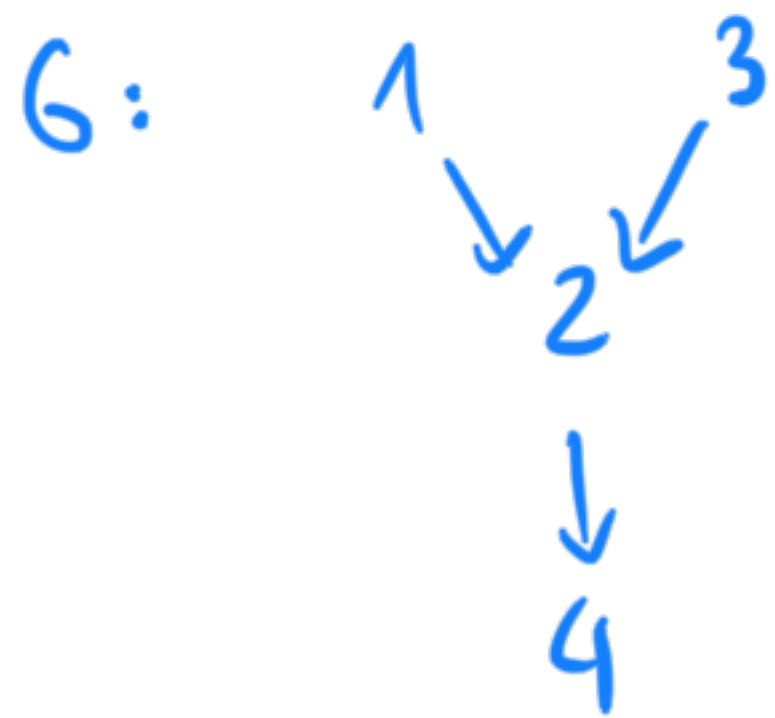
Algorithm output





Step 1: Skeleton learning - example

True causal graph



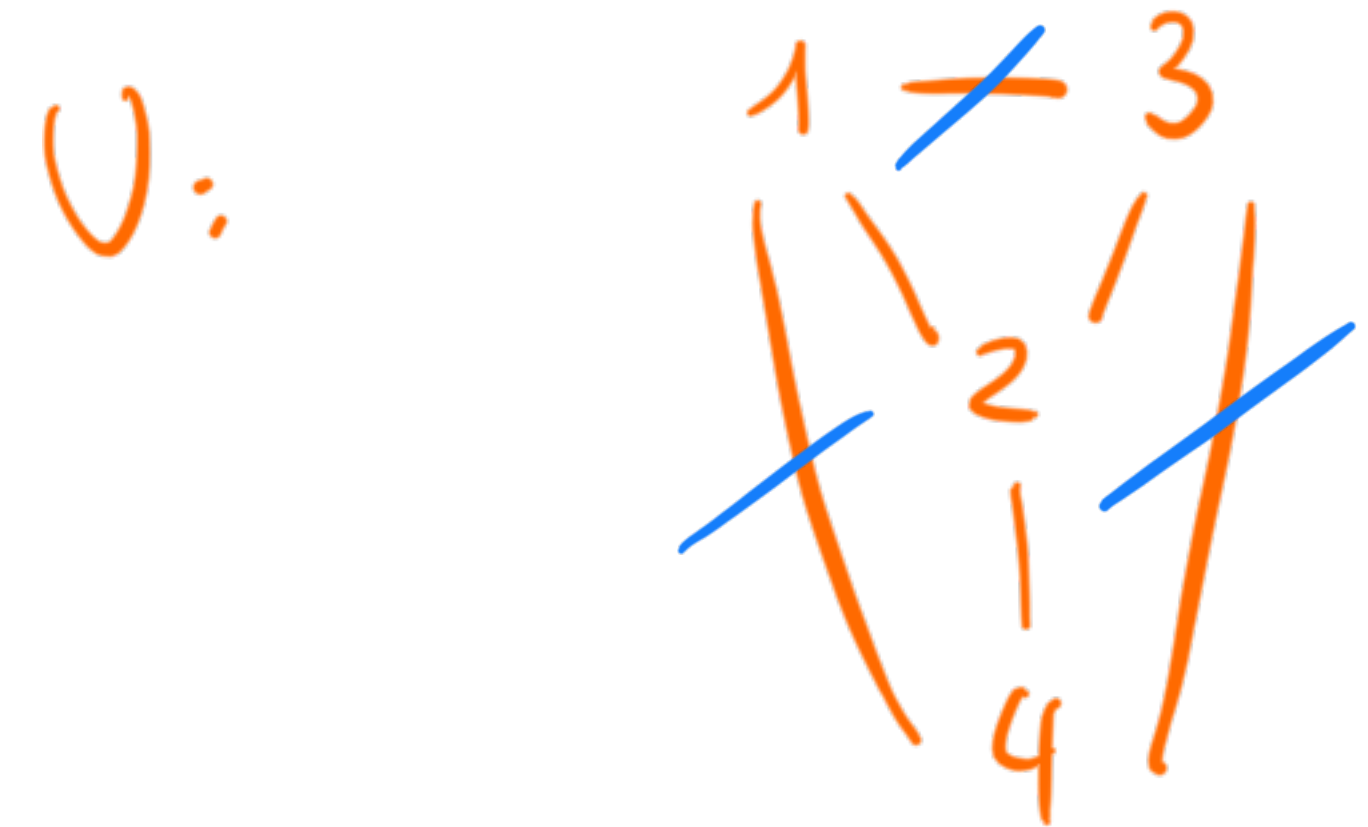
P: $P(x_1) \cdot P(x_3) \cdot P(x_2 | x_1, x_3) P(x_4 | x_2)$

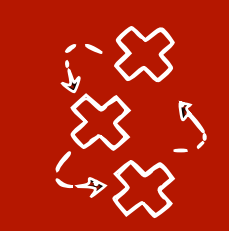
Data

x_1	x_2	x_3	x_4
1	1	1	0
1	0	0	1
0	1	0	1
1	0	1	0
1	0	0	0

$x_3 \perp\!\!\!\perp x_4 | x_2$

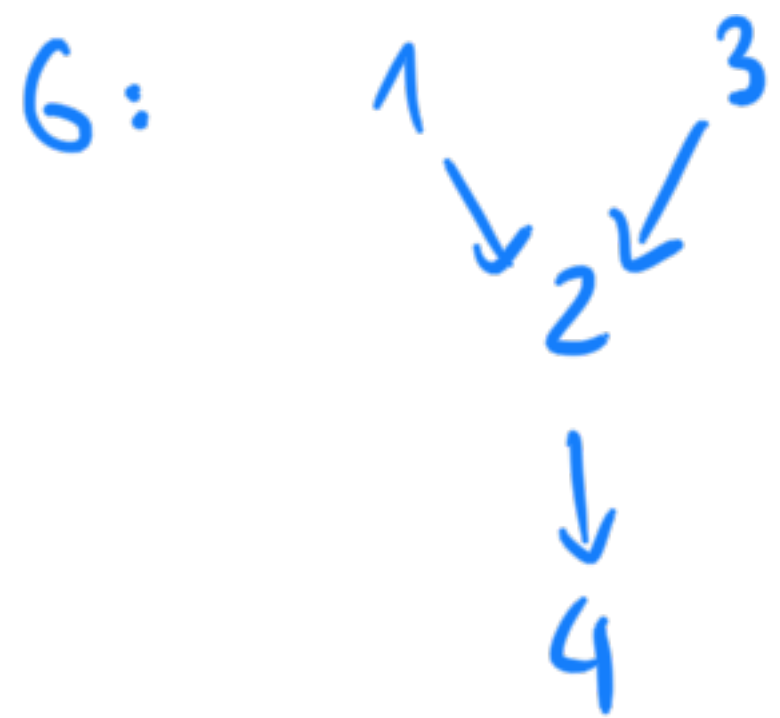
Algorithm output





Step 1: Skeleton learning - example

True causal graph



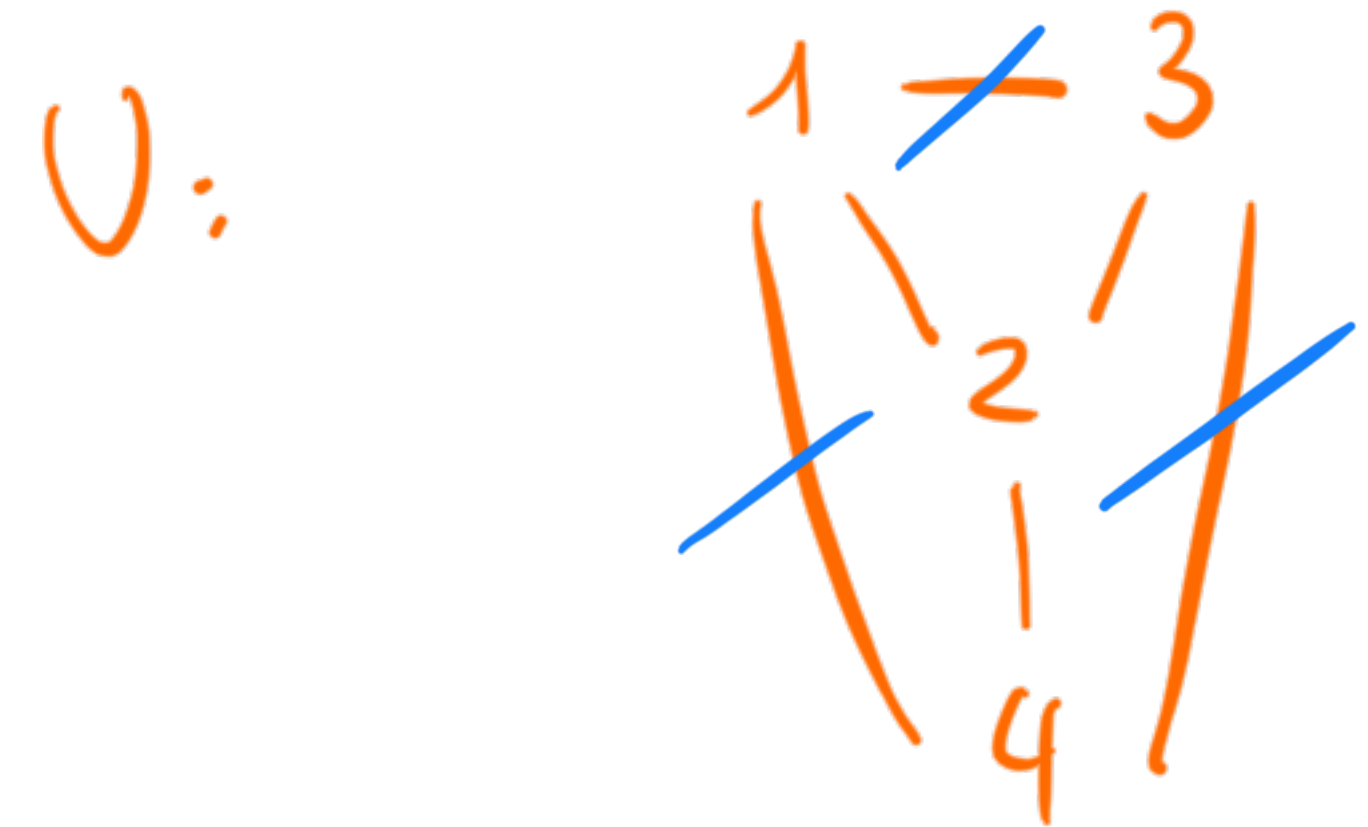
P: $P(x_1) \cdot P(x_3) \cdot P(x_2 | x_1, x_3) P(x_4 | x_2)$

Data

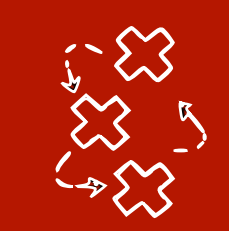
x_1	x_2	x_3	x_4
1	1	1	0
1	0	0	1
0	1	0	1
1	0	1	0
1	0	0	0

$x_1 \perp\!\!\!\perp x_4 \mid x_2, x_3$
 $x_3 \perp\!\!\!\perp x_4 \mid x_2, x_1$

Algorithm output

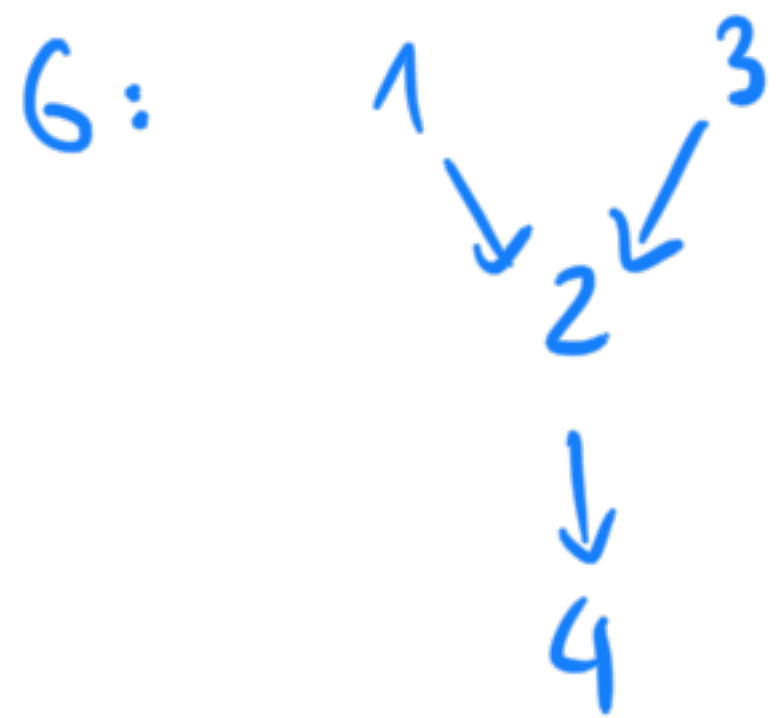


Nothing changes, we can stop



Step 1: Skeleton learning - example

True causal graph

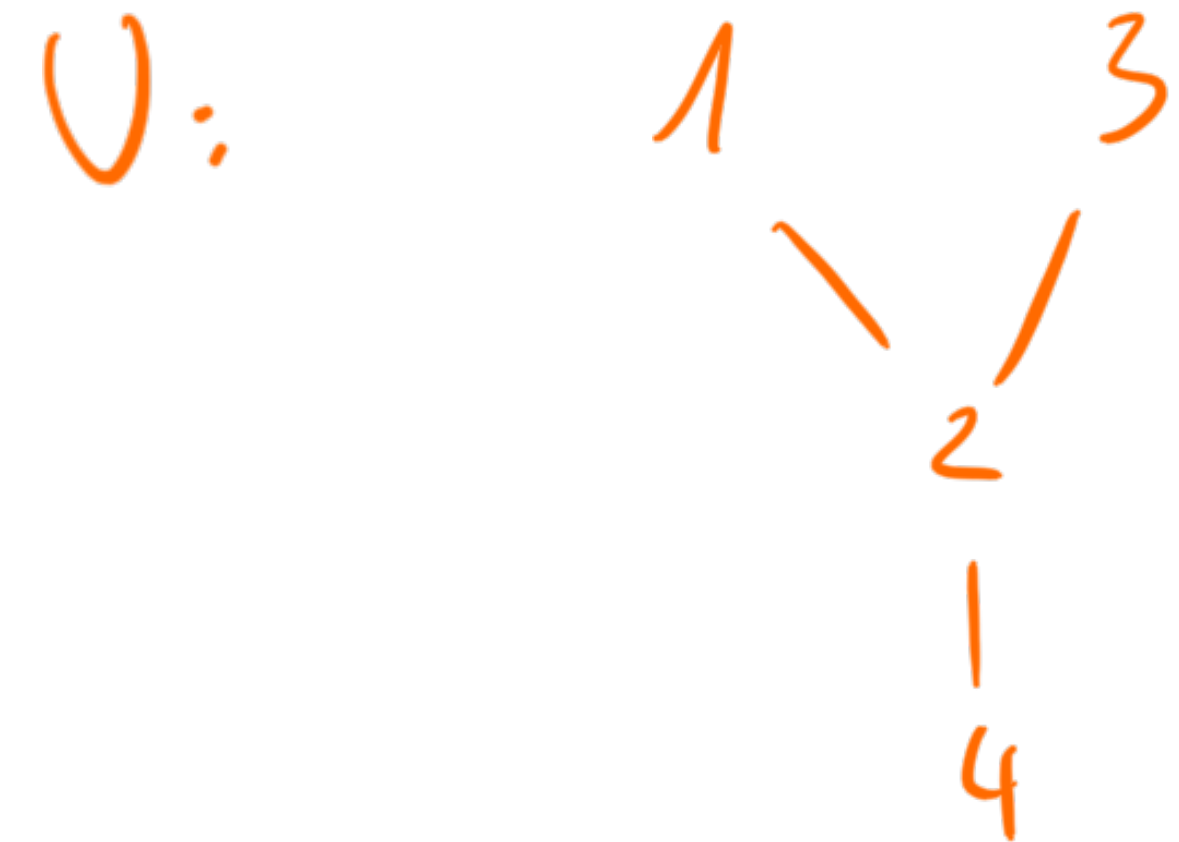


$$P: P(X_1) \cdot P(X_3) \cdot P(X_2 | X_1, X_3) P(X_4 | X_2)$$

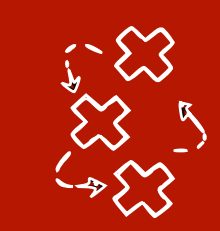
Data

X_1	X_2	X_3	X_4
1	1	1	0
1	0	0	1
0	1	0	1
1	0	1	0
1	0	0	0

Algorithm output

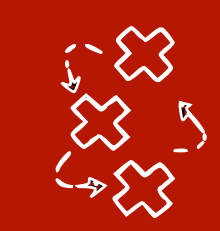


Step 1 finds the correct skeleton



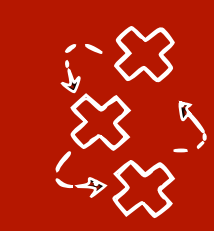
SGS algorithm (Spirtes, Glymour, Scheines)

- Assuming p is Markov and faithful to an unknown graph G
- We can estimate a CPDAG from samples of p in three steps:
 1. Determine the **skeleton**
 2. Determine the **v-structures** (*given the tests in the previous phase*)
 3. Direct as many remaining edges as possible
- **Note:** the directed parts of the CPDAG will agree with G , but some parts might stay undirected



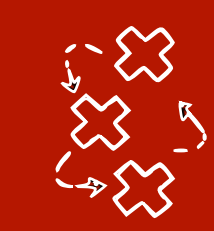
Step 2: Determine v-structures

- A triple of nodes (i, j, k) in a DAG G is a **an unshielded triple** if $i - j, j - k$ and **i is not adjacent to k** , i.e. $i \not\sim k$, in G



Step 2: Determine v-structures

- A triple of nodes (i, j, k) in a DAG G is a **an unshielded triple** if $i - j, j - k$ and **i is not adjacent to k** , i.e. $i \neq k$, in G
 1. Start from the skeleton U from previous step
 2. For each unshielded triple (i, j, k) in U , i.e. $i - j, j - k$ and $i \neq k$ in U
 - For all $S \subseteq V \setminus \{i, j, k\}$ check if $X_i \not\perp\!\!\!\perp X_k \mid X_j \cup X_S$ in data
 - If this is true, $i \rightarrow j \leftarrow k$ is a **v-structure**

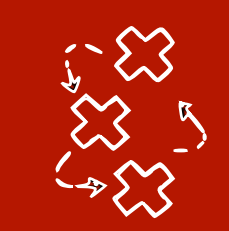


Step 2: Determine v-structures

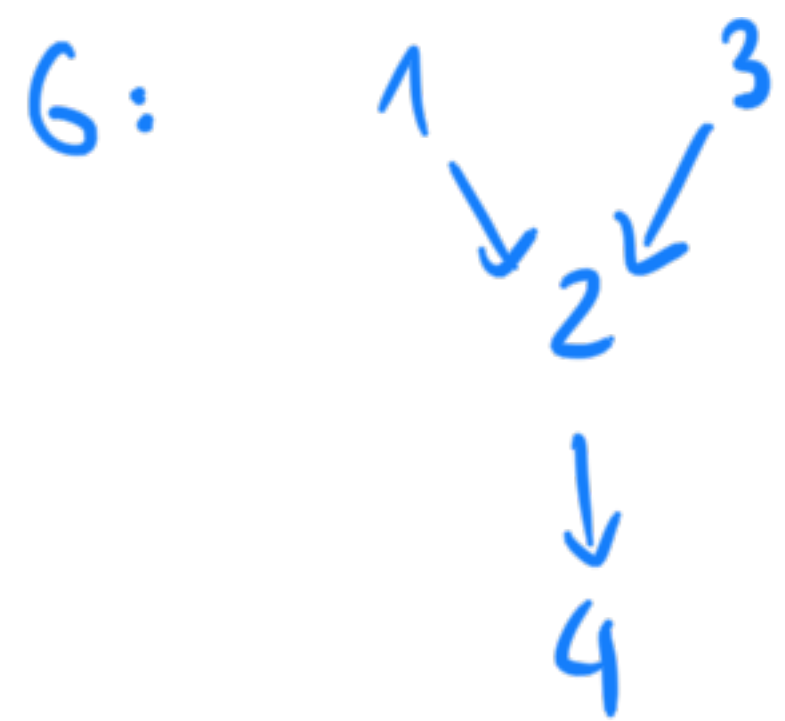
- A triple of nodes (i, j, k) in a DAG G is a **an unshielded triple** if $i - j, j - k$ and **i is not adjacent to k** , i.e. $i \not\sim k$ in G

Keep in mind: for unshielded triples (i, j, k) we check if X_i, X_k are always dependent **given any conditioning set containing X_j**

1. Start from the skeleton
2. For each unshielded triple (i, j, k) in U , i.e. $i - j, j - k$ and $i \not\sim k$ in U
 - For all $S \subseteq V \setminus \{i, j, k\}$ check if $X_i \not\perp\!\!\!\perp X_k \mid X_j \cup X_S$ in data
 - If this is true, $i \rightarrow j \leftarrow k$ is a v-structure

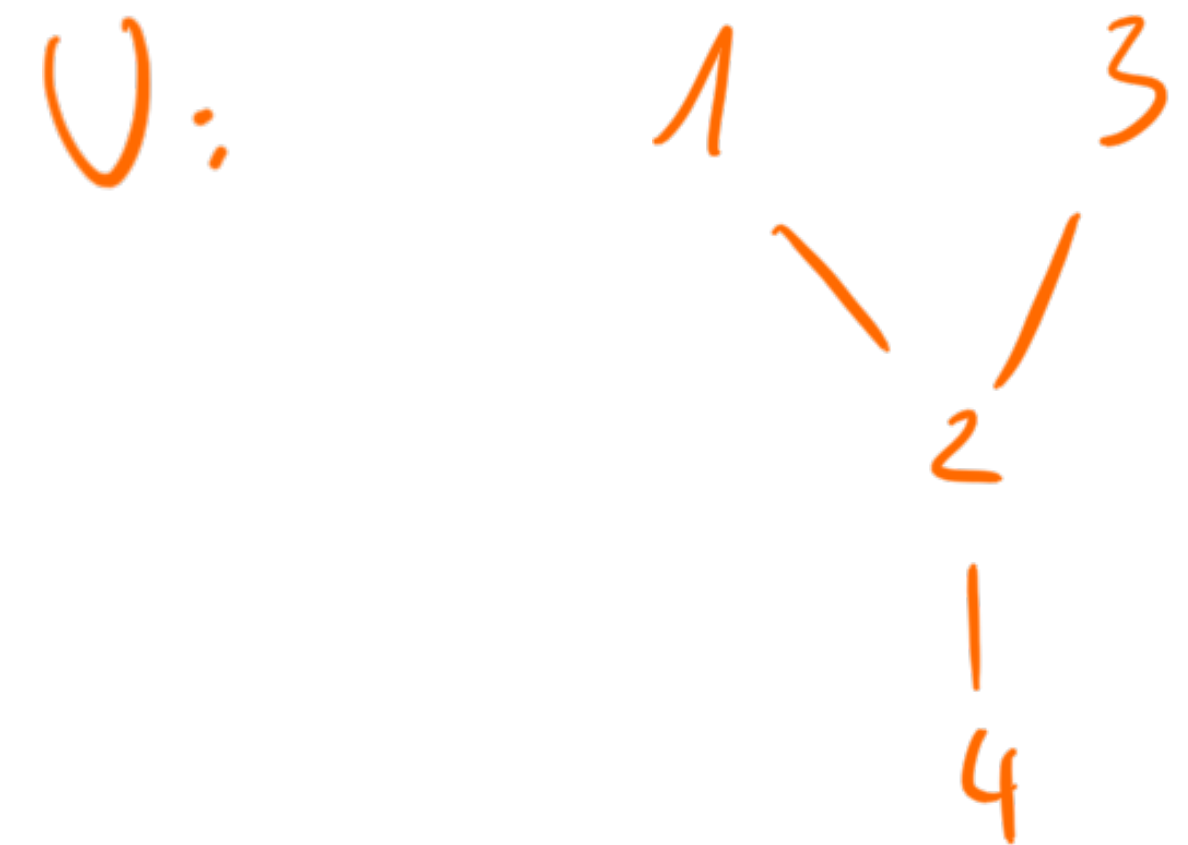


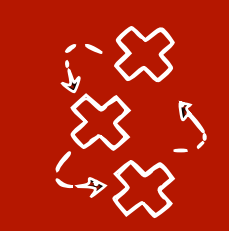
Step 2: Determine v-structures - example



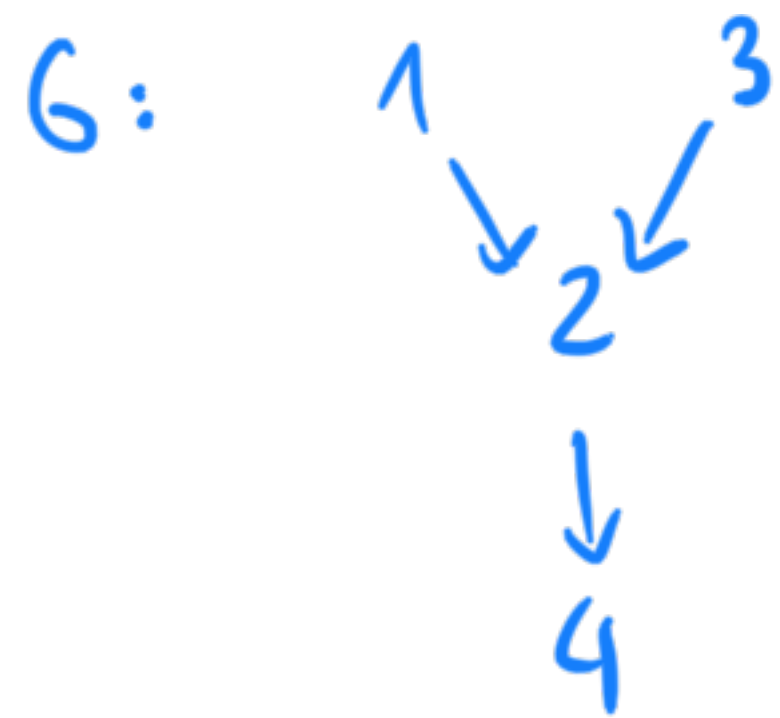
$$P: P(x_1) \cdot P(x_3) \cdot P(x_2 | x_1, x_3) P(x_4 | x_2)$$

x_1	x_2	x_3	x_4
1	1	1	0
1	0	0	1
0	1	0	1
1	0	1	0
1	0	0	0



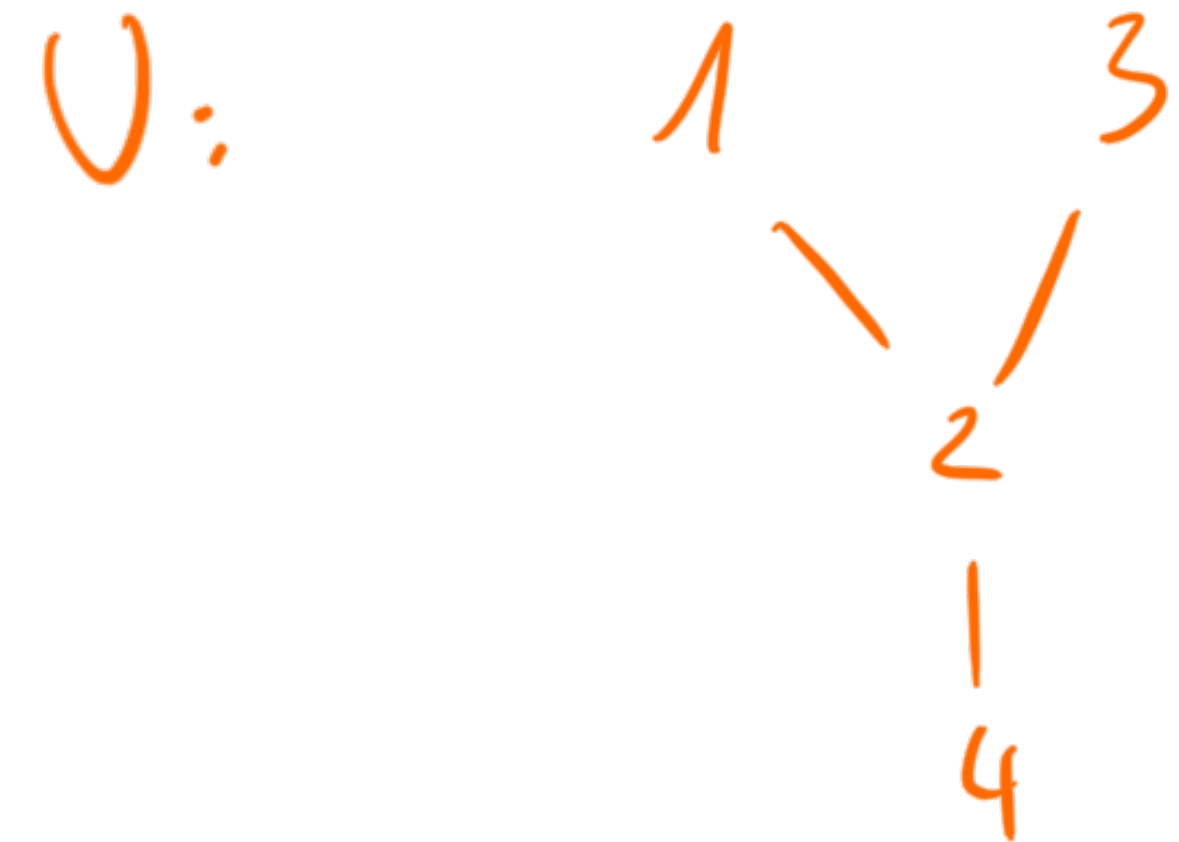


Step 2: Determine v-structures - example



P: $P(x_1) \cdot P(x_3) \cdot P(x_2 | x_1, x_3) P(x_4 | x_2)$

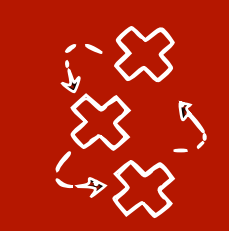
x_1	x_2	x_3	x_4
1	1	1	0
1	0	0	1
0	1	0	1
1	0	1	0
1	0	0	0



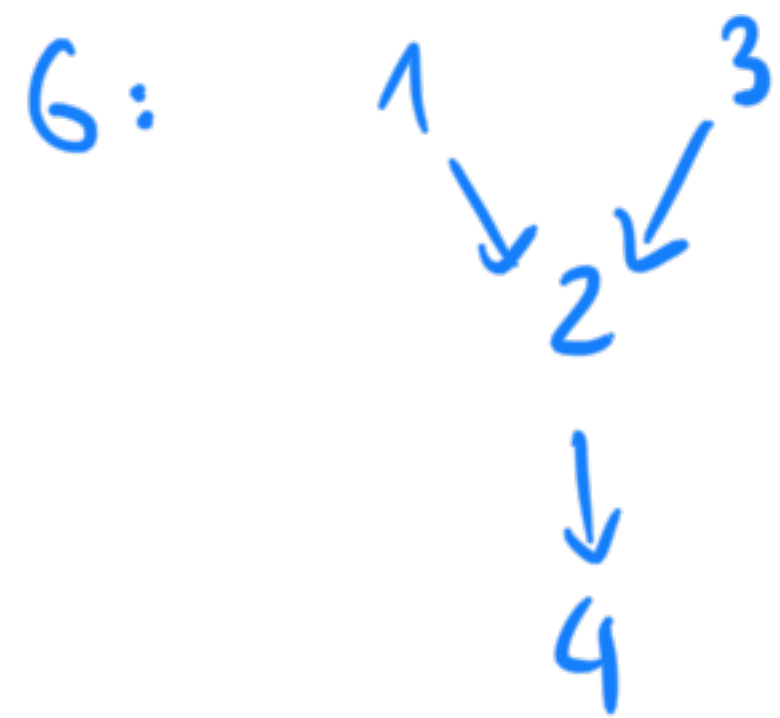
1-2-4: $X_1 \perp\!\!\!\perp X_4 | X_2$
 $X_1 \perp\!\!\!\perp X_4 | X_2, X_3$ ✗

We can reuse the CIs from phase 1

- $X_1 \perp\!\!\!\perp X_3$
- $X_1 \perp\!\!\!\perp X_4 | X_2$
- $X_3 \perp\!\!\!\perp X_4 | X_2$
- $X_1 \perp\!\!\!\perp X_4 | X_2, X_3$
- $X_3 \perp\!\!\!\perp X_4 | X_2, X_1$

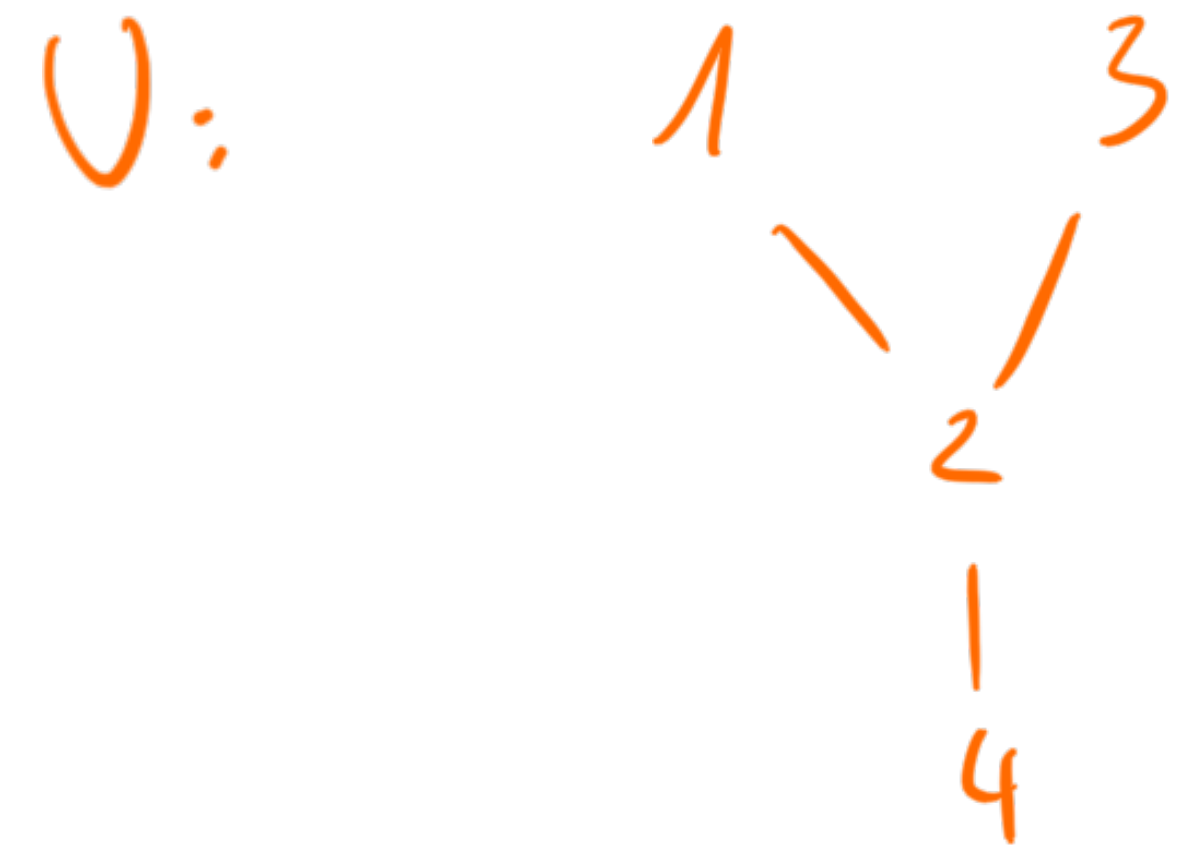


Step 2: Determine v-structures - example



P: $P(x_1) \cdot P(x_3) \cdot P(x_2 | x_1, x_3) P(x_4 | x_2)$

x_1	x_2	x_3	x_4
1	1	1	0
1	0	0	1
0	1	0	1
1	0	1	0
1	0	0	0



3-2-4

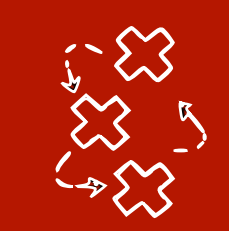
$x_3 \perp\!\!\!\perp x_4 | x_2$

$x_3 \perp\!\!\!\perp x_4 | x_2, x_1$

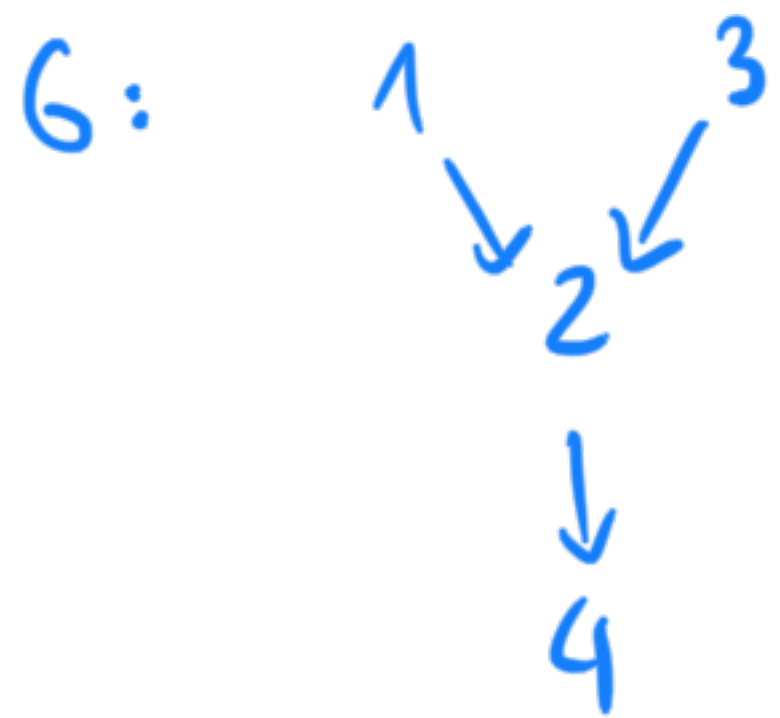
x

We can reuse the CIs from phase 1

- $x_1 \perp\!\!\!\perp x_3$
- $x_1 \perp\!\!\!\perp x_4 | x_2$
- $x_3 \perp\!\!\!\perp x_4 | x_2$
- $x_1 \perp\!\!\!\perp x_4 | x_2, x_3$
- $x_3 \perp\!\!\!\perp x_4 | x_2, x_1$

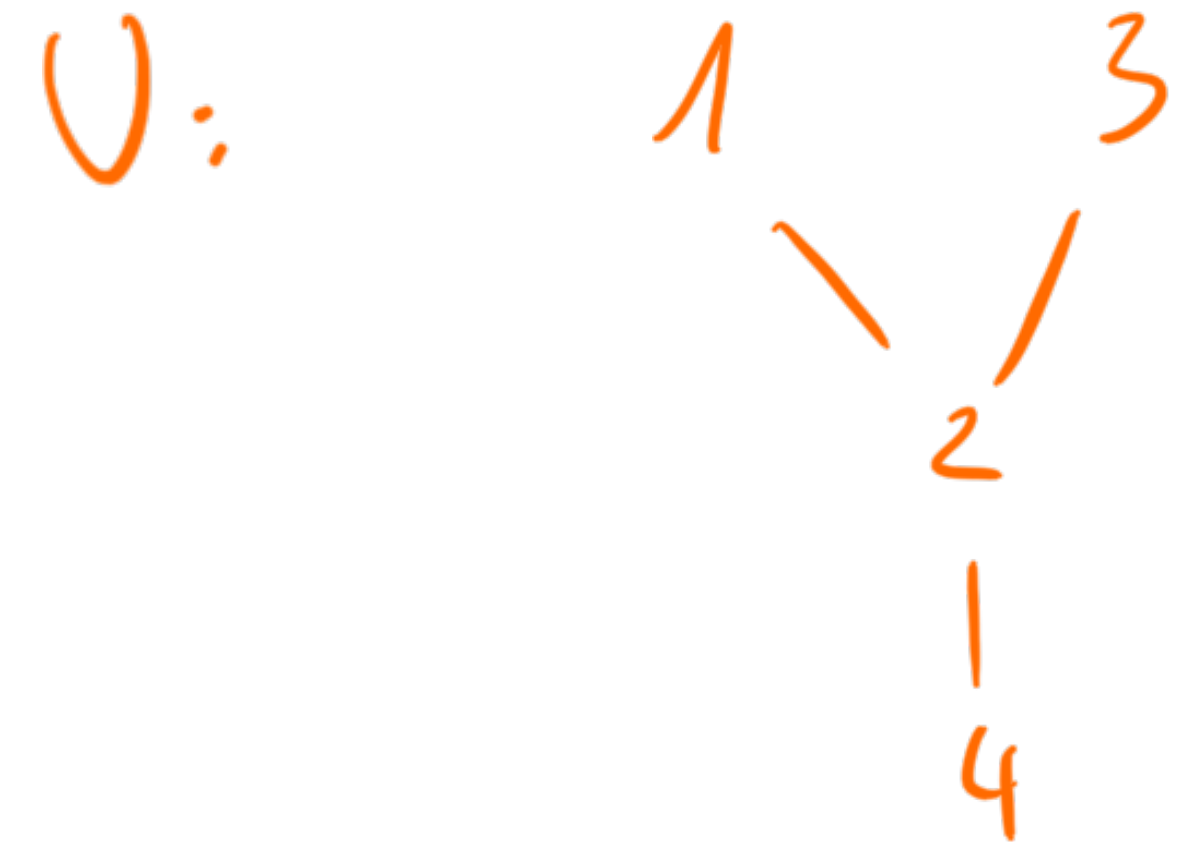


Step 2: Determine v-structures - example



P: $P(x_1) \cdot P(x_3) \cdot P(x_2 | x_1, x_3) P(x_4 | x_2)$

x_1	x_2	x_3	x_4
1	1	1	0
1	0	0	1
0	1	0	1
1	0	1	0
1	0	0	0



1-2-3

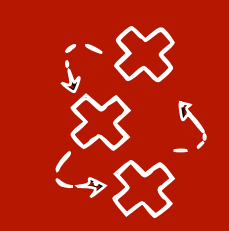
~~$x_1 \perp\!\!\!\perp x_3 | x_2$~~ ✓

~~$x_1 \perp\!\!\!\perp x_3 | x_2, x_4$~~ ✓

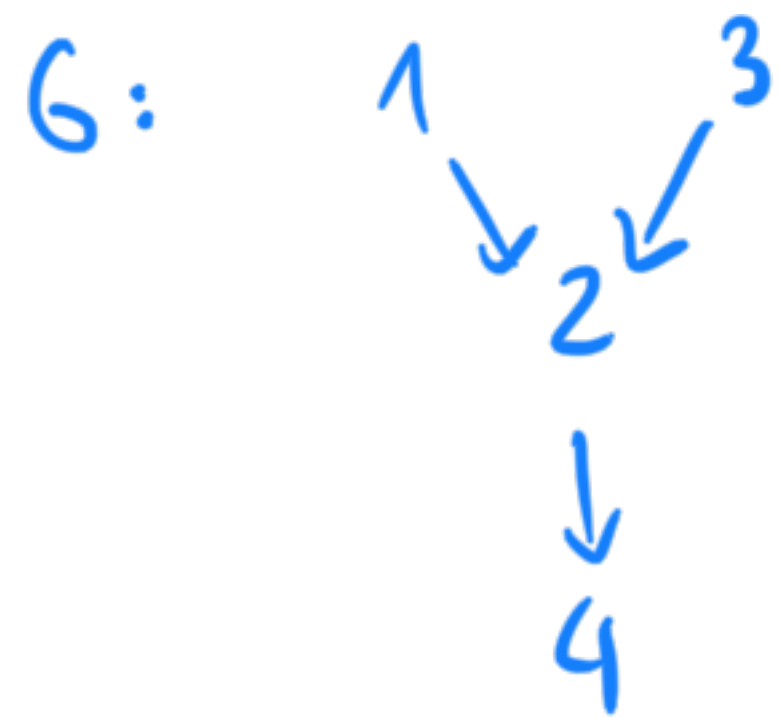
We can reuse the CIs from phase 1

- $x_1 \perp\!\!\!\perp x_3$
- $x_1 \perp\!\!\!\perp x_4 | x_2$
- $x_3 \perp\!\!\!\perp x_4 | x_2$
- $x_1 \perp\!\!\!\perp x_4 | x_2, x_3$
- $x_3 \perp\!\!\!\perp x_4 | x_2, x_1$

$x_1 \perp\!\!\!\perp x_3$ has $Z = \emptyset$, so it does not contain x_2 and we ignore it

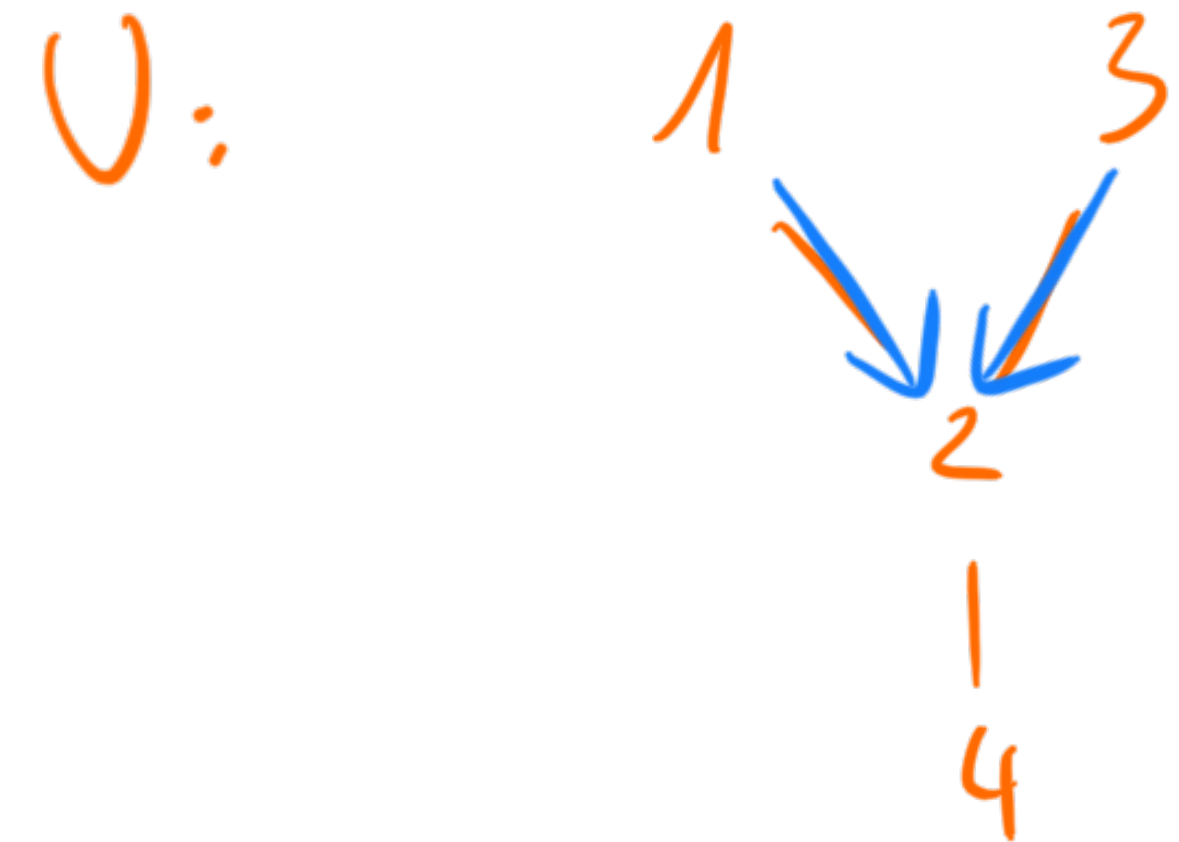


Step 2: Determine v-structures - example



P: $P(x_1) \cdot P(x_3) \cdot P(x_2 | x_1, x_3) P(x_4 | x_2)$

x_1	x_2	x_3	x_4
1	1	1	0
1	0	0	1
0	1	0	1
1	0	1	0
1	0	0	0

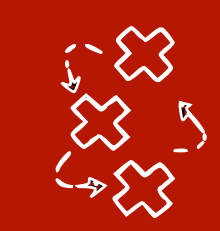


1-2-3

~~$x_1 \perp\!\!\!\perp x_3 | x_2$~~ ✓
 ~~$x_1 \perp\!\!\!\perp x_3 | x_2, x_4$~~ ✓

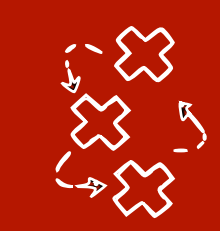
We can reuse the CIs from phase 1

- $x_1 \perp\!\!\!\perp x_3$
- $x_1 \perp\!\!\!\perp x_4 | x_2$
- $x_3 \perp\!\!\!\perp x_4 | x_2$
- $x_1 \perp\!\!\!\perp x_4 | x_2, x_3$
- $x_3 \perp\!\!\!\perp x_4 | x_2, x_1$



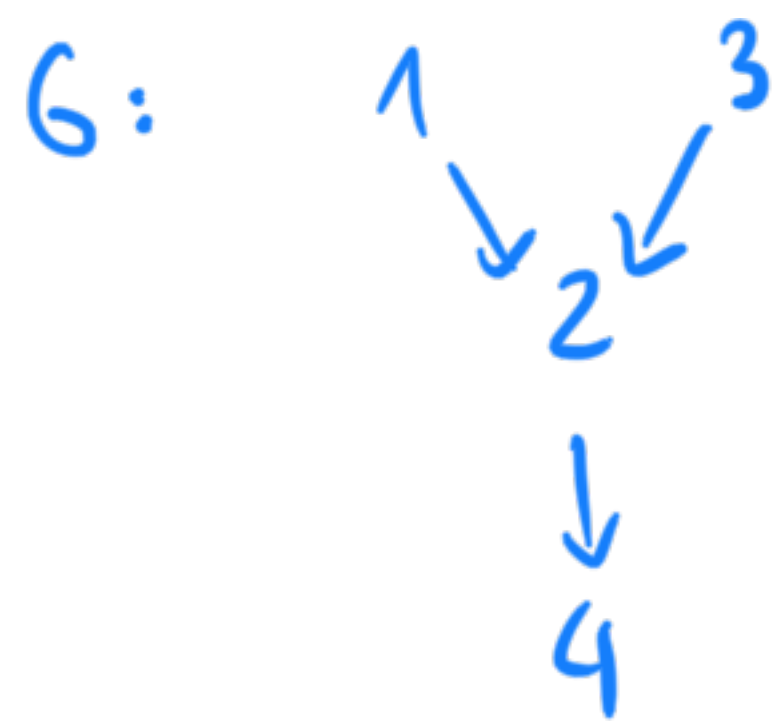
SGS algorithm (Spirtes, Glymour, Scheines)

- Assuming P is Markov and faithful to an unknown graph G
- We can estimate a CPDAG from samples of P in three steps:
 1. Determine the **skeleton**
 2. Determine the **v-structures** (*given the tests in the previous phase*)
 3. Direct as many remaining edges as possible
- **Note**: the directed parts of the CPDAG will agree with G , but some parts might stay undirected



Step 3: Direct as many edges as possible - example

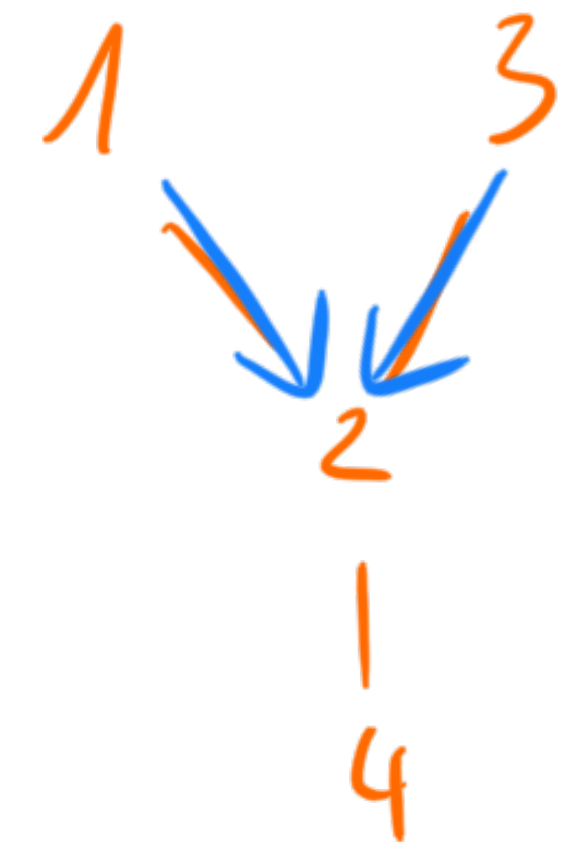
- Cannot create cycles or new v-structures
- Some of the edges can be oriented to disallow these situations to happen

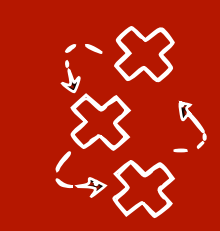


$$P: P(X_1) \cdot P(X_3) \cdot P(X_2 | X_1, X_3) P(X_4 | X_2)$$

X_1	X_2	X_3	X_4
1	1	1	0
1	0	0	1
0	1	0	1
1	0	1	0
1	0	0	0

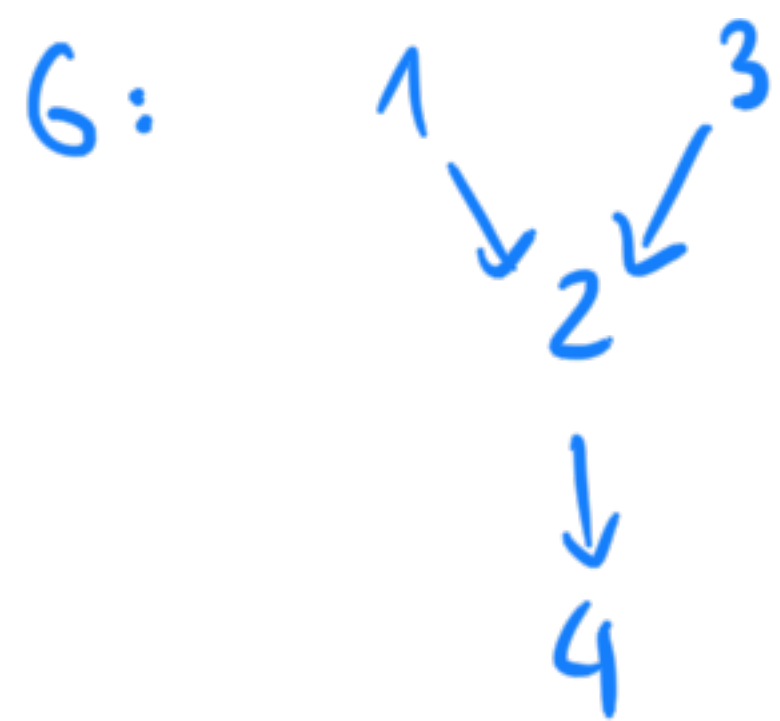
U:





Step 3: Direct as many edges as possible - example

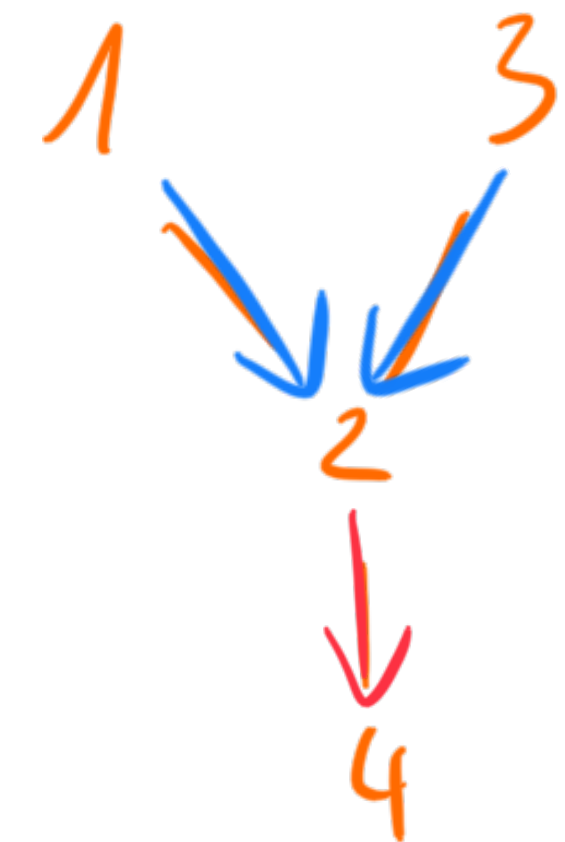
- Cannot create cycles or new v-structures
- Some of the edges can be oriented to disallow these situations to happen

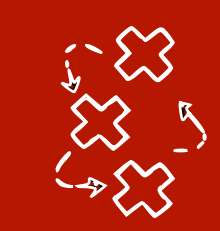


$$P: P(x_1) \cdot P(x_3) \cdot P(x_2 | x_1, x_3) P(x_4 | x_2)$$

x_1	x_2	x_3	x_4
1	1	1	0
1	0	0	1
0	1	0	1
1	0	1	0
1	0	0	0

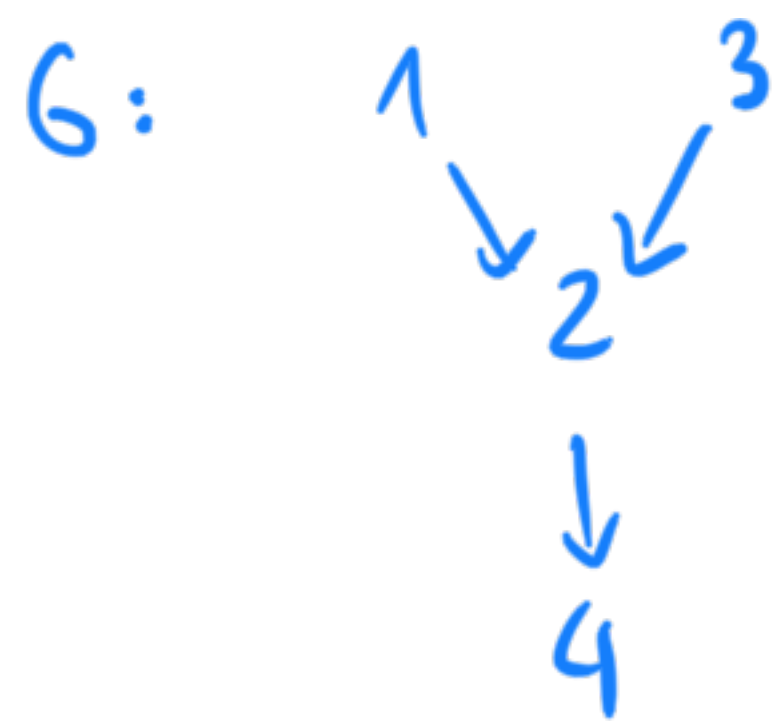
U:





Step 3: Direct as many edges as possible - example

- Cannot create cycles or new v-structures
- Some of the edges can be oriented to disallow these situations to happen



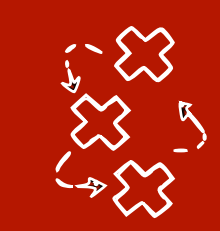
$$P: P(x_1) \cdot P(x_3) \cdot P(x_2 | x_1, x_3) P(x_4 | x_2)$$

x_1	x_2	x_3	x_4
1	1	1	0
1	0	0	1
0	1	0	1
1	0	1	0
1	0	0	0

CPDAG:



Same as G
(special case)



Step 3: Meek's rules (1995)

Sound and complete rules for additional orientations (also with added background knowledge)

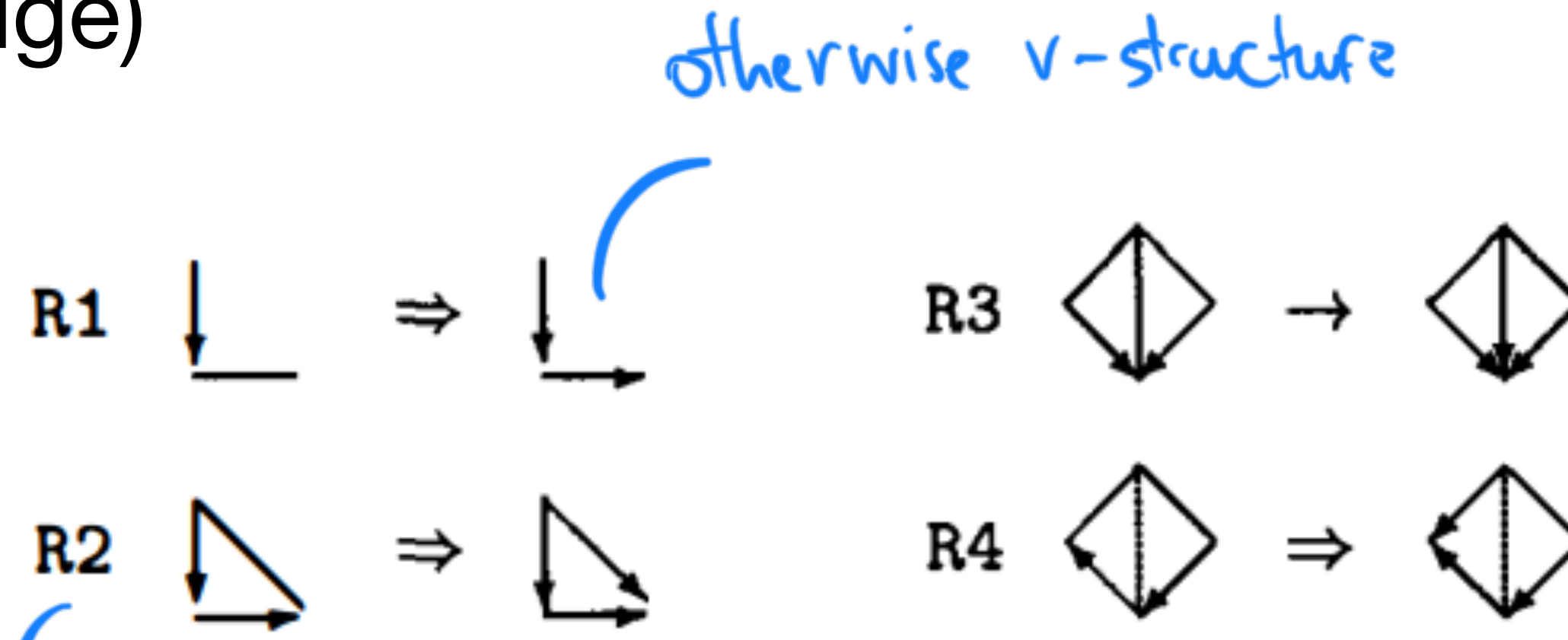
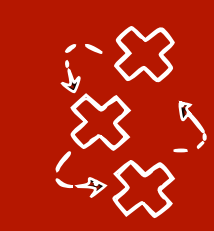
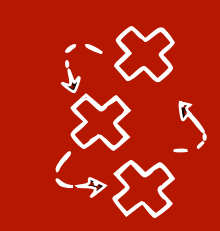


Figure 1: Orientation rules for patterns



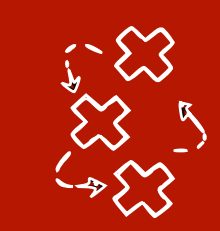
SGS algorithm (Spirtes, Glymour, Scheines)

- Assuming p is Markov and faithful to an unknown graph G
- We can estimate a CPDAG from samples of p in three steps:
 1. Determine the **skeleton**
 2. Determine the **v-structures** (*given the tests in the previous phase*)
 3. Direct as many remaining edges as possible
- **Computationally inefficient:** potentially $O(2^p)$ tests
 - Even if we reuse the test results from skeleton phase in other phases



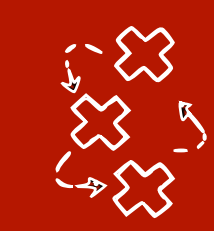
PC algorithm (Peter Spirtes, Clark Glymour)

- Assuming p is Markov and faithful to an unknown graph G
- We can estimate a CPDAG from samples of p in three steps:
 1. Determine the **skeleton in an optimised way**
 2. Determine the **v-structures** (*given the tests in the previous phase*)
 3. Direct as many remaining edges as possible



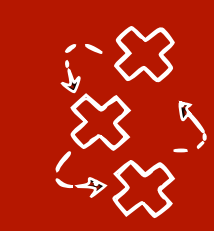
PC algorithm skeleton learning

- If i is not adjacent to j , then they can be d-separated by $\text{Pa}(i)$ or $\text{Pa}(j)$
- Determine the **skeleton in an optimised way**
 - Since we do not know the parents we will use the nodes that are adjacent, $\text{Adj}(i)$ or $\text{Adj}(j)$ in U at a given iteration (*superset*)



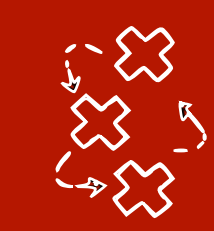
PC algorithm skeleton learning

- If i is not adjacent to j , then they can be d-separated by $\text{Pa}(i)$ or $\text{Pa}(j)$
 - Determine the **skeleton in an optimised way**
 - Since we do not know the parents we will use the nodes that are adjacent, $\text{Adj}(i)$ or $\text{Adj}(j)$ in U at a given iteration (*superset*)
1. Start with **completely connected undirected graph U**



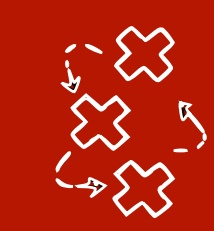
PC algorithm skeleton learning

- If i is not adjacent to j , then they can be d-separated by $\text{Pa}(i)$ or $\text{Pa}(j)$
 - Determine the **skeleton in an optimised way**
 - Since we do not know the parents we will use the nodes that are adjacent, $\text{Adj}(i)$ or $\text{Adj}(j)$ in U at a given iteration (*superset*)
1. Start with **completely connected undirected graph U**
 2. For $k = 0, 1, 2, \dots, p - 2$
 - If $i - j$ in U and **there exists a set $S \subseteq \text{Adj}(i) \setminus \{j\}$ of size at least k**



PC algorithm skeleton learning

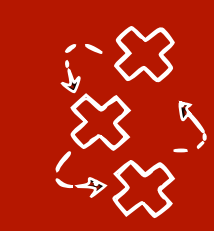
- If i is not adjacent to j , then they can be d-separated by $\text{Pa}(i)$ or $\text{Pa}(j)$
 - Determine the **skeleton in an optimised way**
 - Since we do not know the parents we will use the nodes that are adjacent, $\text{Adj}(i)$ or $\text{Adj}(j)$ in U at a given iteration (*superset*)
1. Start with **completely connected undirected graph U**
 2. For $k = 0, 1, 2, \dots, p - 2$
 - If $i - j$ in U and **there exists a set $S \subseteq \text{Adj}(i) \setminus \{j\}$ of size at least k**
 - If holds $X_i \perp\!\!\!\perp X_j \mid S$, then **remove $i - j$ in U**



PC algorithm skeleton learning

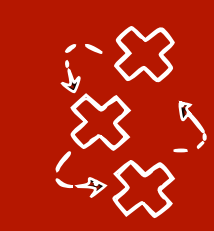
- If i is not adjacent to j , then they can be d-separated by $\text{Pa}(i)$ or $\text{Pa}(j)$
- Determine the **skeleton in an optimised way**
 - Since we do not know the parents we will use the nodes that are adjacent, $\text{Adj}(i)$ or $\text{Adj}(j)$ in U at a given iteration (*superset*)
- 1. Start with **completely connected**
- 2. For $k = 0, 1, 2, \dots, p - 2$
 - If $i - j$ in U and **there exists a set $S \subseteq \text{Adj}(i) \setminus \{j\}$ of size at least k**
 - If holds $X_i \perp\!\!\!\perp X_j \mid S$, then **remove $i - j$ in U**

We stop when there are no more untested sets of size k in the adjacencies of any variable



SGS vs PC

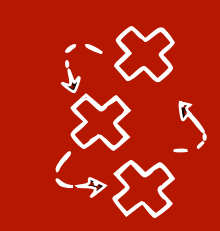
- **SGS:** we can estimate a CPDAG from samples of p in three steps:
 1. Determine the **skeleton**
 2. Determine the **v-structures**
 3. Direct as many remaining edges as possible
- **PC:** we can estimate a CPDAG from samples of p in three steps:
 1. Determine the **skeleton in an optimised way**
 2. Determine the **v-structures**
 3. Direct as many remaining edges as possible



PC algorithm - when does it fail?

- If the conditional independence tests give the wrong result
 - Too few samples
 - A very weak dependence
 - Wrong parametric assumption (e.g. partial correlation on nonlinear data)

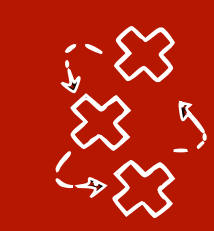
If you're curious, you can see here some variants (like conservative PC) that try to circumvent the problem by doing more tests <https://rdrr.io/cran/pcalg/man/pc.html>



PC algorithm - when does it fail?

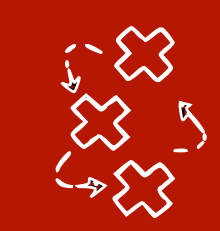
- If the conditional independence tests give the wrong result
 - Too few samples
 - A very weak dependence
 - Wrong parametric assumption (e.g. partial correlation on nonlinear data)
- If there are unmeasured confounders or selection bias
 - For example Chocolate - Nobel prizes





PC algorithm - when does it fail?

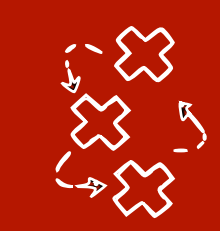
- If the conditional independence tests give the wrong result
 - Too few samples
 - A very weak dependence
 - Wrong parametric assumption (e.g. partial correlation on nonlinear data)
- If there are unmeasured confounders or selection bias
- We use advanced constraint-based algorithms - Fast Causal Inference (FCI)
 - Chapter 6 in [SGS] Causation Prediction and Search
 - <https://www.researchgate.net/publication/242448131> Causation Prediction and Search



Break?

Optional SGS exercise:

<https://drive.google.com/file/d/14IR7BSH2N7ZasRXIO5nXrF3Xn49hnmRn/view?usp=sharing>



Causal discovery - this class

Constraint-based causal discovery

- Conditional independence tests
- Observational data
- Output: MEC
- SGS, PC

Score-based causal discovery

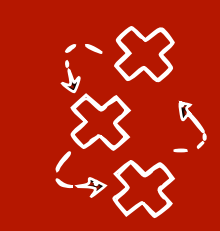
- Penalised likelihood
- Observational data
- Output: MEC
- GES, MMHC

Restricted models

- Nonlinear additive noise, Linear Non-Gaussianity
- Observational data
- Output: DAG
- RESIT, LINGAM

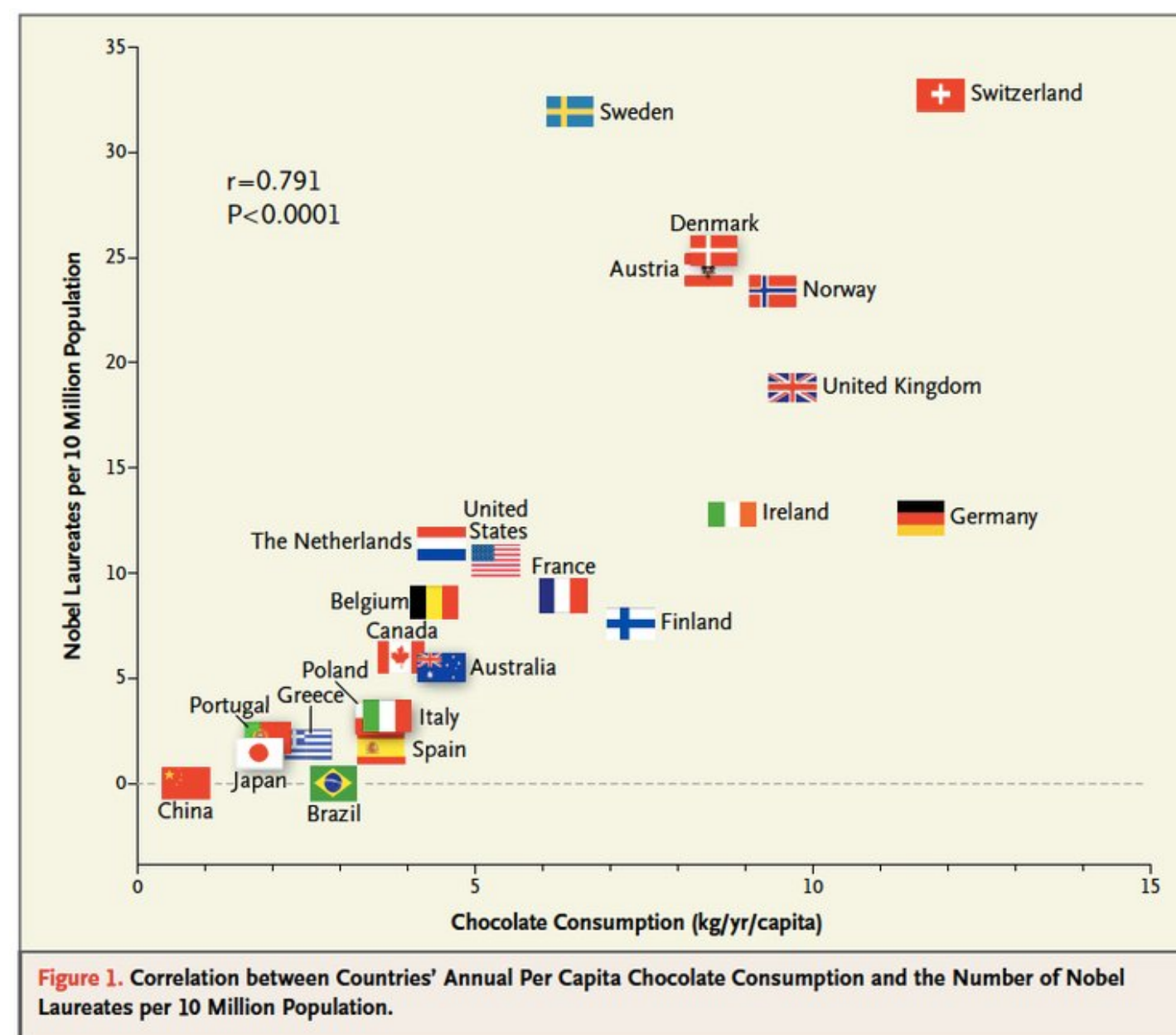
Interventional causal discovery / causal invariance

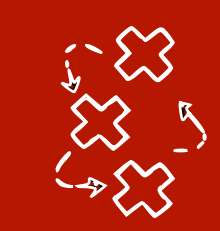
- Observational and Interventional data
- Output: parents of Y, I-MEC
- ICP, JCI



Learning from interventional data - intuition

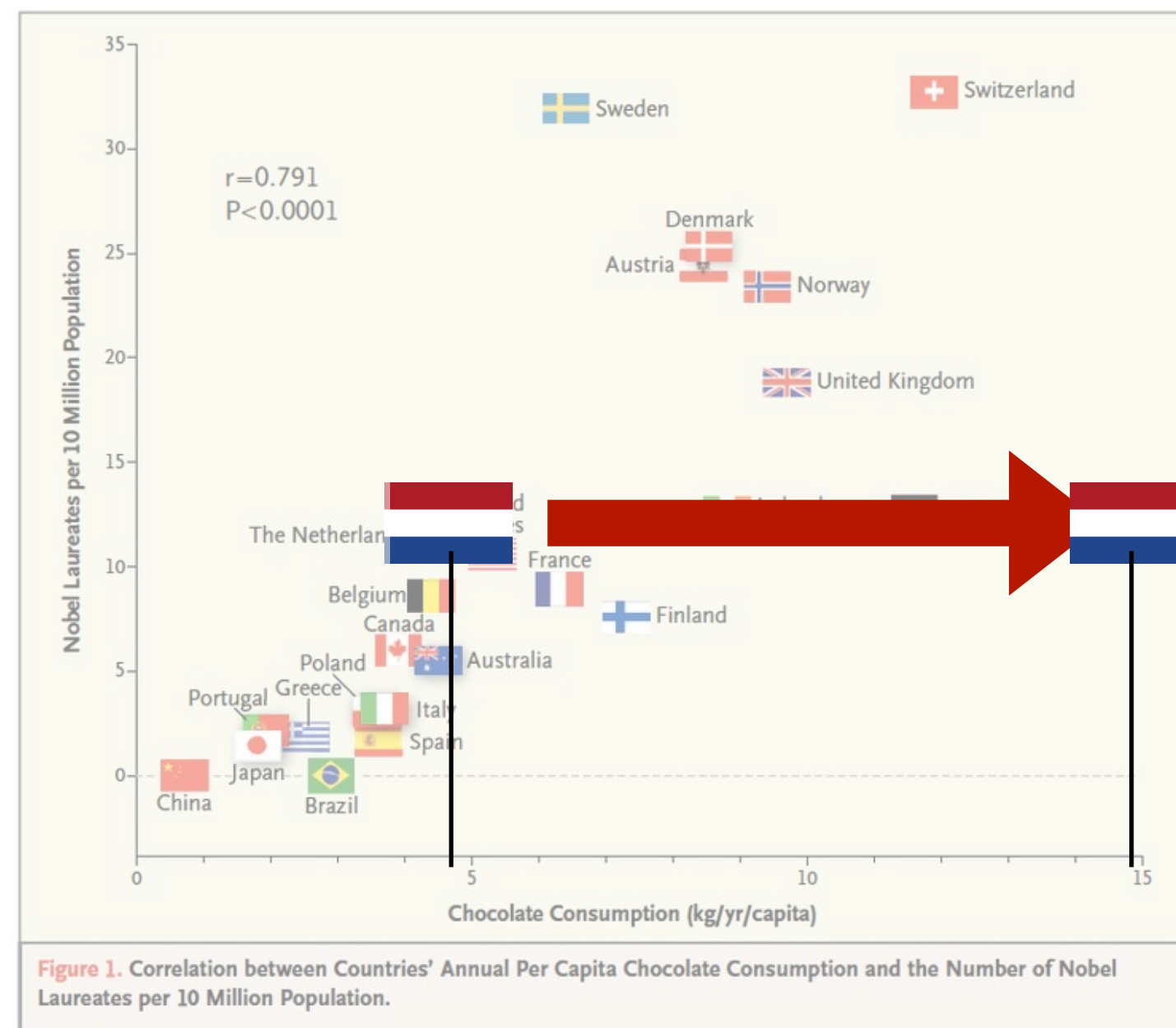
Until now we have only used **observational data**.



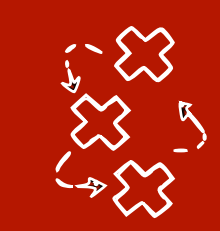


Learning from interventional data - intuition

Hypothetical world: we perform the experiment and see these results:

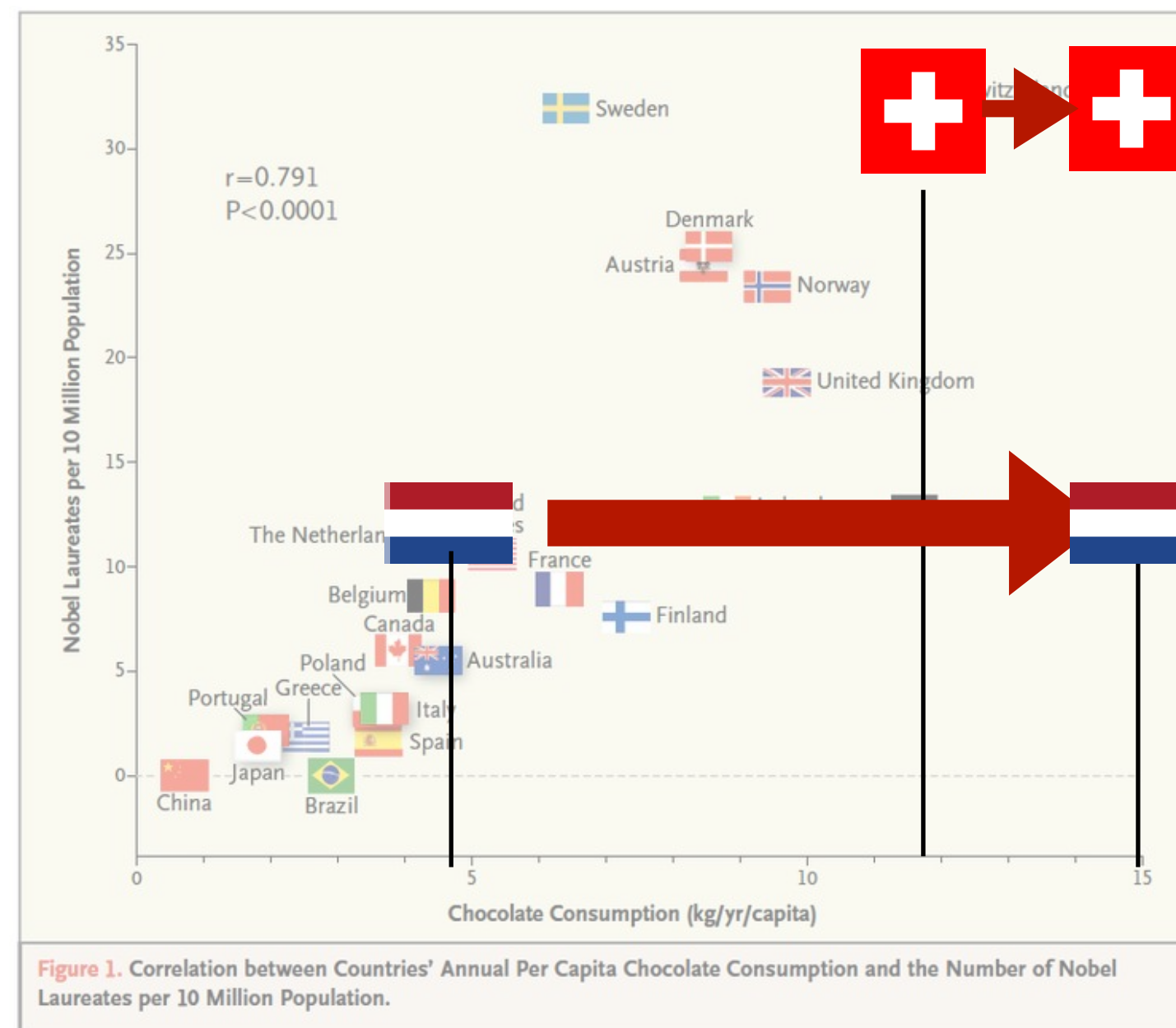


NL eats more chocolate => nothing changes



Learning from interventional data - intuition

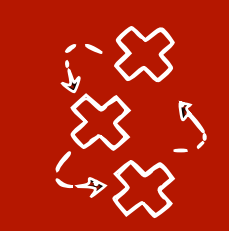
Hypothetical world: we perform the experiment and see these results:



NL eats more chocolate => nothing changes

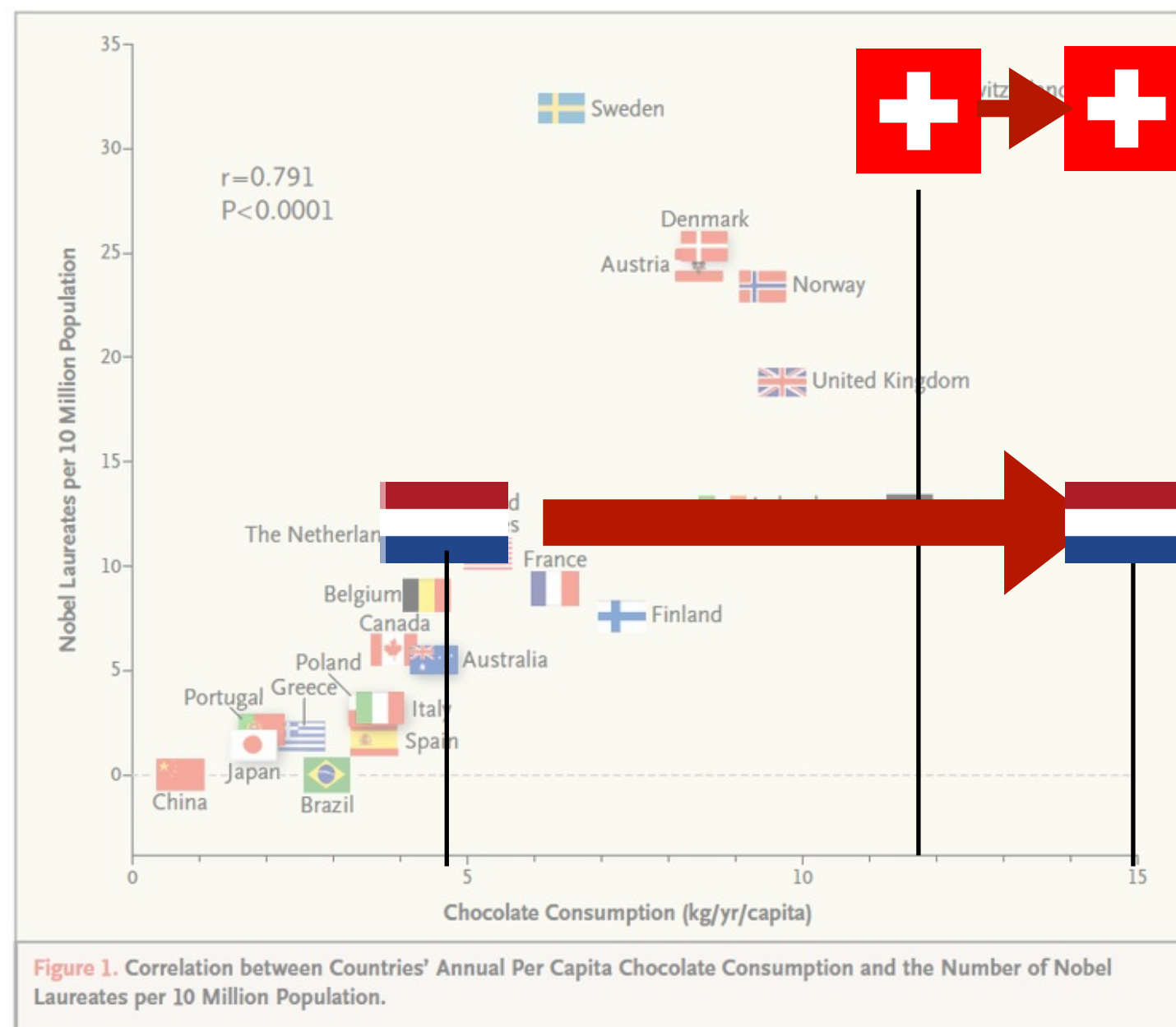
... and similarly for other countries (and other values)

Chocolate does not cause Nobel prizes

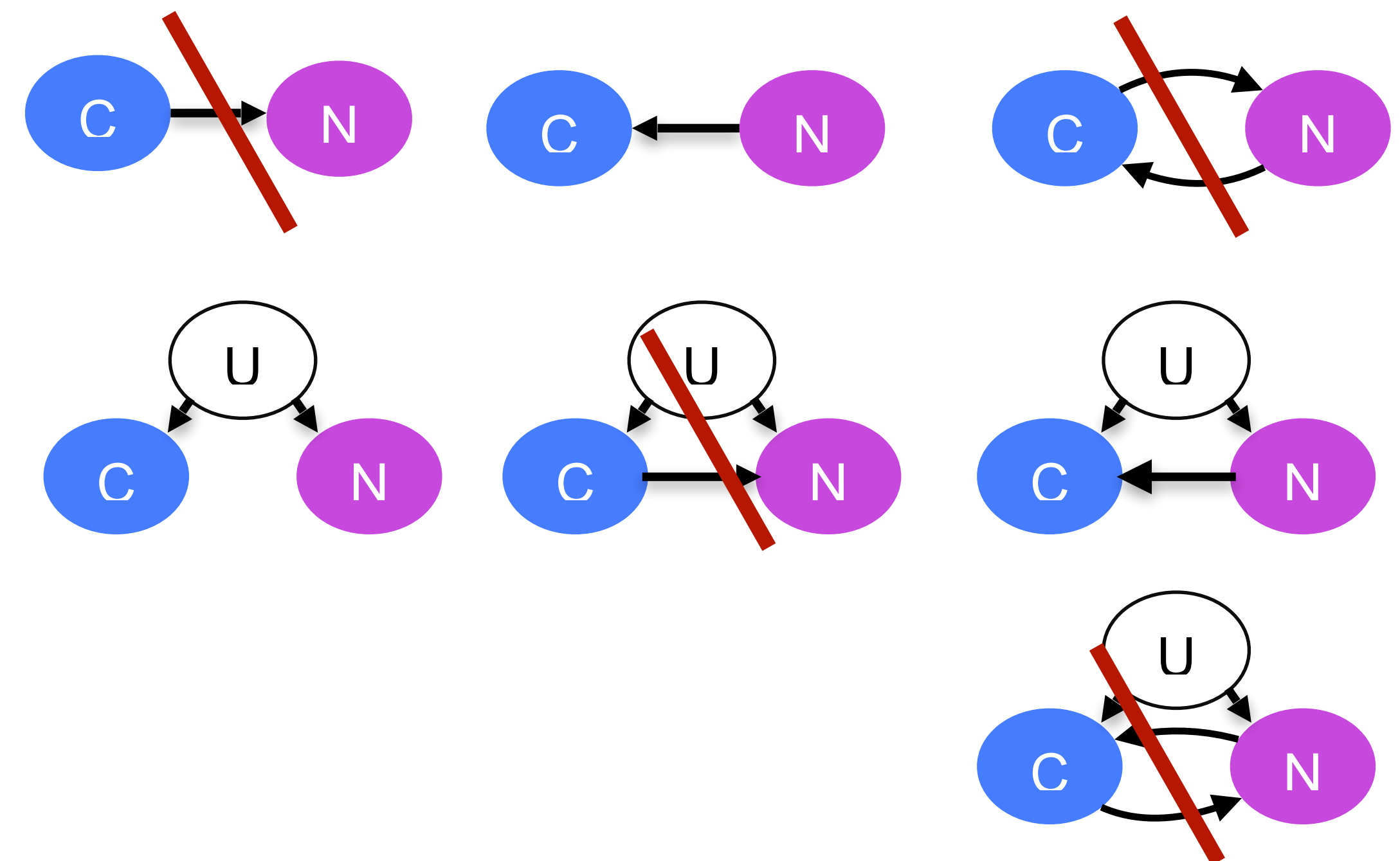


Learning from interventional data - intuition

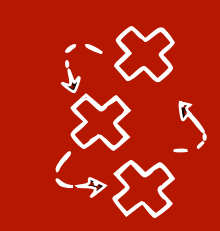
Hypothetical world: we perform the experiment and see these results:



Chocolate does not cause Nobel prizes

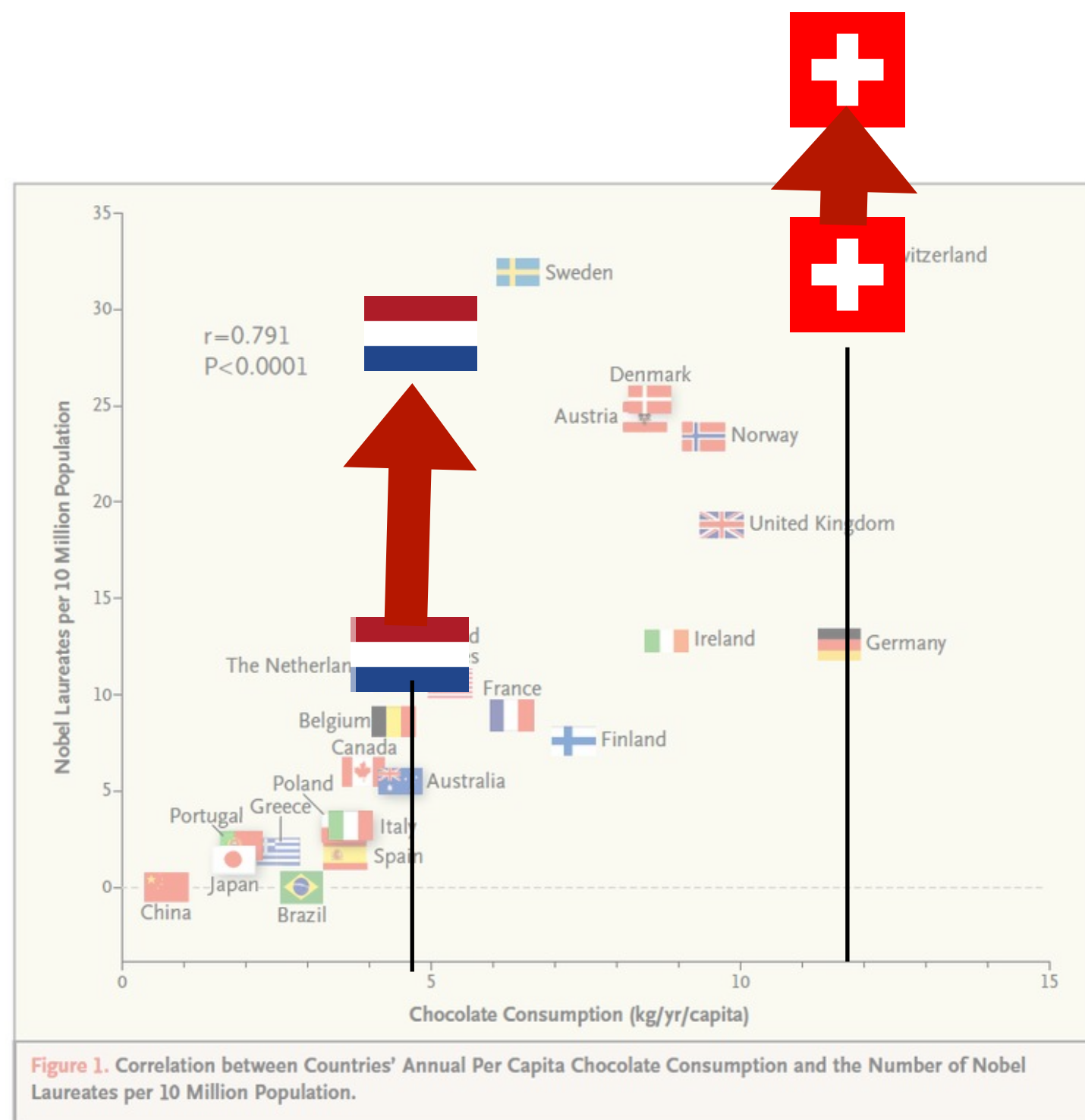


[Messerli, 2012] <https://www.nejm.org/doi/full/10.1056/NEJMon1211064>

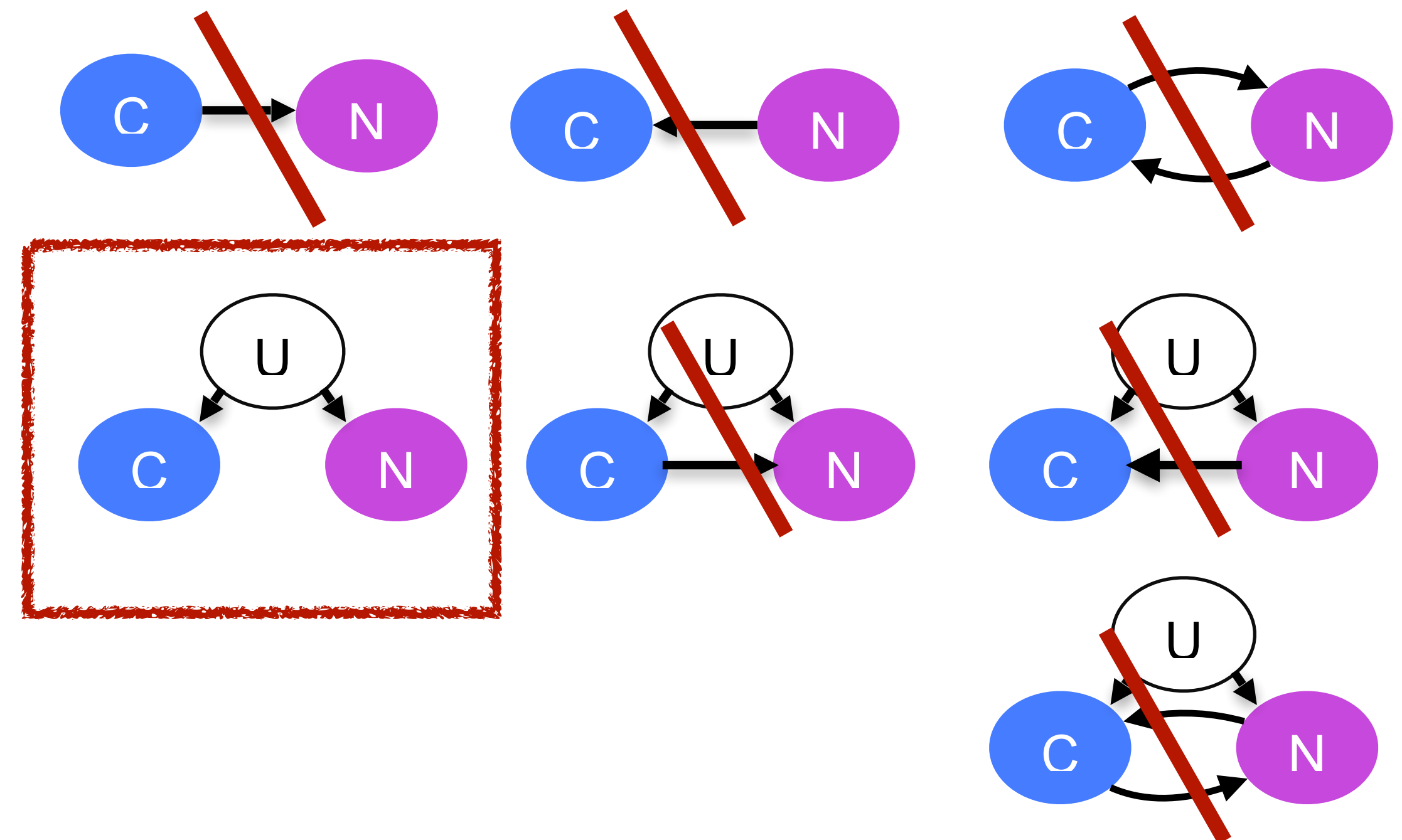


Learning from interventional data - intuition

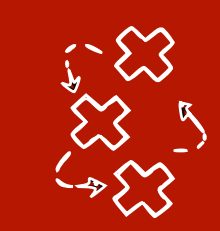
Hypothetical world: we perform another experiment and see these results:



Nobel does not cause chocolate

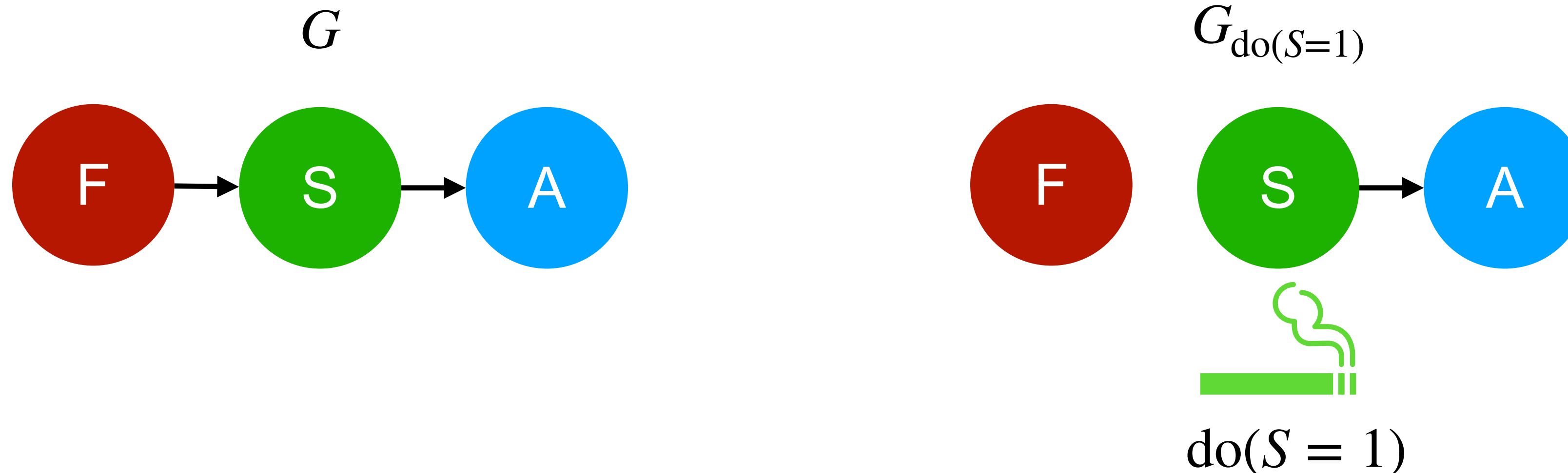


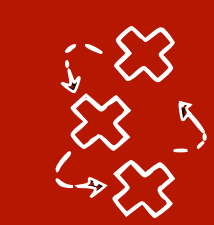
[Messerli, 2012] <https://www.nejm.org/doi/full/10.1056/NEJMon1211064>



Single-node interventions identify parents and children

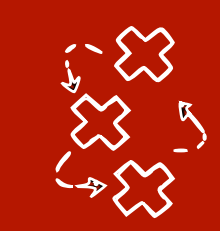
- The skeleton (and v-structures) can be identified from observational data
- Intervening on a node i identifies its parents and children:
 - For all j adjacent i in G :
 - If j is not adjacent i in $G_{do(i)}$ then $j \in Pa(i)$. Otherwise $j \in Ch(i)$





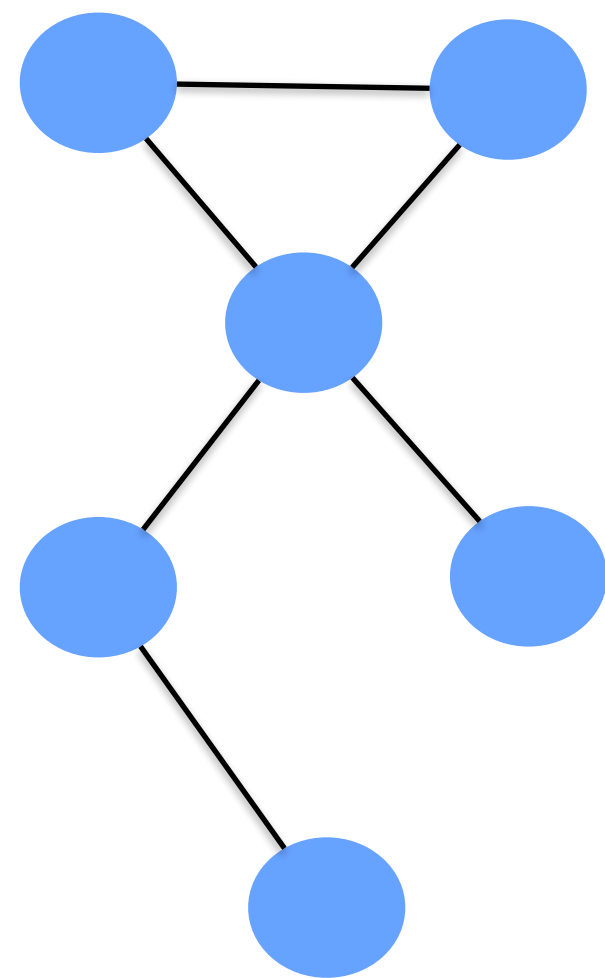
Single-node interventions identify parents and children

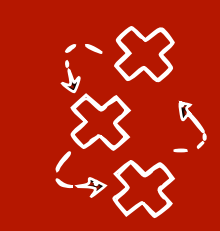
- The skeleton (and v-structures) can be identified from observational data
- Intervening on a node i identifies its parents and children:
 - For all j adjacent i in G :
 - If j is not adjacent i in $G_{do(i)}$ then $j \in \text{Pa}(i)$. Otherwise $j \in \text{Ch}(i)$
- Worst case: need $p - 1$ interventions to fully identify the graph (for $p > 2$)
- **Intervention on multiple nodes:** worst case need $O(\log_2 p)$ experiments to fully orient the graph [\[Hyttinen et al 2013\]](#).



Side note: Intervention design/Experiment selection

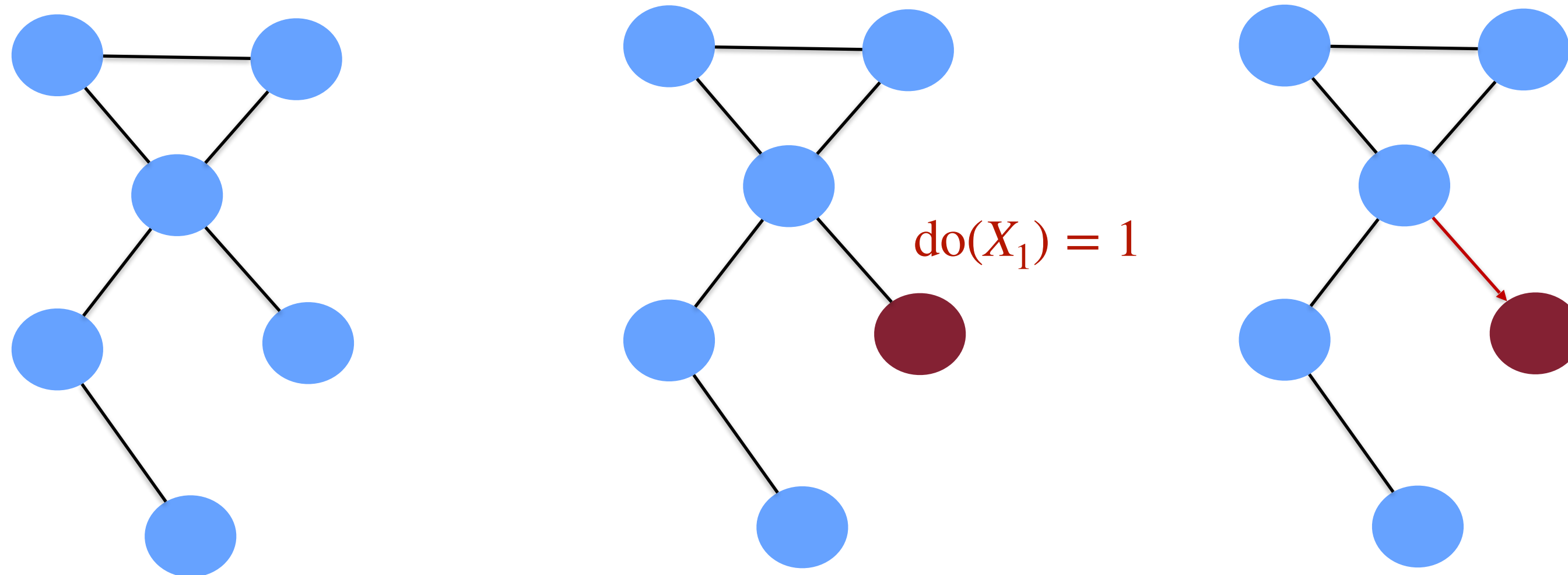
Design a set of interventions, so that we can **accurately** reconstruct **as much as possible** the causal graph **with the least samples**, also when **noisy**

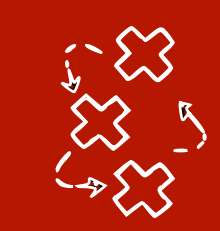




Side note: Intervention design/Experiment selection

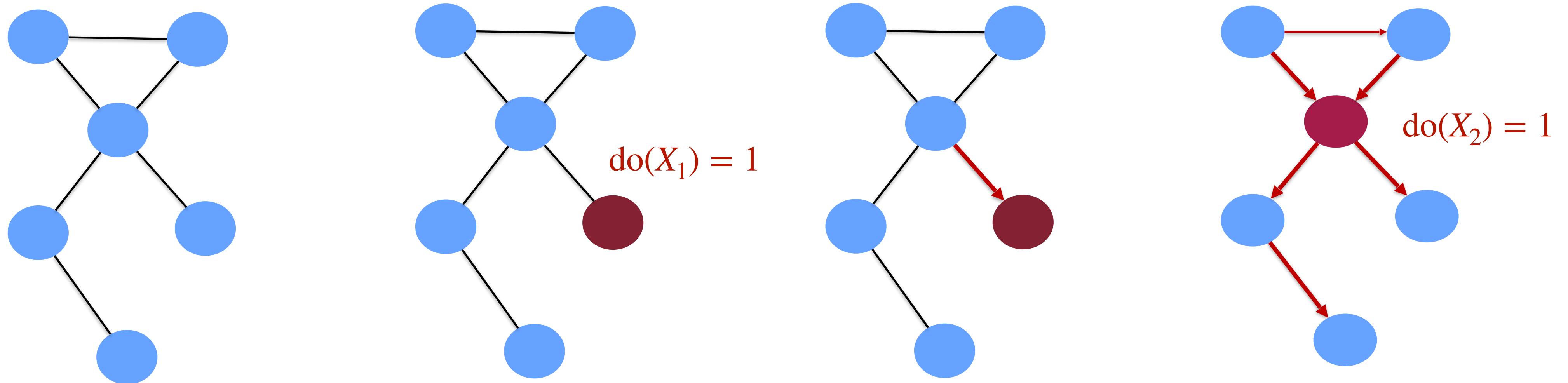
Design a set of interventions, so that we can **accurately** reconstruct **as much as possible** the causal graph **with the least samples**, also when **noisy**

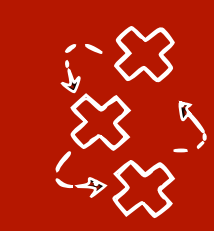




Side note: Intervention design/Experiment selection

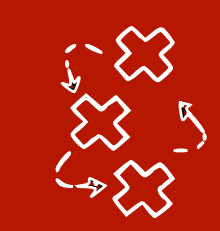
Design a set of interventions, so that we can **accurately** reconstruct **as much as possible** the causal graph **with the least samples**, also when **noisy**





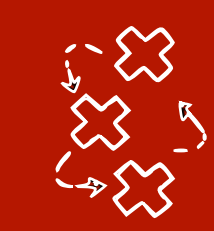
Learning from multiple contexts

- Now we cannot decide which intervention to perform (**intervention design**)
 - In intervention design, we have **known intervention targets**, e.g. $\text{do}(S = 1)$



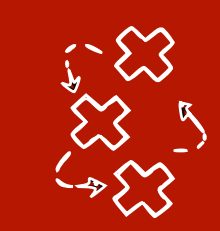
Learning from multiple contexts

- Now we cannot decide which intervention to perform (**intervention design**)
 - In intervention design, we have **known intervention targets**, e.g. $\text{do}(S = 1)$
- Instead, somebody gives us a **set of data from multiple contexts**
 - Possibly **unknown intervention targets**



Learning from multiple contexts

- Now we cannot decide which intervention to perform (**intervention design**)
 - In intervention design, we have **known intervention targets**, e.g. $\text{do}(S = 1)$
- Instead, somebody gives us a **set of data from multiple contexts**
 - Possibly **unknown intervention targets**
 - Possibly **soft interventions** instead of **perfect interventions**

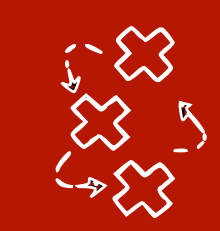


Perfect vs soft interventions

- Recap: we introduced an operator that can represent a **hypothetical intervention** on the whole population, i.e. a perturbation of the system:

$$\text{do}(X_i = x_i) \text{ which changes } P(X_i | X_{\text{Pa}(i)}) \rightarrow \mathbf{1}(X_i = x_i)$$

- This is called a **perfect** (or surgical) **intervention**
- There are also other types of intervention, e.g. **soft interventions which change** $P(X_i | X_{\text{Pa}(i)}) \rightarrow P'(X_i | X_{\text{Pa}(i)})$, which might not change the graph



Example 3.2 in Elements of Causal Inference

$$\begin{cases} X \leftarrow \epsilon_x \\ Y \leftarrow 4 \cdot X + \epsilon_Y \end{cases}$$

$$\epsilon_X, \epsilon_Y \sim \mathcal{N}(0,1)$$

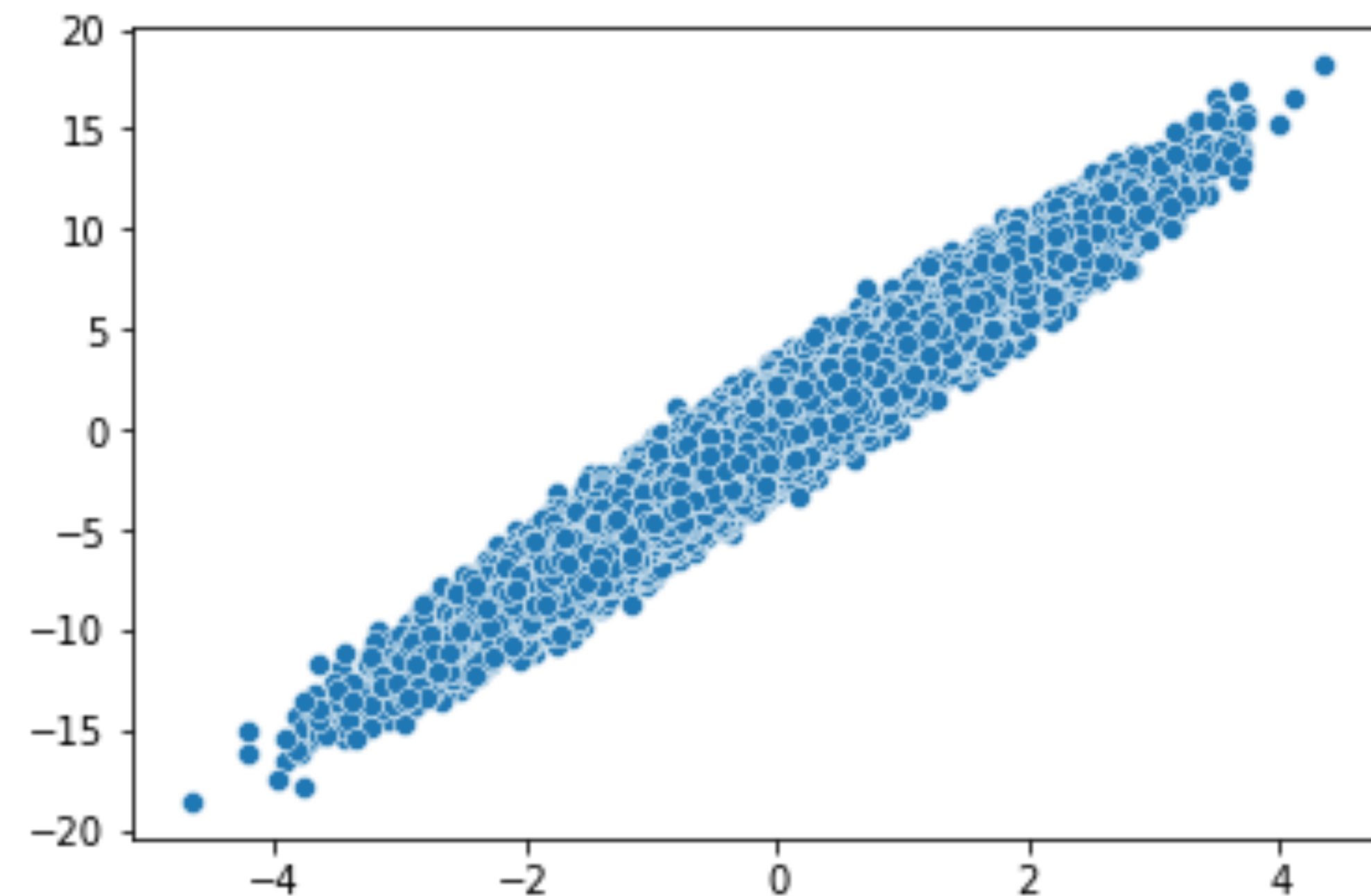
$$P(X) = \mathcal{N}(0,1)$$

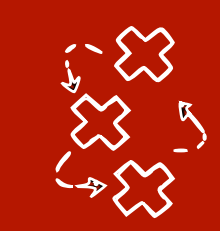
$$P(Y) = \mathcal{N}(0,17)$$

```
# We need a lot of samples to plot the conditional distribution:  
n_samples=100000
```

```
x = randn(n_samples)  
y = 4 * x + randn(n_samples)  
# plot P(X,Y)  
sns.scatterplot(x=x,y=y)
```

<AxesSubplot:>





Example 3.2 in Elements of Causal Inference

$$\begin{cases} X \leftarrow \epsilon_X \\ Y \leftarrow 4 \cdot X + \epsilon_Y \end{cases} \quad Y \leftarrow 4$$

$$\epsilon_X, \epsilon_Y \sim \mathcal{N}(0,1)$$

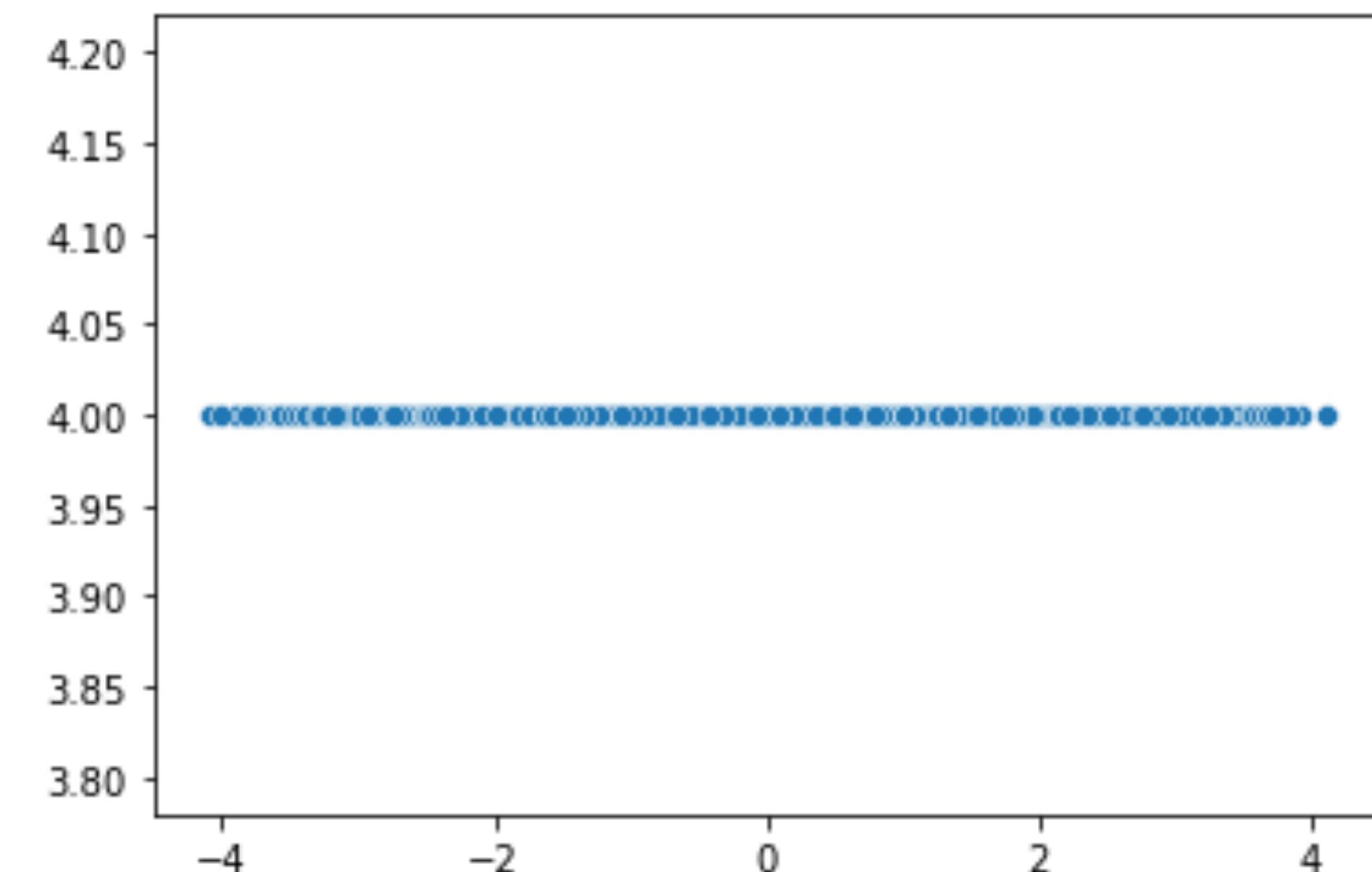
$$P(X \mid \text{do}(Y=4)) = \mathcal{N}(0,1)$$

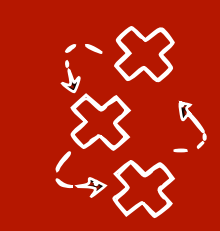
$$P(Y \mid \text{do}(Y=4)) = \begin{cases} 1 & Y=4 \\ 0 & Y \neq 4 \end{cases}$$

```
x_do_y = randn(n_samples)
y_do_y = 4

# plot P(X, Y | do(X=2))
sns.scatterplot(x=x_do_y, y=y_do_y)
```

<AxesSubplot:>





Soft interventions, shift interventions

$$\begin{cases} X \leftarrow \epsilon_x \\ Y \leftarrow 4 \cdot X + \epsilon_Y \end{cases}$$

$$\epsilon_X, \epsilon_Y \sim \mathcal{N}(0,1)$$

$$P(X) = \mathcal{N}(0,1)$$

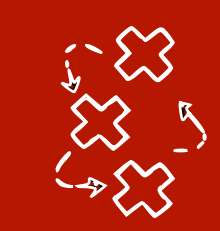
$$P(Y) = \mathcal{N}(0,17)$$

$$\begin{cases} X \leftarrow \epsilon_x \\ Y \leftarrow 4 \cdot X + \epsilon_Y \end{cases}$$

$$Y \leftarrow 2 \cdot X + \epsilon_Y$$

$$\begin{cases} X \leftarrow \epsilon_x \\ Y \leftarrow 4 \cdot X + \epsilon_Y + \epsilon' \end{cases}$$

$$Pa_{G^{do}(Y)}(Y) \subseteq Pa_G(Y)$$

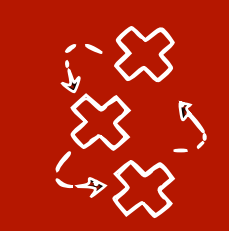


Causal parents as example of robust prediction

$$\begin{cases} X_1 = \varepsilon_1 \\ Y = X_1 + \varepsilon_Y \\ X_2 = Y + \varepsilon_{X_2} \end{cases}$$

$\varepsilon_1, \varepsilon_Y \sim N(0, 1), \varepsilon_{X_2} \sim N(0, 0.01)$

$$X_1 \rightarrow Y \rightarrow X_2$$



Causal parents as example of robust prediction

$$\begin{cases} X_1 = \varepsilon_1 \\ Y = X_1 + \varepsilon_Y \\ X_2 = Y + \varepsilon_{X_2} \end{cases}$$

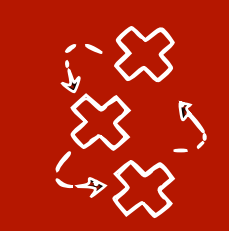
$$\varepsilon_1, \varepsilon_Y \sim N(0, 1), \quad \varepsilon_{X_2} \sim N(0, 0.01)$$

$$M1: Y \sim X_1$$

$$M2: Y \sim X_2$$

$$X_1 \rightarrow Y \rightarrow X_2$$

M2 has smaller error



Causal parents as example of robust prediction

$$\begin{cases} X_1 = \varepsilon_1 \\ Y = X_1 + \varepsilon_Y \\ X_2 = Y + \varepsilon_{X_2} \end{cases}$$

$\varepsilon_1, \varepsilon_Y \sim N(0, 1), \varepsilon_{X_2} \sim N(0, 0.01)$

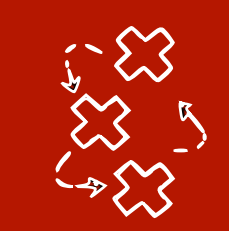
M1: $Y \sim X_1$

M2: $Y \sim X_2$

$$X_1 \rightarrow Y \rightarrow X_2$$

$$X_1 \rightarrow Y \quad X_2 \text{ do}(X_2)$$

M2 has smaller error
but it fails in $\text{do}(X_2)$



Causal parents as example of robust prediction

$$\begin{cases} X_1 = \varepsilon_1 \\ Y = X_1 + \varepsilon_Y \\ X_2 = Y + \varepsilon_{X_2} \end{cases}$$

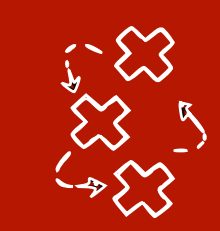
$$\varepsilon_1, \varepsilon_Y \sim N(0, 1), \quad \varepsilon_{X_2} \sim N(0, 0.01)$$

$$M1: Y \sim X_1$$

$$M2: Y \sim X_2$$

$$X_1 \rightarrow Y \rightarrow X_2$$

Using causal parents of Y helps
with DISTRIBUTION SHIFTS



Causal discovery simplified overview

Constraint-based causal discovery

Score-based causal discovery

Restricted models

- Nonlinear additive

Interventional causal discovery / causal invariance

Causal inference using invariant prediction: identification and confidence intervals

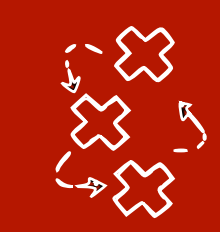
Jonas Peters, Peter Bühlmann, Nicolai Meinshausen

- Observational data
- Output: MEC
- SGS, PC

- Output: MEC
- GES, MMHC

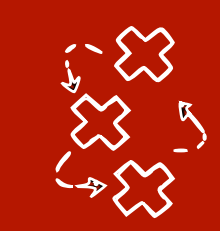
- Observational data
- Output: DAG
- RESIT, LINGAM

- Interventional data
- Output: parents of Y, I-MEC
 - ICP, JCI



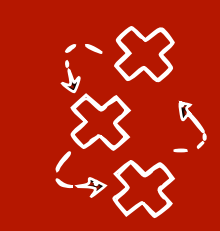
Invariant Causal Prediction (ICP) [Peters et al 2016]

- Given a target variable Y and features (X_1, \dots, X_p) , we want to **find the causal parents of Y , i.e. $\text{Pa}(Y)$**



Invariant Causal Prediction (ICP) [Peters et al 2016]

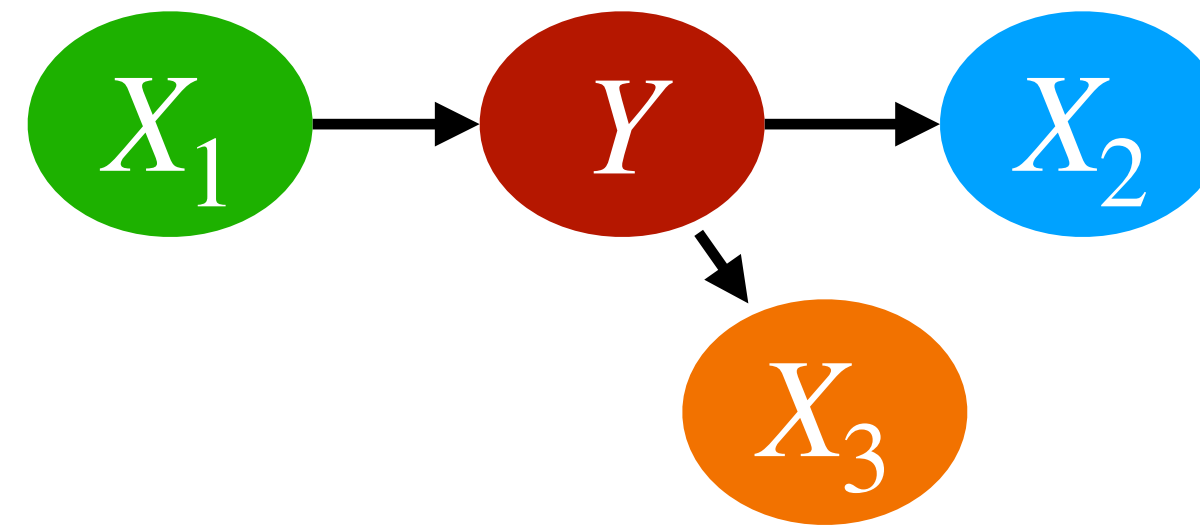
- Given a target variable Y and features (X_1, \dots, X_p) , we want to **find the causal parents of Y , i.e. $\text{Pa}(Y)$**
- We assume we have access to **a set of different environments \mathcal{E}** (e.g. interventional or observational data), s.t. for $e \in \mathcal{E}$, $(X_1^e, \dots, X_p^e, Y^e) \sim P^e$



Multiple environments

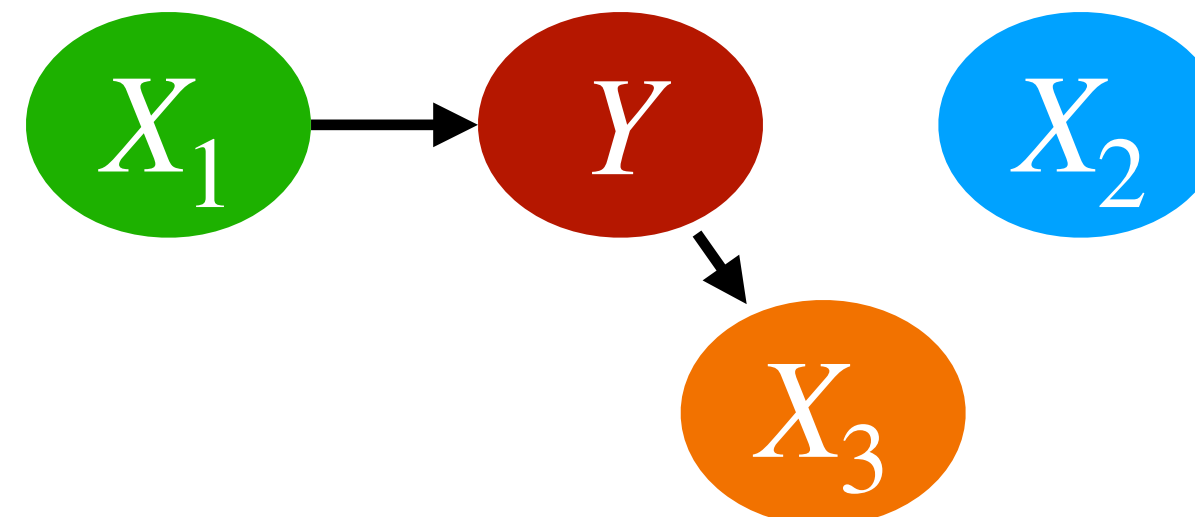
$$\begin{cases} \epsilon_1, \epsilon_2, \epsilon_3, \epsilon_Y \sim \mathcal{N}(0,1) \\ X_1 = 10 + \epsilon_1 \\ Y = 3X_1 + \epsilon_Y \\ X_2 = -2Y + \epsilon_2 \\ X_3 = 2Y + 0.1\epsilon_3 \end{cases}$$

$E = 0$



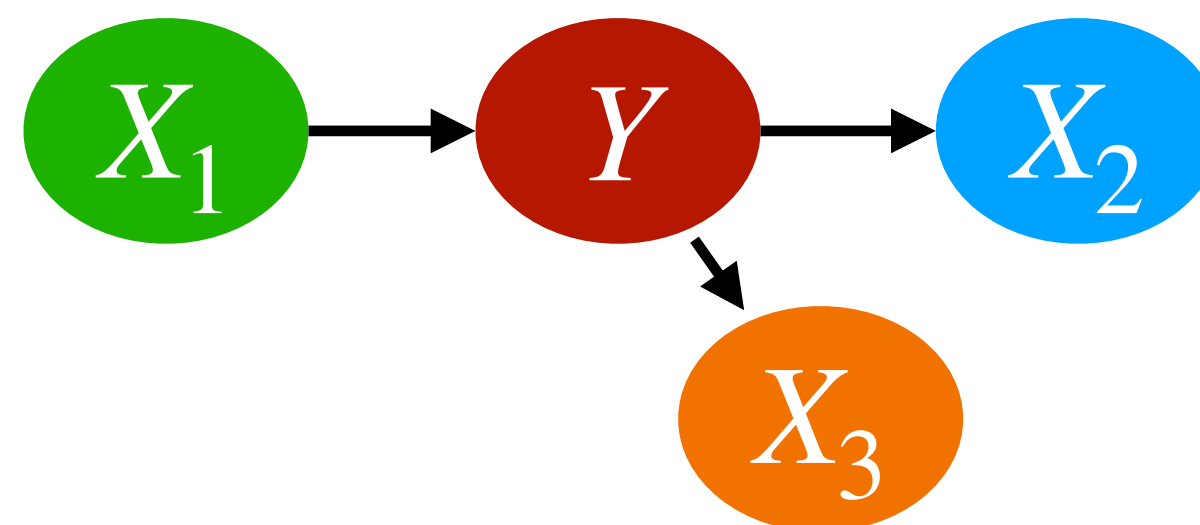
$$\begin{cases} \epsilon_1, \epsilon_2, \epsilon_3, \epsilon_Y \sim \mathcal{N}(0,1) \\ X_1 = 10 + \epsilon_1 \\ Y = 3X_1 + \epsilon_Y \\ X_2 = 1 \\ X_3 = 2Y + 0.1\epsilon_3 \end{cases}$$

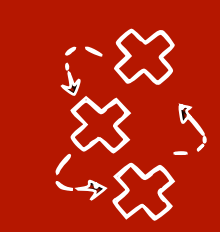
$E = 1$



$$\begin{cases} \epsilon_1, \epsilon_2, \epsilon_3, \epsilon_Y \sim \mathcal{N}(0,1) \\ X_1 = 10 + \epsilon_1 \\ Y = 3X_1 + \epsilon_Y \\ X_2 = 10Y + \epsilon_2 \\ X_3 = 2Y + 0.1\epsilon_3 \end{cases}$$

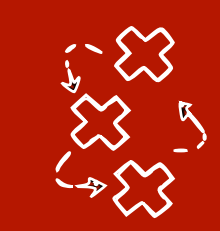
$E = 2$





Invariant Causal Prediction (ICP) [Peters et al 2016]

- Given a target variable Y and features (X_1, \dots, X_p) , we want to **find the causal parents of Y , i.e. $\text{Pa}(Y)$**
- We assume we have access to **a set of different environments \mathcal{E}** (e.g. interventional or observational data), s.t. for $e \in \mathcal{E}$, $(X_1^e, \dots, X_p^e, Y^e) \sim P^e$
- We further assume that in **none of the environments Y is intervened upon**
 - We can then show that $e, f \in \mathcal{E} : P^e(Y^e | \text{Pa}(Y^e)) = P^f(Y^f | \text{Pa}(Y^f))$
- We represent the environment index with E



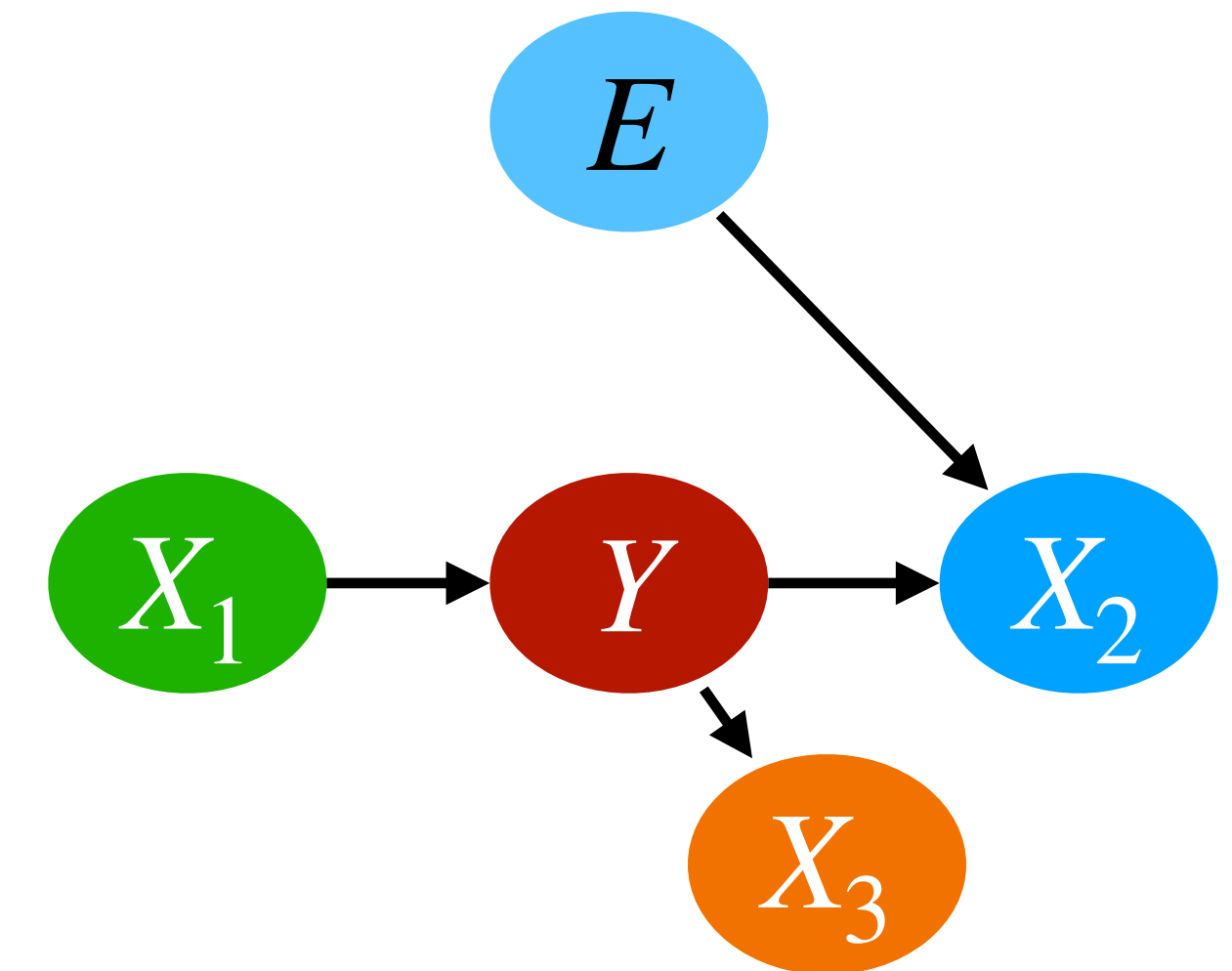
Multiple environments in a single SCM

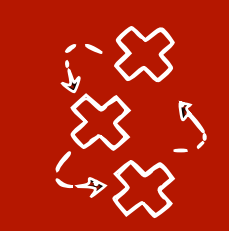
$$\begin{cases} \epsilon_1, \epsilon_2, \epsilon_3, \epsilon_Y \sim \mathcal{N}(0,1) \\ X_1 = 10 + \epsilon_1 \\ Y = 3X_1 + \epsilon_Y \\ X_2 = -2Y + \epsilon_2 \\ X_3 = 2Y + 0.1\epsilon_3 \end{cases} \quad E = 0$$

$$\begin{cases} \epsilon_1, \epsilon_2, \epsilon_3, \epsilon_Y \sim \mathcal{N}(0,1) \\ X_1 = 10 + \epsilon_1 \\ Y = 3X_1 + \epsilon_Y \\ X_2 = 1 \\ X_3 = 2Y + 0.1\epsilon_3 \end{cases} \quad E = 1$$

$$\begin{cases} \epsilon_1, \epsilon_2, \epsilon_3, \epsilon_Y \sim \mathcal{N}(0,1) \\ X_1 = 10 + \epsilon_1 \\ Y = 3X_1 + \epsilon_Y \\ X_2 = 10Y + \epsilon_2 \\ X_3 = 2Y + 0.1\epsilon_3 \end{cases} \quad E = 2$$

$$\begin{cases} \epsilon_1, \epsilon_2, \epsilon_3, \epsilon_Y \sim \mathcal{N}(0,1) \\ X_1 = 10 + \epsilon_1 \\ Y = 3X_1 + \epsilon_Y \\ X_2 = \begin{cases} -2Y + \epsilon_2 & \text{if } E = 0 \\ 1 & \text{if } E = 1 \\ 10Y + \epsilon_2 & \text{if } E = 2 \end{cases} \\ X_3 = 2Y + 0.1\epsilon_3 \end{cases}$$





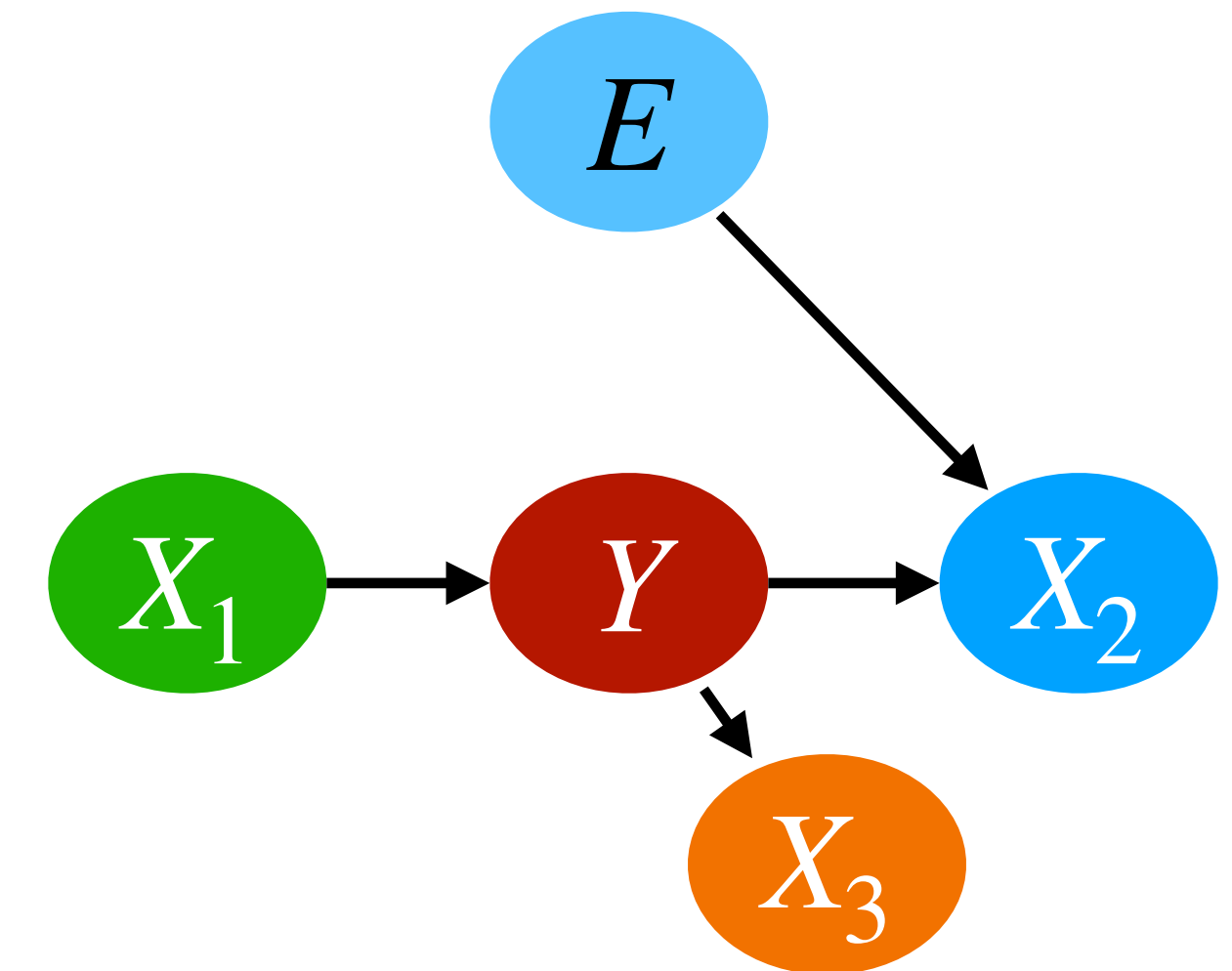
Multiple environments in a single SCM

$$\begin{cases} \epsilon_1, \epsilon_2, \epsilon_3, \epsilon_Y \sim \mathcal{N}(0,1) \\ X_1 = 10 + \epsilon_1 \\ Y = 3X_1 + \epsilon_Y \\ X_2 = -2Y + \epsilon_2 \\ X_3 = 2Y + 0.1\epsilon_3 \end{cases} \quad E = 0$$

$$\begin{cases} \epsilon_1, \epsilon_2, \epsilon_3, \epsilon_Y \sim \mathcal{N}(0,1) \\ X_1 = 10 + \epsilon_1 \\ Y = 3X_1 + \epsilon_Y \\ X_2 = 1 \\ X_3 = 2Y + 0.1\epsilon_3 \end{cases} \quad E = 1$$

$$\begin{cases} \epsilon_1, \epsilon_2, \epsilon_3, \epsilon_Y \sim \mathcal{N}(0,1) \\ X_1 = 10 + \epsilon_1 \\ Y = 3X_1 + \epsilon_Y \\ X_2 = 10Y + \epsilon_2 \\ X_3 = 2Y + 0.1\epsilon_3 \end{cases} \quad E = 2$$

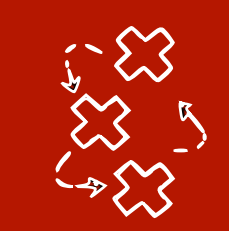
$$\begin{cases} \epsilon_1, \epsilon_2, \epsilon_3, \epsilon_Y \sim \mathcal{N}(0,1) \\ X_1 = 10 + \epsilon_1 \\ Y = 3X_1 + \epsilon_Y \\ X_2 = \begin{cases} -2Y + \epsilon_2 & \text{if } E = 0 \\ 1 & \text{if } E = 1 \\ 10Y + \epsilon_2 & \text{if } E = 2 \end{cases} \\ X_3 = 2Y + 0.1\epsilon_3 \end{cases}$$



$$P^e(Y^e | \text{Pa}(Y^e)) = P^f(Y^f | \text{Pa}(Y^f))$$

$$E \perp\!\!\!\perp Y | \text{Pa}(Y) \iff E \perp_d Y | \text{Pa}(Y)$$

$$E \perp_d Y | X_1$$

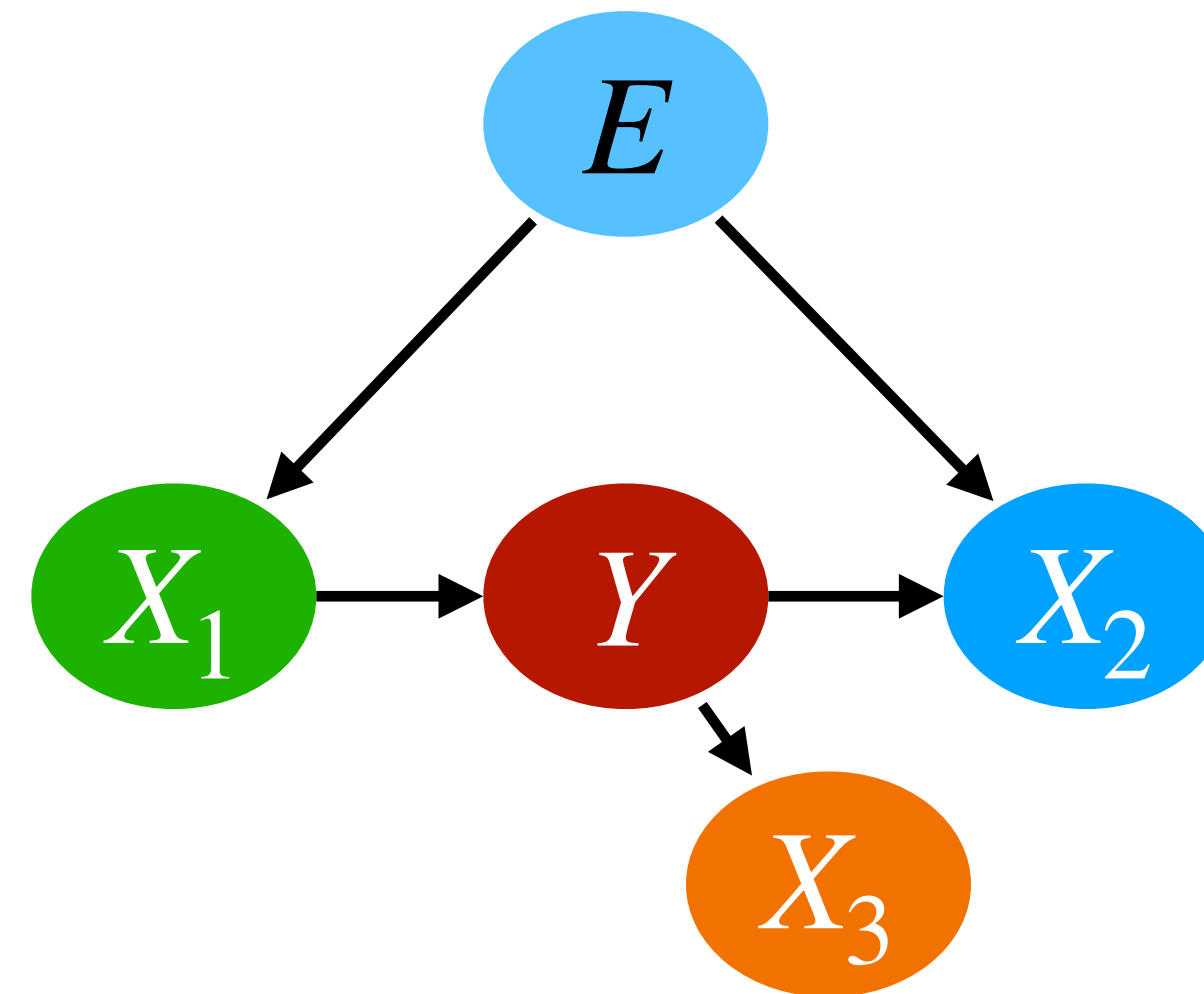


Multiple environments in a single SCM

$$\begin{cases} \epsilon_1, \epsilon_2, \epsilon_3, \epsilon_Y \sim \mathcal{N}(0,1) \\ X_1 = 10 + \epsilon_1 \\ Y = 3X_1 + \epsilon_Y \\ X_2 = -2Y + \epsilon_2 \\ X_3 = 2Y + 0.1\epsilon_3 \end{cases} \quad E = 0$$

$$\begin{cases} \epsilon_1, \epsilon_2, \epsilon_3, \epsilon_Y \sim \mathcal{N}(0,1) \\ X_1 = 100 + \epsilon_1 \\ Y = 3X_1 + \epsilon_Y \\ X_2 = 1 \\ X_3 = 2Y + 0.1\epsilon_3 \end{cases} \quad E = 1$$

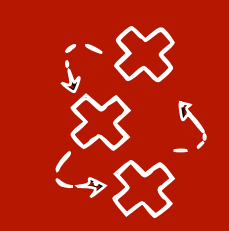
$$\begin{cases} \epsilon_1, \epsilon_2, \epsilon_3, \epsilon_Y \sim \mathcal{N}(0,1) \\ X_1 = 10 + \epsilon_1 \\ Y = 3X_1 + \epsilon_Y \\ X_2 = 10Y + \epsilon_2 \\ X_3 = 2Y + 0.1\epsilon_3 \end{cases} \quad E = 2$$



$$\mathbf{S} = \text{Pa}(Y) \implies E \perp_d Y | \mathbf{S}$$

$$E \perp_d Y | \mathbf{S} \implies \mathbf{S} = \text{Pa}(Y)?$$

$$E \perp_d Y | X_1$$

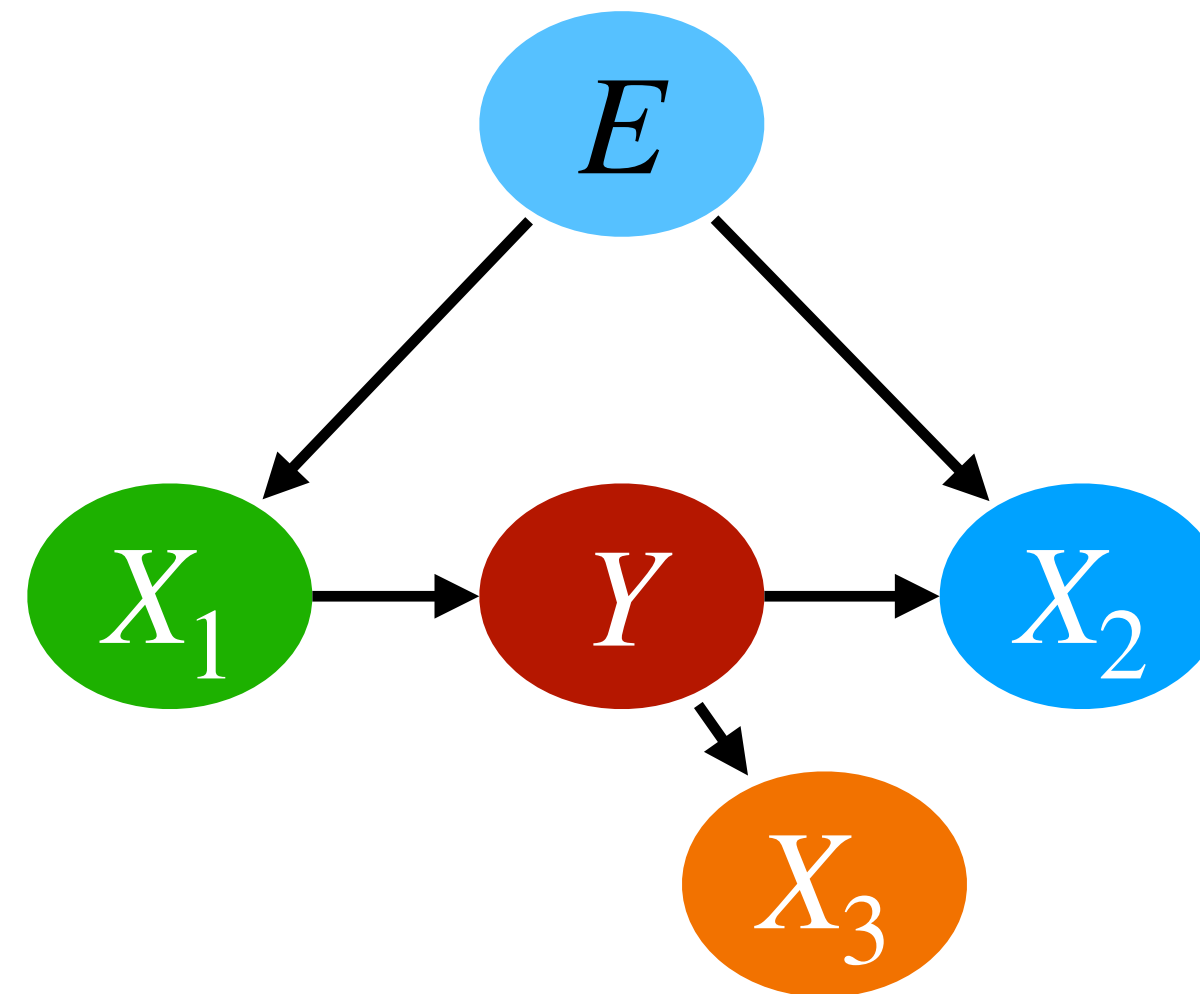


Multiple environments in a single SCM

$$\begin{cases} \epsilon_1, \epsilon_2, \epsilon_3, \epsilon_Y \sim \mathcal{N}(0,1) \\ X_1 = 10 + \epsilon_1 \\ Y = 3X_1 + \epsilon_Y \\ X_2 = -2Y + \epsilon_2 \\ X_3 = 2Y + 0.1\epsilon_3 \end{cases} \quad E = 0$$

$$\begin{cases} \epsilon_1, \epsilon_2, \epsilon_3, \epsilon_Y \sim \mathcal{N}(0,1) \\ X_1 = 100 + \epsilon_1 \\ Y = 3X_1 + \epsilon_Y \\ X_2 = 1 \\ X_3 = 2Y + 0.1\epsilon_3 \end{cases} \quad E = 1$$

$$\begin{cases} \epsilon_1, \epsilon_2, \epsilon_3, \epsilon_Y \sim \mathcal{N}(0,1) \\ X_1 = 10 + \epsilon_1 \\ Y = 3X_1 + \epsilon_Y \\ X_2 = 10Y + \epsilon_2 \\ X_3 = 2Y + 0.1\epsilon_3 \end{cases} \quad E = 2$$



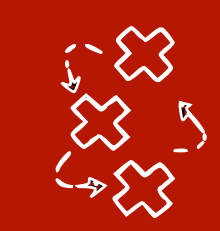
Other subset of nodes also satisfy this relationship

$$\mathbf{S} = \text{Pa}(Y) \implies E \perp_d Y | \mathbf{S}$$

$$E \perp_d Y | \mathbf{S} \not\implies \mathbf{S} = \text{Pa}(Y)$$

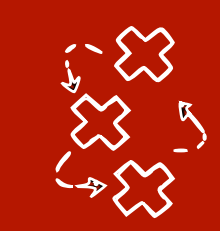
$$E \perp_d Y | X_1$$

$$E \perp_d Y | \{X_1, X_3\}$$



Invariant Causal Prediction (ICP)

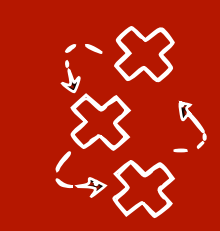
- We assume we have access to **a set of different environments** \mathcal{E} (e.g. interventional or observational data), s.t. for $e \in \mathcal{E}$, $(X_1^e, \dots, X_p^e, Y^e) \sim P^e$
- We further assume that in **none of the environments Y is intervened upon**



Invariant Causal Prediction (ICP)

- We assume we have access to **a set of different environments** \mathcal{E} (e.g. interventional or observational data), s.t. for $e \in \mathcal{E}$, $(X_1^e, \dots, X_p^e, Y^e) \sim P^e$
- We further assume that in **none of the environments** Y is intervened upon
- We represent the environment index with E
- If there are **no latent confounders**, one can prove that:

$$\bigcap_{S \subseteq \{1, \dots, p\} \text{ s.t. } E \perp\!\!\!\perp Y | S} S \subseteq \text{Pa}(Y)$$

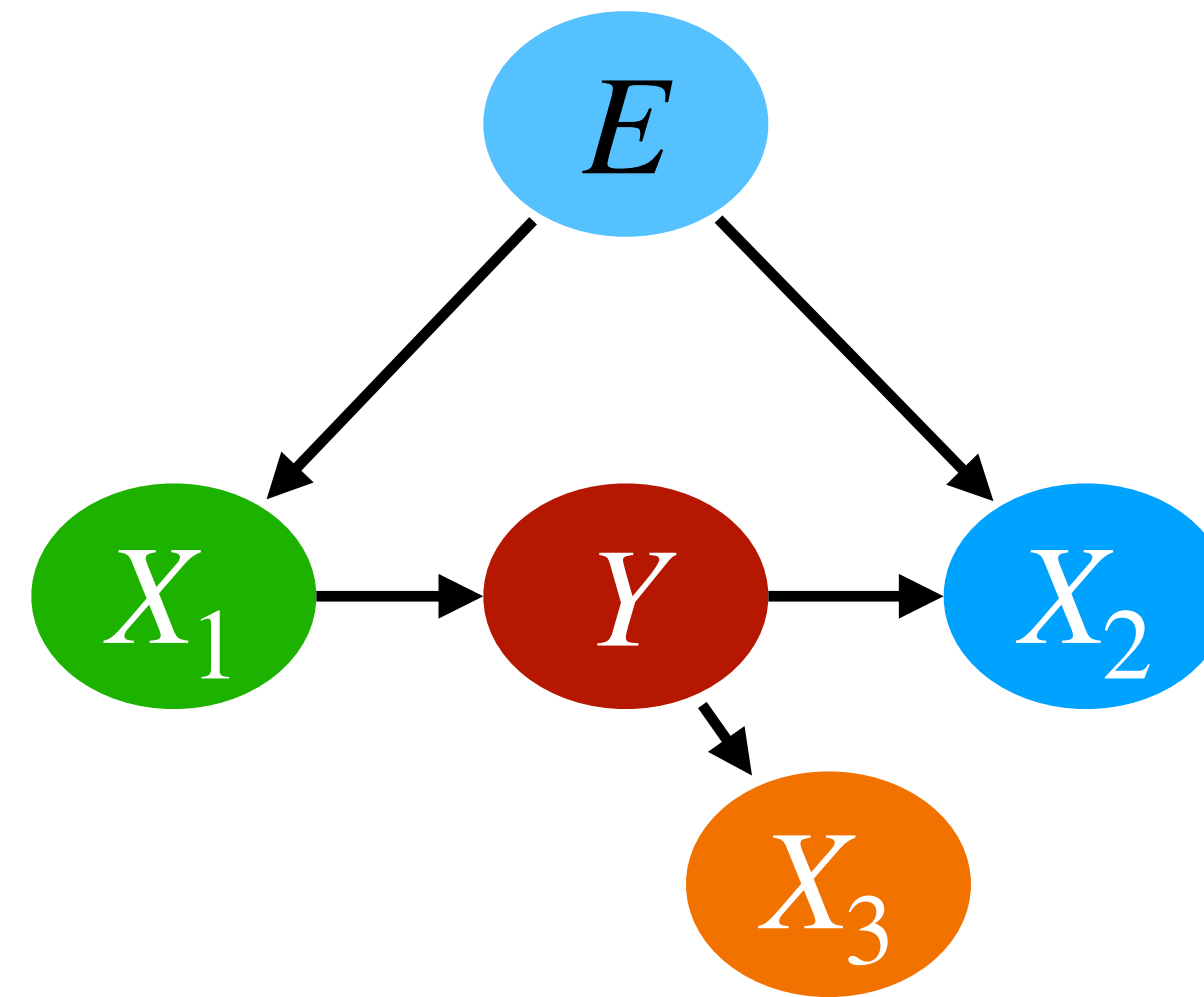


Invariant Causal Prediction example

$$\begin{cases} \epsilon_1, \epsilon_2, \epsilon_3, \epsilon_Y \sim \mathcal{N}(0,1) \\ X_1 = 10 + \epsilon_1 \\ Y = 3X_1 + \epsilon_Y \\ X_2 = -2Y + \epsilon_2 \\ X_3 = 2Y + 0.1\epsilon_3 \end{cases} \quad E = 0$$

$$\begin{cases} \epsilon_1, \epsilon_2, \epsilon_3, \epsilon_Y \sim \mathcal{N}(0,1) \\ X_1 = 100 + \epsilon_1 \\ Y = 3X_1 + \epsilon_Y \\ X_2 = 1 \\ X_3 = 2Y + 0.1\epsilon_3 \end{cases} \quad E = 1$$

$$\begin{cases} \epsilon_1, \epsilon_2, \epsilon_3, \epsilon_Y \sim \mathcal{N}(0,1) \\ X_1 = 10 + \epsilon_1 \\ Y = 3X_1 + \epsilon_Y \\ X_2 = 10Y + \epsilon_2 \\ X_3 = 2Y + 0.1\epsilon_3 \end{cases} \quad E = 2$$

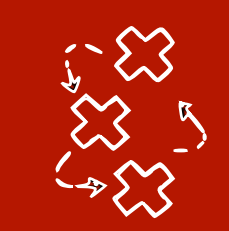


$$\bigcap_{S \subseteq \{1, \dots, p\} \text{ s.t. } E \perp\!\!\!\perp Y | S} S \subseteq \text{Pa}(Y)$$

$$E \perp_d Y | X_1$$

$$E \perp_d Y | \{X_1, X_3\}$$

$$\{X_1\} \cap \{X_1, X_3\} = \{X_1\} = \text{Pa}(Y)$$

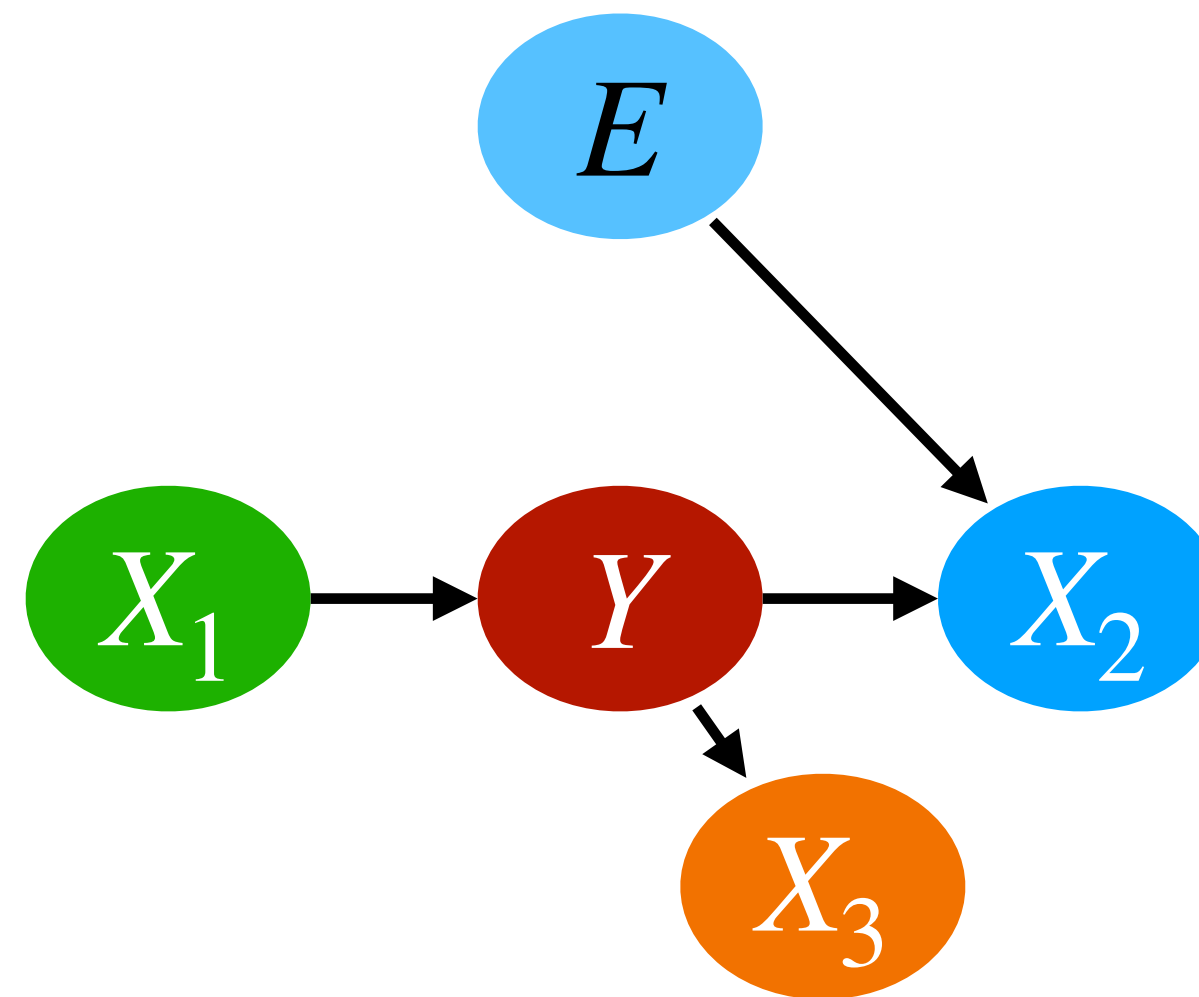


Invariant Causal Prediction example 2

$$\begin{cases} \epsilon_1, \epsilon_2, \epsilon_3, \epsilon_Y \sim \mathcal{N}(0,1) \\ X_1 = 10 + \epsilon_1 \\ Y = 3X_1 + \epsilon_Y \\ X_2 = -2Y + \epsilon_2 \\ X_3 = 2Y + 0.1\epsilon_3 \end{cases} \quad E = 0$$

$$\begin{cases} \epsilon_1, \epsilon_2, \epsilon_3, \epsilon_Y \sim \mathcal{N}(0,1) \\ X_1 = 10 + \epsilon_1 \\ Y = 3X_1 + \epsilon_Y \\ X_2 = 1 \\ X_3 = 2Y + 0.1\epsilon_3 \end{cases} \quad E = 1$$

$$\begin{cases} \epsilon_1, \epsilon_2, \epsilon_3, \epsilon_Y \sim \mathcal{N}(0,1) \\ X_1 = 10 + \epsilon_1 \\ Y = 3X_1 + \epsilon_Y \\ X_2 = 10Y + \epsilon_2 \\ X_3 = 2Y + 0.1\epsilon_3 \end{cases} \quad E = 2$$



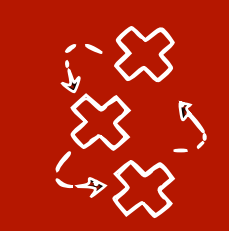
$$\bigcap_{S \subseteq \{1, \dots, p\} \text{ s.t. } E \perp\!\!\!\perp Y | S} S \subseteq \text{Pa}(Y)$$

$$E \perp_d Y | X_1 \quad \{X_1\} \cap \{X_1, X_3\} \cap \emptyset = \emptyset \subseteq \text{Pa}(Y)$$

$$E \perp_d Y | \{X_1, X_3\}$$

$$E \perp_d Y$$

ICP finds SUBSETS of parents



ICP improves with more environments

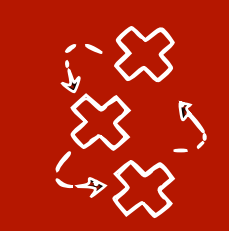
$$E \rightarrow X_1 \rightarrow X_2 \rightarrow Y$$

+ new environment e_3

$$X_2 = f^{\text{new}}(X_1, \varepsilon_2)$$

$$\left. \begin{array}{l} E \perp\!\!\!\perp Y \mid X_1 \\ E \perp\!\!\!\perp Y \mid X_2 \\ E \perp\!\!\!\perp Y \mid X_1, X_2 \end{array} \right\} \cap \emptyset$$

ICP on these environments finds empty set



ICP improves with more environments

$$E \rightarrow X_1 \rightarrow X_2 \rightarrow Y$$

$$\left. \begin{array}{l} E \perp\!\!\!\perp Y \mid X_1 \\ E \perp\!\!\!\perp Y \mid X_2 \\ E \perp\!\!\!\perp Y \mid X_1, X_2 \end{array} \right\} \cap \emptyset$$

ICP on these environments finds empty set

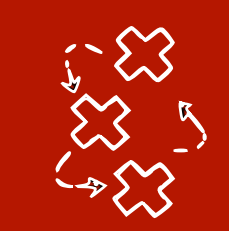
+ new environment e_3

$$X_2 = f^{\text{new}}(X_1, \varepsilon_2)$$



ICP on these environments (including new) finds X_2

$$E \perp\!\!\!\perp Y \mid X_2 \cap X_2 \in Pa(Y) \\ E \perp\!\!\!\perp Y \mid X_2, X_1$$



ICP improves with more environments

$$E \rightarrow X_1 \rightarrow X_2 \rightarrow Y$$

$$\left. \begin{array}{l} E \perp\!\!\!\perp Y \mid X_1 \\ E \perp\!\!\!\perp Y \mid X_2 \\ E \perp\!\!\!\perp Y \mid X_1, X_2 \end{array} \right\} \cap \emptyset$$

ICP on these environments finds empty set

+ new environment e_3

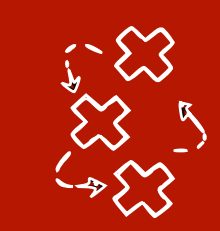
$$X_2 = f^{\text{new}}(X_1, \varepsilon_2)$$



ICP on these environments (including new) finds X_2

$$E \perp\!\!\!\perp Y \mid X_2 \cap X_2 \in \text{Pa}(Y) \\ E \perp\!\!\!\perp Y \mid X_2, X_1$$

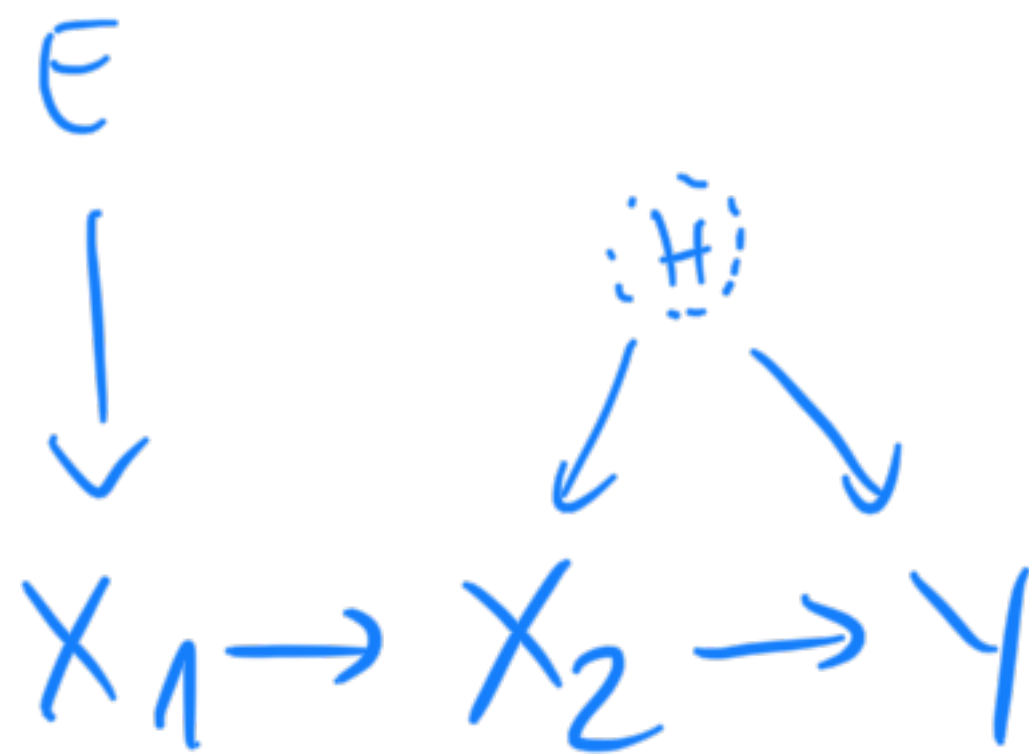
If all variables are caused by E (so we see enough environments), then we find ALL parents



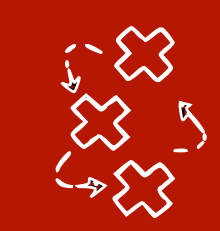
Invariant Causal Prediction - latent confounders

- If there are latent confounders, one can prove that:

$$\bigcap_{S \subseteq \{1, \dots, p\} \text{ s.t. } Y \perp\!\!\!\perp E | S} S \subseteq \text{Anc}(Y)$$



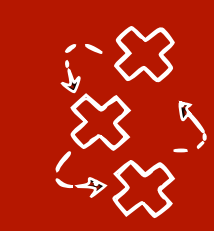
$$\begin{aligned} & Y \not\perp\!\!\!\perp E | X_2 \\ & Y \perp\!\!\!\perp E | X_1 \Rightarrow X_1 \in \text{Anc}(Y) \\ & Y \perp\!\!\!\perp E | X_1, X_2 \end{aligned}$$



Learning from multiple contexts

- Now we cannot decide which intervention to perform (**intervention design**)
 - We then also have **known intervention targets**, e.g. $\text{do}(S = 1)$
- Instead, somebody gives us a **set of data from multiple contexts**
 - Possibly **unknown intervention targets**
 - Possibly **soft interventions** instead of **perfect interventions**

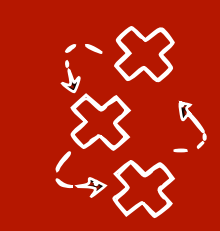
**ICP finds subsets of parents,
what about finding (an
equivalence class of) the
causal graph?**



Joint Causal Inference from Multiple Contexts

Joris M. Mooij, Sara Magliacane, Tom Claassen

- We represent different distributions (including interventional) as an **unknown joint causal graph** (*possibly cyclic or with latent confounders*)
- We **add context variables** so we can **disentangle** changes in distribution across the datasets

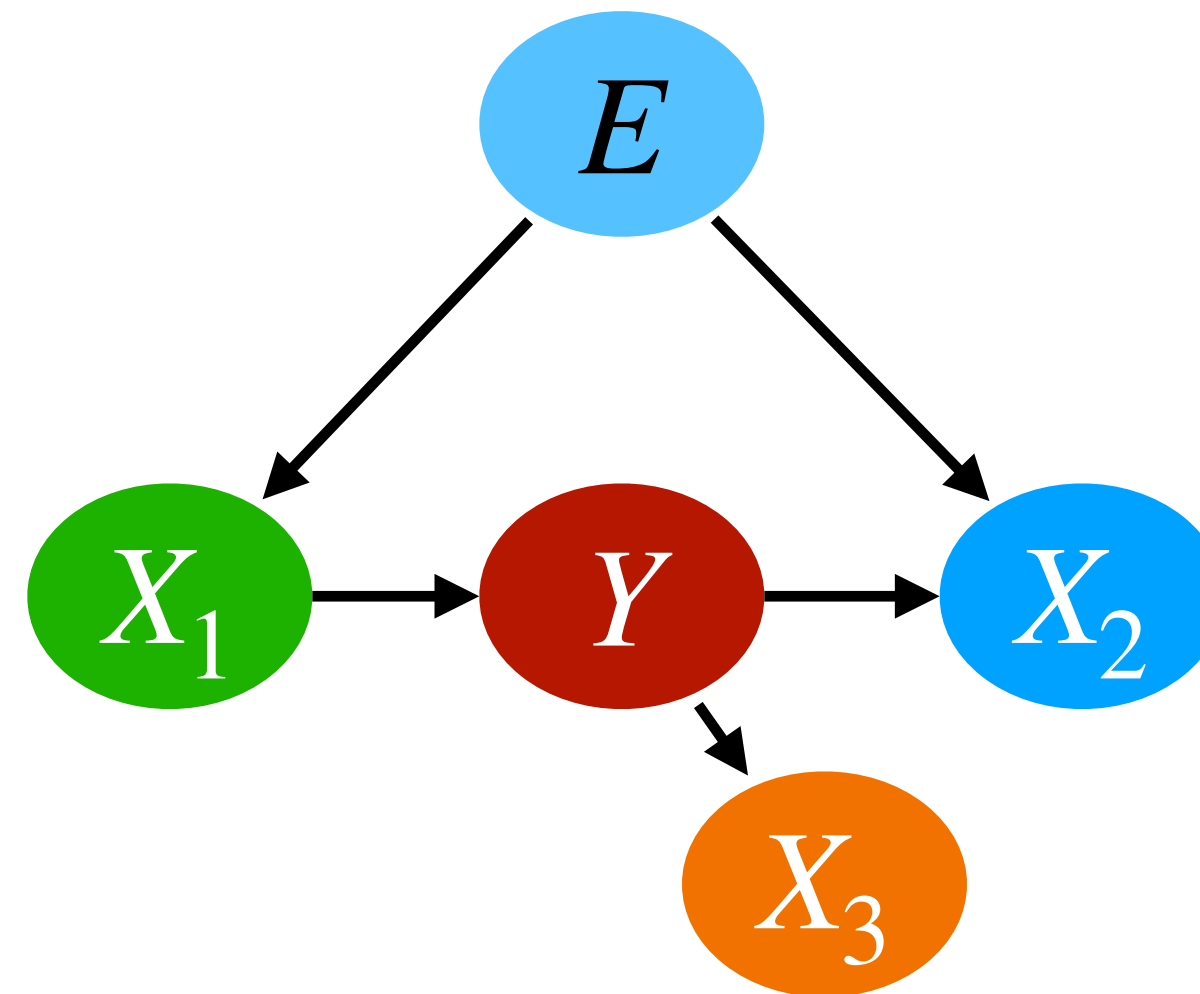


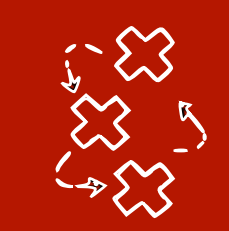
Joint Causal Inference intuition

$$\begin{cases} \epsilon_1, \epsilon_2, \epsilon_3, \epsilon_Y \sim \mathcal{N}(0,1) \\ X_1 = 10 + \epsilon_1 \\ Y = 3X_1 + \epsilon_Y \\ X_2 = -2Y + \epsilon_2 \\ X_3 = 2Y + 0.1\epsilon_3 \end{cases} \quad E = 0$$

$$\begin{cases} \epsilon_1, \epsilon_2, \epsilon_3, \epsilon_Y \sim \mathcal{N}(0,1) \\ X_1 = 100 + \epsilon_1 \\ Y = 3X_1 + \epsilon_Y \\ X_2 = -2Y + \epsilon_2 \\ X_3 = 2Y + 0.1\epsilon_3 \end{cases} \quad E = 1$$

$$\begin{cases} \epsilon_1, \epsilon_2, \epsilon_3, \epsilon_Y \sim \mathcal{N}(0,1) \\ X_1 = 10 + \epsilon_1 \\ Y = 3X_1 + \epsilon_Y \\ X_2 = 10Y + \epsilon_2 \\ X_3 = 2Y + 0.1\epsilon_3 \end{cases} \quad E = 2$$



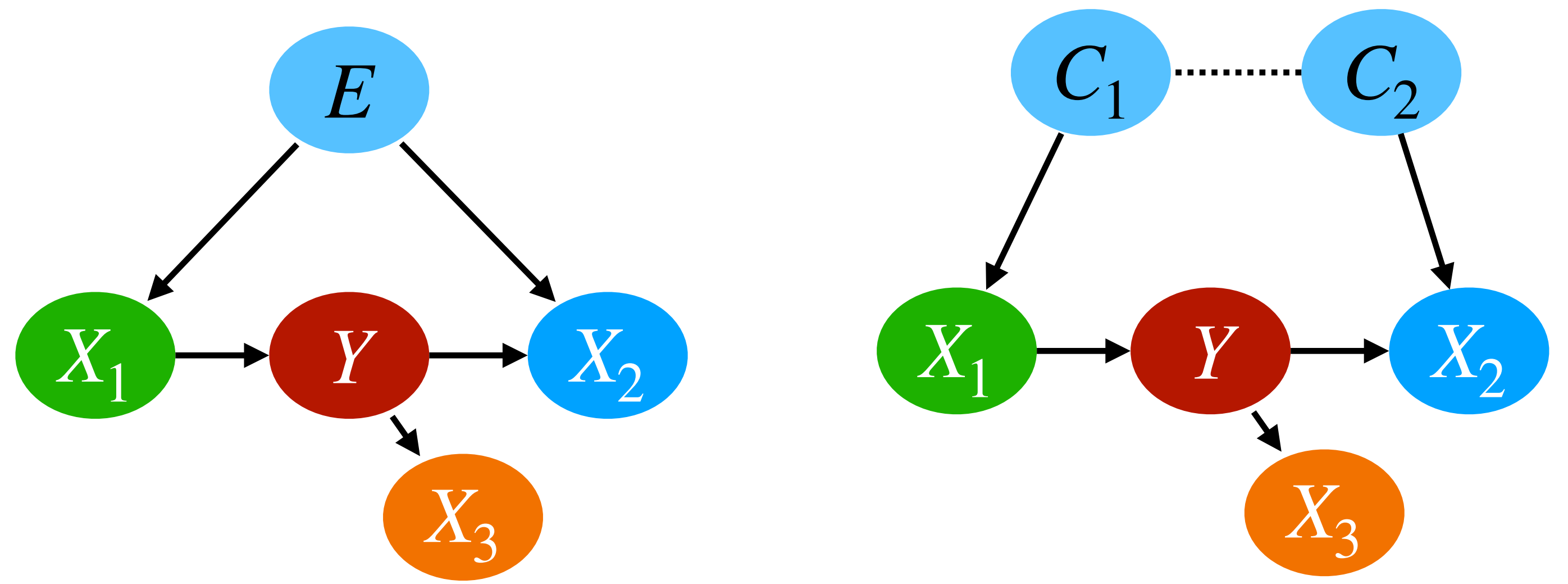


Joint Causal Inference intuition

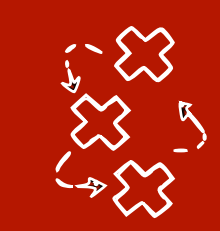
$$\begin{cases} \epsilon_1, \epsilon_2, \epsilon_3, \epsilon_Y \sim \mathcal{N}(0,1) \\ X_1 = 10 + \epsilon_1 \\ Y = 3X_1 + \epsilon_Y \\ X_2 = -2Y + \epsilon_2 \\ X_3 = 2Y + 0.1\epsilon_3 \end{cases} \quad \begin{matrix} E = 0 \\ C_1 = 0 \\ C_2 = 0 \end{matrix}$$

$$\begin{cases} \epsilon_1, \epsilon_2, \epsilon_3, \epsilon_Y \sim \mathcal{N}(0,1) \\ X_1 = 100 + \epsilon_1 \\ Y = 3X_1 + \epsilon_Y \\ X_2 = -2Y + \epsilon_2 \\ X_3 = 2Y + 0.1\epsilon_3 \end{cases} \quad \begin{matrix} E = 1 \\ C_1 = 1 \\ C_2 = 0 \end{matrix}$$

$$\begin{cases} \epsilon_1, \epsilon_2, \epsilon_3, \epsilon_Y \sim \mathcal{N}(0,1) \\ X_1 = 10 + \epsilon_1 \\ Y = 3X_1 + \epsilon_Y \\ X_2 = 10Y + \epsilon_2 \\ X_3 = 2Y + 0.1\epsilon_3 \end{cases} \quad \begin{matrix} E = 2 \\ C_1 = 0 \\ C_2 = 1 \end{matrix}$$



Adding context variables C_1 and C_2 helps disentangle the changes in each environment



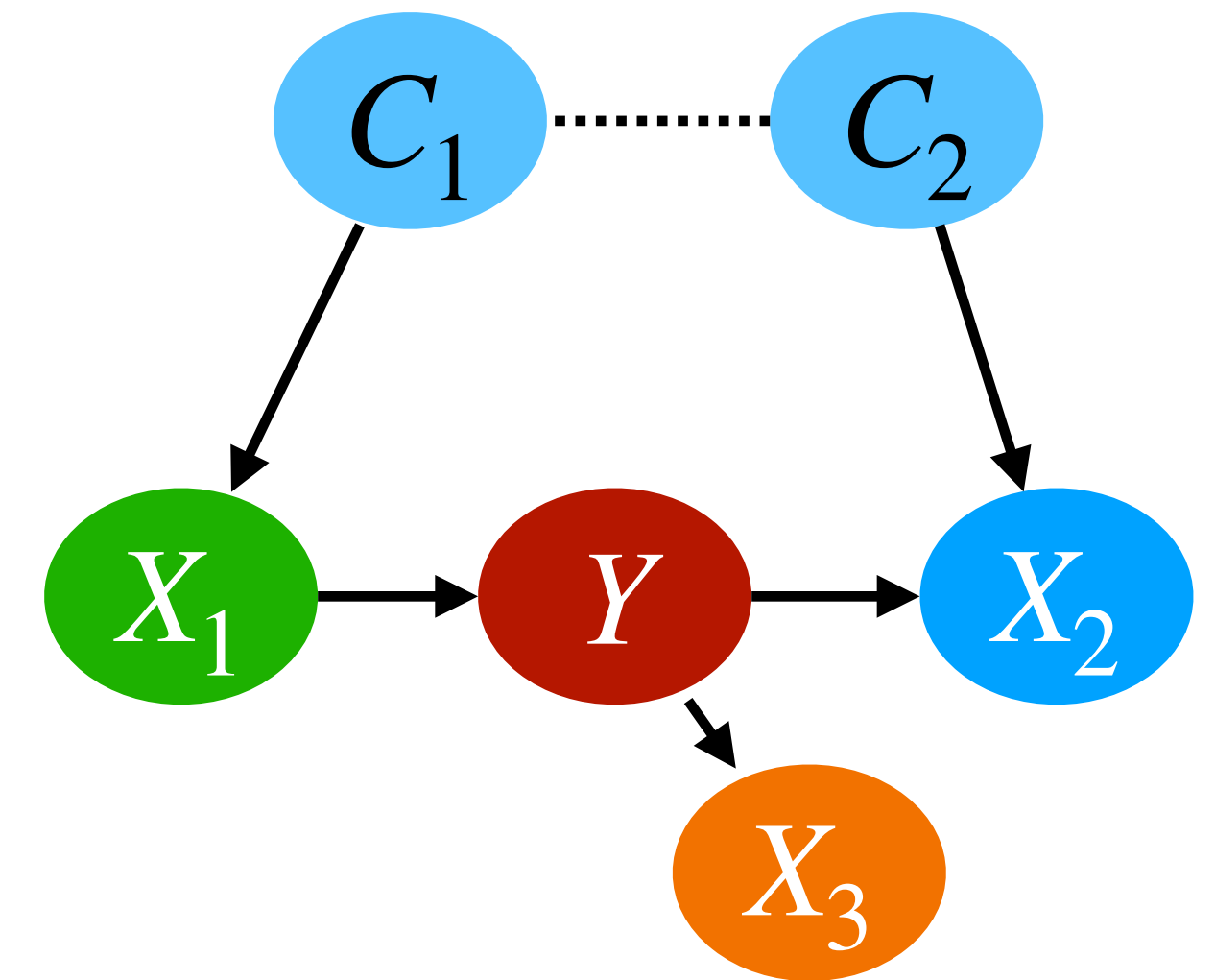
Joint Causal Inference intuition

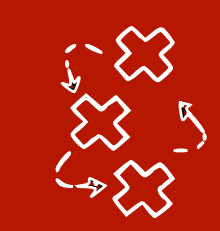
$$\begin{cases} \epsilon_1, \epsilon_2, \epsilon_3, \epsilon_Y \sim \mathcal{N}(0,1) \\ X_1 = 10 + \epsilon_1 \\ Y = 3X_1 + \epsilon_Y \\ X_2 = -2Y + \epsilon_2 \\ X_3 = 2Y + 0.1\epsilon_3 \end{cases} \quad \begin{matrix} E = 0 \\ C_1 = 0 \\ C_2 = 0 \end{matrix}$$

$$\begin{cases} \epsilon_1, \epsilon_2, \epsilon_3, \epsilon_Y \sim \mathcal{N}(0,1) \\ X_1 = 100 + \epsilon_1 \\ Y = 3X_1 + \epsilon_Y \\ X_2 = -2Y + \epsilon_2 \\ X_3 = 2Y + 0.1\epsilon_3 \end{cases} \quad \begin{matrix} E = 1 \\ C_1 = 1 \\ C_2 = 0 \end{matrix}$$

$$\begin{cases} \epsilon_1, \epsilon_2, \epsilon_3, \epsilon_Y \sim \mathcal{N}(0,1) \\ X_1 = 10 + \epsilon_1 \\ Y = 3X_1 + \epsilon_Y \\ X_2 = 10Y + \epsilon_2 \\ X_3 = 2Y + 0.1\epsilon_3 \end{cases} \quad \begin{matrix} E = 2 \\ C_1 = 0 \\ C_2 = 1 \end{matrix}$$

$$\begin{cases} \epsilon_1, \epsilon_2, \epsilon_3, \epsilon_Y \sim \mathcal{N}(0,1) \\ X_1 = \begin{cases} 10 + \epsilon_1 & \text{if } C_1 = 0 \\ 100 + \epsilon_1 & \text{if } C_1 = 1 \end{cases} \\ Y = 3X_1 + \epsilon_Y \\ X_2 = \begin{cases} -2Y + \epsilon_2 & \text{if } C_2 = 0 \\ 10Y + \epsilon_2 & \text{if } C_2 = 1 \end{cases} \\ X_3 = 2Y + 0.1\epsilon_3 \end{cases}$$



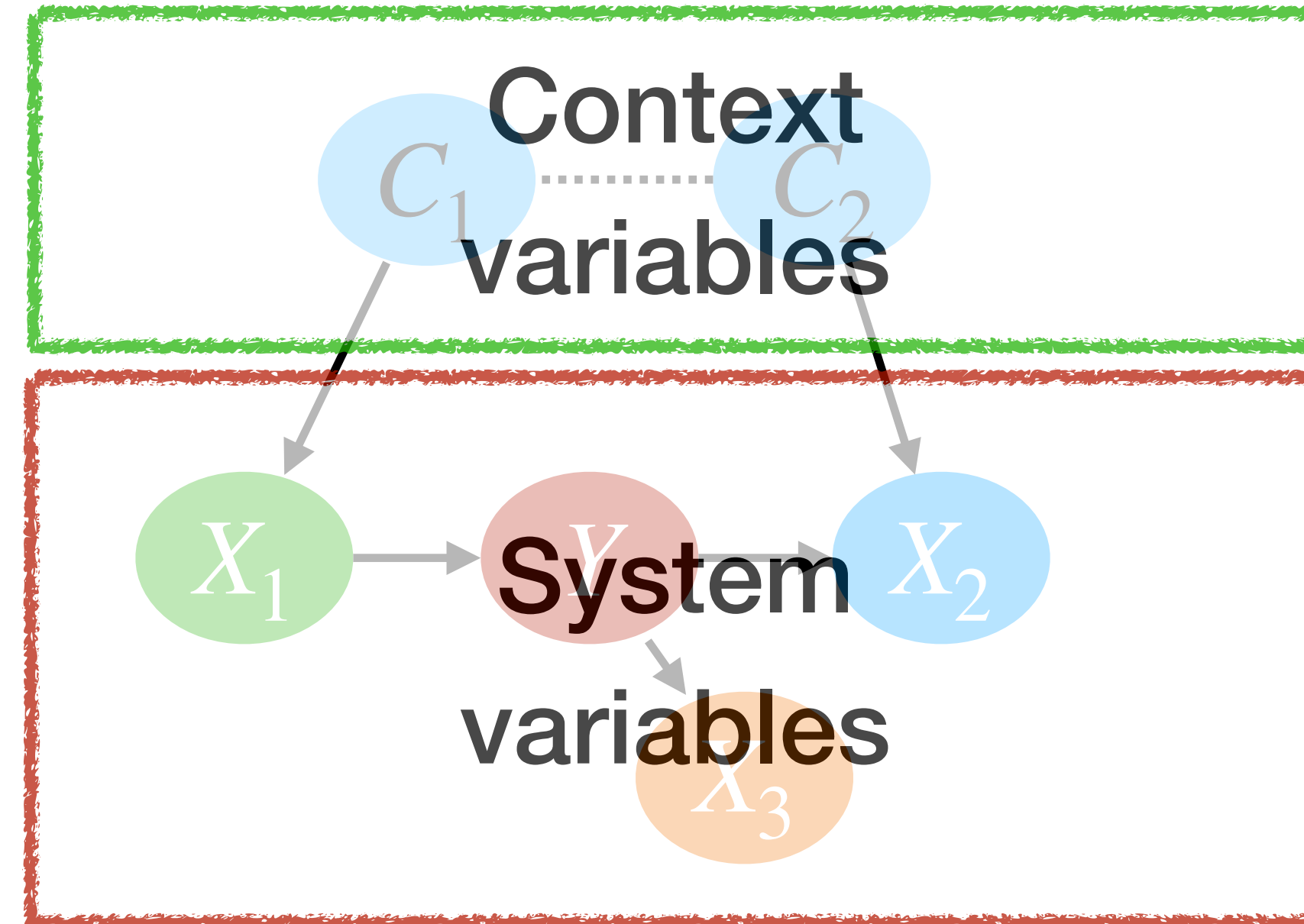


Joint Causal Inference intuition

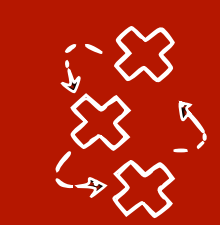
$$\begin{cases} \epsilon_1, \epsilon_2, \epsilon_3, \epsilon_Y \sim \mathcal{N}(0,1) \\ X_1 = 10 + \epsilon_1 \\ Y = 3X_1 + \epsilon_Y \\ X_2 = -2Y + \epsilon_2 \\ X_3 = 2Y + 0.1\epsilon_3 \end{cases} \quad \begin{matrix} E = 0 \\ C_1 = 10 \\ C_2 = -2 \end{matrix}$$

$$\begin{cases} \epsilon_1, \epsilon_2, \epsilon_3, \epsilon_Y \sim \mathcal{N}(0,1) \\ X_1 = 100 + \epsilon_1 \\ Y = 3X_1 + \epsilon_Y \\ X_2 = -2Y + \epsilon_2 \\ X_3 = 2Y + 0.1\epsilon_3 \end{cases} \quad \begin{matrix} E = 1 \\ C_1 = 100 \\ C_2 = -2 \end{matrix}$$

$$\begin{cases} \epsilon_1, \epsilon_2, \epsilon_3, \epsilon_Y \sim \mathcal{N}(0,1) \\ X_1 = 10 + \epsilon_1 \\ Y = 3X_1 + \epsilon_Y \\ X_2 = 10Y + \epsilon_2 \\ X_3 = 2Y + 0.1\epsilon_3 \end{cases} \quad \begin{matrix} E = 2 \\ C_1 = 10 \\ C_2 = 10 \end{matrix}$$



The context variables C_1 and C_2 can be also descriptive of the intervention in each environment

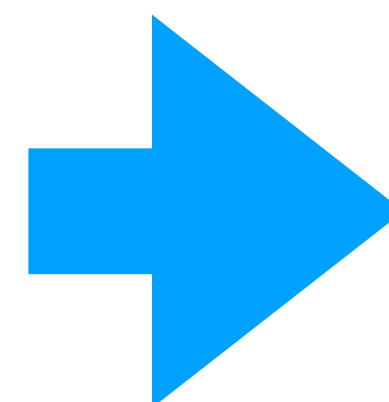


Joint Causal Inference from Multiple Contexts

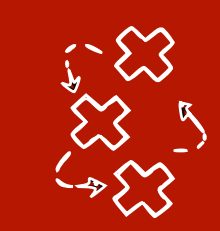
Joris M. Mooij, Sara Magliacane, Tom Claassen

- We represent different distributions (including interventional) as an **unknown joint causal graph** (possibly cyclic or with latent confounders)
- We **add context variables** so we can **disentangle** changes in distribution across the datasets

	X1	X2	X3
Normal	0,1	2	0
Normal	0,2	3	0
	X1	X2	X3
Gene A	3,1	2	1
Gene A	3,2	3	1
	X1	X2	X3
Gene B	0.2	1	0
Gene B	0.3	1	1
Gene B	0.3	2	1
Gene B	0.4	1	1



C1	C2	X1	X2	X3
0	0	0,1	2	0
0	0	0,2	3	0
0	0	1,1	2	1
0	0	0,1	3	0
1	0	3,1	2	1
1	0	3,2	3	1
1	0	4	1	1
1	0	3,2	3	1
0	1	0,2	1	0
0	1	0,3	1	1
0	1	0,3	2	1
0	1	0,4	1	1

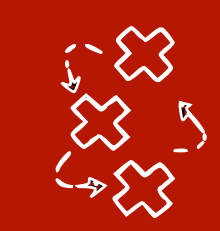


Joint Causal Inference from Multiple Contexts

Joris M. Mooij, Sara Magliacane, Tom Claassen

- We **add context variables** so we can **disentangle** changes in distribution across the datasets (and optionally background knowledge, e.g. context variables are uncaused)
- We can reuse **any standard method for observational data** that fits any chosen assumptions

C1	C2	X1	X2	X3
0	0	0,1	2	0
0	0	0,2	3	0
0	0	1,1	2	1
0	0	0,1	3	0
1	0	3,1	2	1
1	0	3,2	3	1
1	0	4	1	1
1	0	3,2	3	1
0	1	0,2	1	0
0	1	0,3	1	1
0	1	0,3	2	1
0	1	0,4	1	1

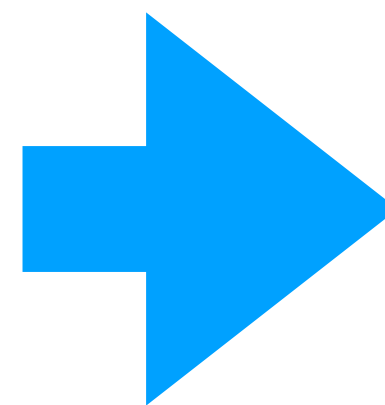


Joint Causal Inference from Multiple Contexts

Joris M. Mooij, Sara Magliacane, Tom Claassen

- We **add context variables** so we can **disentangle** changes in distribution across the datasets (and optionally background knowledge, e.g. context variables are uncaused)
- We can reuse **any standard method for observational data** that fits any chosen assumptions

C1	C2	X1	X2	X3
0	0	0,1	2	0
0	0	0,2	3	0
0	0	1,1	2	1
0	0	0,1	3	0
1	0	3,1	2	1
1	0	3,2	3	1
1	0	4	1	1
1	0	3,2	3	1
0	1	0,2	1	0
0	1	0,3	1	1
0	1	0,3	2	1
0	1	0,4	1	1

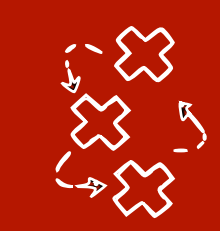


$$X_2 \perp\!\!\!\perp C_2$$

$$X_1 \perp\!\!\!\perp C_2 \mid C_1$$

$$X_2 \perp\!\!\!\perp C_1 \mid X_3$$

...

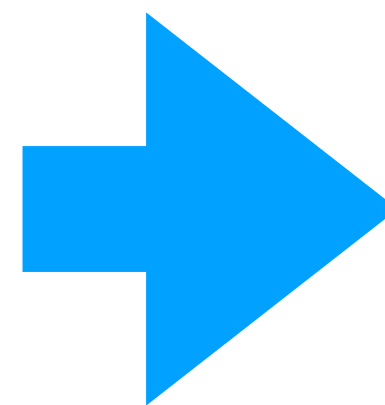


Joint Causal Inference from Multiple Contexts

Joris M. Mooij, Sara Magliacane, Tom Claassen

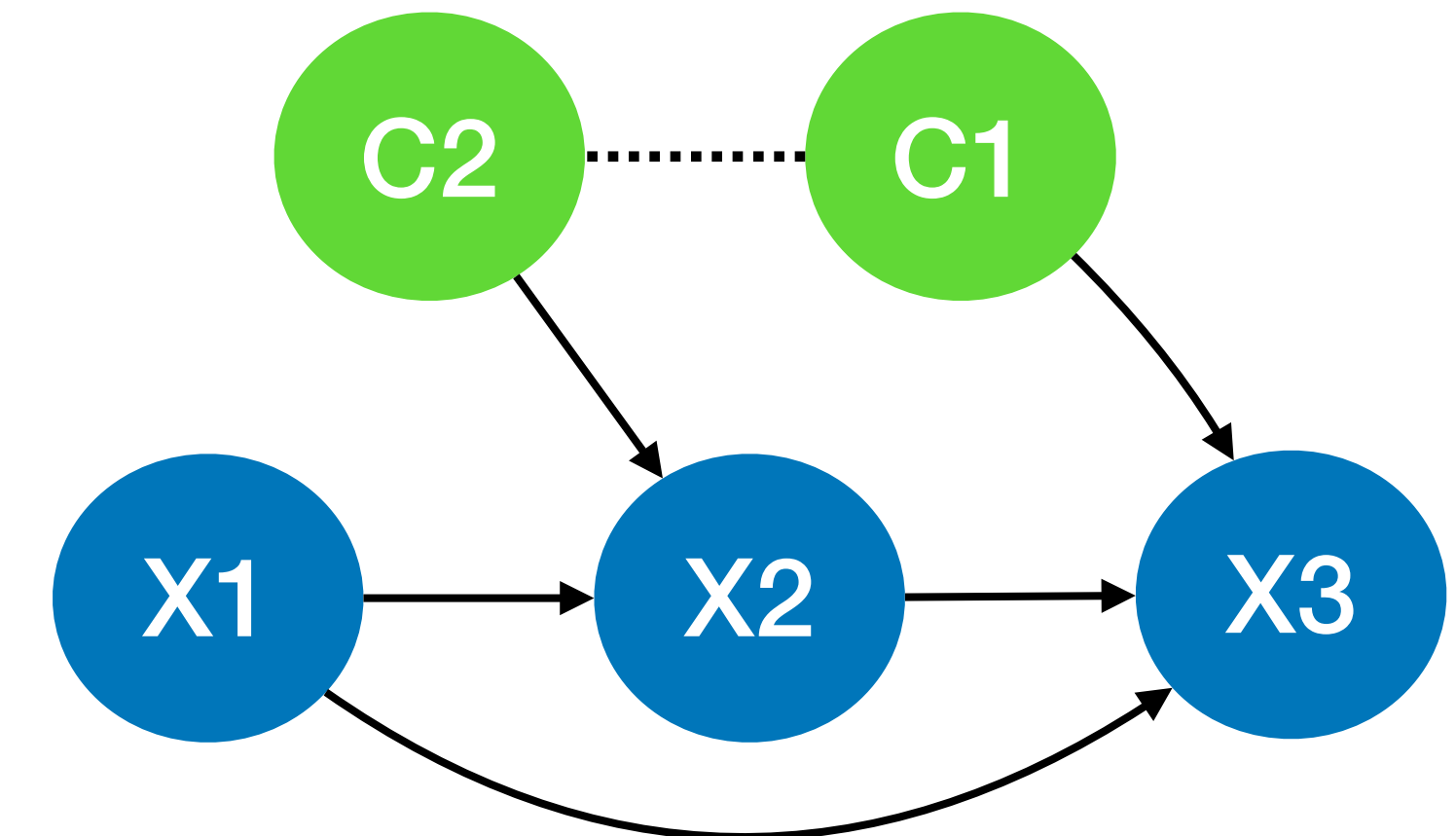
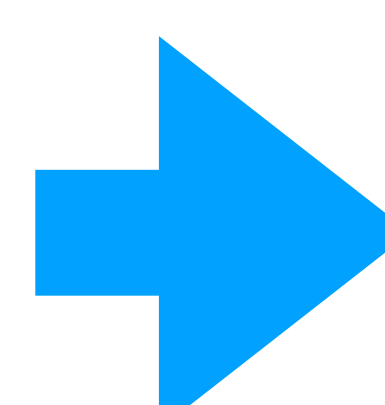
- We **add context variables** so we can **disentangle** changes in distribution across the datasets (and optionally background knowledge, e.g. context variables are uncaused)
- We can reuse **any standard method for observational data** that fits any chosen assumptions

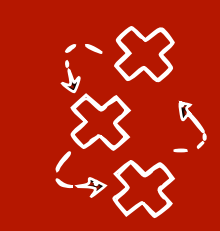
C1	C2	X1	X2	X3
0	0	0,1	2	0
0	0	0,2	3	0
0	0	1,1	2	1
0	0	0,1	3	0
1	0	3,1	2	1
1	0	3,2	3	1
1	0	4	1	1
1	0	3,2	3	1
0	1	0,2	1	0
0	1	0,3	1	1
0	1	0,3	2	1
0	1	0,4	1	1



$X_2 \perp\!\!\!\perp C_2$
 $X_1 \perp\!\!\!\perp C_2 \mid C_1$
 $X_2 \perp\!\!\!\perp C_1 \mid X_3$
 ...

PC-JCI
FCI-JCI



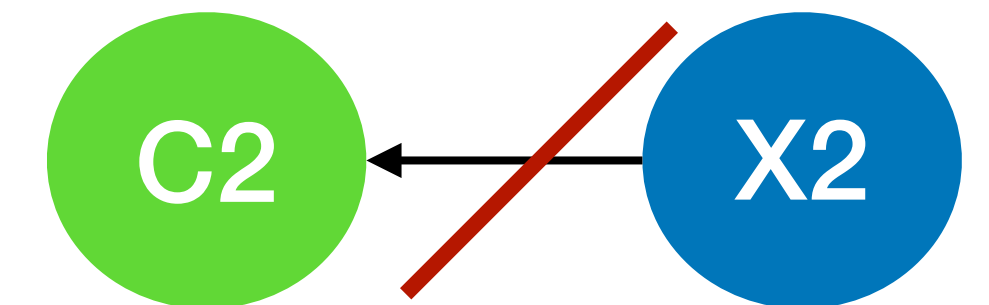


Joint Causal Inference from Multiple Contexts

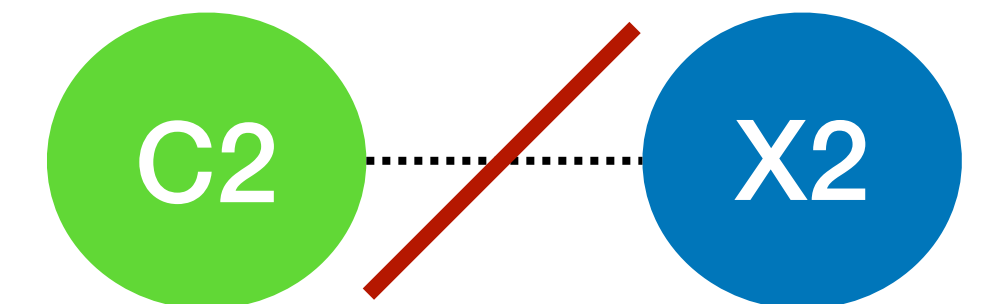
Joris M. Mooij, Sara Magliacane, Tom Claassen

- Additional optional background knowledge based on assumptions:

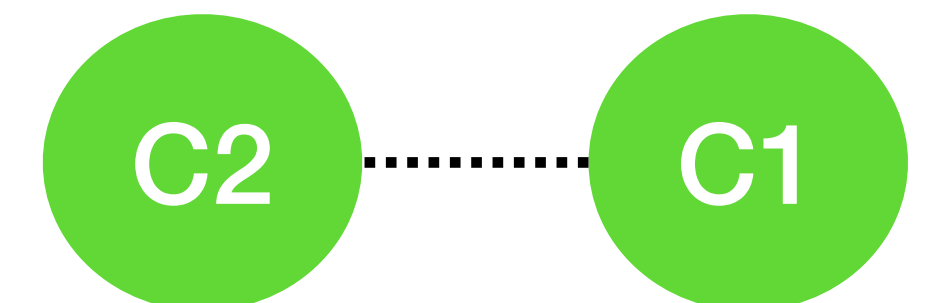
1. No system variable causes any context variable.



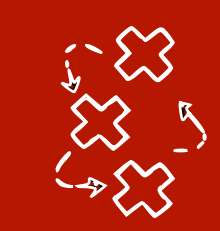
2. No context variable is confounded with a system variable.



3. The context variables do not cause each other and they are assumed to be confounded.



<https://arxiv.org/abs/1611.10351>



Joint Causal Inference from Multiple Contexts

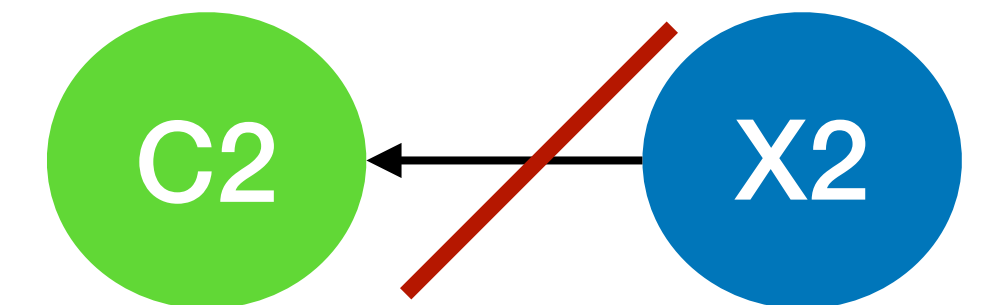
Joris M. Mooij, Sara Magliacane, Tom Claassen

In this talk we assume that there are no latent confounders (except some dependence between the context variables)

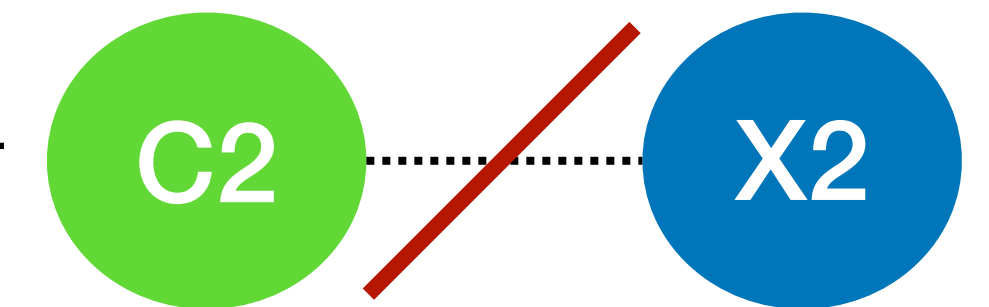
• Additional

assumptions:

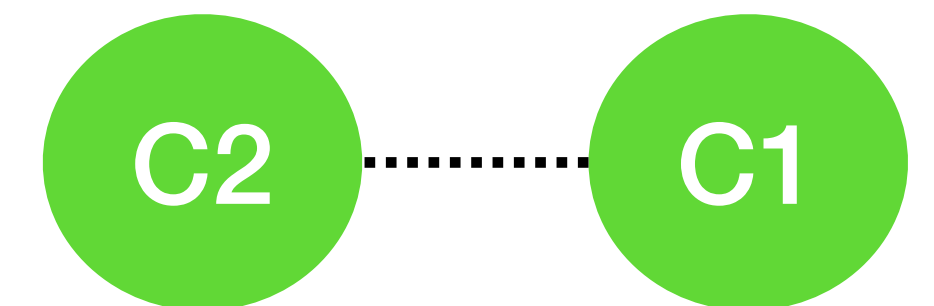
1. No system variable causes any context variable.



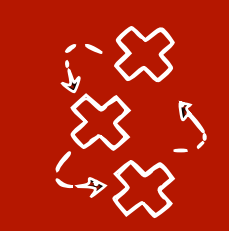
2. ~~No context variable is confounded with a system variable.~~



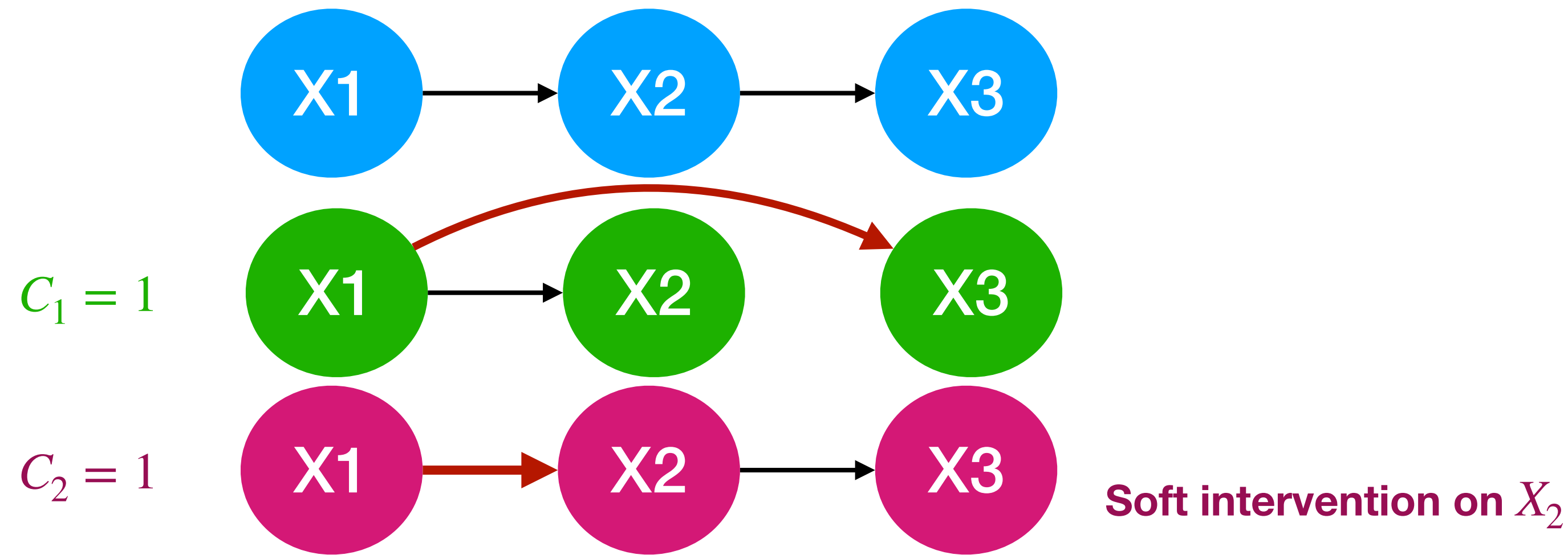
3. The context variables do not cause each other and they are assumed to be confounded.



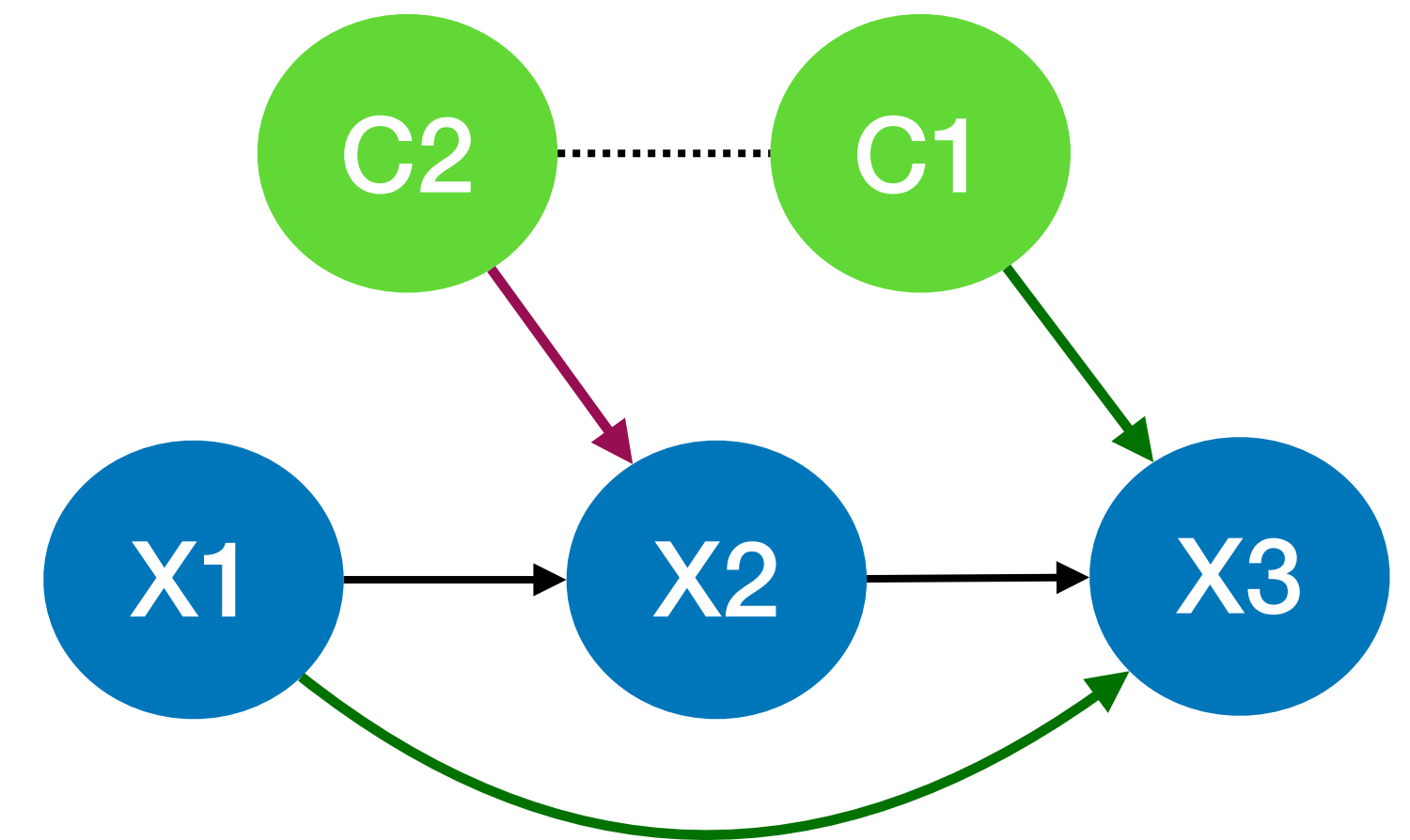
<https://arxiv.org/abs/1611.10351>



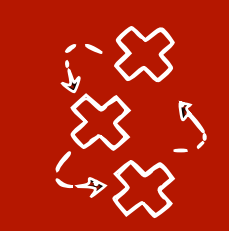
Joint Causal Inference example - setting



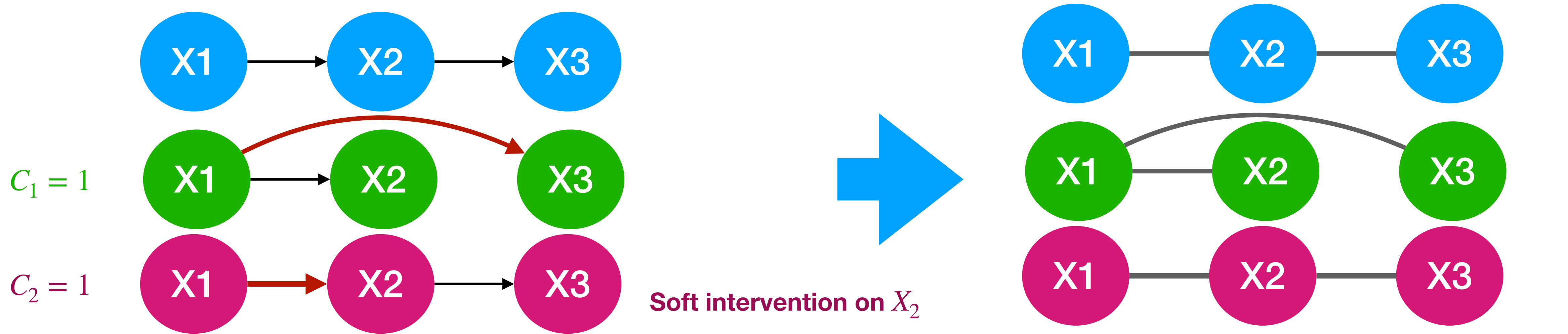
Single graphs in each environment



The joint graph is the union of the single graphs + edges from context variables for the causal mechanisms that change

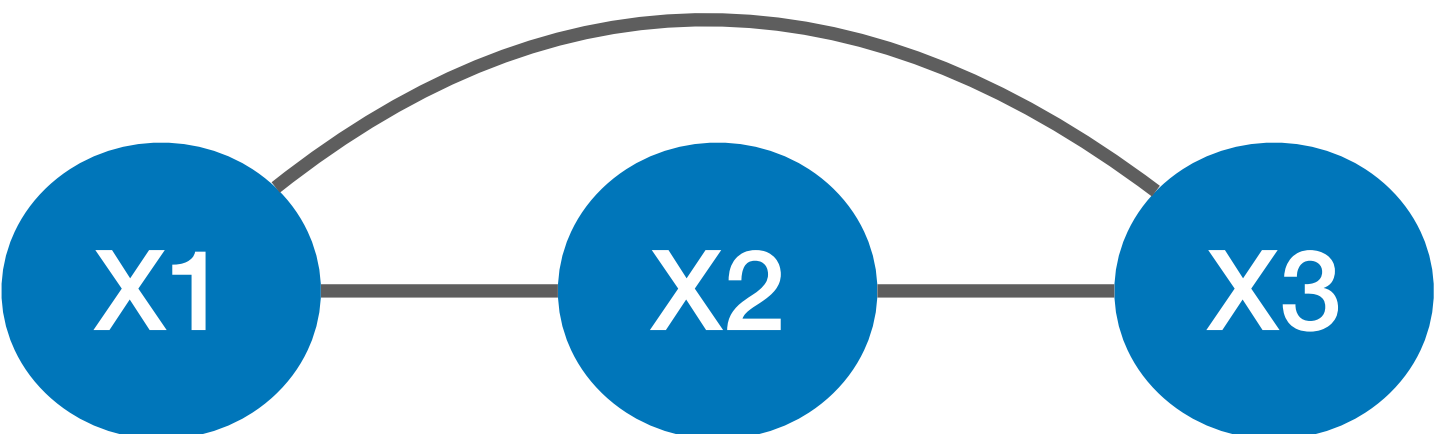


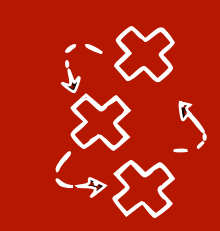
Learning graphs separately in each environment with PC



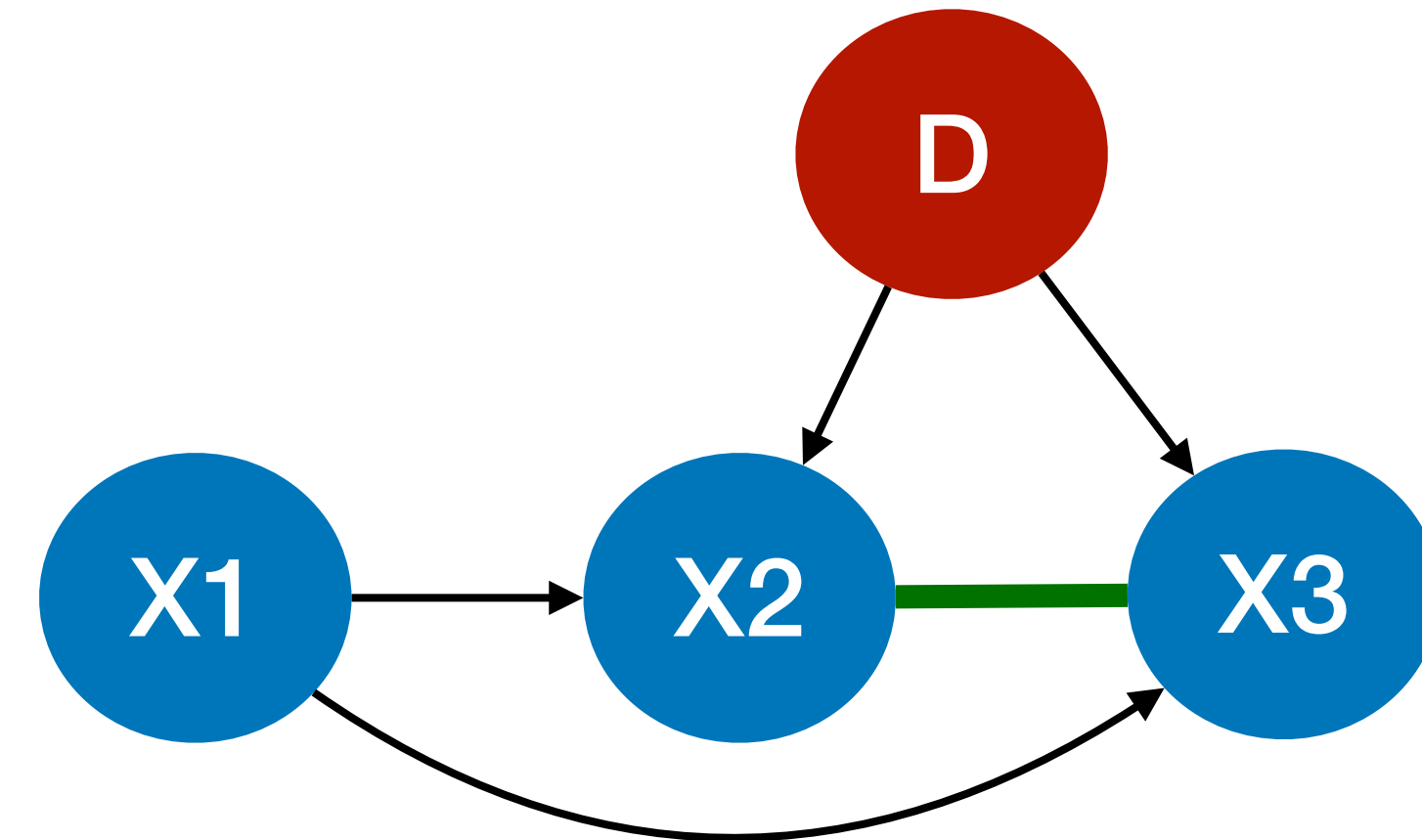
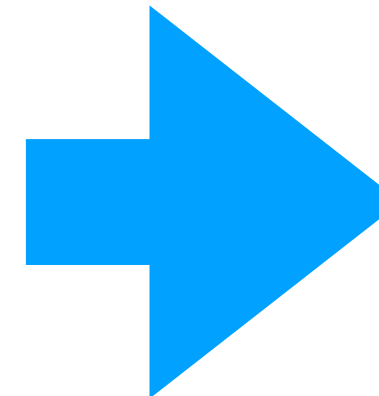
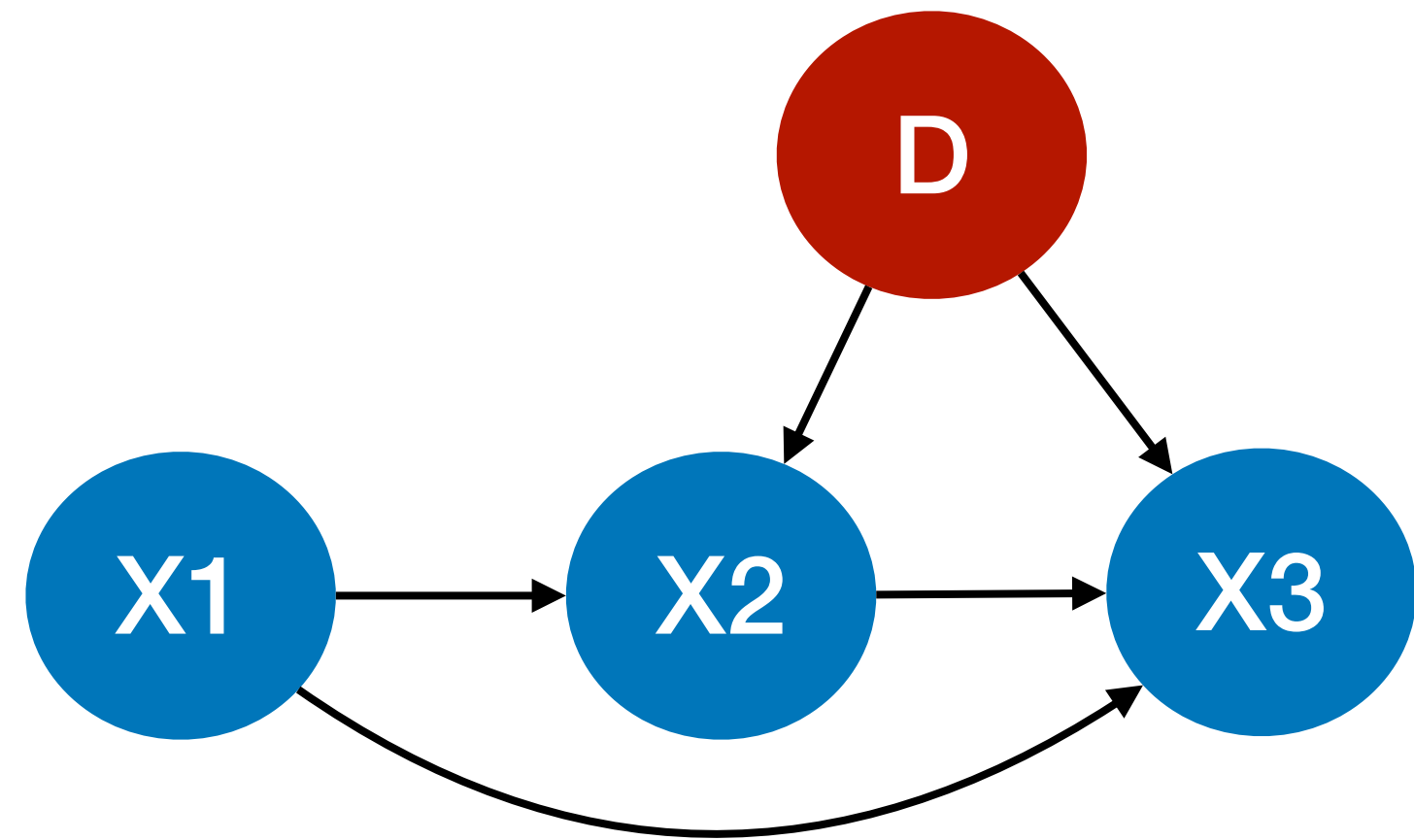
Single graphs in each environment

CPDAGs from each environment



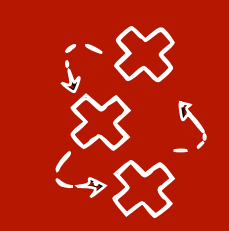


Joint causal inference + PC with a single domain variable

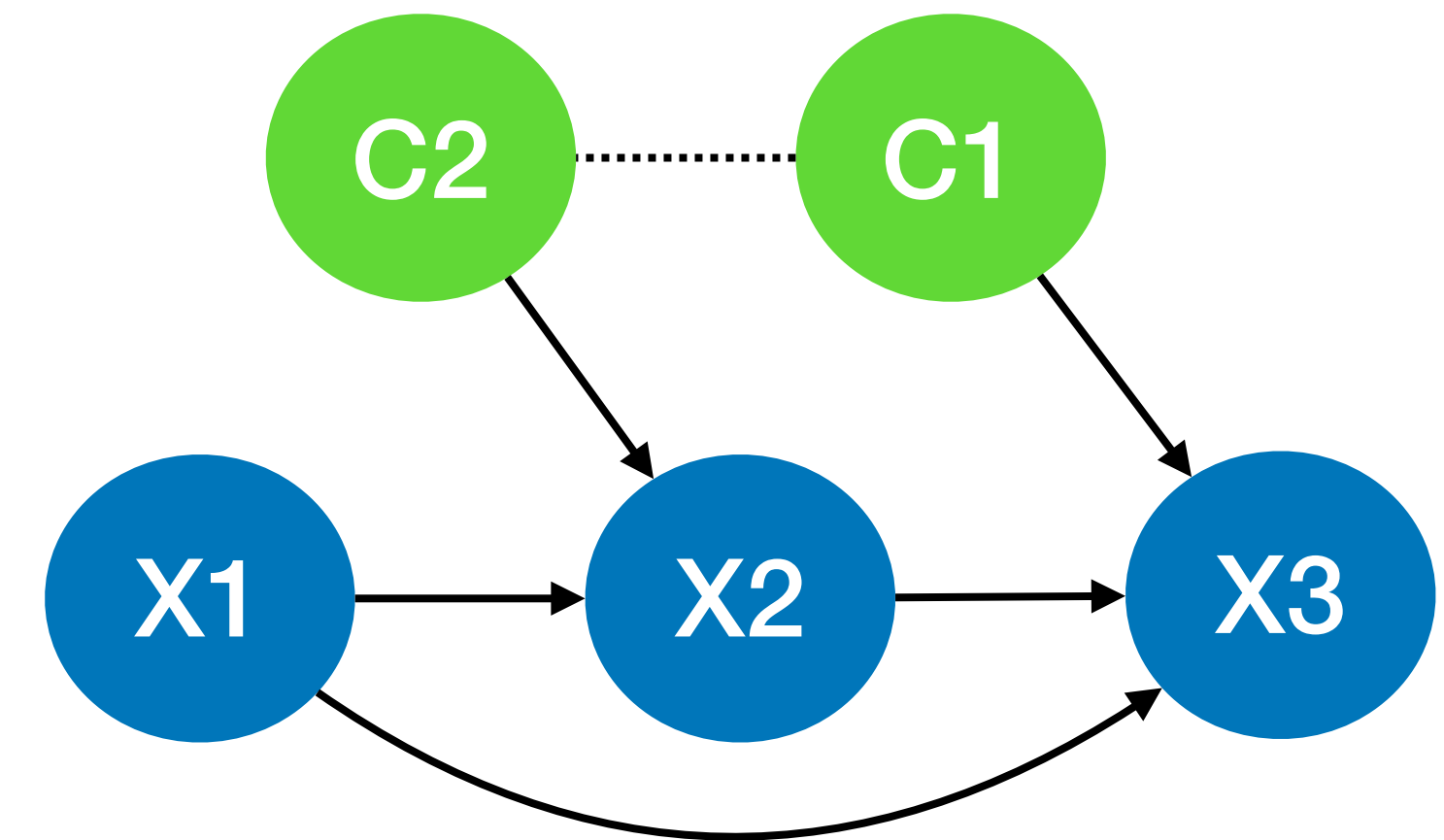
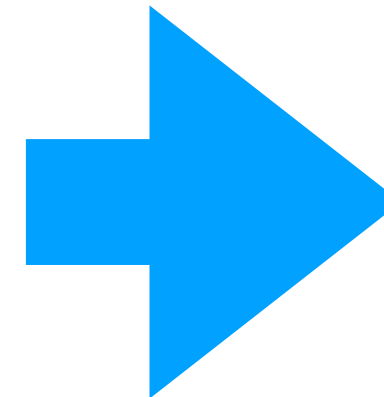
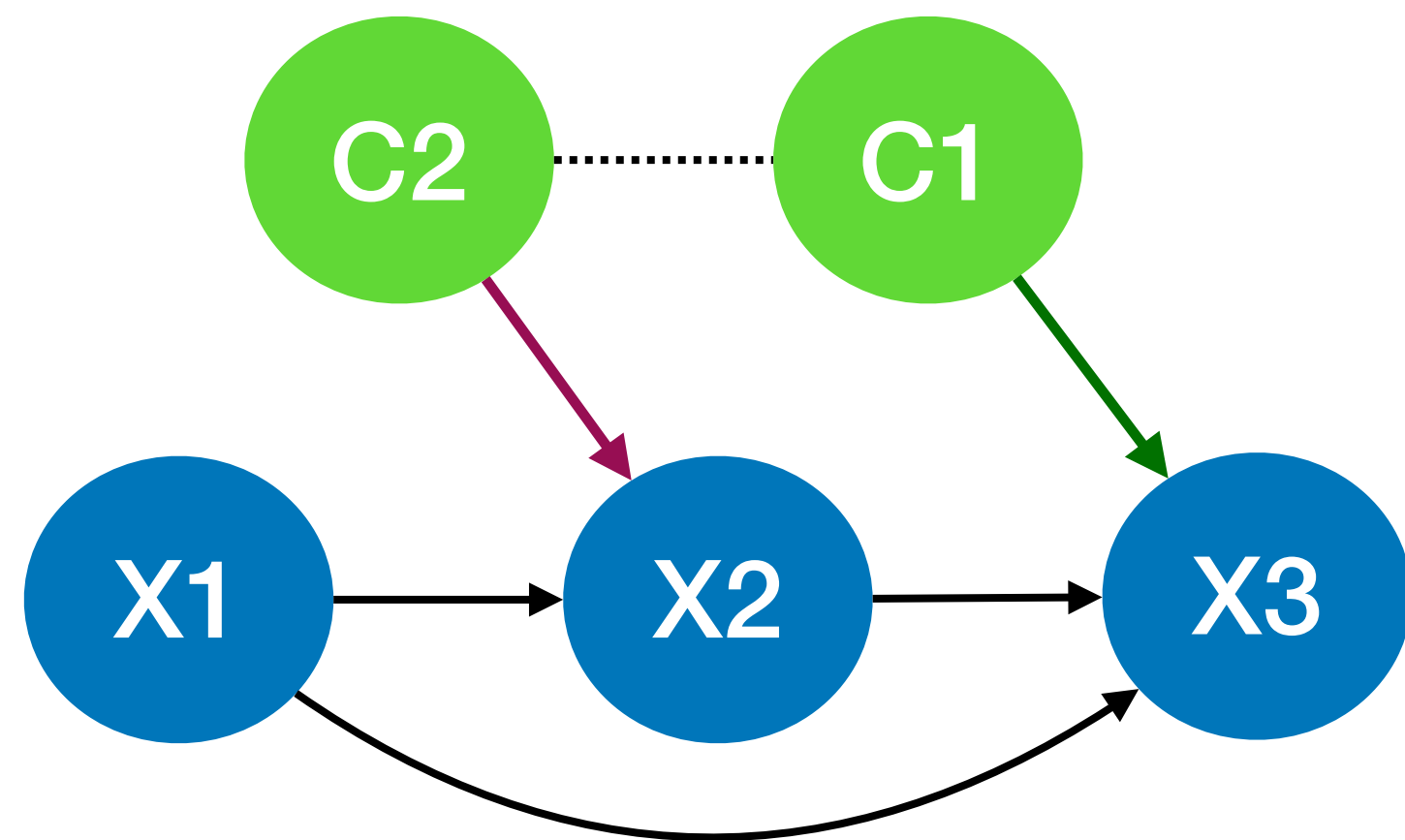


True underlying joint graph with domain variable

CPDAG with Joint Causal Inference and PC with a single domain variable

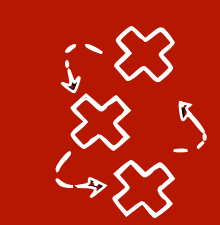


Joint causal inference + PC with multiple context variables



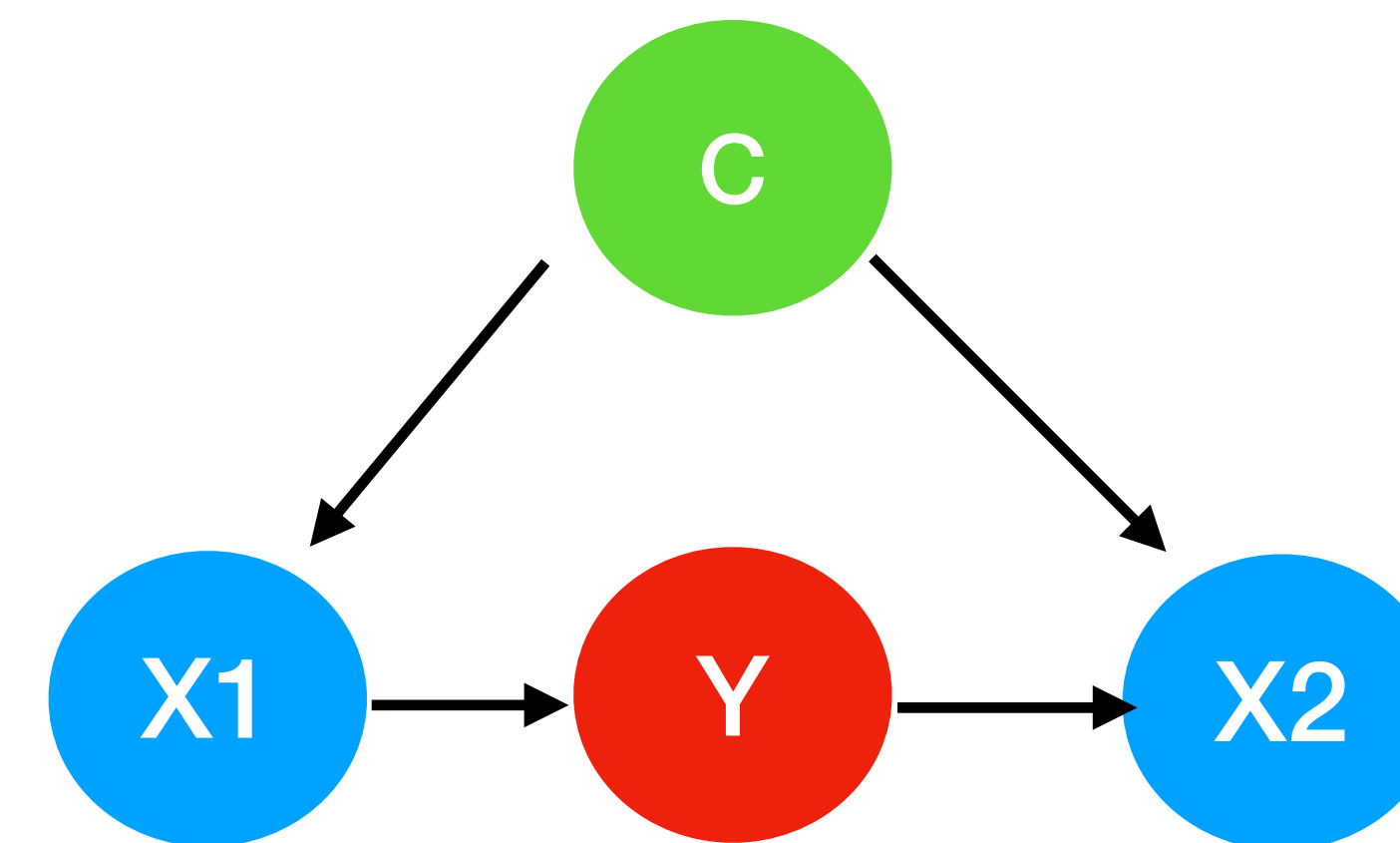
True underlying joint graph with multiple context variable

CPDAG with Joint Causal Inference and PC with multiple context variables



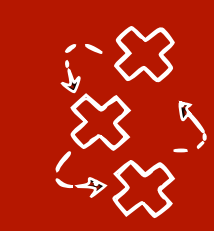
Application - Domain adaptation

	C	X1	X2	Y
Source domain	0	0,1	2	0
	0	0,2	3	0
	0	1,1	2	1
	0	0,1	3	0
Target domain	1	3,1	2	?
	1	3,2	3	?
	1	4	2	?
	1	3,2	3	?



We can represent $P(X1, X2, Y, C)$ with an **(unknown)** causal graph

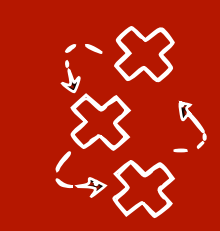
- We can represent the different datasets jointly with Joint Causal Inference
- We can use it to reason about **features that offer a robust prediction of Y**



This class

- Introduction to causal discovery
 - Common assumptions: causal sufficiency, acyclicity, faithfulness
 - Constraint-based causal discovery on observational data (causal sufficiency)
 - SGS, PC
-
- Learning from multiple contexts or interventional data
 - Invariant Causal Prediction
 - Joint Causal Inference

Inspired by <https://stat.ethz.ch/lectures/ss21/causality.php>



Questions?

