

# Causality

Evgeny Burnaev

Rodrigo Rivera

Skoltech

RESEARCH

CHRISTMAS 2011: DEATH'S DOMINION

**How fast does the Grim Reaper walk? Receiver operating characteristics curve analysis in healthy men aged 70 and over**

 OPEN ACCESS

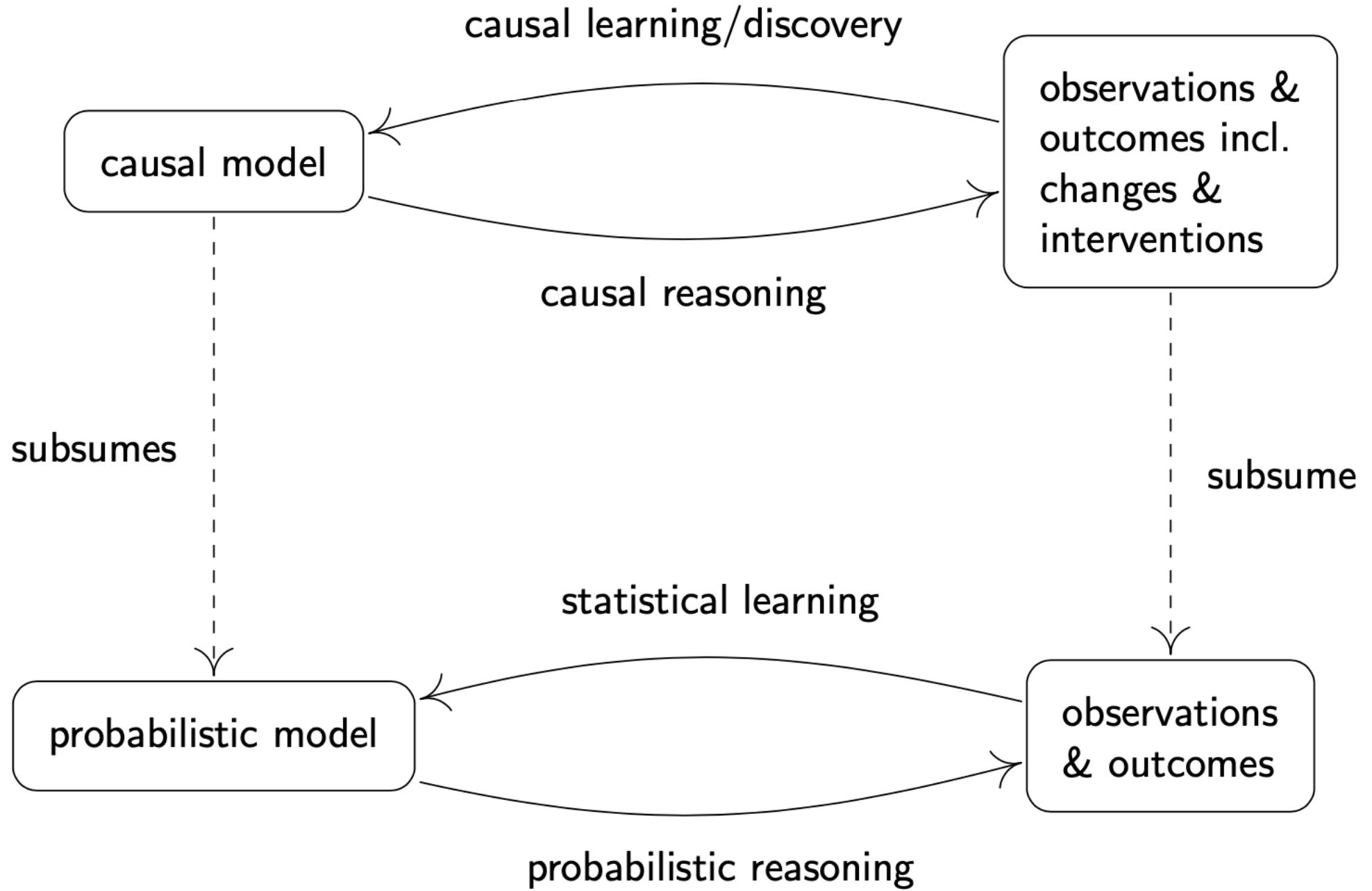
Fiona F Stanaway *research fellow*<sup>1</sup>, Danijela Gnjidic *research fellow*<sup>2,3,4</sup>, Fiona M Blyth *deputy*

answer: 0.82m/s

(picture by Belle Mellor)



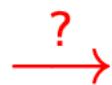
## **Part II: Causal Discovery**



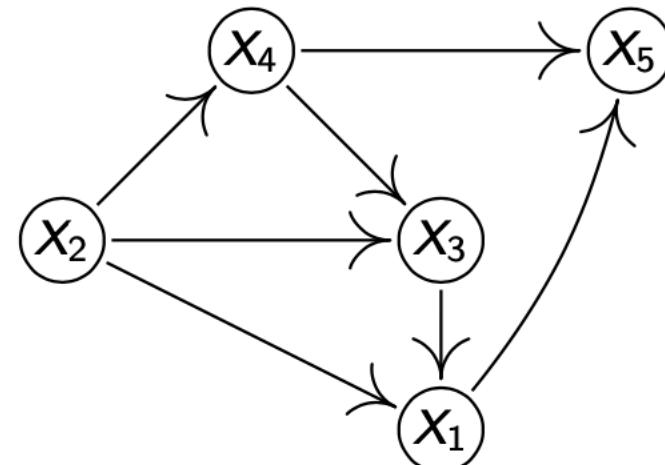
## The Problem of Causal Discovery:

observed iid data  
from  $P(X_1, \dots, X_5)$

$X_1$	$X_2$	$X_3$	$X_4$	$X_5$
3.4	-0.3	5.8	-2.1	2.2
1.7	-0.2	7.0	-1.2	0.4
-2.4	-0.1	4.3	-0.7	3.5
2.3	-0.3	5.5	-1.1	-4.4
3.5	-0.2	3.9	-0.9	-3.9
:	:	:	:	:



causal model, e.g. DAG  $\mathcal{G}$



I USED TO THINK  
CORRELATION IMPLIED  
CAUSATION.



THEN I TOOK A  
STATISTICS CLASS.  
NOW I DON'T.



SOUNDS LIKE THE  
CLASS HELPED.

WELL, MAYBE.



Correlation (Dependence) does not imply causation

Correlation (Dependence) does not imply causation ... but:

Correlation (Dependence) does not imply causation ... but:

**Reichenbach's common cause principle.**

Assume that  $X \perp\!\!\!\perp Y$ . Then

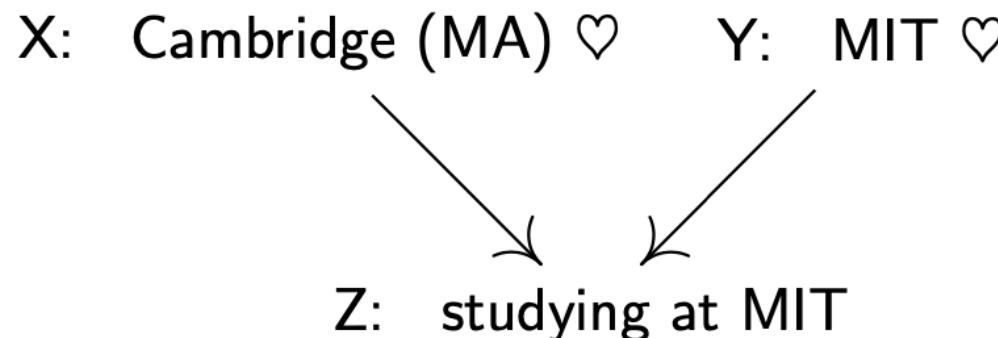
- $X$  “causes”  $Y$ ,
- $Y$  “causes”  $X$ ,
- there is a hidden common “cause” or
- combination of the above.

Correlation (Dependence) does not imply causation ... but:

### **Reichenbach's common cause principle.**

Assume that  $X \perp\!\!\!\perp Y$ . Then

- $X$  “causes”  $Y$ ,
- $Y$  “causes”  $X$ ,
- there is a hidden common “cause” or
- combination of the above.
- (In practice implicit conditioning also happens:



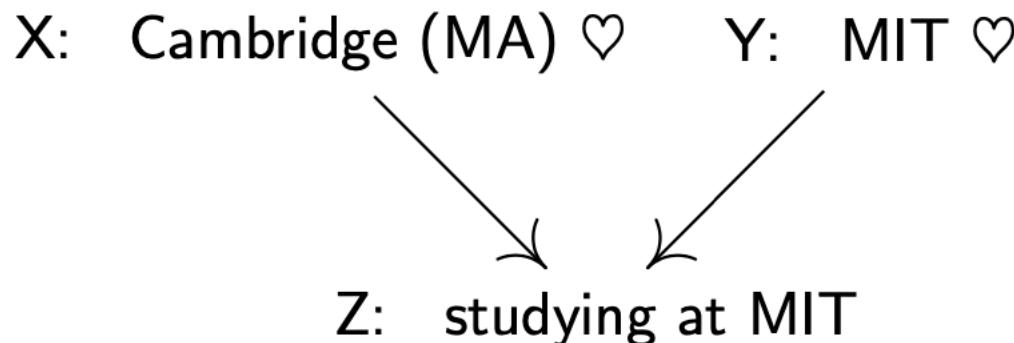
aka “selection bias”).

Correlation (Dependence) does not imply causation ... but:

**Reichenbach's common cause principle.**

Assume that  $X \not\perp\!\!\!\perp Y$ . Then

- $X$  “causes”  $Y$ ,
- $Y$  “causes”  $X$ ,
- there is a hidden common “cause” or
- combination of the above.
- (In practice implicit conditioning also happens:

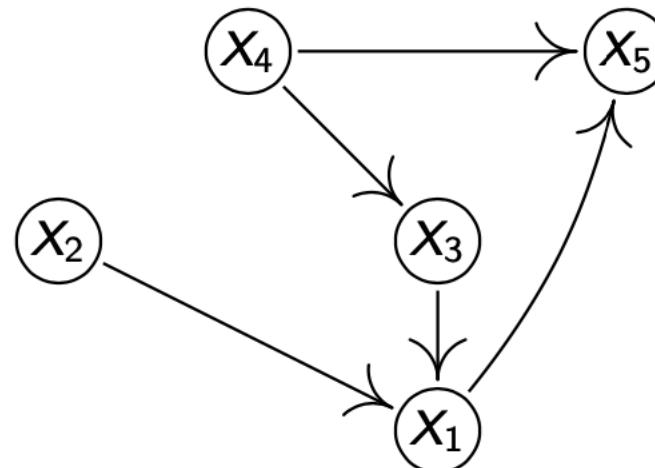


aka “selection bias”). Formalization of this idea...

# Definition: graphs

$G = (V, E)$  with  $E \subseteq V \times V$ . The rest is as in real life!

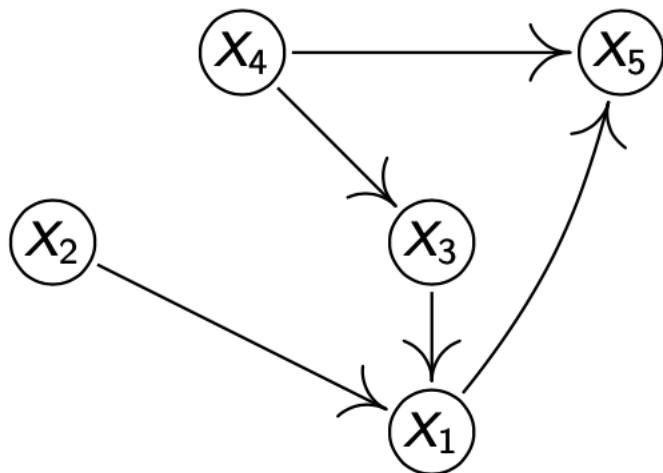
- parents, children, descendants, ancestors, ...
- paths, directed paths
- immoralities (or v-structures)
- $d$ -separation (see next)
- ...



# Definition: d-separation

$X_i$  and  $X_j$  are  $d$ -separated by  $S$  if all paths between  $X_i$  and  $X_j$  are blocked by  $S$ .

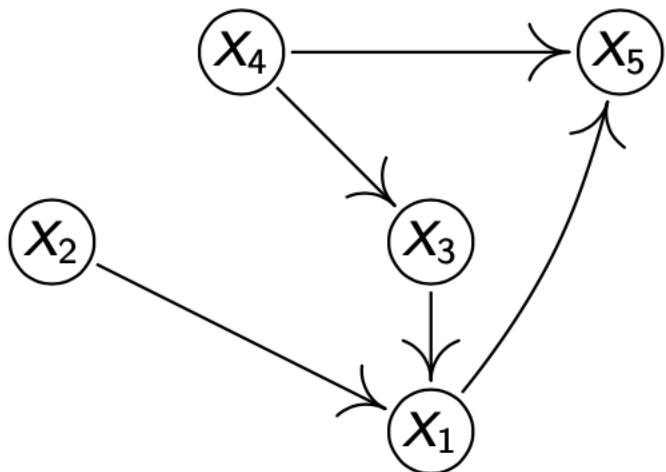
Check, whether all paths blocked!!



# Definition: d-separation

$X_i$  and  $X_j$  are  $d$ -separated by  $\mathcal{S}$  if all paths between  $X_i$  and  $X_j$  are blocked by  $\mathcal{S}$ .

Check, whether all paths blocked!!

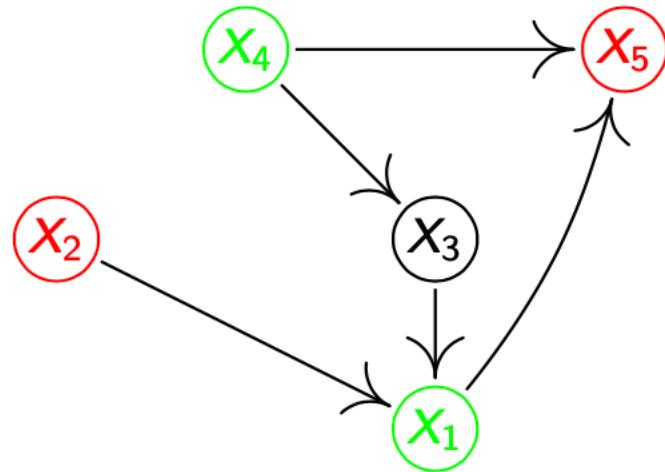


- $\dots \rightarrow \circ \rightarrow \dots$  ○ blocks a path.
- $\dots \leftarrow \circ \rightarrow \dots$  ○ blocks a path.
- $\dots \rightarrow \circ \leftarrow \dots$  ○ blocks a path.

# Definition: d-separation

$X_i$  and  $X_j$  are  $d$ -separated by  $\mathcal{S}$  if all paths between  $X_i$  and  $X_j$  are blocked by  $\mathcal{S}$ .

Check, whether all paths blocked!!



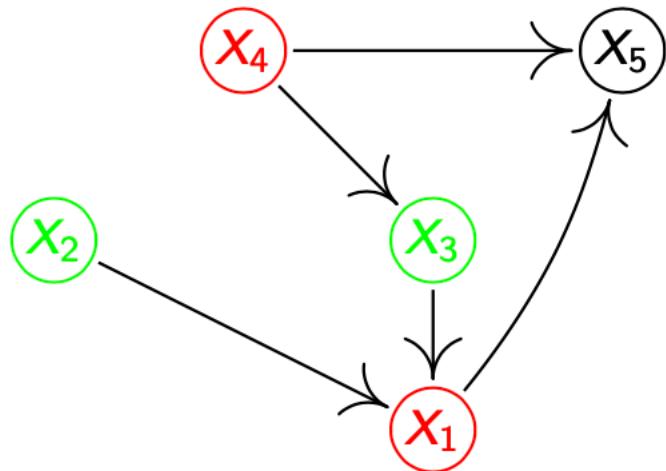
- $\dots \rightarrow \circ \rightarrow \dots$  ○ blocks a path.
- $\dots \leftarrow \circ \rightarrow \dots$  ○ blocks a path.
- $\dots \rightarrow \circ \leftarrow \dots$  ○ blocks a path.

$X_2$  and  $X_5$  are  $d$ -sep. by  $\{X_1, X_4\}$

# Definition: d-separation

$X_i$  and  $X_j$  are  $d$ -separated by  $\mathcal{S}$  if all paths between  $X_i$  and  $X_j$  are blocked by  $\mathcal{S}$ .

Check, whether all paths blocked!!



- ... → ○ → ... ○ blocks a path.
- ... ← ○ → ... ○ blocks a path.
- ... → ○ ← ... ○ blocks a path.

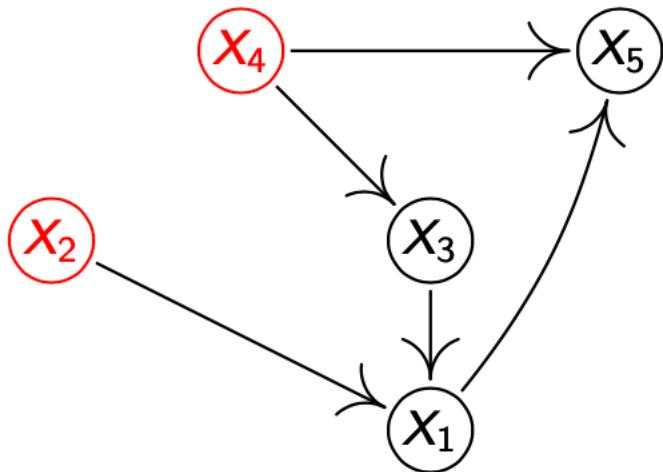
$X_2$  and  $X_5$  are  $d$ -sep. by  $\{X_1, X_4\}$

$X_4$  and  $X_1$  are  $d$ -sep. by  $\{X_2, X_3\}$

# Definition: d-separation

$X_i$  and  $X_j$  are  $d$ -separated by  $S$  if all paths between  $X_i$  and  $X_j$  are blocked by  $S$ .

Check, whether all paths blocked!!



- ... → ○ → ... ○ blocks a path.
- ... ← ○ → ... ○ blocks a path.
- ... → ○ ← ... ○ blocks a path.

$X_2$  and  $X_5$  are  $d$ -sep. by  $\{X_1, X_4\}$

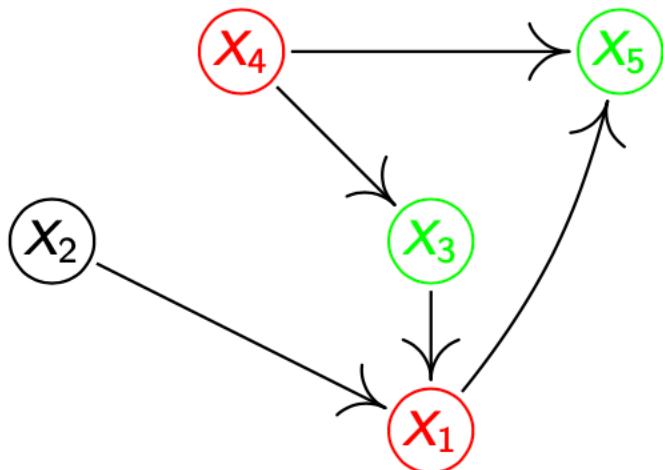
$X_4$  and  $X_1$  are  $d$ -sep. by  $\{X_2, X_3\}$

$X_2$  and  $X_4$  are  $d$ -sep. by  $\{\}$

# Definition: d-separation

$X_i$  and  $X_j$  are  $d$ -separated by  $S$  if all paths between  $X_i$  and  $X_j$  are blocked by  $S$ .

Check, whether all paths blocked!!



- ... → ○ → ... ○ blocks a path.
- ... ← ○ → ... ○ blocks a path.
- ... → ○ ← ... ○ blocks a path.

$X_2$  and  $X_5$  are  $d$ -sep. by  $\{X_1, X_4\}$

$X_4$  and  $X_1$  are  $d$ -sep. by  $\{X_2, X_3\}$

$X_2$  and  $X_4$  are  $d$ -sep. by  $\{\}$

$X_4$  and  $X_1$  are NOT  $d$ -sep. by  $\{X_3, X_5\}$

## Definition

$P$  is Markov w.r.t.  $G$  if

$$X_i \text{ and } X_j \text{ are } d\text{-separated by } \mathcal{S} \text{ in } G \Rightarrow X_i \perp\!\!\!\perp X_j \mid \mathcal{S}$$

## Definition

$P$  is Markov w.r.t.  $G$  if

$$X_i \text{ and } X_j \text{ are } d\text{-separated by } \mathcal{S} \text{ in } G \Rightarrow X_i \perp\!\!\!\perp X_j \mid \mathcal{S}$$

## Proposition

*Let the distribution  $P$  be Markov wrt a causal graph  $G$ . Then, Reichenbach's common cause principle is satisfied.*

Proof: dependent variables must be  $d$ -connected.

There are three equivalent formulations of the Markov condition.

- (i) **global Markov property:**

$$\mathbf{A} \text{ d-sep } \mathbf{B} \mid \mathbf{C} \text{ in } G \Rightarrow \mathbf{A} \perp\!\!\!\perp \mathbf{B} \mid \mathbf{C}$$

- (ii) **local Markov property:** each variable is independent of its non-descendants given its parents, and
- (iii) **Markov factorization property:**

$$p(\mathbf{x}) = p(x_1, \dots, x_d) = \prod_{j=1}^d p(x_j \mid \mathbf{pa}_j^G).$$

(assume existence of density)

## Definition

$P$  is Markov w.r.t.  $G$  if

$X_i$  and  $X_j$  are  $d$ -separated by  $S$  in  $G$   $\Rightarrow$   $X_i \perp\!\!\!\perp X_j | S$

## Definition

$P$  is Markov w.r.t.  $G$  if

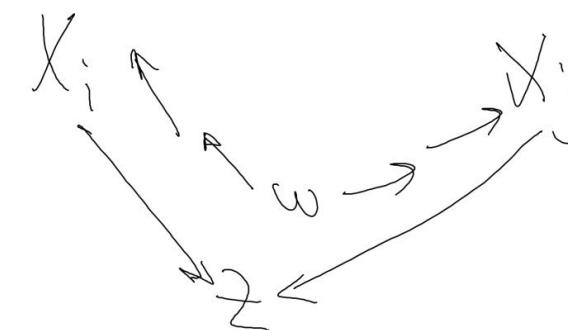
$$X_i \text{ and } X_j \text{ are } d\text{-separated by } \mathcal{S} \text{ in } G \Rightarrow X_i \perp\!\!\!\perp X_j | \mathcal{S}$$

## Definition

$P$  is faithful w.r.t.  $G$  if

$$X_i \text{ and } X_j \text{ are } d\text{-separated by } \mathcal{S} \text{ in } G \Leftarrow X_i \perp\!\!\!\perp X_j | \mathcal{S}$$

Examples...



$X_i \& X_j$  d-sep by  $\emptyset$   
 $\Rightarrow X_i \perp\!\!\!\perp X_j$

Ex. Assume P is MARKOV &  
faithful w.r.t.  $\mathcal{C}$ . The following  
lists are complete

(i)  $X \perp\!\!\! \perp Z$  (Var: X, Y, Z)

(ii)  $X \perp\!\!\! \perp Y \mid Z$  (- || -)

(iii)  $X \perp\!\!\! \perp Y$ ,  $X \perp\!\!\! \perp W \mid Z$ ,  $X \perp\!\!\! \perp W \mid Z, Y$

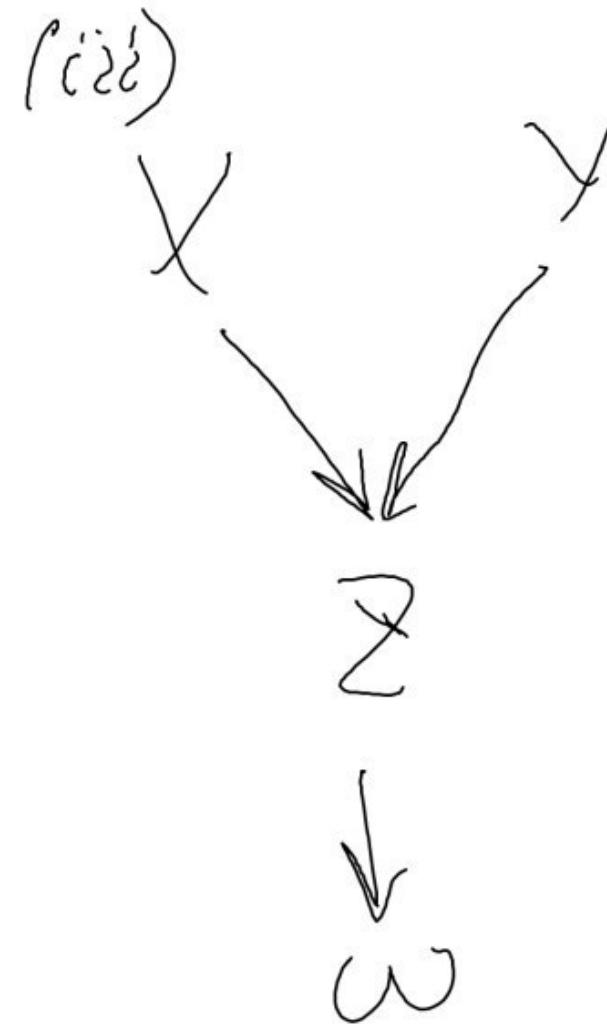
$Y \perp\!\!\! \perp W \mid Z, X$ ,  $Y \perp\!\!\! \perp W \mid Z$

(i)  $X \rightarrow Y \rightarrow Z$

(ii)  $X \rightarrow Z \rightarrow Y$

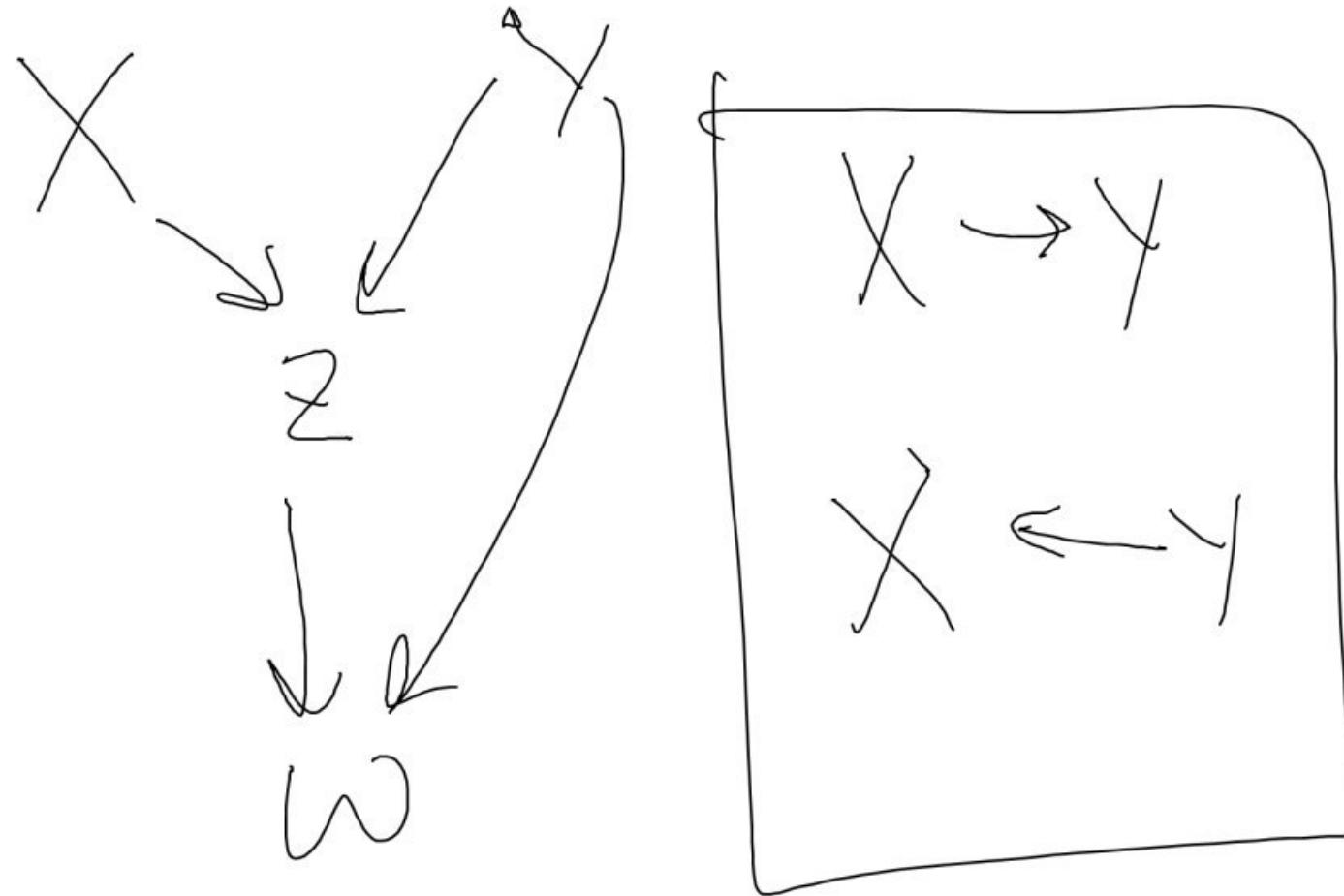
OR  $X \leftarrow Z \leftarrow Y$

OR  $X \leftarrow Z \rightarrow Y$

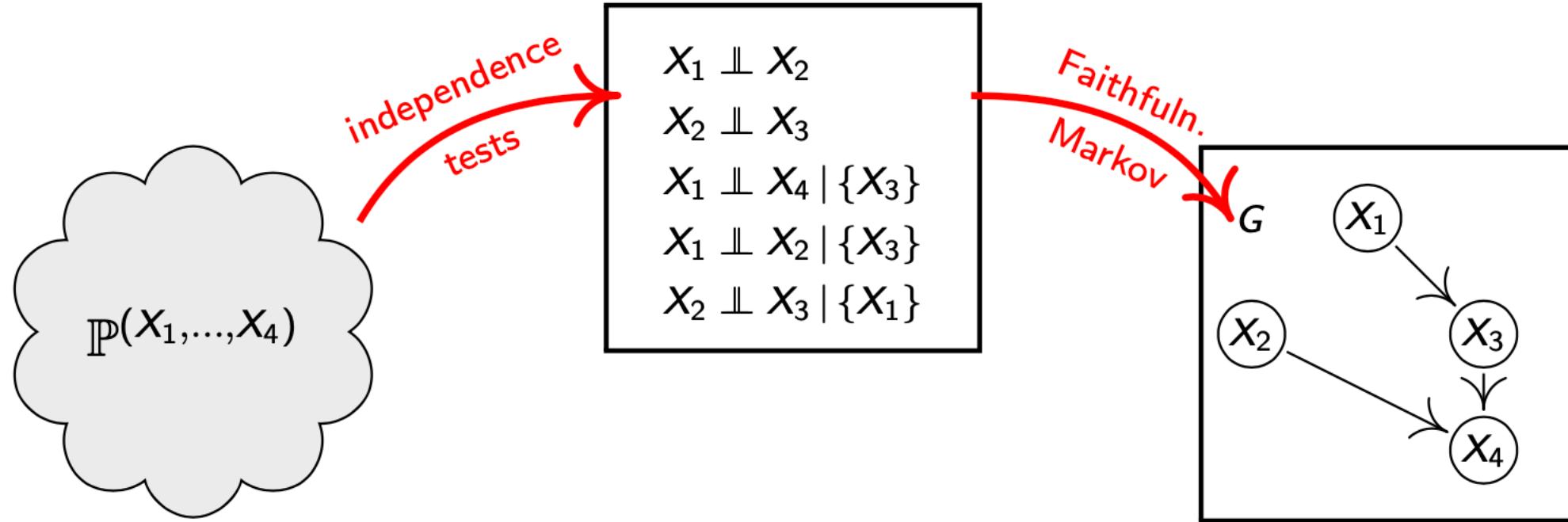


def  $G, H$  encoding same set is  
d-sep  $S \rightsquigarrow$  "Markov equivalence"  
members of the same Markov  
equivalence class

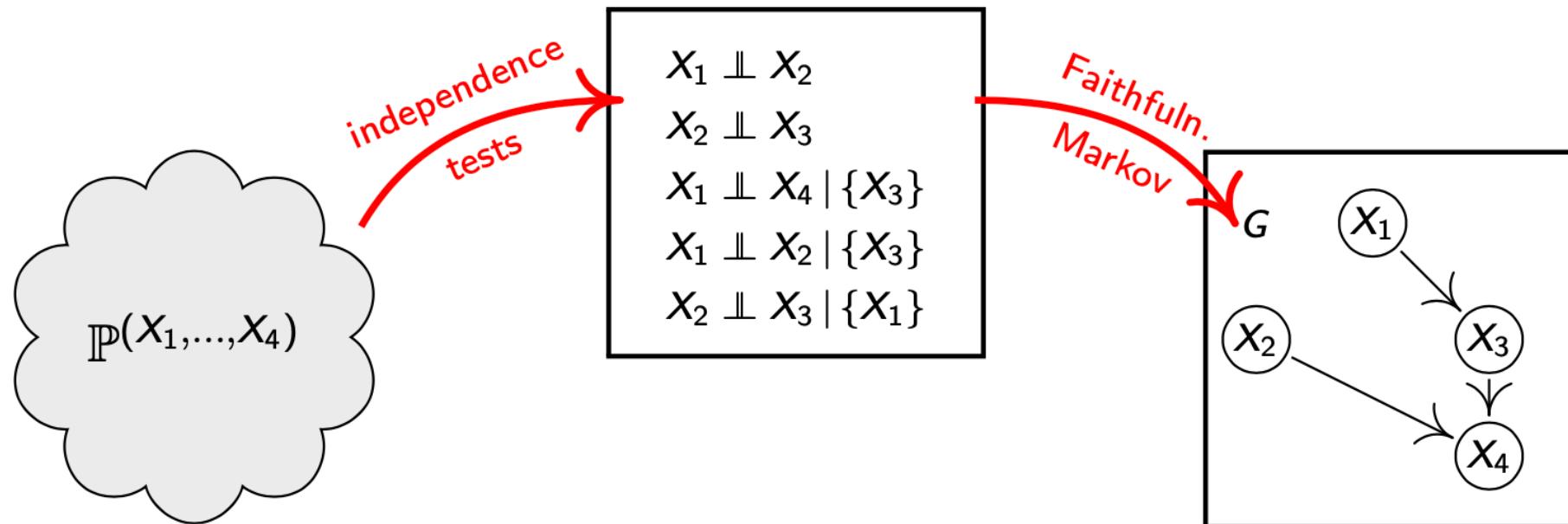
Prop:  $G, H$  Markov equivalent  
 $\Leftrightarrow$  they have the same skeletons  
of the same V-structures



# Idea 1: independence-based methods



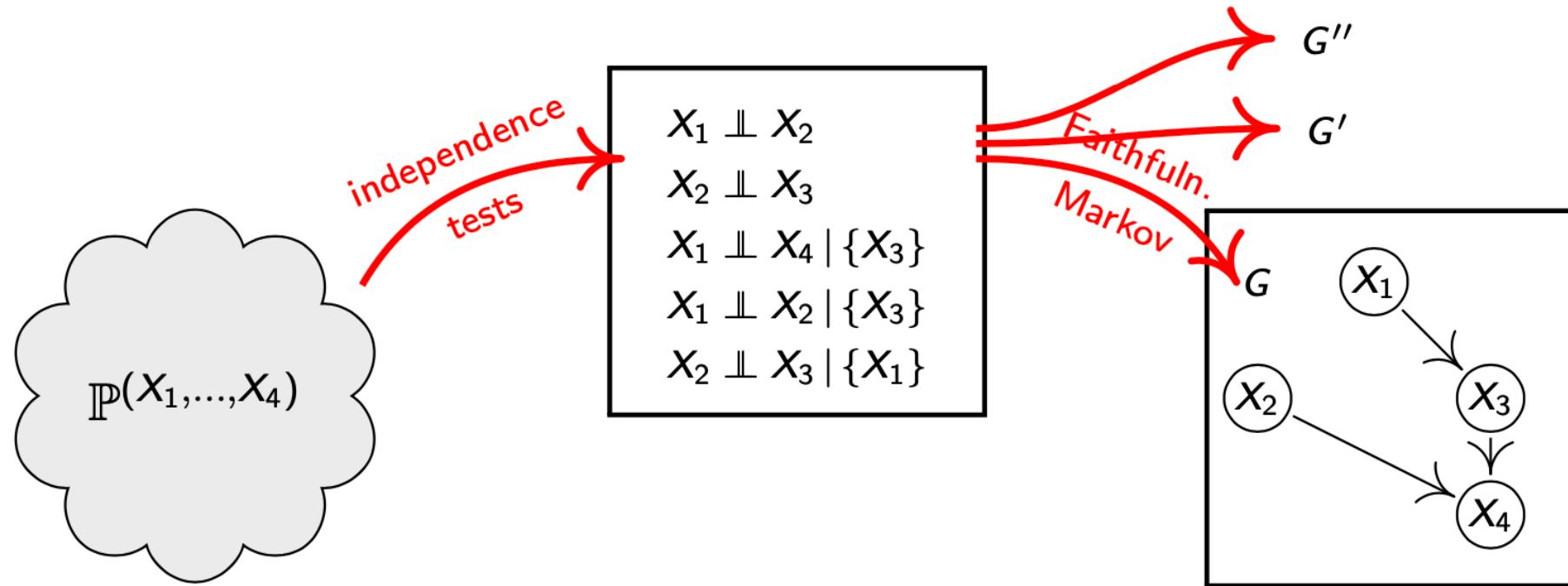
# Idea 1: independence-based methods



Method: IC (Pearl 2009); PC, FCI (Spirtes et al., 2000)

- ① Find all (cond.) independences from the data.
- ② Select the DAG(s) that corresponds to these independences.

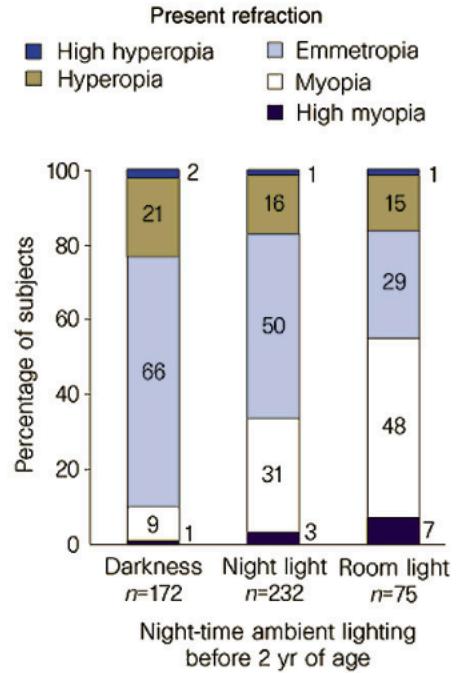
# Idea 1: independence-based methods



Method: IC (Pearl 2009); PC, FCI (Spirtes et al., 2000)

- ① Find all (cond.) independences from the data.
- ② Select the DAG(s) that corresponds to these independences.

# Example: myopia



We have

- night light  $\not\perp$  child myopia
- night light  $\perp\!\!\!\perp$  child myopia | parent myopia
- no other independences

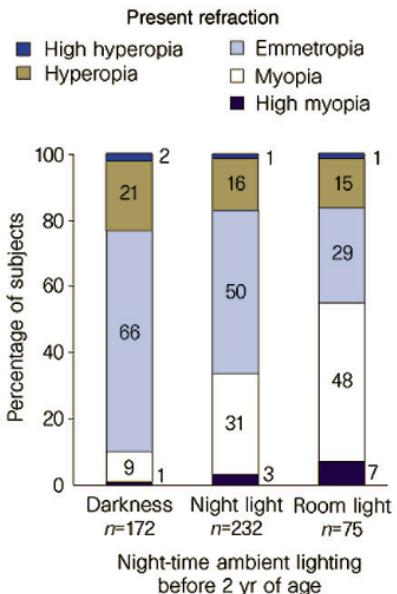
Quinn et al.: *Myopia and ambient lighting at night*, Nature 1999

Zadnik et al.: *Vision: Myopia and ambient night-time light.*, Nature 2000

Gwiazda et al.: *Vision: Myopia and ambient night-time light.*, Nature 2000

and therefore ...

# Example: myopia



We have

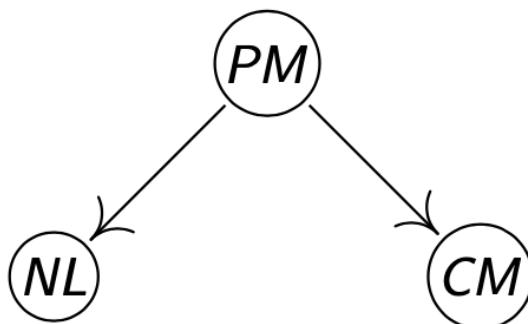
- night light  $\perp\!\!\!\perp$  child myopia
- night light  $\perp\!\!\!\perp$  child myopia | parent myopia
- no other independences

Quinn et al.: *Myopia and ambient lighting at night*, Nature 1999

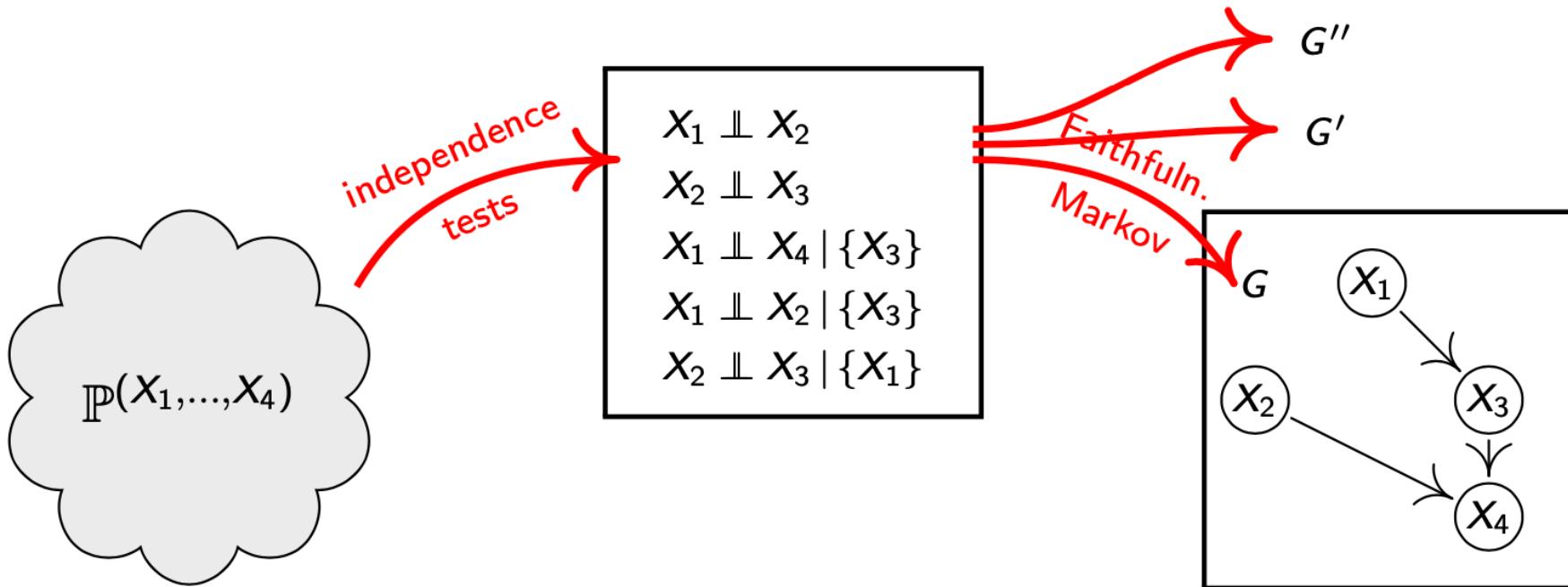
Zadnik et al.: *Vision: Myopia and ambient night-time light.*, Nature 2000

Gwiazda et al.: *Vision: Myopia and ambient night-time light.*, Nature 2000

and therefore ...



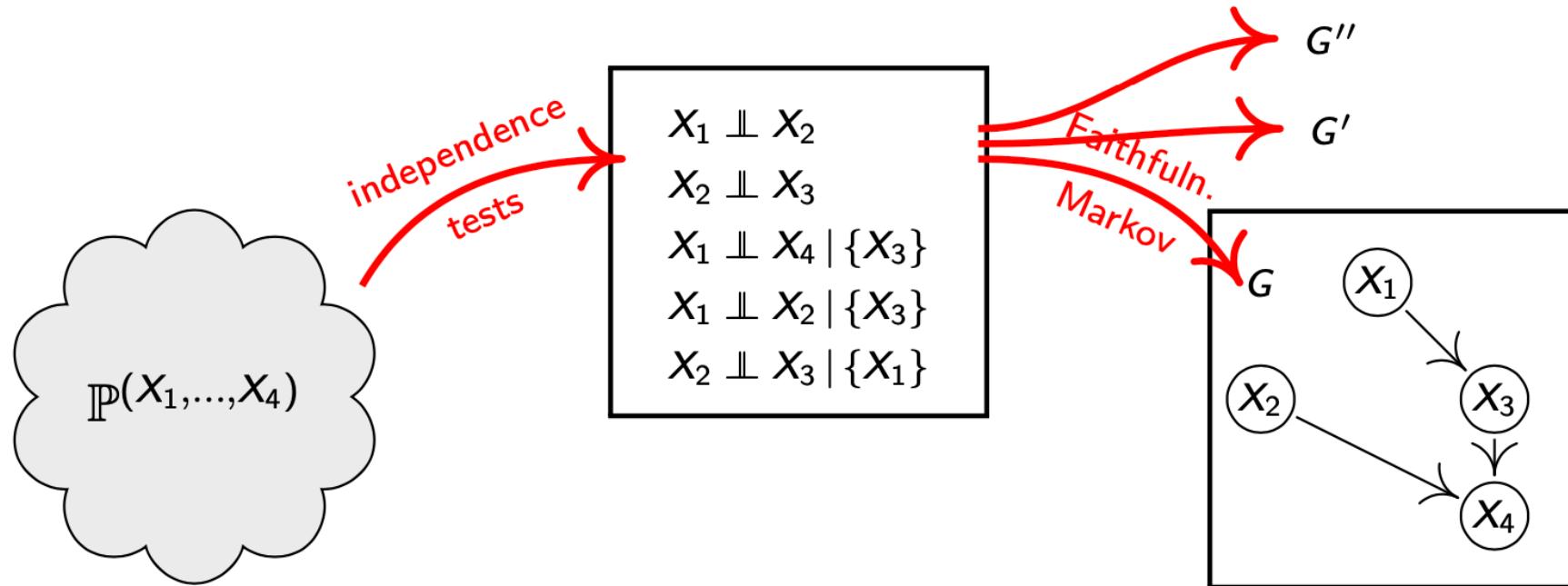
# Idea 1: independence-based methods



Method: IC (Pearl 2009); PC, FCI (Spirtes et al., 2000)

- ① Find all (cond.) independences from the data.
- ② Select the DAG(s) that corresponds to these independences.

# Idea 1: independence-based methods



Method: IC (Pearl 2009); PC, FCI (Spirtes et al., 2000)

- ① Find all (cond.) independences from the data. Be smart.
- ② Select the DAG(s) that corresponds to these independences.

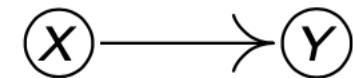
What do we do with two variables? (Nothing is possible in general.)

# Idea 2: restricted structural causal models

Consider a distribution generated by

$$Y = \alpha X + N_Y$$

with  $N_Y, X$  ind.



# Idea 2: restricted structural causal models

Consider a distribution generated by

$$\boxed{Y = \alpha X + N_Y}$$

with  $N_Y, X$  ind.

$$(X) \rightarrow (Y)$$

Then, if  $(X, N_Y)$  is non-Gaussian, there is no

$$\cancel{\boxed{X = \beta Y + M_X}}$$

with  $M_X, Y$  ind.

$$(X) \leftarrow (Y)$$

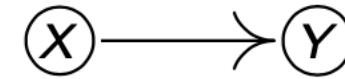
Shimizu et al. 2006

# Idea 2: restricted structural causal models

Consider a distribution corresponding to

$$\boxed{Y = 2X + N_Y}$$

with  $N_Y, X \stackrel{\text{ind}}{\sim} \mathcal{U}$

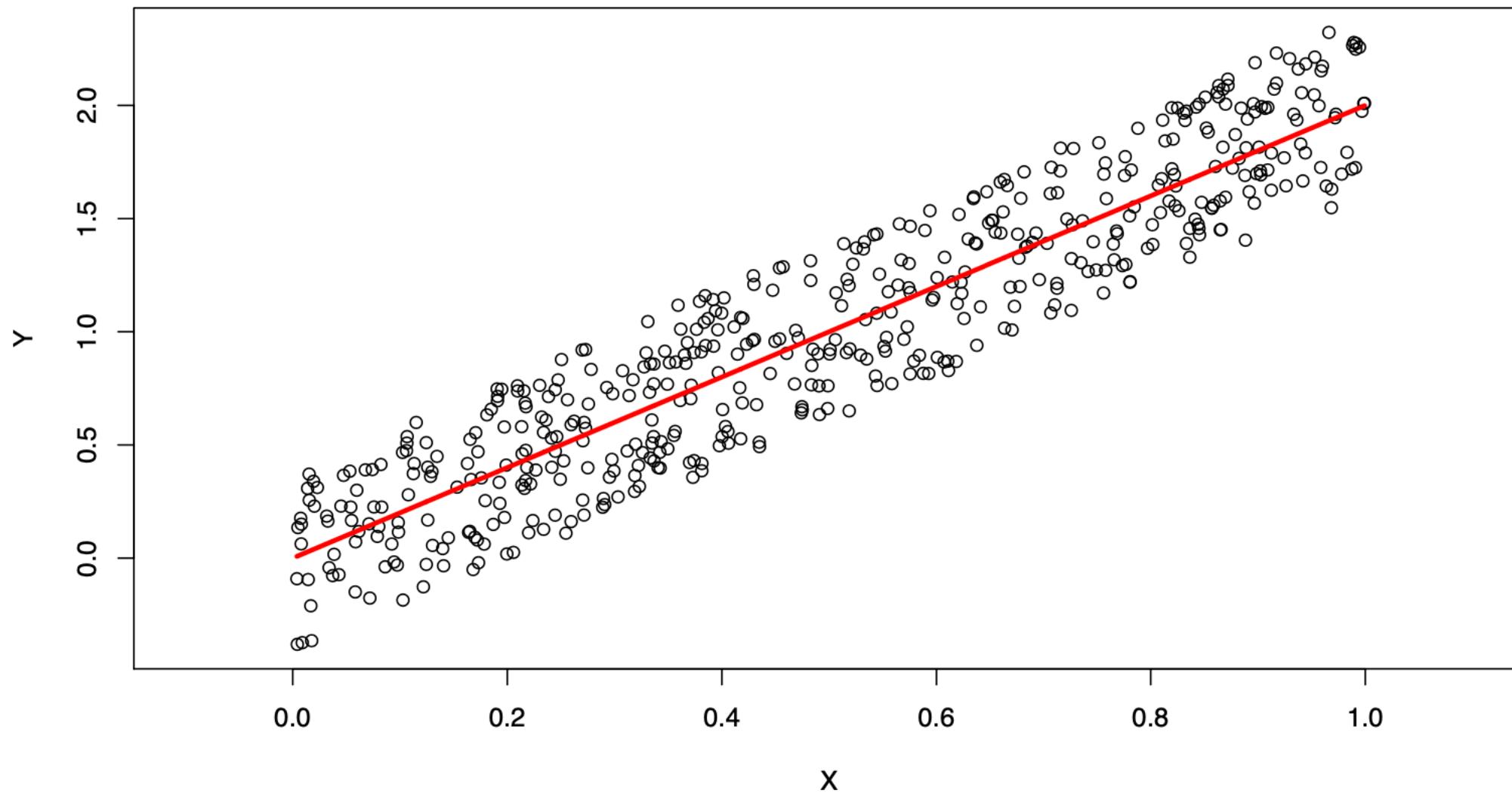


with

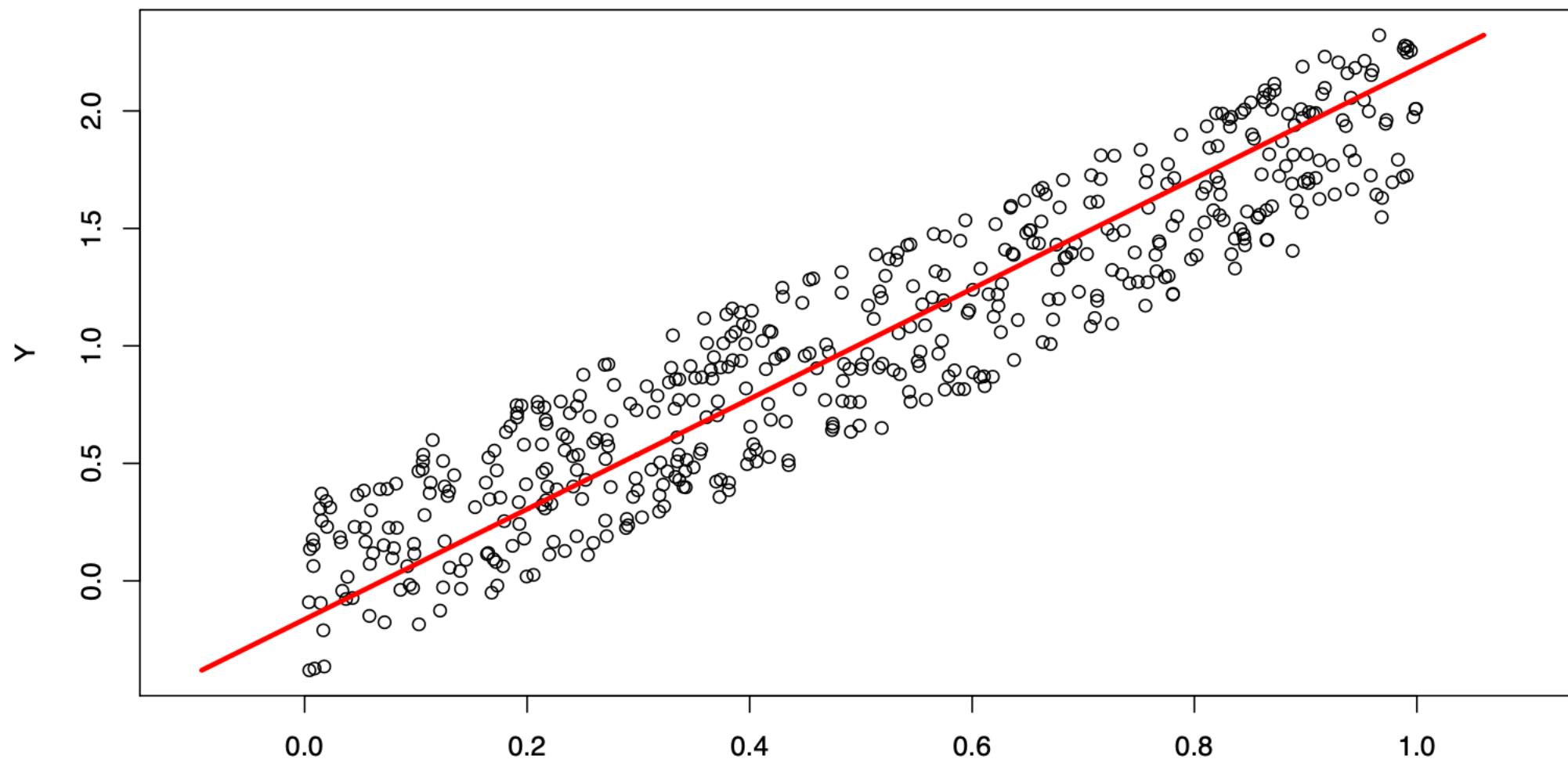
$$X \sim \mathcal{U}[-1, 1]$$

$$N_Y \sim \mathcal{U}[-0.4, 0.4]$$

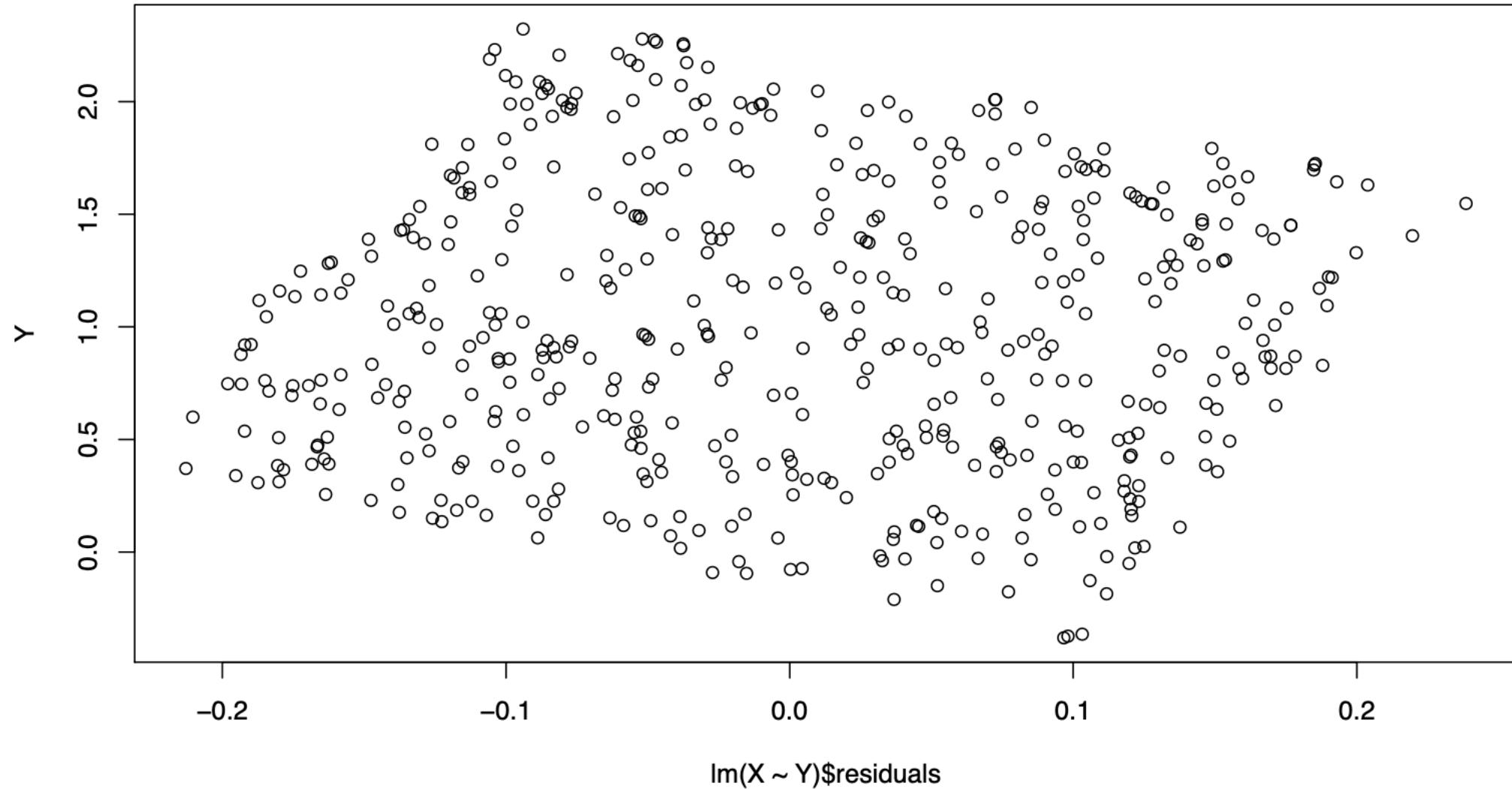
# Idea 2: restricted structural causal models



# Idea 2: restricted structural causal models



# Idea 2: restricted structural causal models



# Idea 2: restricted structural causal models

Method...

Prop.

$\forall P \text{ on } (X, Y)$

$\exists g, N_Y, \quad Y = g(X, N_Y)$

$N_Y \perp\!\!\!\perp X$

$\exists f, N_X, \quad X = f(Y, N_X)$

$N_X \perp\!\!\!\perp Y$

Prop.

$\forall \rho$  on  $(X, Y)$

$\exists g, N_Y, Y = g(X, N_Y)$

$N_Y \perp\!\!\!\perp X$

$\exists f, N_X, X = f(Y, N_X)$

$N_X \perp\!\!\!\perp Y$

inverse  
cdf of  
 $X|Y$

uniform

Method

1) Regress  $Y_{\text{on}} X$   
to check ind. between  
residuals and  $X$

2)

—||—  $X$  on  $Y$   
—||—  $H$   
—||— and  $Y$

$X \sim N_X$

$$Y = X^3 + \alpha_Y$$

HSIC test:

	p-value
residuals vs. X	0.77
residuals vs. Y	0.004

## Idea 2: restricted structural causal models

Theory...

Th.: Assume

$$Y = \alpha X + \nu_Y$$

$$\exists \beta, M_X \quad X = \beta Y + \mu_X$$

$\Rightarrow (X, \nu_Y)$  is Gaussian

# Idea 2: restricted structural causal models

Peters et al ICML 2009 (univariate), Bauer et al ICML 2016 (multivariate)

## Theorem

Let  $(X_t)_t$  be a causal<sup>a</sup> solution of an ARMA( $p, q$ ) process:

$$X_t - \phi_1 X_{t-1} - \dots - \phi_p X_{t-p} = Z_t + \theta_1 Z_{t-1} + \dots + \theta_q Z_{t-q}.$$

Then,  $X_t$  is time reversible, i.e., a causal solution of an ARMA( $\tilde{p}, \tilde{q}$ ) process with reversed time, if and only if  $(Z_t)_t$  is Gaussian.

---

<sup>a</sup> $(X_t)_t$  causal iff  $Z_t \perp\!\!\!\perp X_{t-k}$ ,  $k > 0$ .

# Idea 2: restricted structural causal models

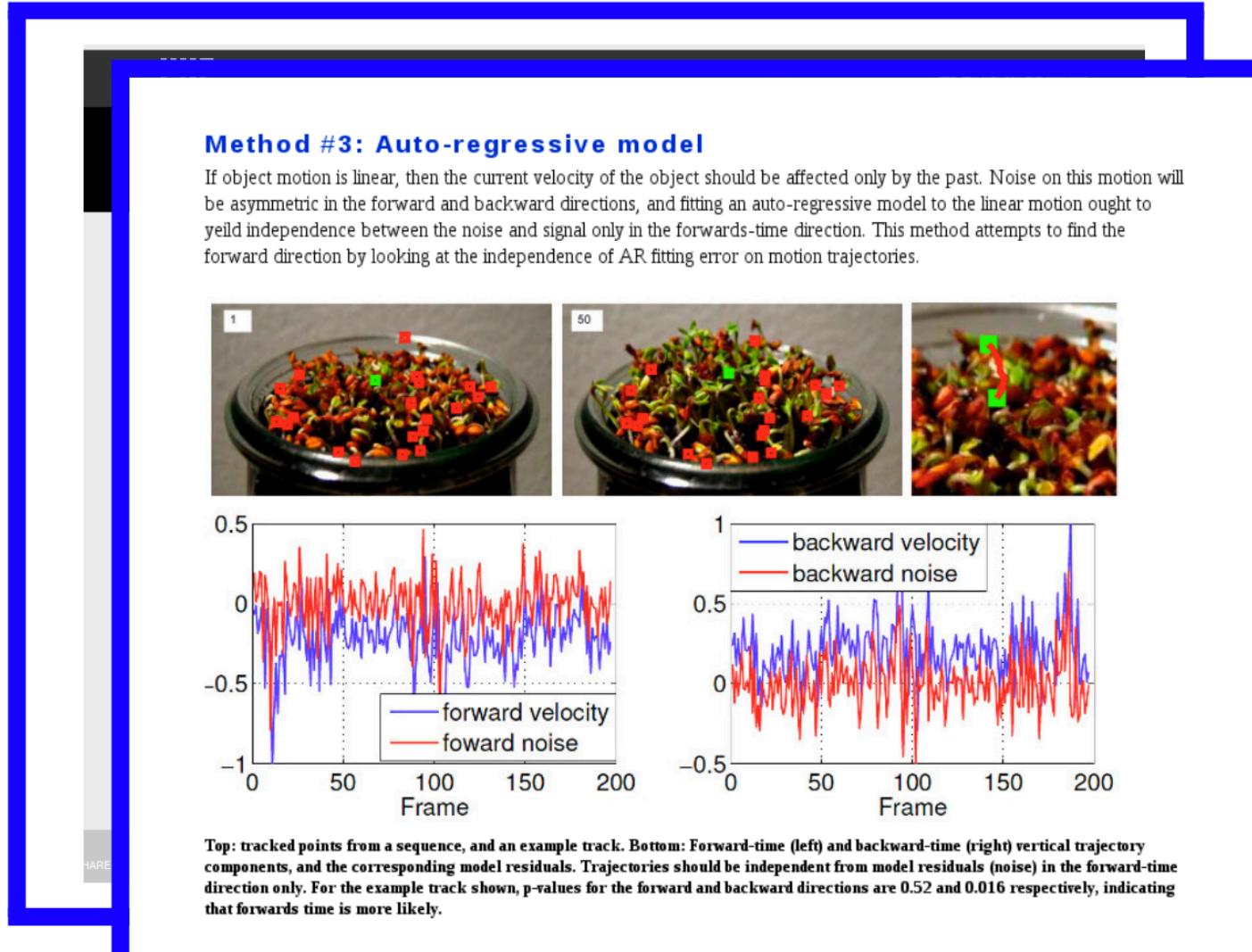
Pickup et al. 2014:



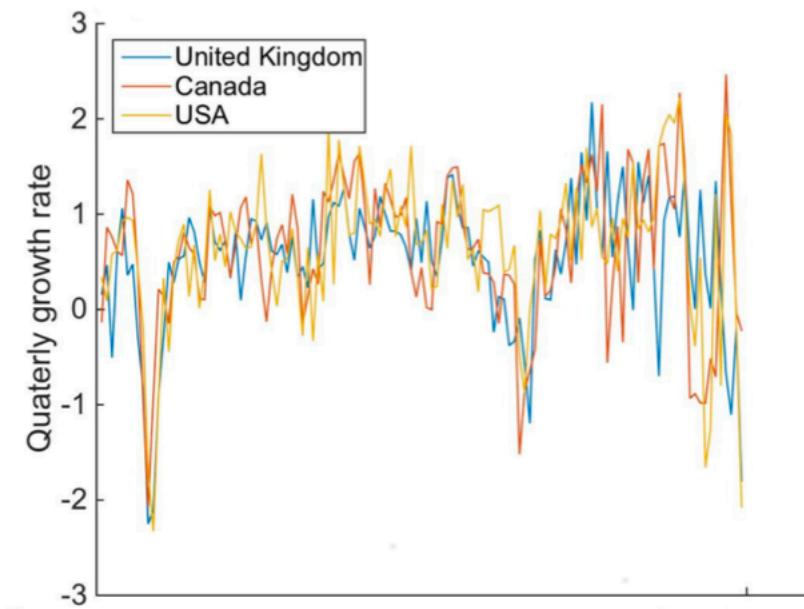
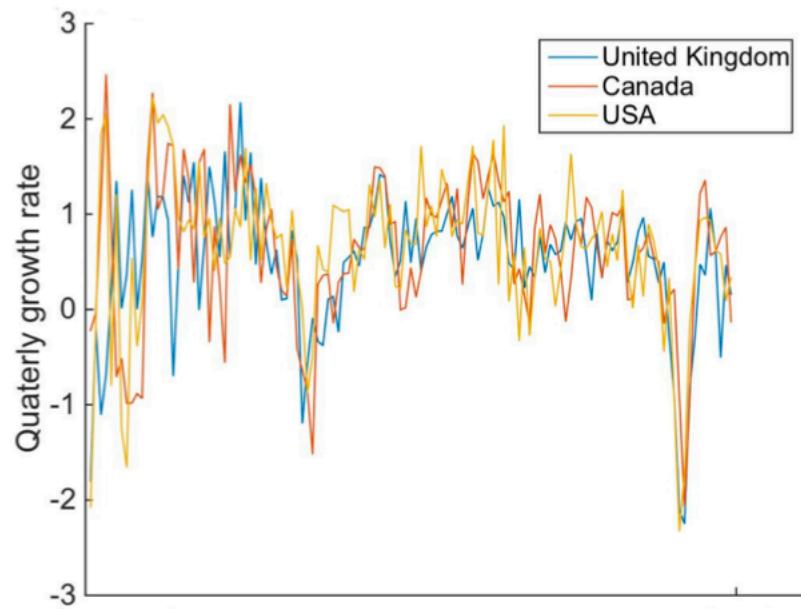
The image shows a screenshot of the MIT News website. At the top, the MIT logo and "Massachusetts Institute of Technology" are visible, along with navigation links for NEWS, VIDEO, SOCIAL, and FOLLOW MIT (with icons for RSS, Twitter, Facebook, Google+, and others). Below the header is a large banner with the text "MIT News" and "ON CAMPUS AND AROUND THE WORLD". To the right of the banner are "Browse" and "Search" buttons. A search bar with a magnifying glass icon is also present. The main content area features a large illustration of several circular arrows in blue, purple, and red, some pointing left and some pointing right, set against a light blue background. Below the illustration, the text "Illustration: Jose-Luis Olivares/MIT" is visible. The main headline reads "Can we see the arrow of time?". Below the headline, a subtext states: "Algorithm can determine, with 80 percent accuracy, whether video is running forward or backward." At the bottom left, the author is listed as "Larry Hardesty | MIT News Office" and the date as "June 20, 2014". On the far left, there is a "SHARE" button. At the bottom right, there is a "RELATED" section and a link to a paper titled "PAPER: 'Seeing the arrow of time.'".

# Idea 2: restricted structural causal models

Pickup et al. 2014:

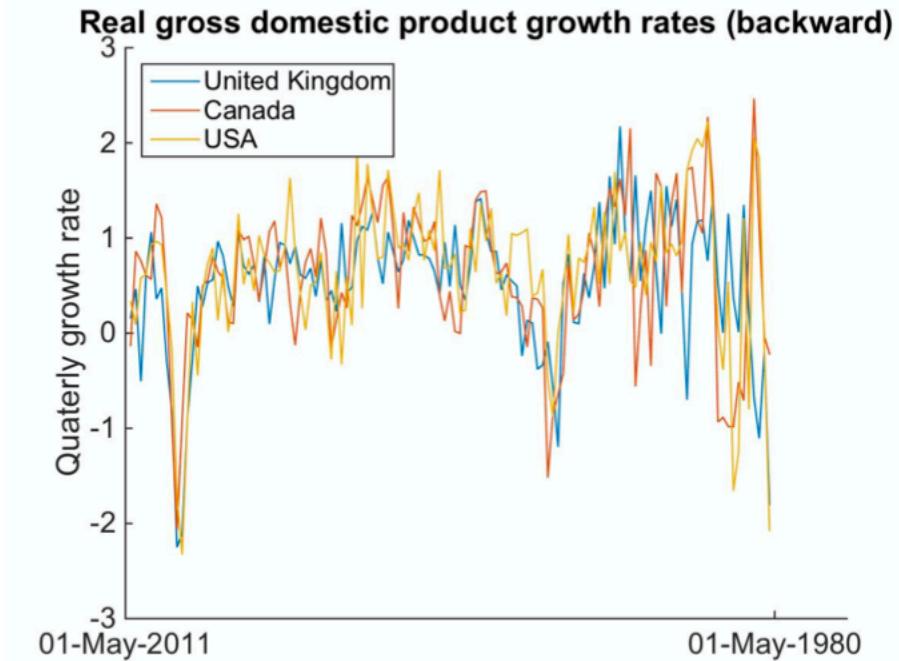
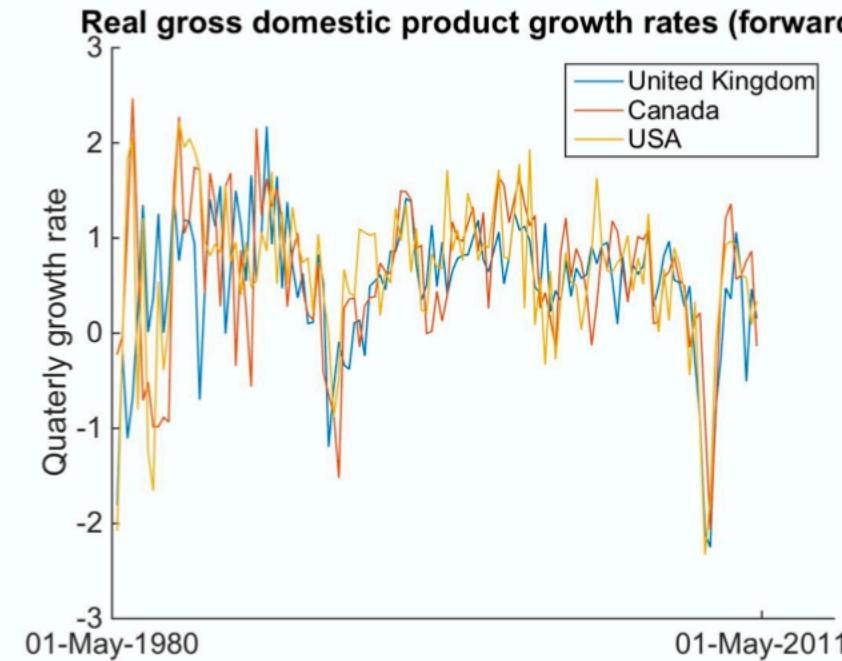


# Idea 2: restricted structural causal models



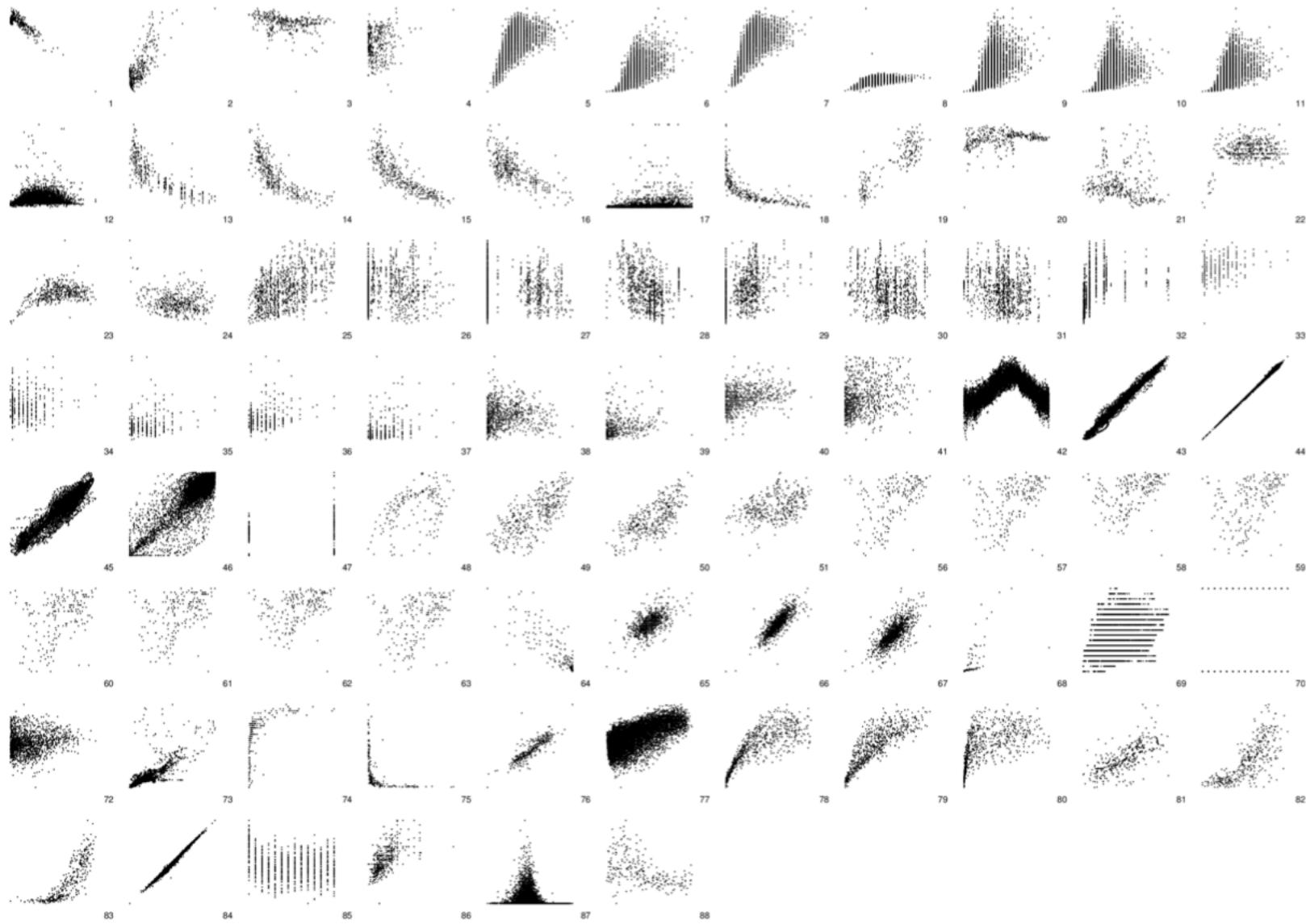
Quarterly growth rates in percentage of GDP for the UK, Canada and USA (Tsay et al 2014).

# Idea 2: restricted structural causal models



Quarterly growth rates in percentage of GDP for the UK, Canada and USA (Tsay et al 2014).

# Idea 2: restricted structural causal models

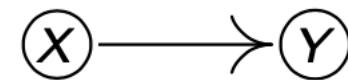


# Idea 2: restricted structural causal models

Consider a distribution entailed by

$$Y = f(X) + N_Y$$

with  $N_Y, X \stackrel{ind}{\sim} \mathcal{N}$

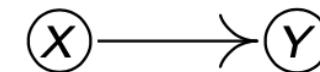


# Idea 2: restricted structural causal models

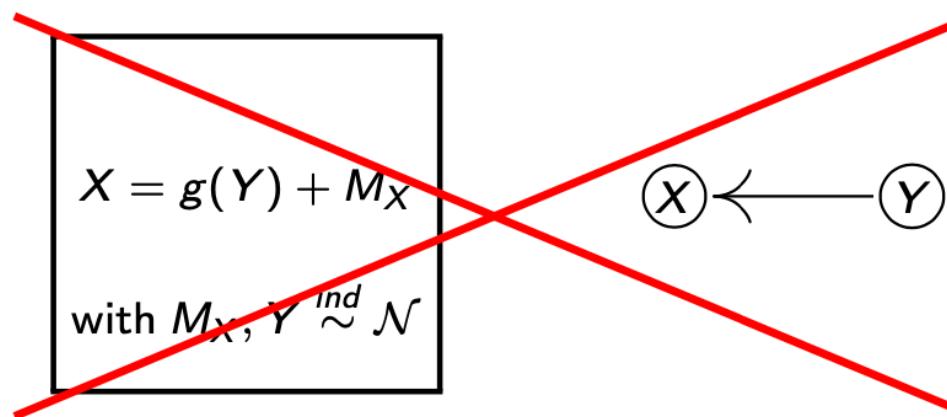
Consider a distribution entailed by

$$\boxed{Y = f(X) + N_Y}$$

with  $N_Y, X \stackrel{ind}{\sim} \mathcal{N}$



Then, if  $f$  is nonlinear, there is no


$$\boxed{X = g(Y) + Mx}$$

with  $Mx, Y \stackrel{ind}{\sim} \mathcal{N}$

# Idea 2: restricted structural causal models

Consider a distribution corresponding to

$$\boxed{Y = \textcolor{red}{X}^3 + N_Y}$$

with  $N_Y, X \stackrel{\textit{ind}}{\sim} \mathcal{N}$

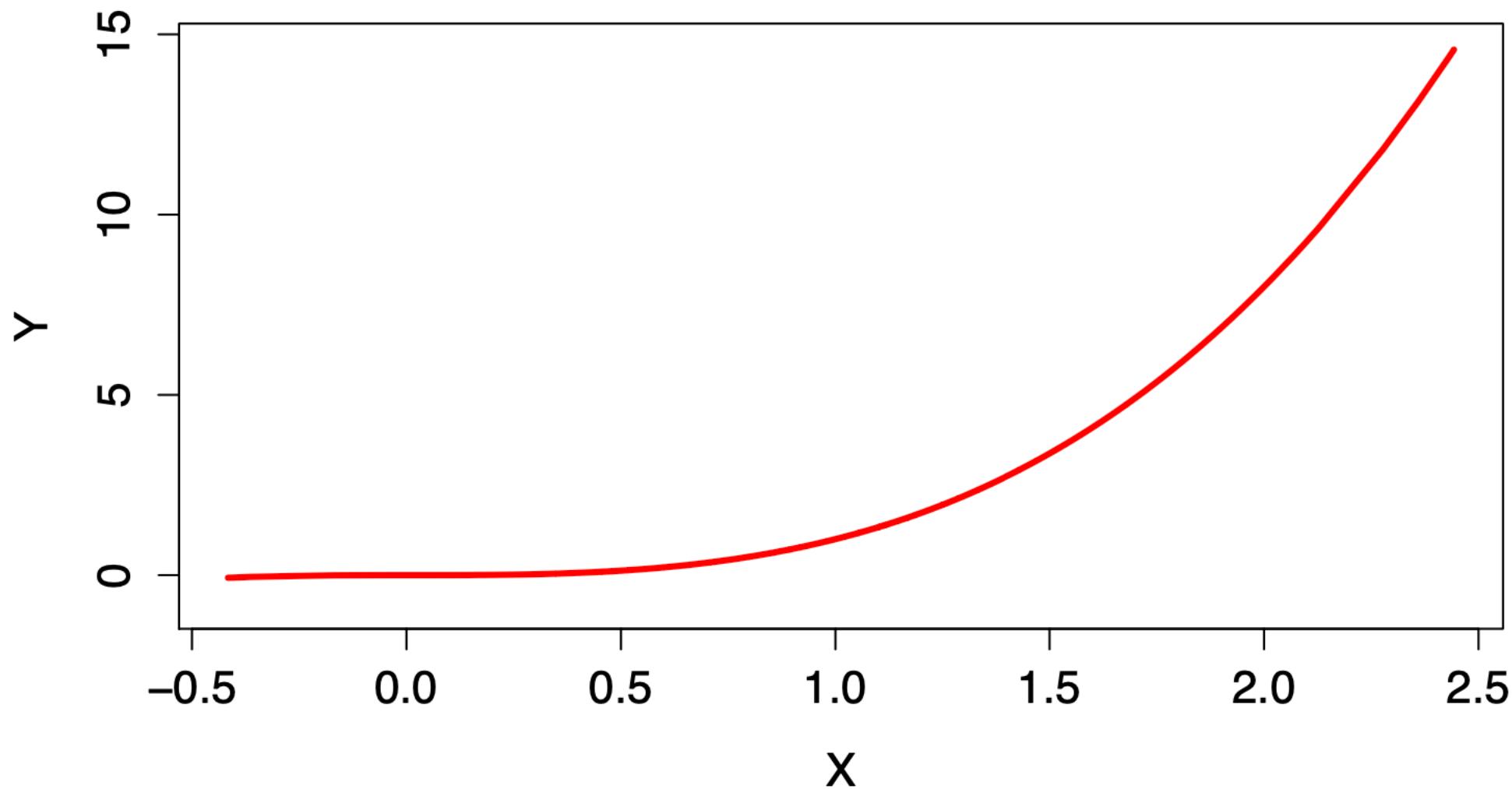


with

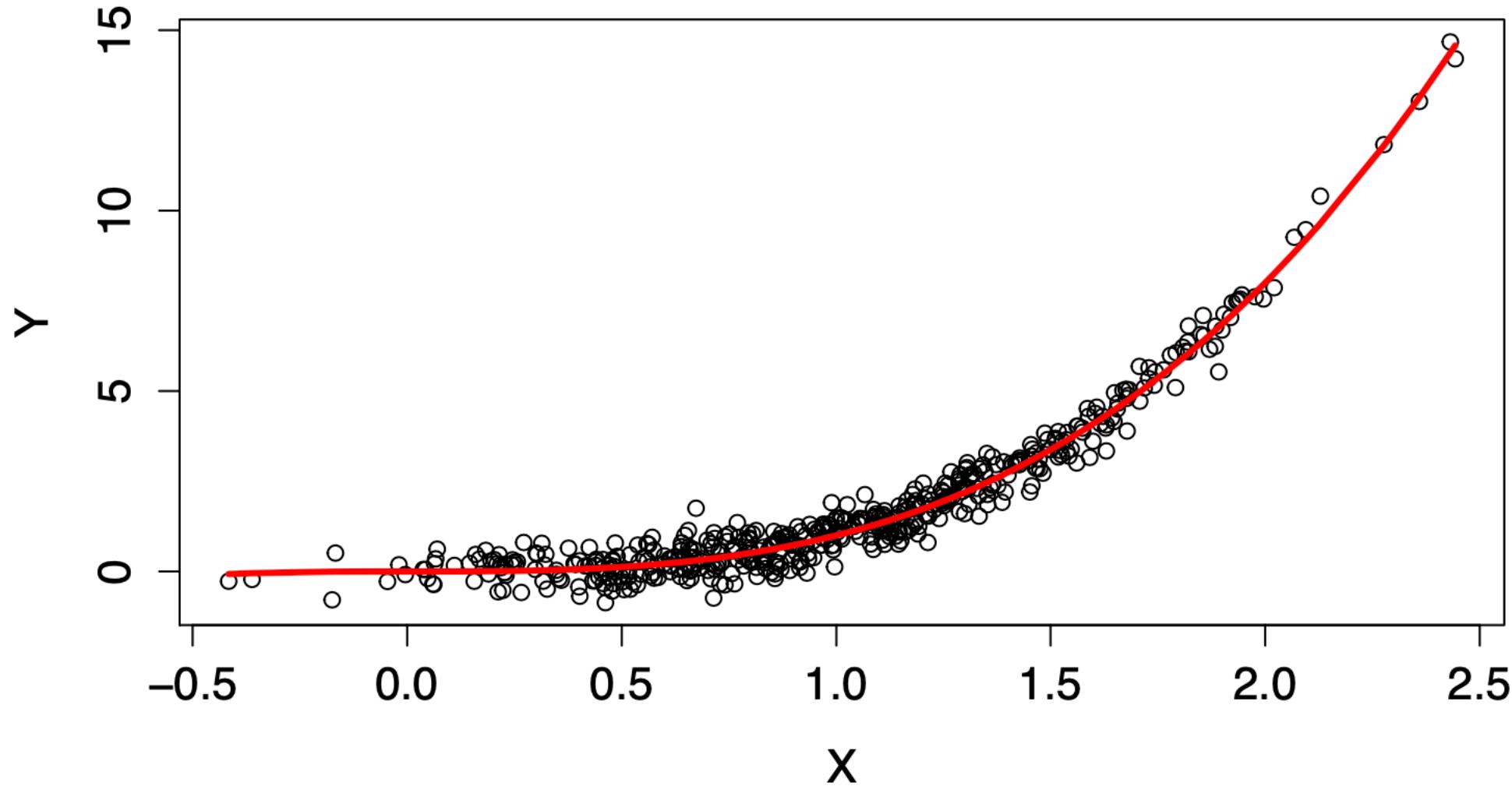
$$X \sim \mathcal{N}(1, 0.5^2)$$

$$N_Y \sim \mathcal{N}(0, 0.4^2)$$

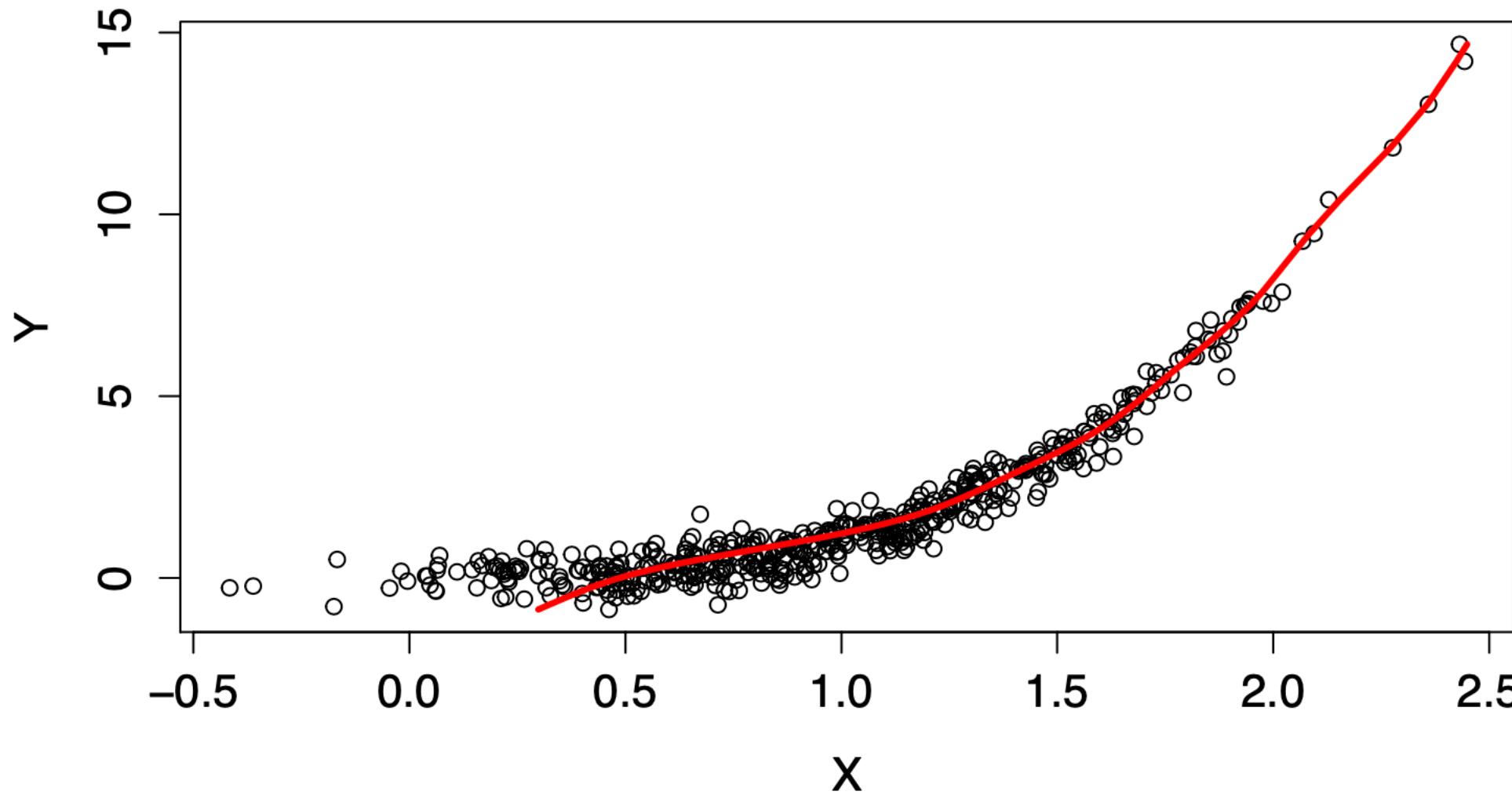
## Idea 2: restricted structural causal models



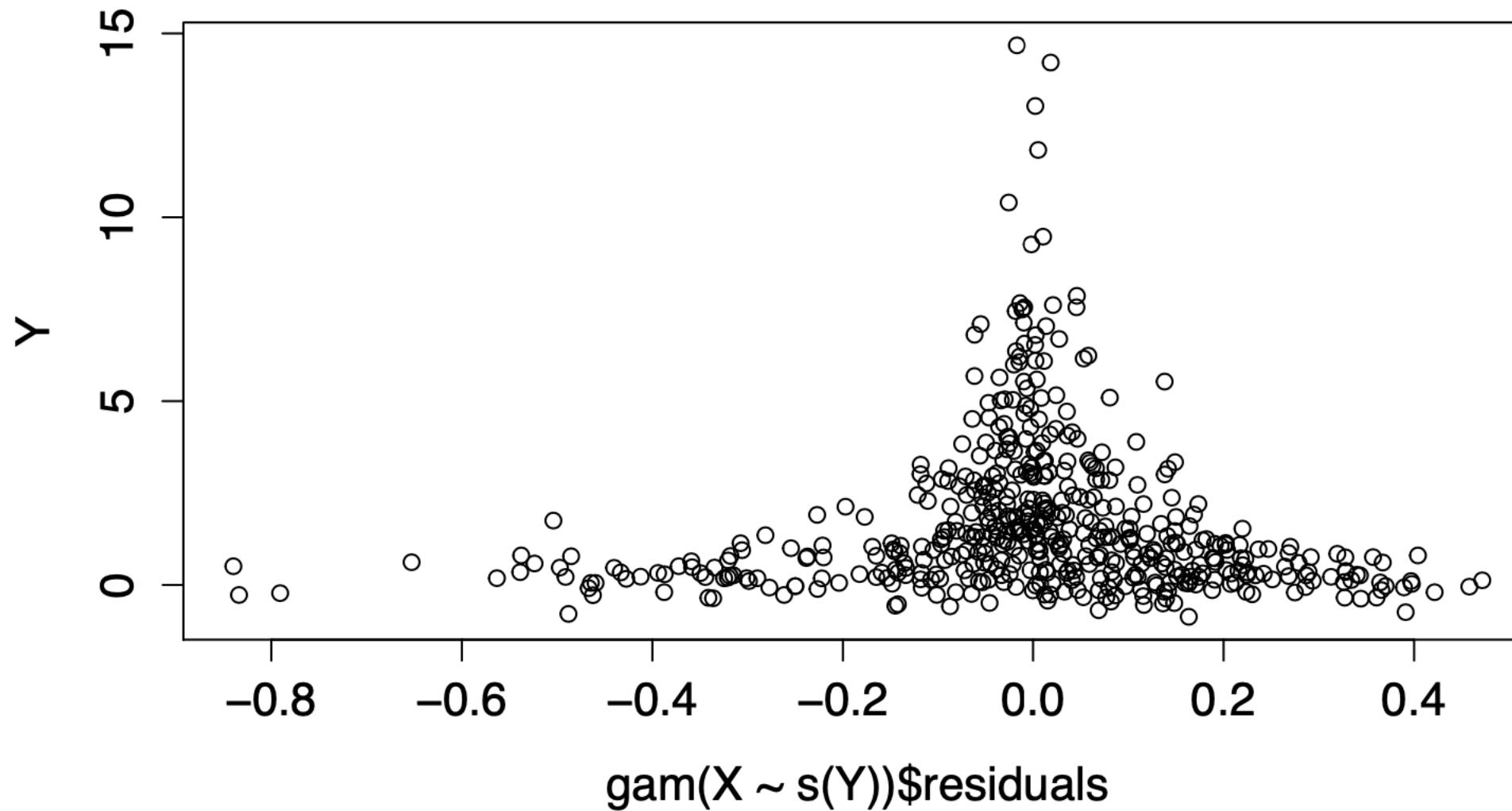
## Idea 2: restricted structural causal models



# Idea 2: restricted structural causal models

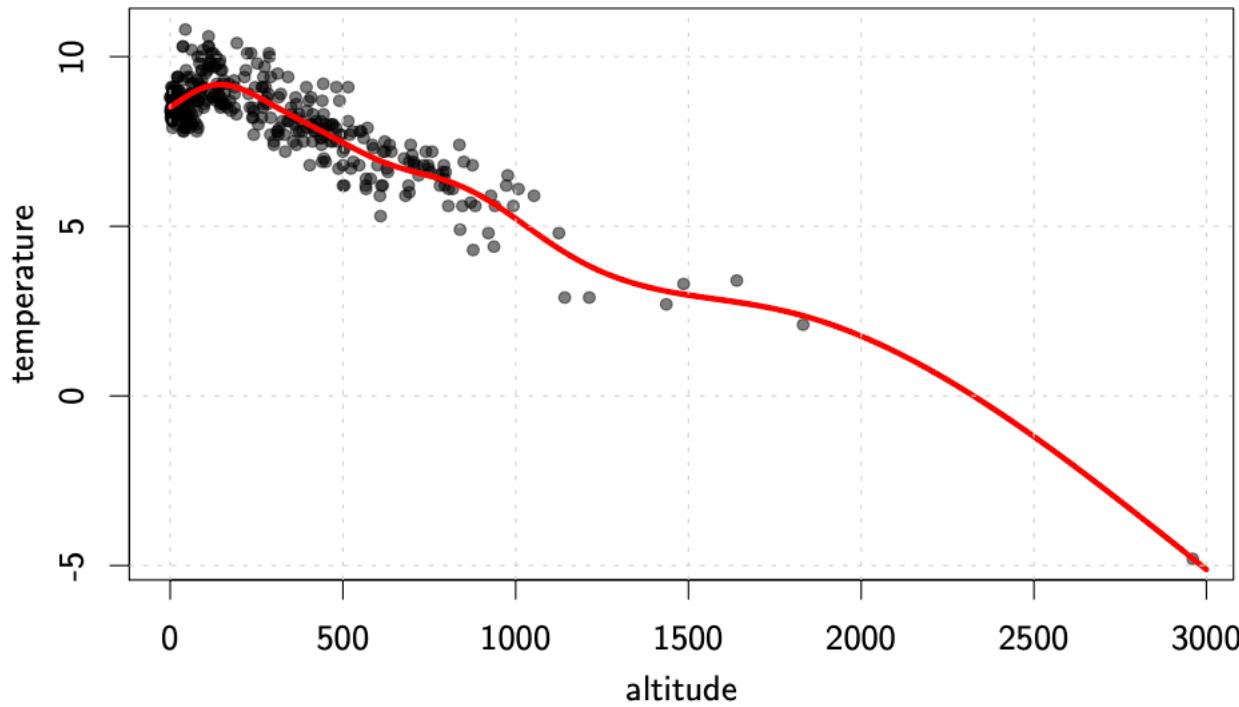


# Idea 2: restricted structural causal models

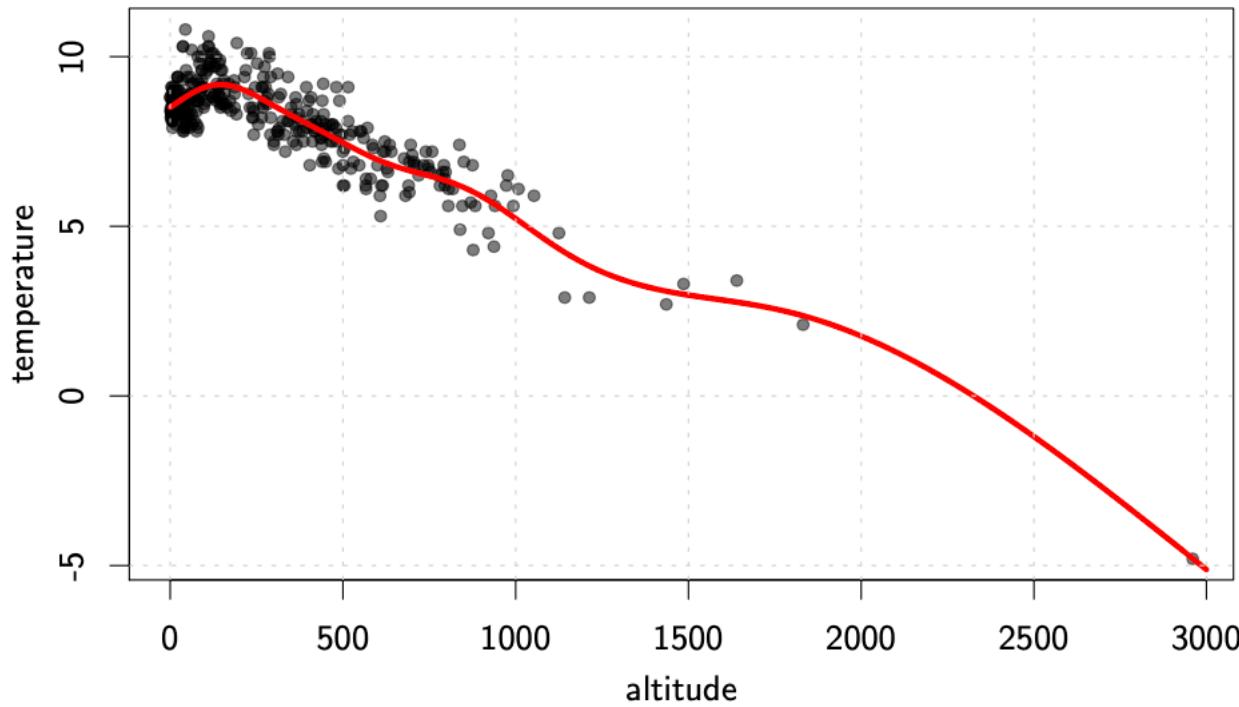


Method...

# Real Data: altitude and temperature



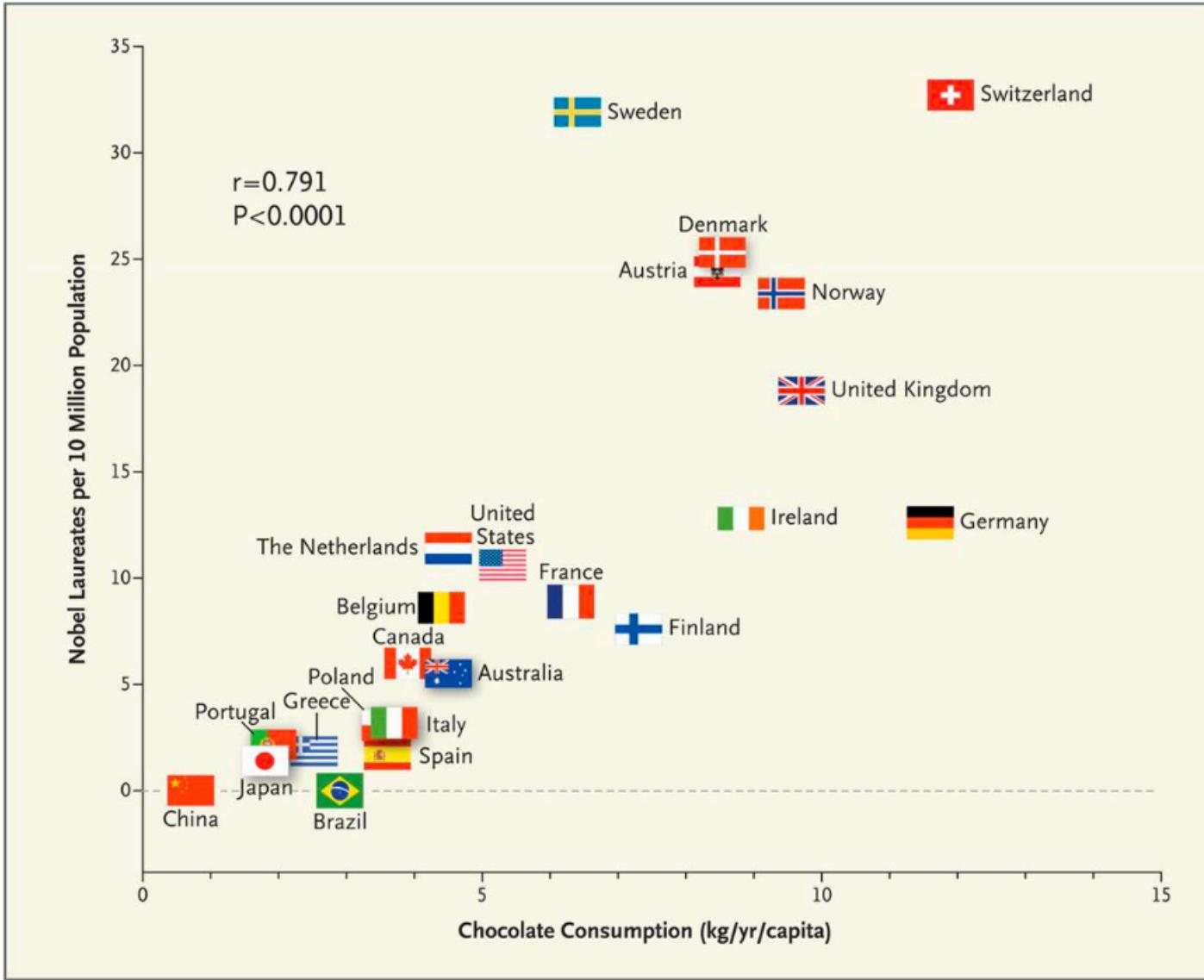
# Real Data: altitude and temperature



p-value forward: 0.024

p-value backward: 0.000000000019

# Example: chocolate

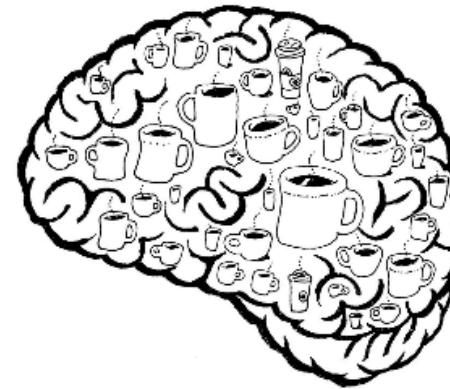


F. H. Messerli: *Chocolate Consumption, Cognitive Function, and Nobel Laureates*, N Engl J Med 2012

# Example: chocolate

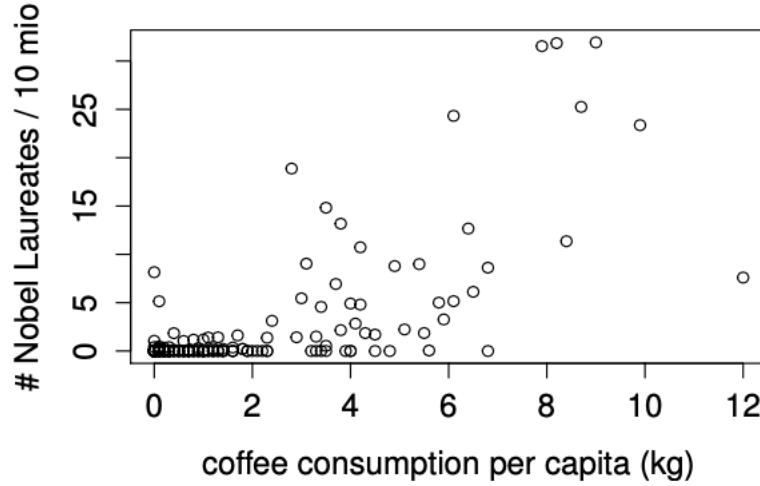


No (not enough) data for chocolate

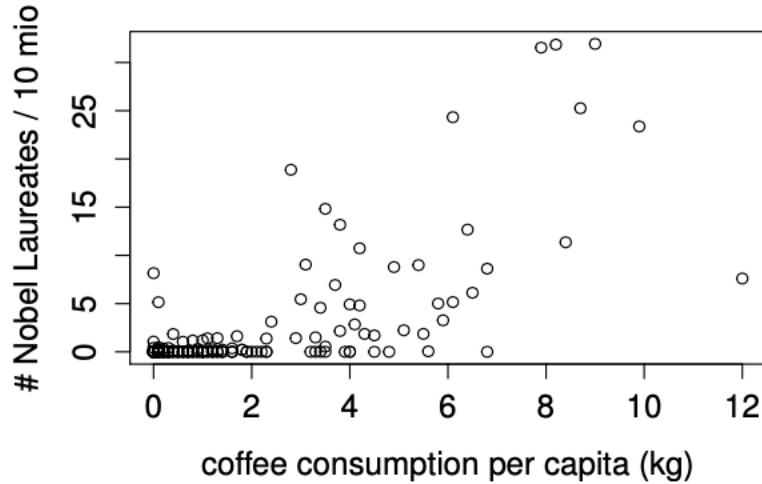


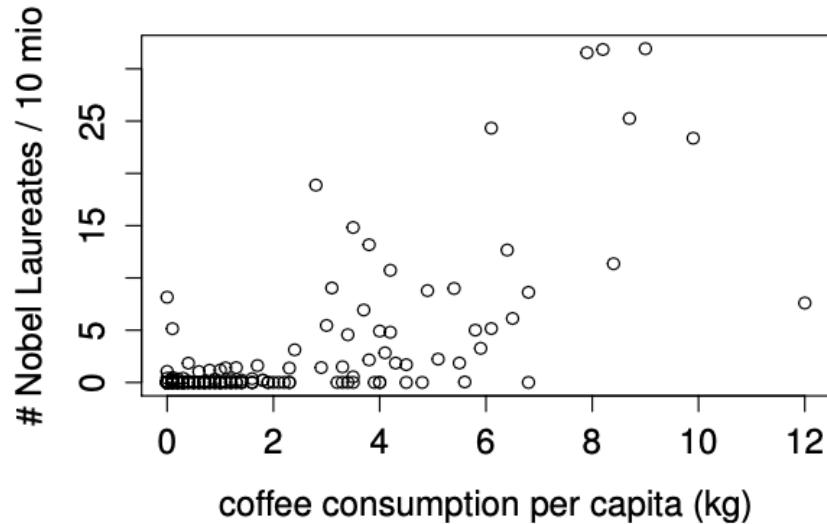
... but we have data for coffee!

# Example: chocolate



# Example: chocolate





Correlation: 0.698  
 $p$ -value:  $< 2.2 \cdot 10^{-16}$

Coffee → Nobel Prize: Dependent residuals ( $p$ -value of  $5.1 \cdot 10^{-78}$ ).  
Nobel Prize → Coffee: Dependent residuals ( $p$ -value of  $3.1 \cdot 10^{-12}$ ).

⇒ Model class too small? Causally insufficient?

Question: When is a  $p$ -value too small?