

Causality

Evgeny Burnaev

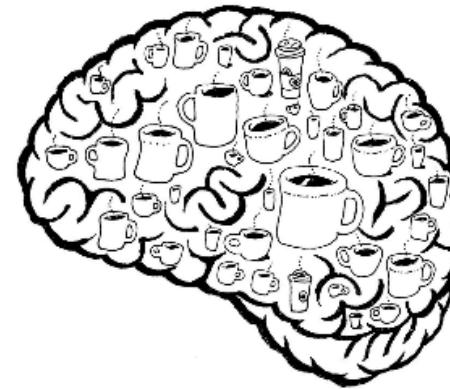
Rodrigo Rivera

Skoltech

Example: chocolate

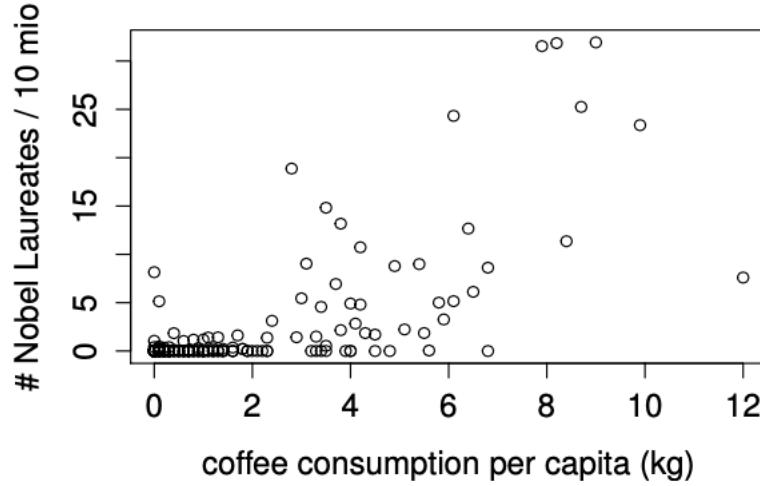


No (not enough) data for chocolate

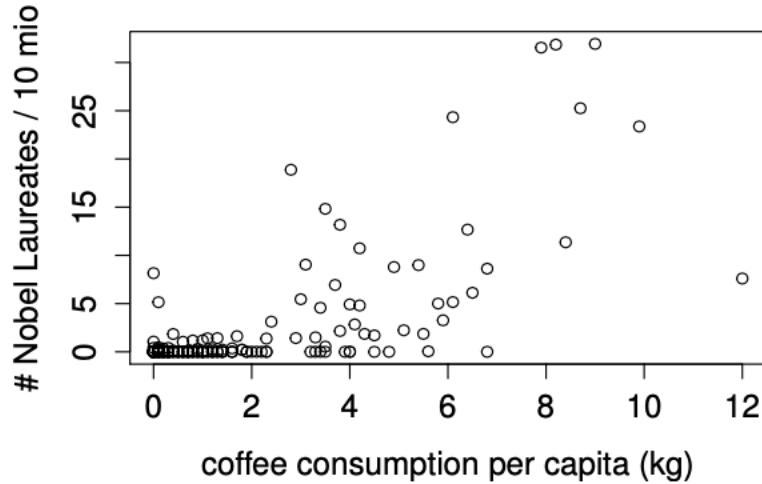


... but we have data for coffee!

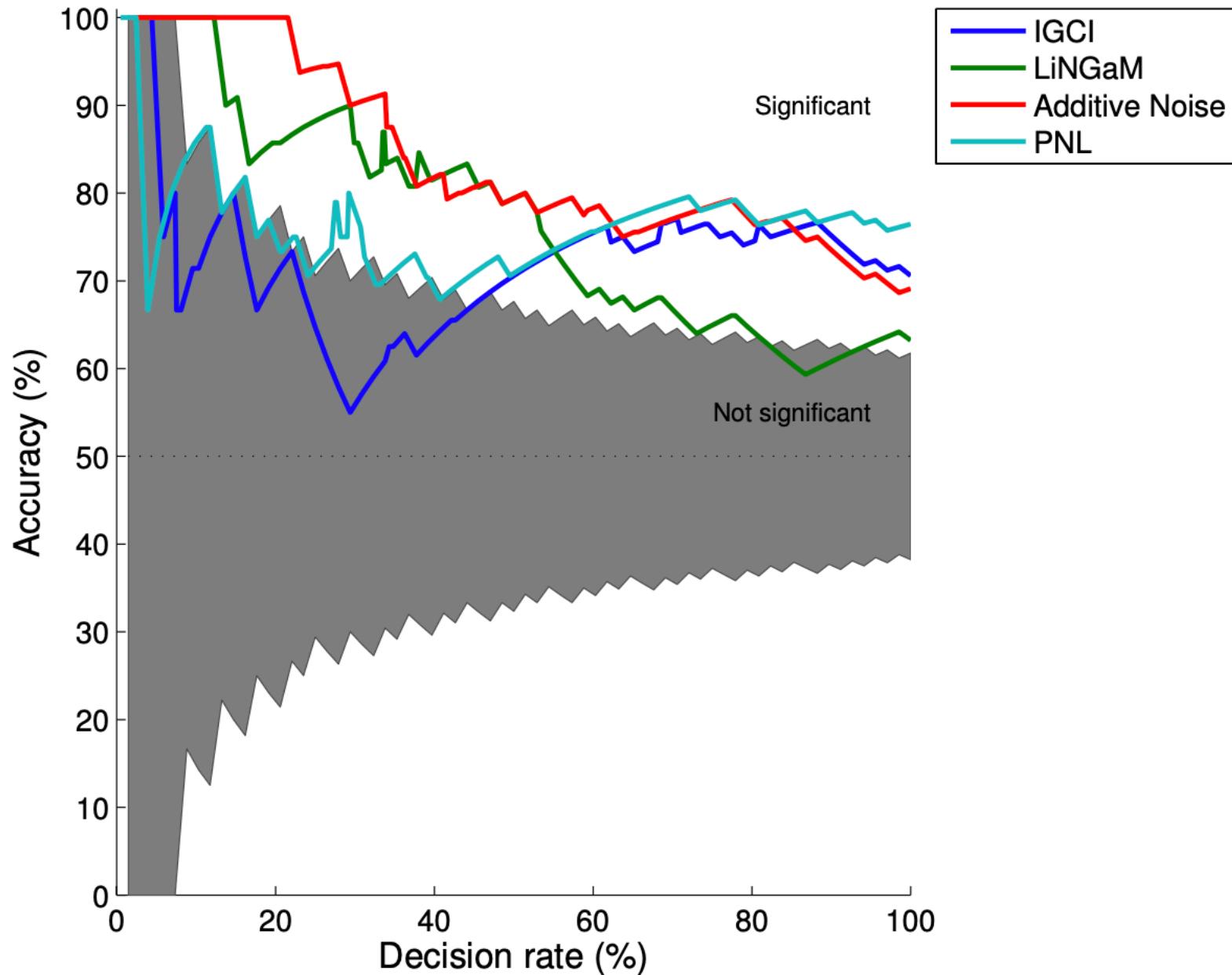
Example: chocolate



Example: chocolate



Real Data: cause-effect pairs



Idea 2: restricted structural causal models

Slightly surprising:

identifiability for two variables \rightsquigarrow identifiability for d variables

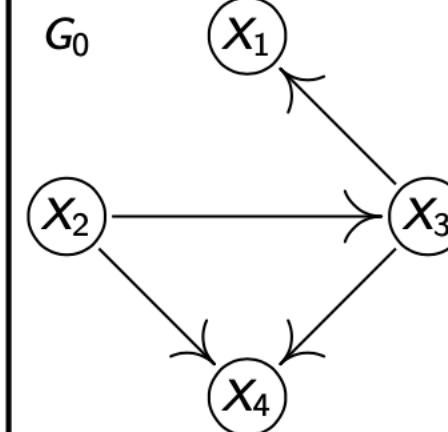
Peters et al.: *Identifiability of Causal Graphs using Functional Models*, UAI 2011

Idea 2: restricted structural causal models

Assume $P(X_1, \dots, X_4)$ has been entailed by

$$\begin{aligned}X_1 &= f_1(X_3, N_1) \\X_2 &= N_2 \\X_3 &= f_3(X_2, N_3) \\X_4 &= f_4(X_2, X_3, N_4)\end{aligned}$$

- N_i jointly independent
- G_0 has no cycles



Structural causal model.

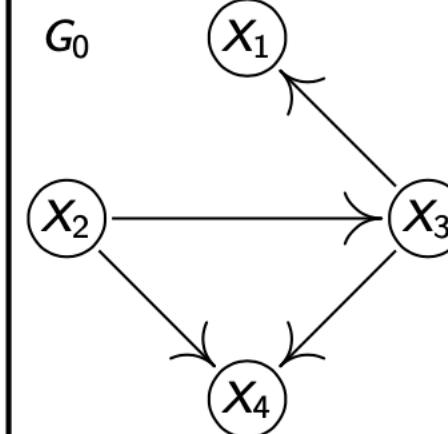
Can the DAG be recovered from $P(X_1, \dots, X_4)$?

Idea 2: restricted structural causal models

Assume $P(X_1, \dots, X_4)$ has been entailed by

$$\begin{aligned}X_1 &= f_1(X_3, N_1) \\X_2 &= N_2 \\X_3 &= f_3(X_2, N_3) \\X_4 &= f_4(X_2, X_3, N_4)\end{aligned}$$

- N_i jointly independent
- G_0 has no cycles



Structural causal model.

Can the DAG be recovered from $P(X_1, \dots, X_4)$?

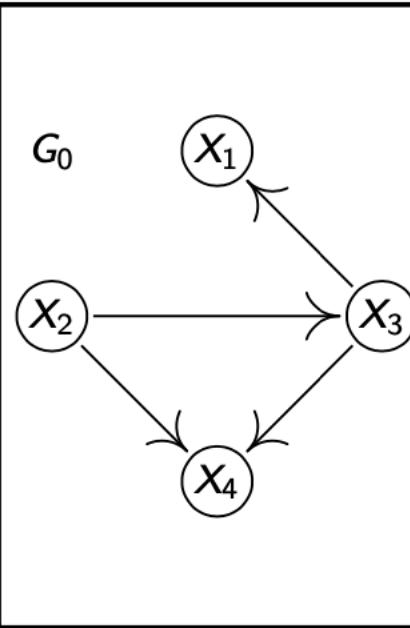
No.

Idea 2: restricted structural causal models

Assume $P(X_1, \dots, X_4)$ has been entailed by

$$\begin{aligned}X_1 &= f_1(X_3) + N_1 \\X_2 &= N_2 \\X_3 &= f_3(X_2) + N_3 \\X_4 &= f_4(X_2, X_3) + N_4\end{aligned}$$

- $N_i \sim \mathcal{N}(0, \sigma_i^2)$ jointly independent
- G_0 has no cycles



Additive noise model with Gaussian noise.

Can the DAG be recovered from $P(X_1, \dots, X_4)$? Yes iff f_i nonlinear.

JP, J. Mooij, D. Janzing and B. Schölkopf: *Causal Discovery with Continuous Additive Noise Models*, JMLR 2014

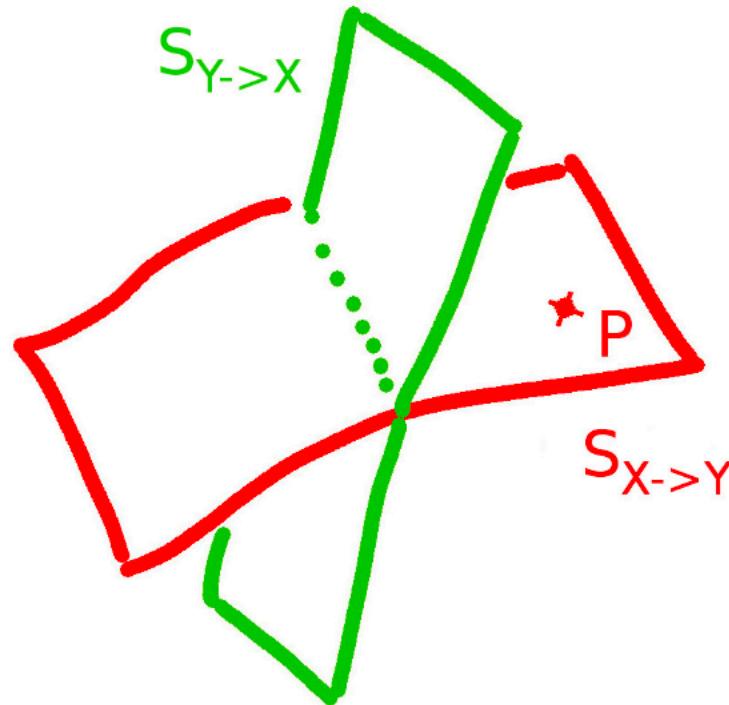
P. Bühlmann, JP, J. Ernest: *CAM: Causal add. models, high-dim. order search and penalized regr.*, Annals of Statistics 2014

Let $P(X_1, \dots, X_d)$ be entailed by an ...

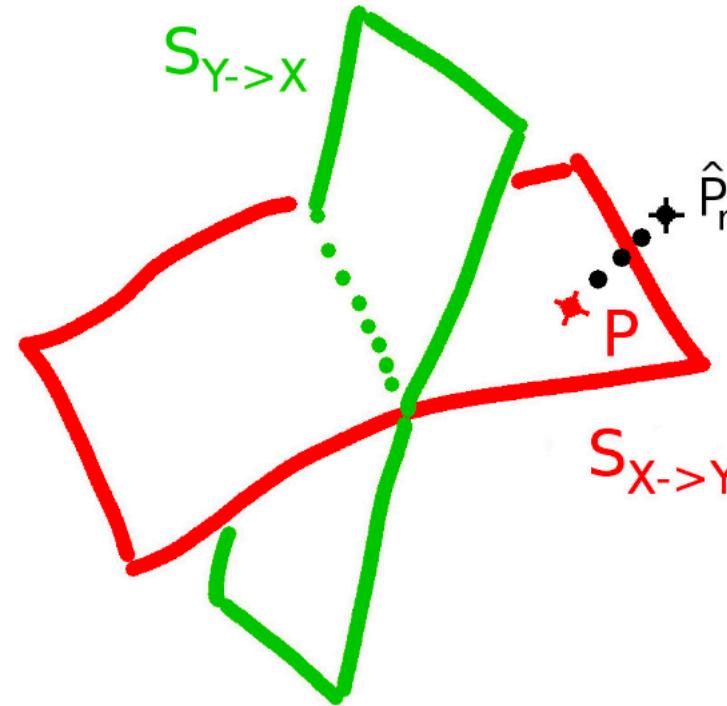
		conditions	identif.
structural causal model:	$X_i = f_i(X_{\text{PA}_i}, N_i)$	-	✗
additive noise model:	$X_i = f_i(X_{\text{PA}_i}) + N_i$	nonlin. fct.	✓
causal additive model:	$X_i = \sum_{k \in \text{PA}_i} f_{ik}(X_k) + N_i$	nonlin. fct.	✓
linear Gaussian model:	$X_i = \sum_{k \in \text{PA}_i} \beta_{ik} X_k + N_i$	linear fct.	✗

(results hold for Gaussian noise)

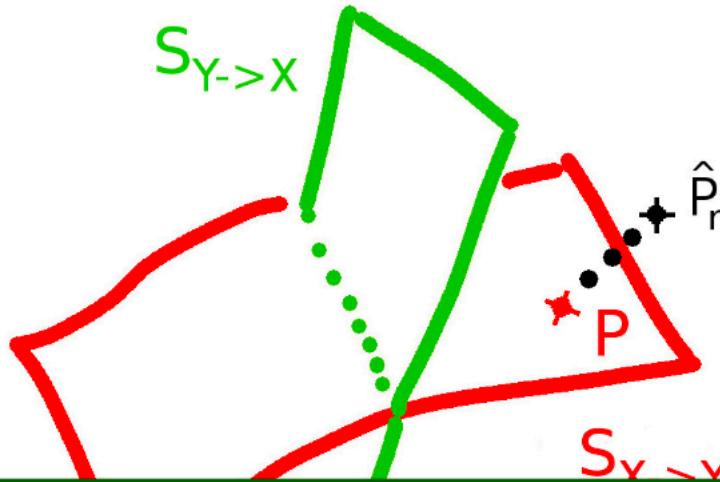
Idea 2: restricted structural causal models



Idea 2: restricted structural causal models



Idea 2: restricted structural causal models



Method: Minimizing KL

Choose the direction that corresponds to the closest subspace...



Idea 2: restricted structural causal models

Consider model classes

$$\mathcal{S}_G := \{Q : Q \text{ entailed by a causal additive model (CAM) with DAG } G\}$$

Define

$$\hat{G}_n := \underset{\text{DAG } G}{\operatorname{argmin}} \inf_{Q \in \mathcal{S}_G} \text{KL}(\hat{P}_n \parallel Q)$$

Idea 2: restricted structural causal models

Consider model classes

$$\mathcal{S}_G := \{Q : Q \text{ entailed by a causal additive model (CAM) with DAG } G\}$$

Define

$$\hat{G}_n := \underset{\substack{\text{DAG } G}}{\operatorname{argmin}} \inf_{Q \in \mathcal{S}_G} \text{KL}(\hat{P}_n \parallel Q)$$

$$\max_{\substack{\text{likelihood}}} \underset{\substack{\text{DAG } G}}{\operatorname{argmin}} \sum_{i=1}^d \log \text{var}(\text{residuals}_{\text{PA}_i^G \rightarrow X_i})$$

Idea 2: restricted structural causal models

Consider model classes

$$\mathcal{S}_G := \{Q : Q \text{ entailed by a causal additive model (CAM) with DAG } G\}$$

Define

$$\hat{G}_n := \underset{\text{DAG } G}{\operatorname{argmin}} \inf_{Q \in \mathcal{S}_G} \text{KL}(\hat{P}_n \parallel Q)$$

$$\max_{\substack{= \\ \text{likelihood}}} \underset{\text{DAG } G}{\operatorname{argmin}} \sum_{i=1}^d \log \text{var}(\text{residuals}_{\text{PA}_i^G \rightarrow X_i})$$

Wait, there is no penalization on the number of edges!

Idea 2: restricted structural causal models

Consider model classes

$$\mathcal{S}_G := \{Q : Q \text{ entailed by a causal additive model (CAM) with DAG } G\}$$

Define

$$\hat{G}_n := \underset{\substack{\text{DAG } G}}{\operatorname{argmin}} \inf_{Q \in \mathcal{S}_G} \text{KL}(\hat{P}_n \parallel Q)$$

$$\max_{\substack{\text{likelihood}}} \underset{\substack{\text{DAG } G}}{\operatorname{argmin}} \sum_{i=1}^d \log \text{var}(\text{residuals}_{\text{PA}_i^G \rightarrow X_i})$$

Wait, there is no penalization on the number of edges!

Wait again, there are too many DAGs!

Idea 2: restricted structural causal models

p		number of DAGs with p nodes
1		1
2		3
3		25
4		543
5		29281
6		3781503
7		1138779265
8		783702329343
9		1213442454842881
10		4175098976430598143
11		31603459396418917607425
12		521939651343829405020504063
13		18676600744432035186664816926721
14		1439428141044398334941790719839535103
15		237725265553410354992180218286376719253505
16		83756670773733320287699303047996412235223138303
17		62707921196923889899446452602494921906963551482675201
18		99421195322159515895228914592354524516555026878588305014783
19		332771901227107591736177573311261125883583076258421902583546773505
20		2344880451051088988152559855229099188899081192234291298795803236068491263
21		3469876828358875002875932843018108822313944540438601719027559113446586077675521
22		1075822921725761493652956179327624326573727662809185218104090000500559527511693495107583
23		69743329837281492647141549700245804876504274990515985894109106401549811985510951501377122074625

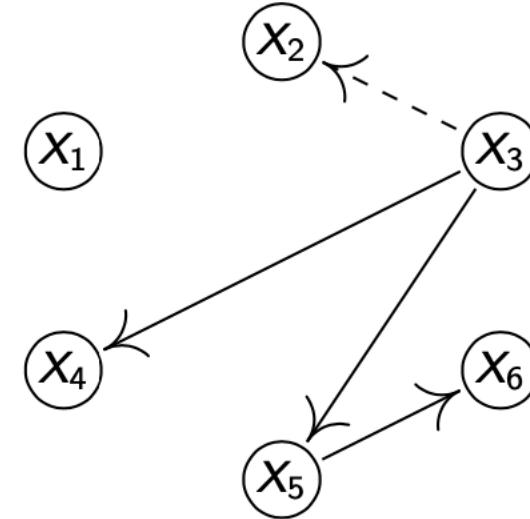
<https://oeis.org/A003024/b003024.txt>

Idea 2: restricted structural causal models

E.g. greedy search!

-	0.2	0.1	0.1	0.1	0.3
0.4	-	0.1	0.1	0.1	0.1
0.1	0.6	-	-	-	0.4
0.1	0.1	-	-	0.1	0.1
0.1	0.1	-	0.1	-	-
0.3	0.1	-	0.1	-	-

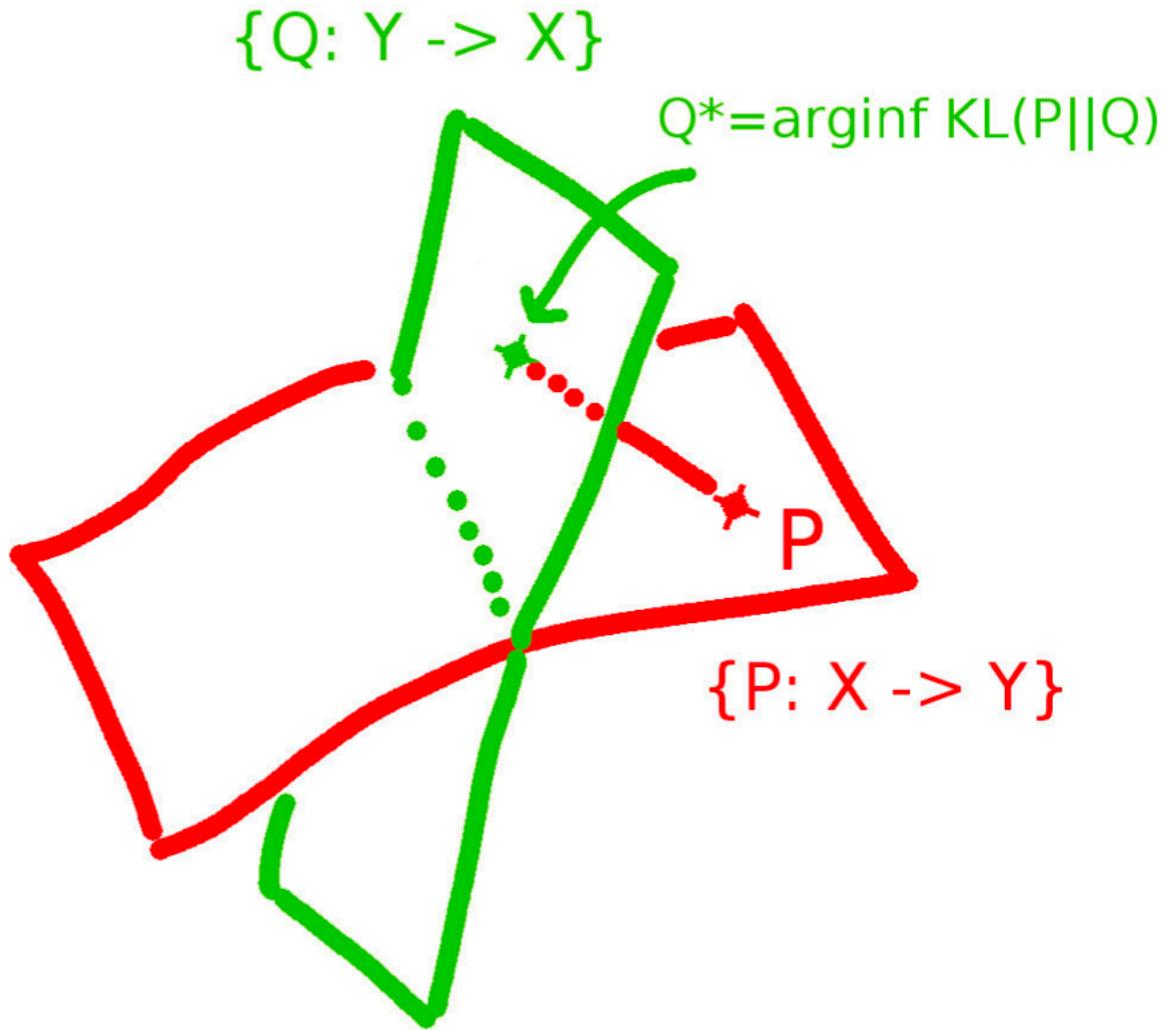
include best edge
→
recompute column



Greedy Addition (e.g. Chickering 2002). Include the edge that leads to the largest increase of the log-likelihood.

Bühlmann, JP, Ernest: *CAM: Causal add. models, high-dim. order search and penalized regr.*, Annals of Statistics 2014

Can we characterize identifiability?



Can we characterize identifiability?

Proposition

Assume $P(X, Y)$ is generated by

$$Y = \beta X^2 + N_Y$$

with independent $X \sim \mathcal{N}(0, \sigma_X^2)$ and $N_Y \sim \mathcal{N}(0, \sigma_{N_Y}^2)$.

Can we characterize identifiability?

Proposition

Assume $P(X, Y)$ is generated by

$$Y = \beta X^2 + N_Y$$

with independent $X \sim \mathcal{N}(0, \sigma_X^2)$ and $N_Y \sim \mathcal{N}(0, \sigma_{N_Y}^2)$.

Then

$$\inf_{Q \in \{Q: Y \rightarrow X\}} \text{KL}(P \parallel Q) > 0 \quad \text{if } \beta \neq 0.$$

Can we characterize identifiability?

Proposition

Assume $P(X, Y)$ is generated by

$$Y = \beta X^2 + N_Y$$

with independent $X \sim \mathcal{N}(0, \sigma_X^2)$ and $N_Y \sim \mathcal{N}(0, \sigma_{N_Y}^2)$.

Then

$$\inf_{Q \in \{Q: Y \rightarrow X\}} \text{KL}(P \parallel Q) = \frac{1}{2} \log \left(1 + 2\beta^2 \frac{\sigma_X^4}{\sigma_{N_Y}^2} \right)$$



Leonardo da Vinci: Mould of the Horses Head

Given an original **drawing** (left) and a copy. How good is the copy?



Leonardo da Vinci: Mould of the Horses Head

Given an original **drawing** (left) and a copy. How good is the copy?

Given a true **causal graph** G and an estimate \hat{G} . How good is the estimate \hat{G} ?

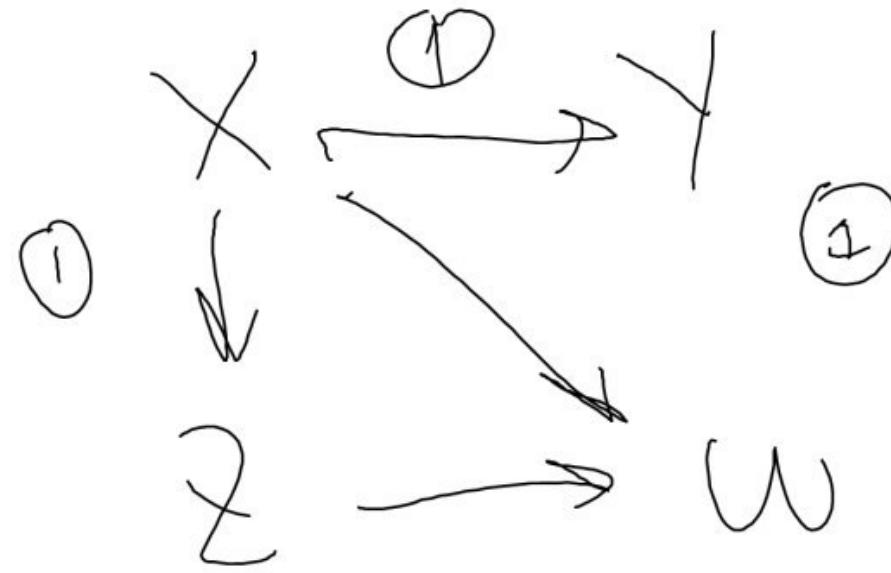
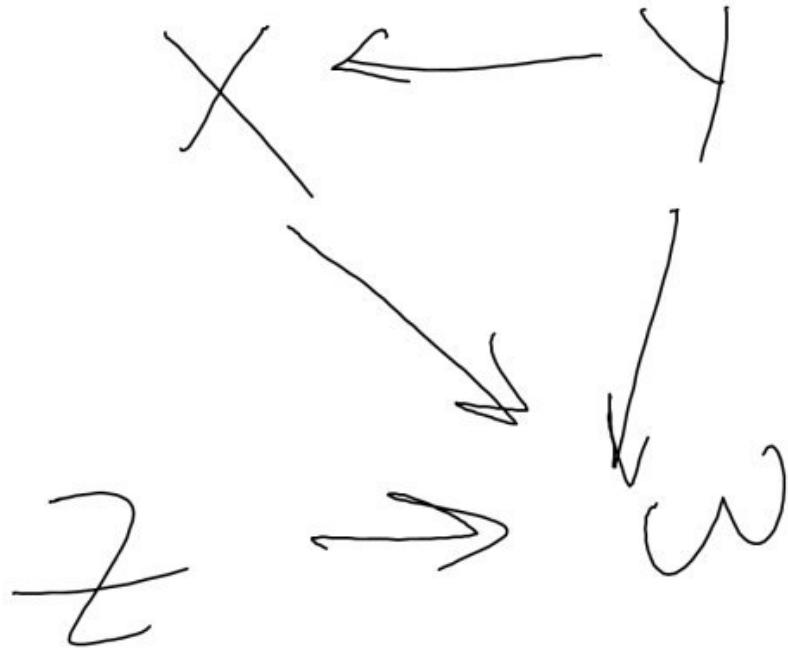


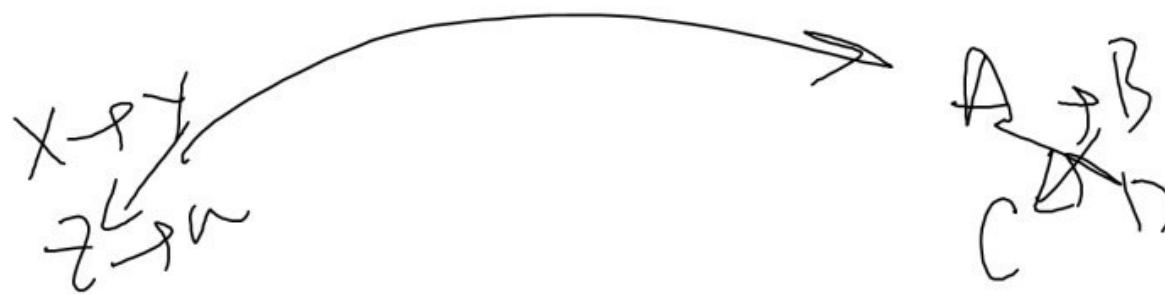
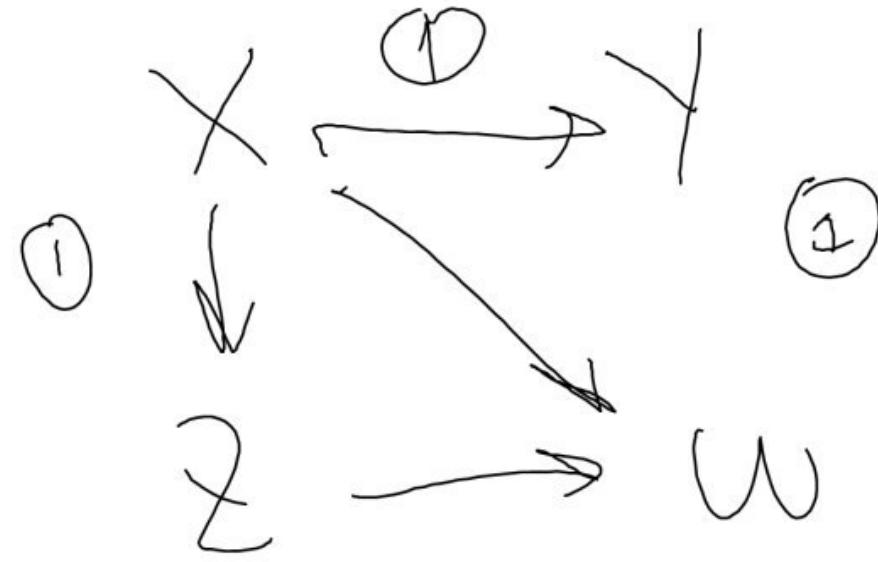
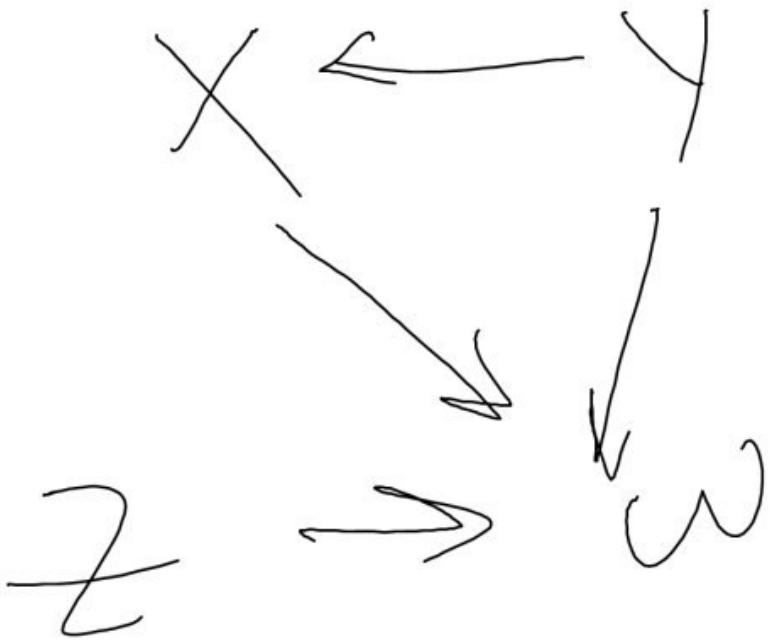
Leonardo da Vinci: Mould of the Horses Head

Given an original **drawing** (left) and a copy. How good is the copy?

Given a true **causal graph** G and an estimate \hat{G} . How good is the estimate \hat{G} ?

What do we want do with it?





Definition: Structural Intervention Distance

For each pair (X, Y) check whether $\text{PA}_X^{\hat{G}}$ is a valid adjustment set for (X, Y) in G **for all distributions** Markov w.r.t. G .

Definition: Structural Intervention Distance

For each pair (X, Y) check whether $\text{PA}_X^{\hat{G}}$ is a valid adjustment set for (X, Y) in G **for all distributions Markov w.r.t. G .** - does not depend on P

$\text{SID}(G, \hat{G})$ equals the number of pairs, for which this is not the case.

Definition: Structural Intervention Distance

For each pair (X, Y) check whether $\text{PA}_X^{\hat{G}}$ is a valid adjustment set for (X, Y) in G **for all distributions Markov w.r.t. G .** - does not depend on P

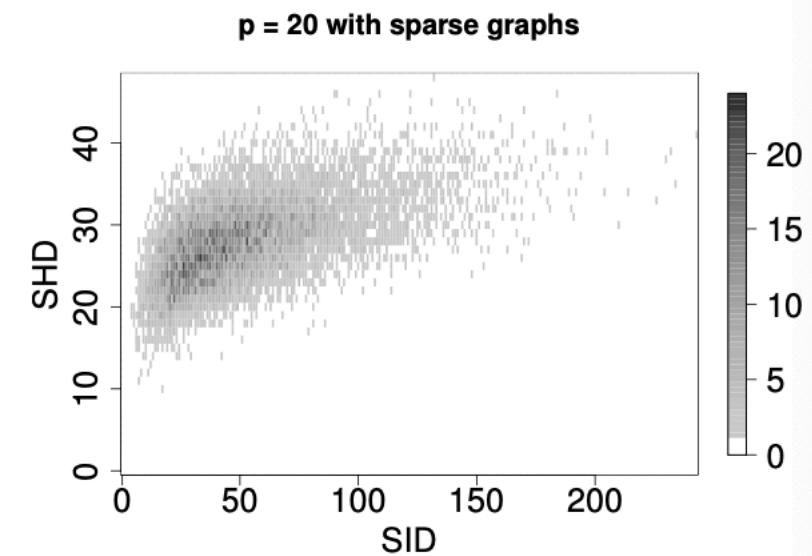
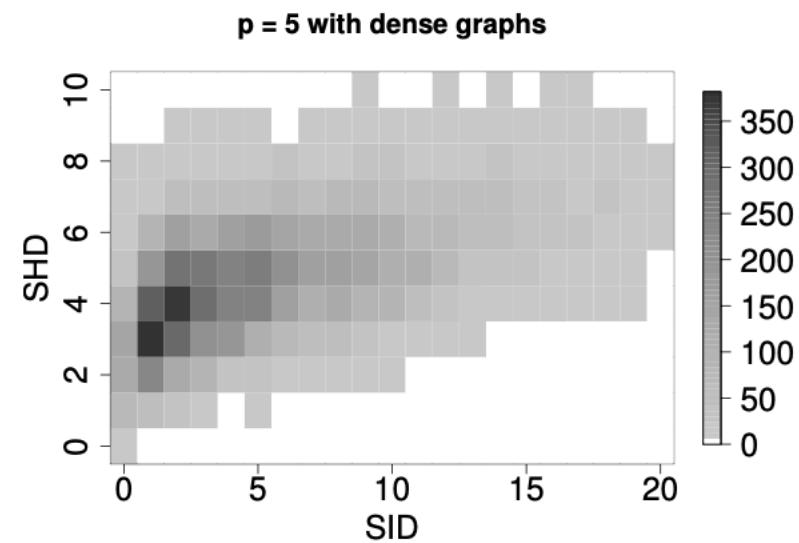
$\text{SID}(G, \hat{G})$ equals the number of pairs, for which this is not the case.

Graphical representation of the SID available!!

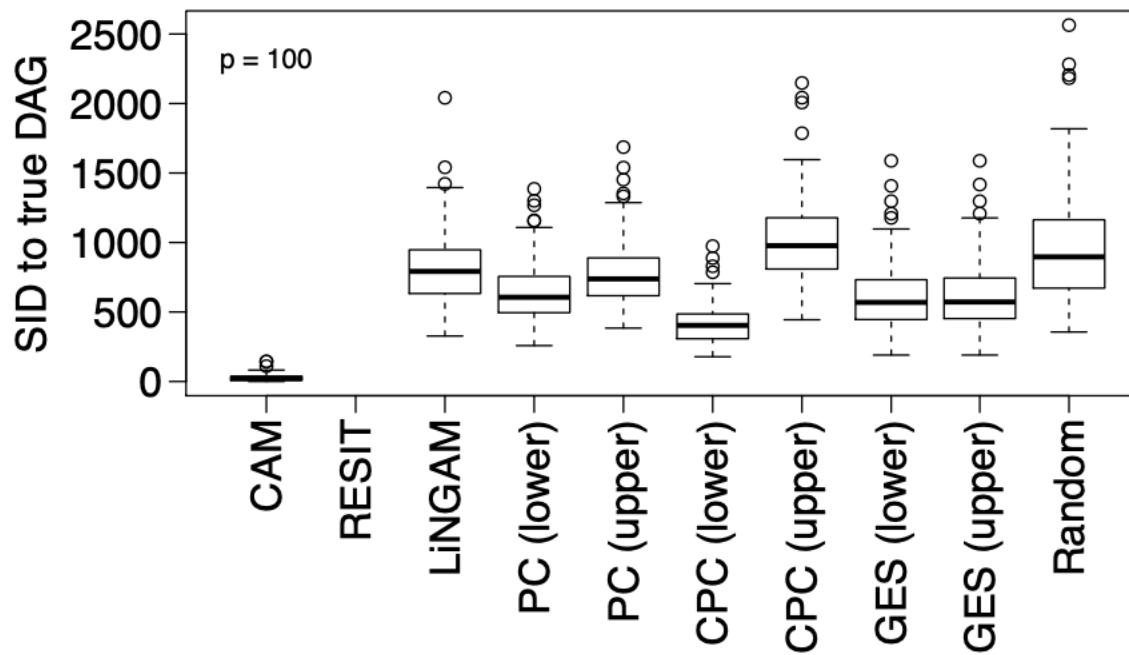
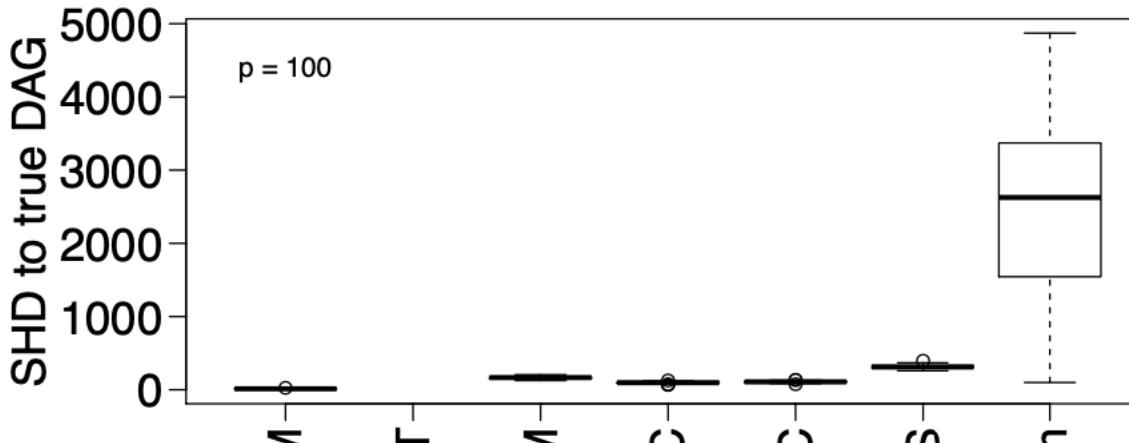
based on Shpitser et al: "On the validity of covariate adjustment for estimating causal effects", UAI 2010.

SHD and SID are quite different!

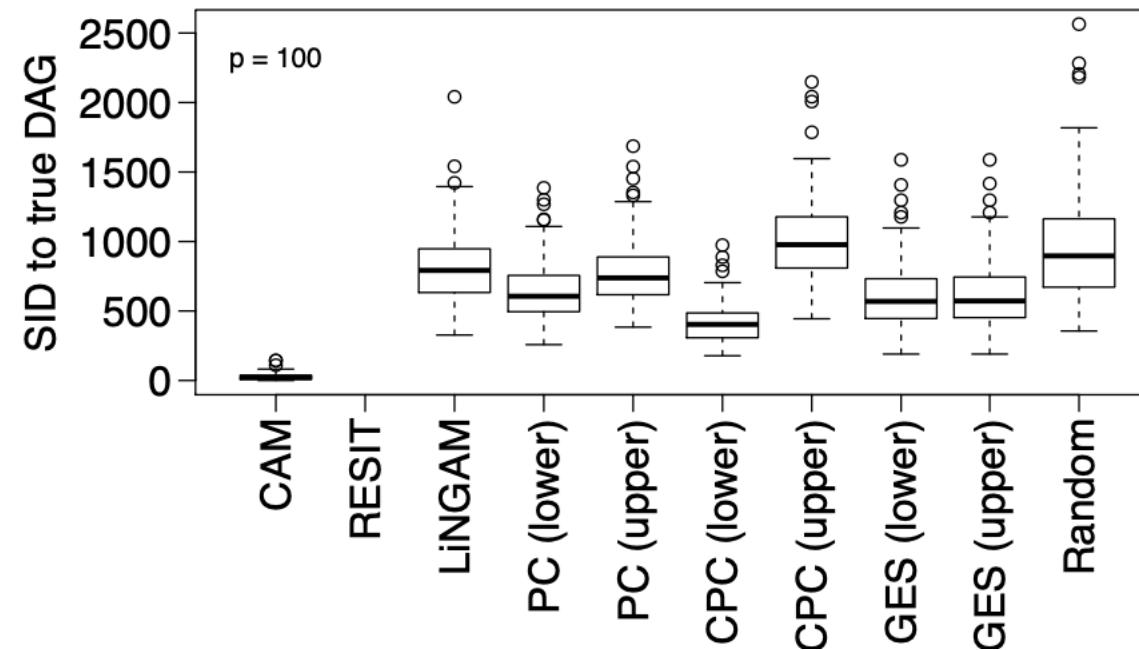
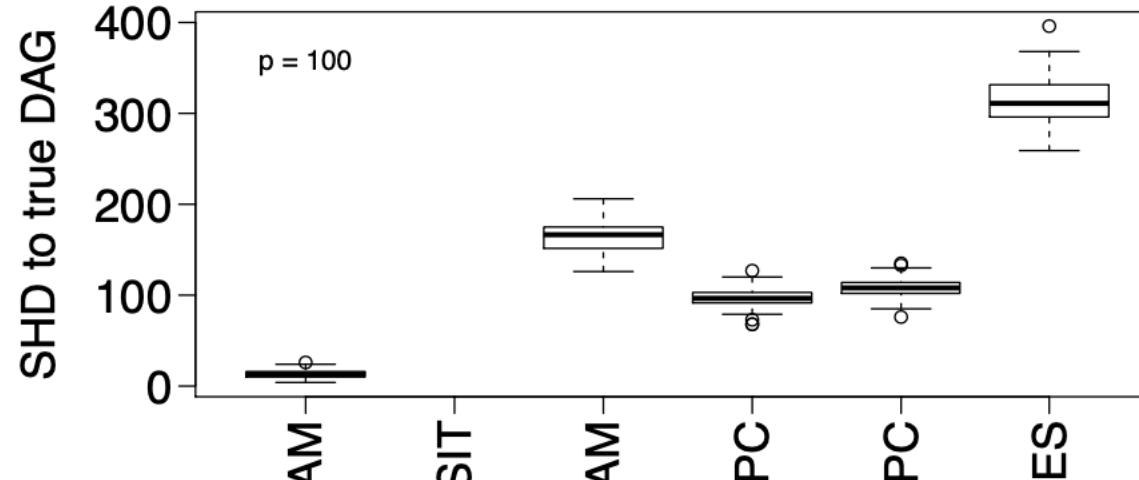
10,000 random DAGs



SHD and SID may lead to different conclusions (nonlinear, Gaussian).



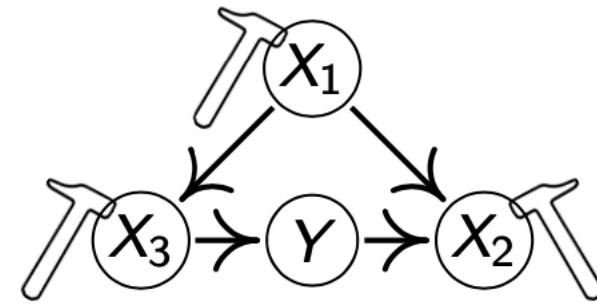
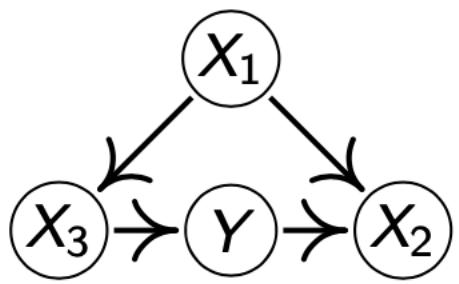
SHD and SID may lead to different conclusions (nonlinear, Gaussian).



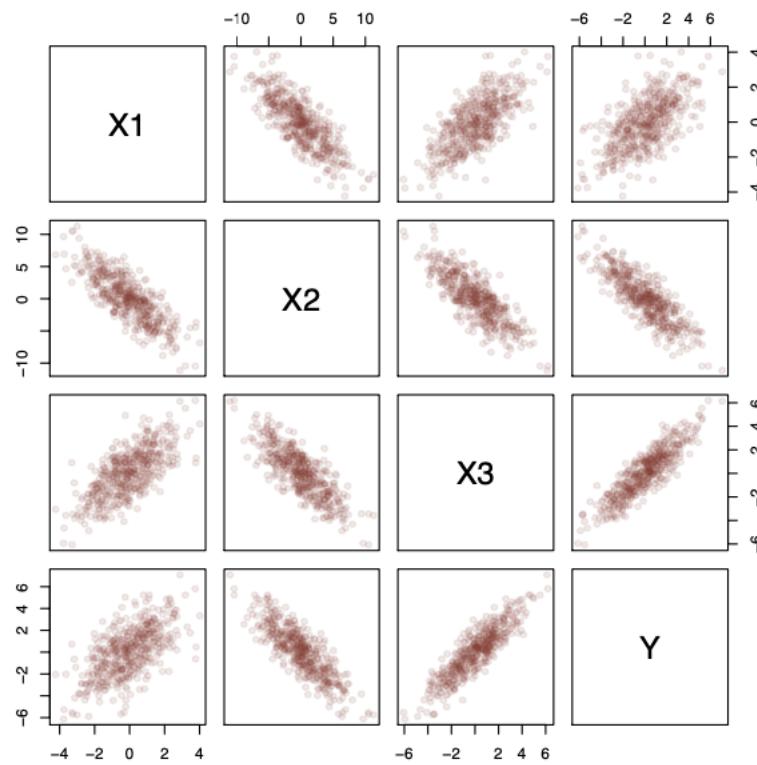
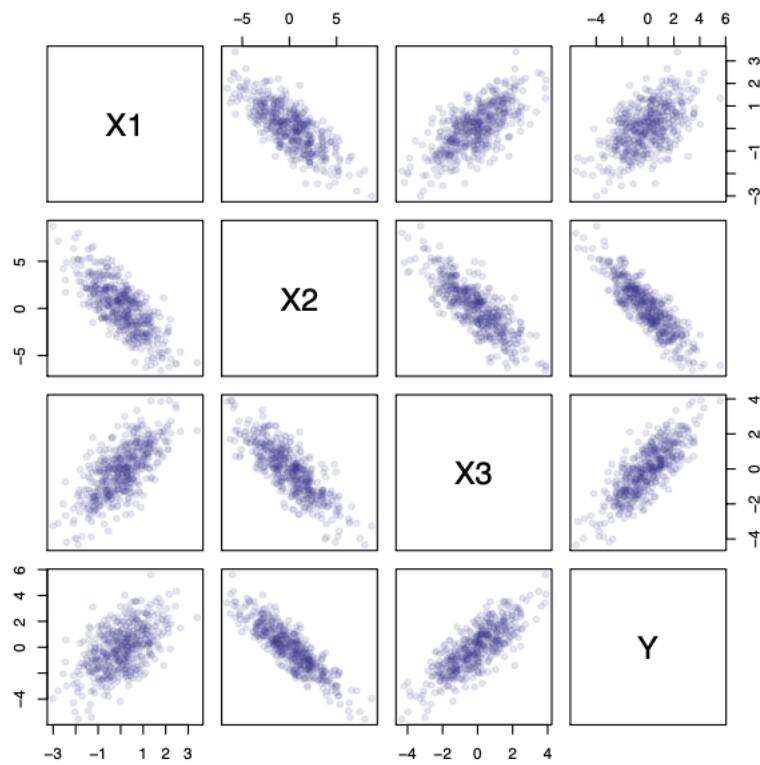
Idea 3: invariant causal prediction

Concentrate on one target variable.

unknown:



known:



linear model

```
> linmod <- lm(Y~X)
> summary(linmod)
```

Coefficients:

	Estimate	Std. Error	t value	Pr(> t)
(Intercept)	0.000322	0.025858	0.012	0.99
X1	-0.444534	0.034306	-12.958	<2e-16 ***
X2	-0.402398	0.016471	-24.430	<2e-16 ***
X3	0.603502	0.025642	23.536	<2e-16 ***

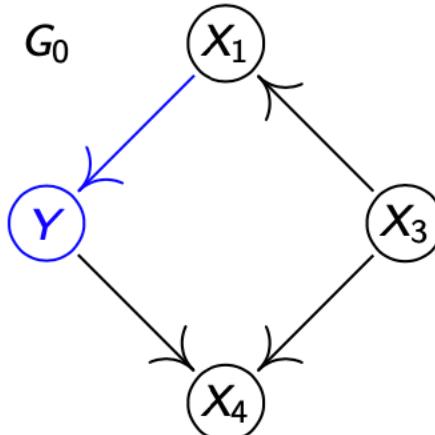
ICP (R-package InvariantCausalPrediction)

Key idea: MUTE ... or:

$P(Y | \text{PA}_Y)$ remains invariant if the struct. equ. for Y does not change.

$$\begin{aligned}X_1 &:= f_1(X_3, N_1) \\Y &:= f_2(X_1, N_2) \\X_3 &:= f_3(N_3) \\X_4 &:= f_4(Y, X_3, N_4)\end{aligned}$$

- N_i jointly independent
- G_0 has no cycles



IMPORTANT: modularity, autonomy

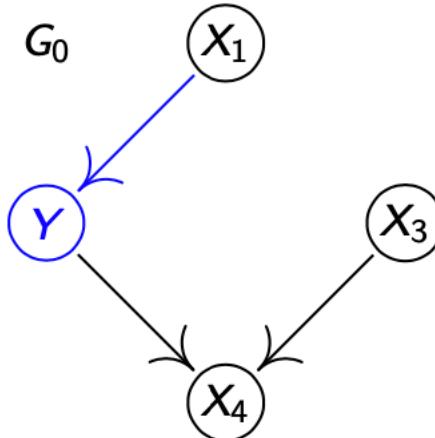
Haavelmo 1944, Aldrich 1989, Pearl 2009, Schölkopf et al. 2012, Bareinboim et al. 2014, Hauser et al. 2014, ...

Key idea:

$P(Y | \text{PA}_Y)$ remains invariant if the struct. equ. for Y does not change.

$$\begin{aligned}X_1 &:= \tilde{f}_1(\tilde{N}_1) \\Y &:= f_2(X_1, N_2) \\X_3 &:= f_3(N_3) \\X_4 &:= f_4(Y, X_3, N_4)\end{aligned}$$

- N_i jointly independent
- G_0 has no cycles



IMPORTANT: modularity, autonomy

Haavelmo 1944, Aldrich 1989, Pearl 2009, Schölkopf et al. 2012, Bareinboim et al. 2014, Hauser et al. 2014, ...

Key idea:

$P(Y | \text{PA}_Y)$ remains invariant if the struct. equ. for Y does not change.

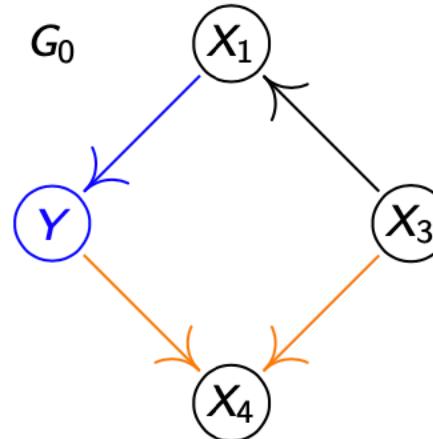
$$X_1 := f_1(X_3, N_1)$$

$$Y := f_2(X_1, N_2)$$

$$X_3 := f_3(N_3)$$

$$X_4 := \tilde{f}_4(Y, X_3, \tilde{N}_4)$$

- N_i jointly independent
- G_0 has no cycles



IMPORTANT: modularity, autonomy

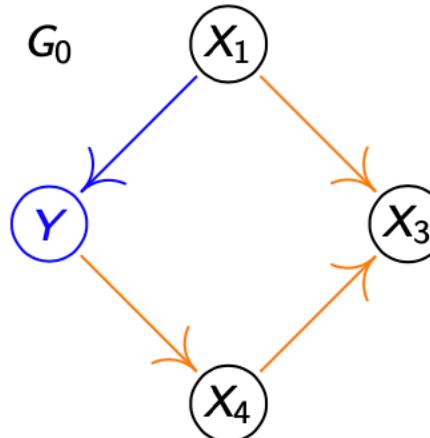
Haavelmo 1944, Aldrich 1989, Pearl 2009, Schölkopf et al. 2012, Bareinboim et al. 2014, Hauser et al. 2014, ...

Key idea:

$P(Y | \mathbf{PA}_Y)$ remains invariant if the struct. equ. for Y does not change.

$$\begin{aligned}X_1 &:= \tilde{f}_1(\tilde{N}_1) \\Y &:= f_2(X_1, N_2) \\X_3 &:= \tilde{f}_3(X_1, X_4, \tilde{N}_3) \\X_4 &:= \tilde{f}_4(Y, \tilde{N}_4)\end{aligned}$$

- N_i jointly independent
- G_0 has no cycles



IMPORTANT: modularity, autonomy

Haavelmo 1944, Aldrich 1989, Pearl 2009, Schölkopf et al. 2012, Bareinboim et al. 2014, Hauser et al. 2014, ...

Given: Data from different environments $e \in \mathcal{E}$, e.g. interventions.

Given: Data from different environments $e \in \mathcal{E}$, e.g. interventions.

Proposition

Let $S^* = \mathbf{PA}_Y$. Then, $H_{0,S^*}(\mathcal{E})$ is true, i.e.,

for all $e \in \mathcal{E}$: X^e has an arbitrary distribution and
 $Y^e | X_{S^*}^e = x$ invariant.

Given: Data from different environments $e \in \mathcal{E}$, e.g. interventions.

Proposition

Let $S^* = \mathbf{PA}_Y$. Then, $H_{0,S^*}(\mathcal{E})$ is true, i.e., there exists γ^* with support S^* that satisfies

for all $e \in \mathcal{E}$: X^e has an arbitrary distribution and

~~$Y^e | X_{S^*}^e = x$ invariant.~~

$$Y^e = X^e \gamma^* + \varepsilon^e, \quad \varepsilon^e \sim F_\varepsilon \text{ and } \varepsilon^e \perp\!\!\!\perp X_{S^*}^e.$$

Given: Data from different environments $e \in \mathcal{E}$, e.g. interventions.

Proposition

Let $S^* = \mathbf{PA}_Y$. Then, $H_{0,S^*}(\mathcal{E})$ is true, i.e., there exists γ^* with support S^* that satisfies

for all $e \in \mathcal{E}$: X^e has an arbitrary distribution and

$\cancel{Y^e | X_{S^*}^e = x}$ invariant.

$$Y^e = X^e \gamma^* + \varepsilon^e, \quad \varepsilon^e \sim F_\varepsilon \text{ and } \varepsilon^e \perp\!\!\!\perp X_{S^*}^e.$$

Goal: Find S^* .

Idea: Check $H_{0,S}(\mathcal{E})$ for several candidates S .

$$S(\mathcal{E}) := \bigcap_{S : H_{0,S}(\mathcal{E}) \text{ is true}} S$$

set	$\{3, 5\}$	$\{3, 7\}$	$S^* = \{1, 3, 6\}$	$\{2\}$	$\{3, 8\}$	\dots
inv. pred.	✓	✗	✓	✗	✓	...

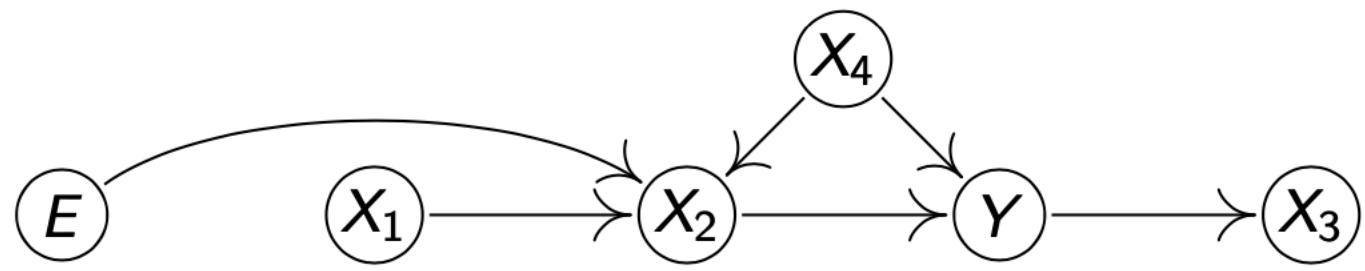
Theorem (PBM 2016)

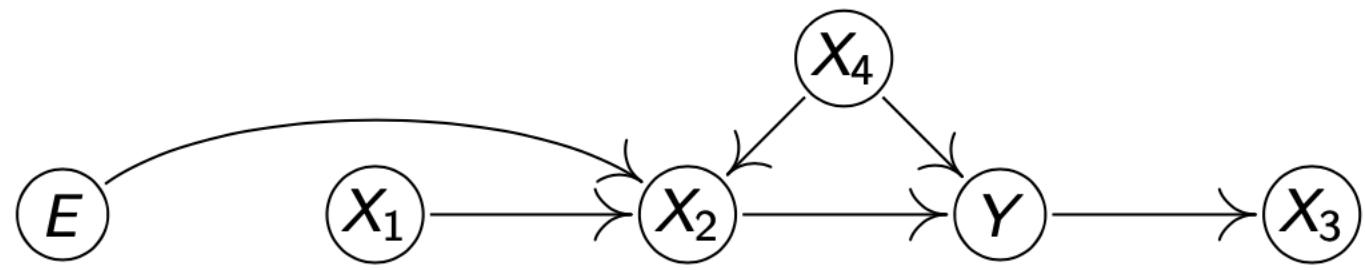
$$P(\hat{S}(\mathcal{E}) \subseteq S^*) \geq 1 - \alpha .$$

Theorem (PBM 2016)

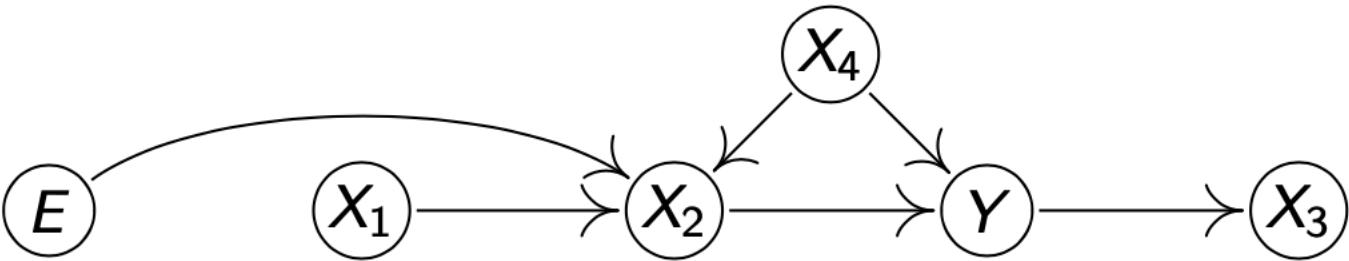
$$P(\hat{S}(\mathcal{E}) \subseteq S^*) \geq 1 - \alpha.$$

Identifiability improves if we have more and stronger interventions, at better places, more heterogeneity in the data.





```
> Y <- X[,2] + X[,4] + noise  
> ICP(X,Y,ExpInd)
```



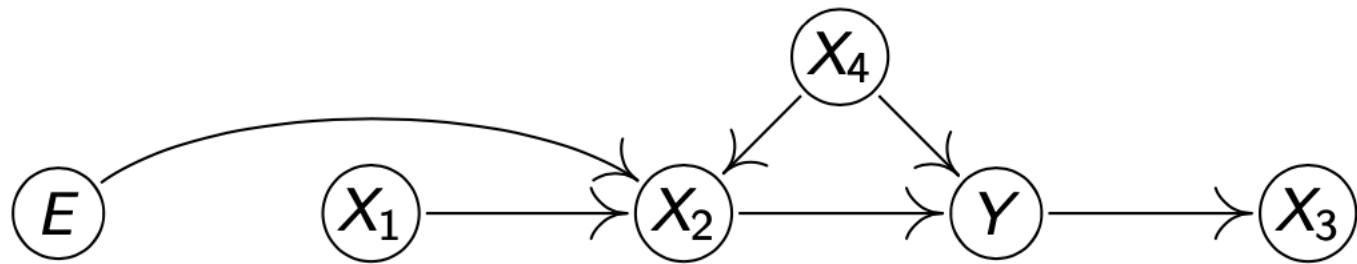
```

> Y <- X[,2] + X[,4] + noise
> ICP(X,Y,ExpInd)

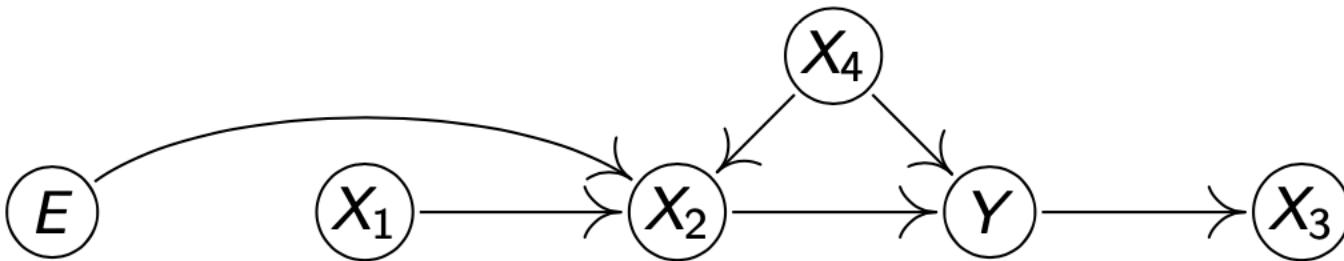
```

accepted set of variables: 2,4
 accepted set of variables: 1,2,4
 accepted set of variables: 2,3,4
 accepted set of variables: 1,2,3,4

	LOWER BOUND	UPPER BOUND	MAXIMIN	EFFECT	P-VALUE
X1	-0.03	0.01		0.00	0.48
X2	0.98	1.01	0.98	< 1e-09 ***	
X3	-0.07	0.00	0.00		0.48
X4	0.95	1.01	0.95	2.6e-05 ***	



```
> Y <- X[,2]^2 + X[,4] + noise  
> ICP(X,Y,ExpInd)
```

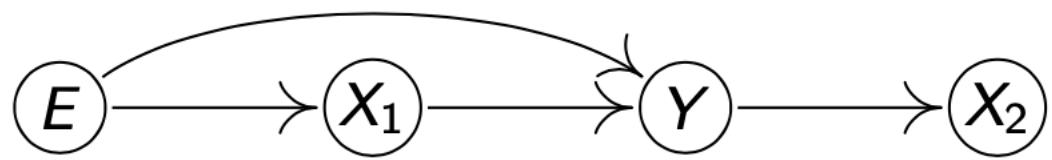


```
> Y <- X[,2]^2 + X[,4] + noise  
> ICP(X, Y, ExpInd)
```

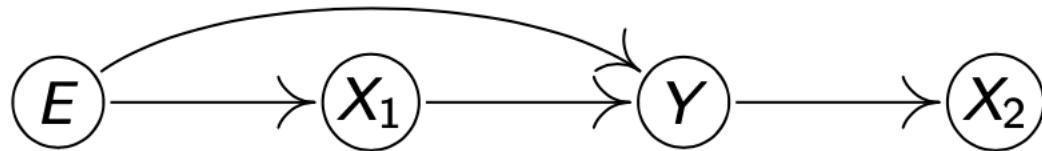
empty set
(all models rejected)

Model violation: nonlinear models

~> usually leads to loss of power, not coverage



```
> Y <- X[,1] + E + noise  
> ICP(X,Y,ExpInd)
```

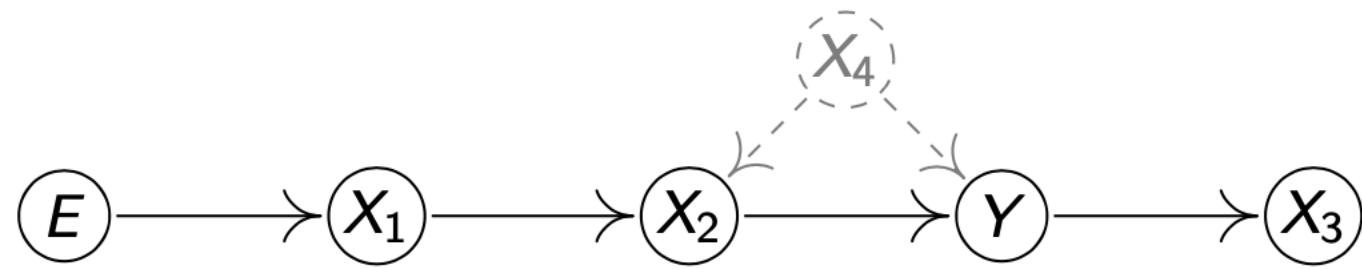


```
> Y <- X[,1] + E + noise  
> ICP(X,Y,ExpInd)
```

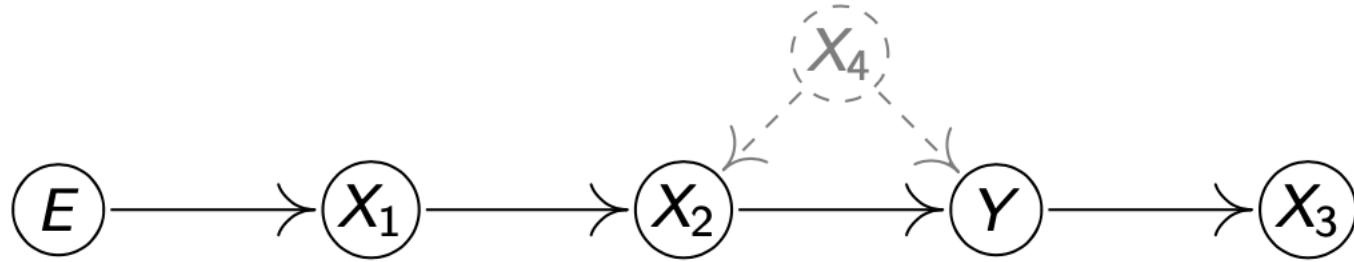
empty set
(all models rejected)

Model violation: intervention on Y

~~ usually leads to loss of power, not coverage



```
> Y <- X[,2] + X[,4] + noise  
> ICP(X[,1:3], Y, ExpInd)
```



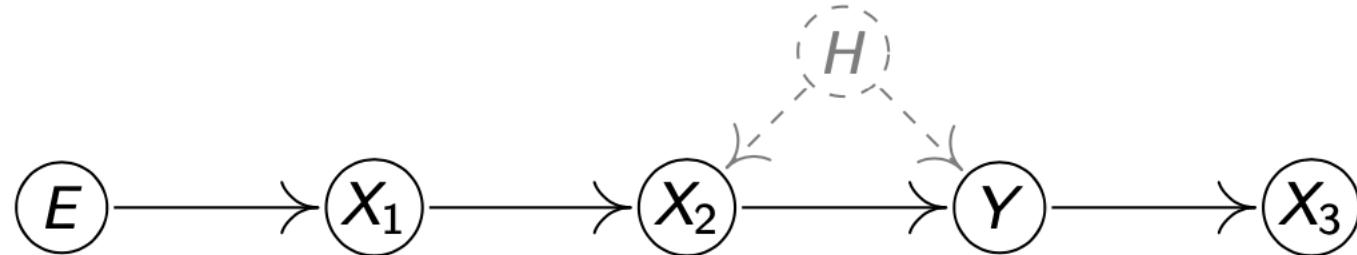
```
> Y <- X[,2] + X[,4] + noise
> ICP(X[,1:3],Y,ExpInd)
```

accepted set of variables: 1
 accepted set of variables: 1,2
 accepted set of variables: 1,3
 accepted set of variables: 1,2,3

	LOWER BOUND	UPPER BOUND	MAXIMIN EFFECT	P-VALUE
X1	-0.87	1.05	0.00	<1e-09 ***
X2	0.00	1.86	0.00	1.00
X3	-1.61	0.00	0.00	0.73

Model violation: hidden variables

~~ coverage still holds if we consider ancestors instead of parents



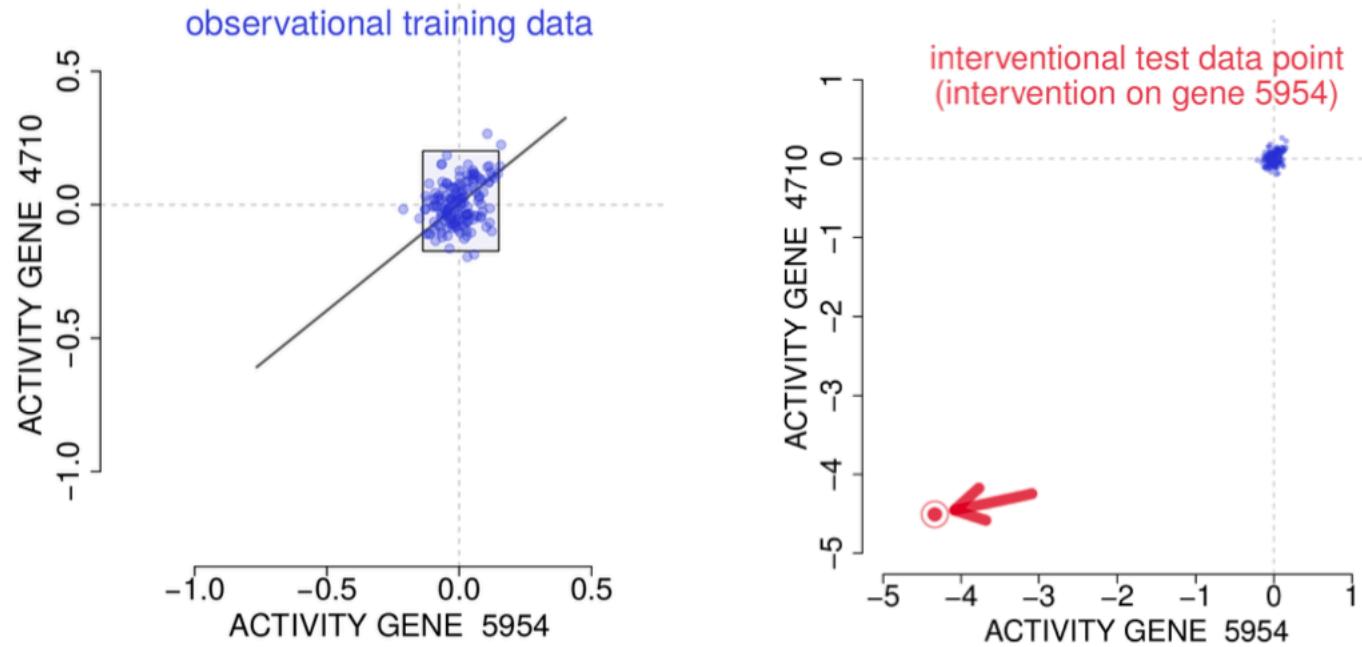
Theorem (PBM 2016)

Assume that the joint distribution over $(Y, X_1, \dots, X_p, H_1, \dots, H_q, E)$ is faithful w.r.t. the augmented graph. Then

$$S(\mathcal{E}) := \bigcap_{\substack{S : H_{0,S}(\mathcal{E}) \text{ is true}}} S \subseteq \mathbf{AN}(Y) \cap \{X_1, \dots, X_p\}.$$

Real data: genetic perturbation experiments for yeast (Kemmeren et al., 2014)

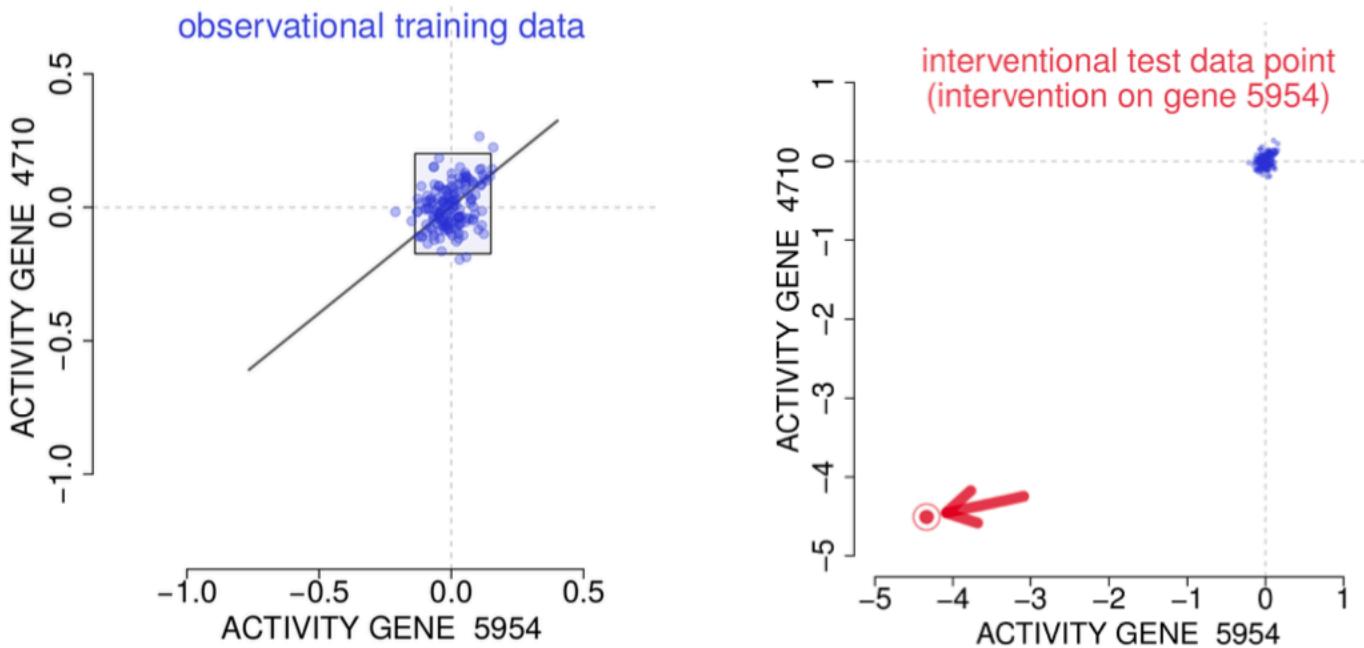
- $p = 6170$ genes
- $n_{obs} = 160$ wild-types
- $n_{int} = 1479$ gene deletions (targets known)



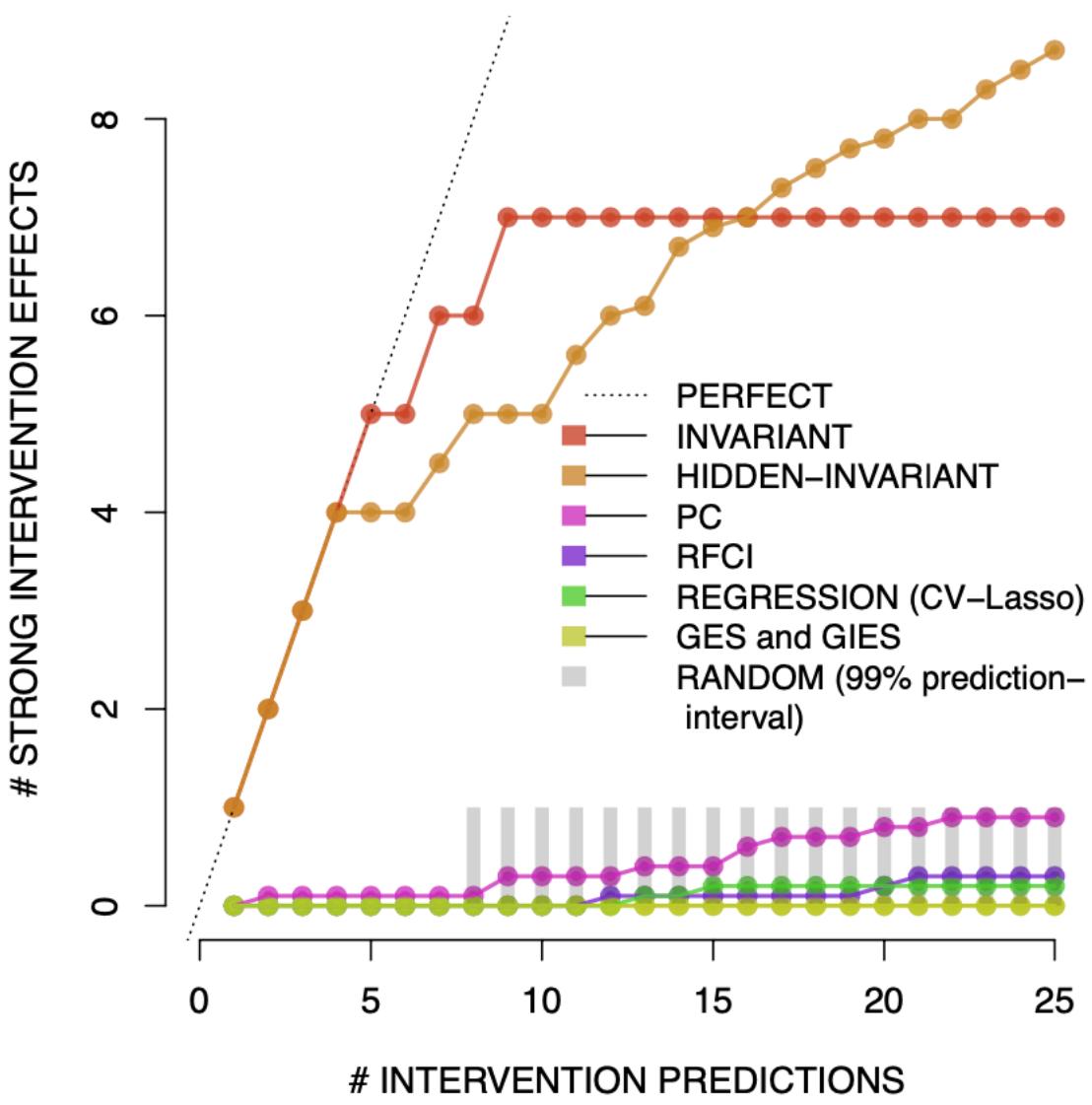
- true hits: $\approx 0.1\%$ of pairs

Real data: genetic perturbation experiments for yeast (Kemmeren et al., 2014)

- $p = 6170$ genes
- $n_{obs} = 160$ wild-types
- $n_{int} = 1479$ gene deletions (targets known)

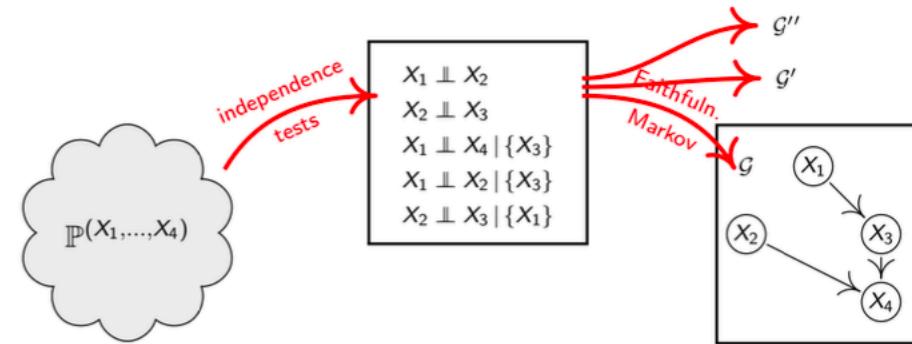


- true hits: $\approx 0.1\%$ of pairs
- our method: $\mathcal{E} = \{obs, int\}$



Summary Part II:

- Idea 1: independence-based methods (single environment)



- Idea 2: additive noise (single environment)

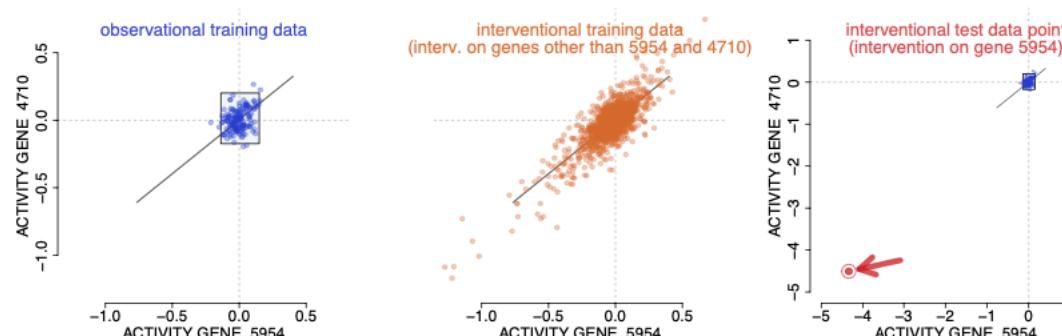
$$X_1 = f_1(X_3) + N_1$$

$$X_2 = N_2$$

$$X_3 = f_3(X_2) + N_3$$

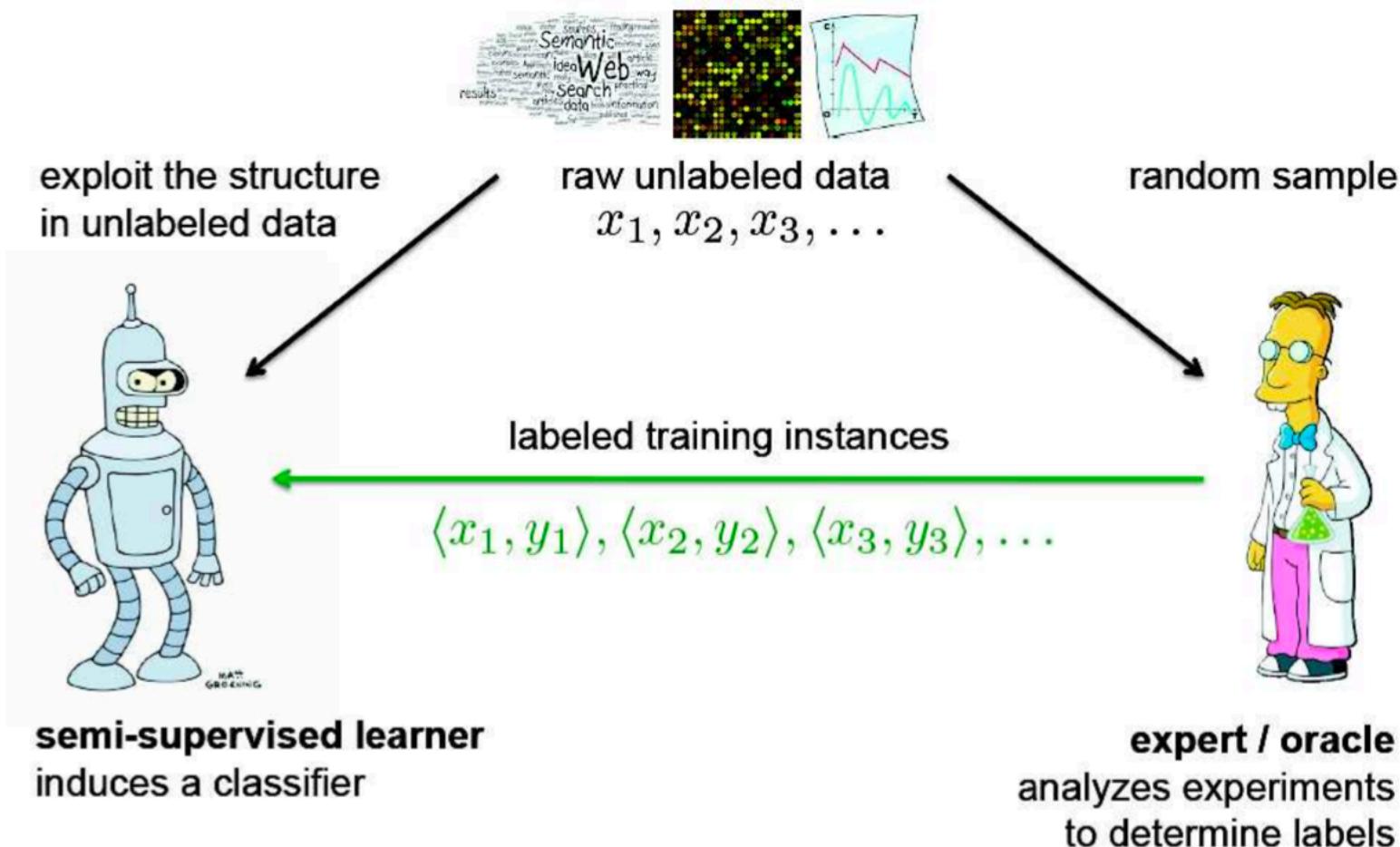
$$X_4 = f_4(X_2, X_3) + N_4$$

- Idea 3: invariant prediction (the more heterogeneity the better!)

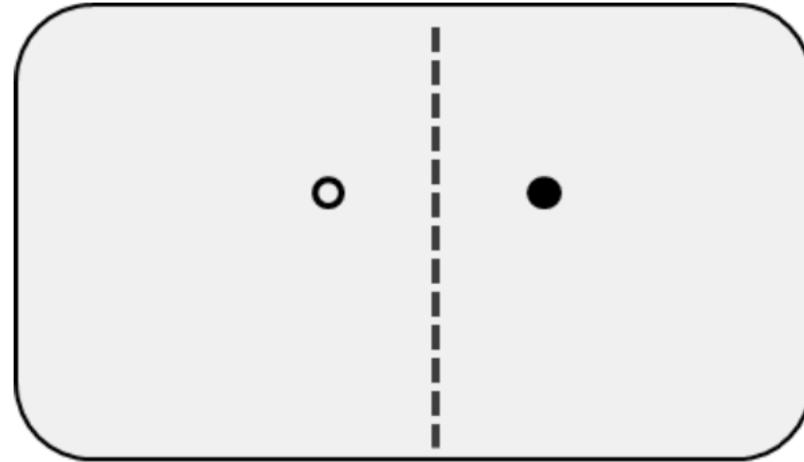


Part III: Applications to Machine Learning

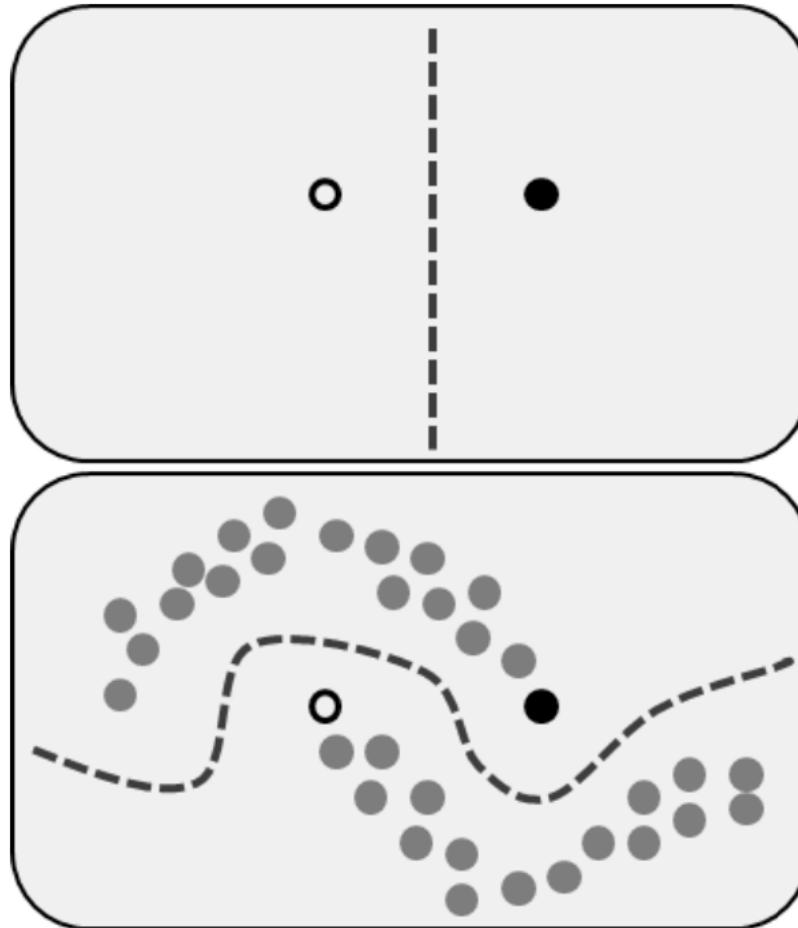
Semi-supervised learning



Semi-supervised learning



Semi-supervised learning



Idea 1: Semi-supervised learning

Consider a Markov factorization w.r.t. causal DAG:

$$p(x_1, \dots, x_d) = \prod_{i=1}^d p(x_i | x_{pa(i)})$$

Idea 1: Semi-supervised learning

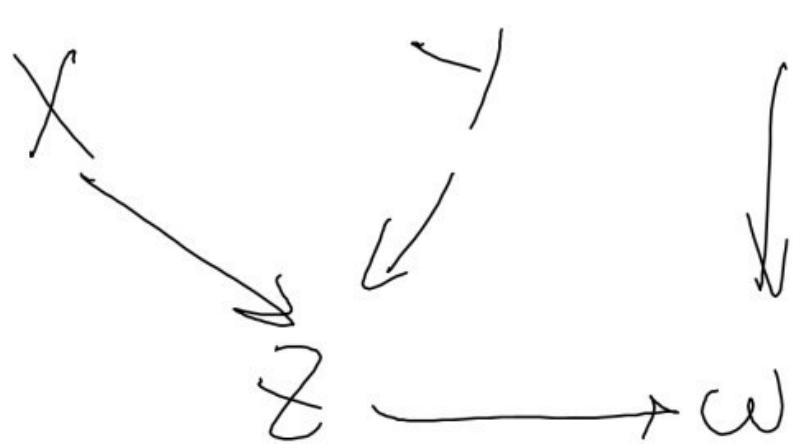
Consider a Markov factorization w.r.t. causal DAG:

$$p(x_1, \dots, x_d) = \prod_{i=1}^d p(x_i | x_{pa(i)})$$

Modularity suggests:

$p(x_1 | x_{pa(1)}), \dots, p(x_d | x_{pa(d)})$ are “independent”

Idea 1: Semi-supervised learning



$$\begin{aligned} p(A, X, Y, Z, W) &= \\ &= p(X) p(Y) \underbrace{p(Z)}_{\sim p(a)} \underbrace{x_{Y|Z}}_{\sim p(w|a, z)}, \end{aligned}$$

$$X = N_x$$

$$Z = f(X, \gamma) + N_z$$

$$W = g(A, Z) + N_w$$

Idea 1: Semi-supervised learning

Consider a Markov factorization w.r.t. causal DAG:

$$p(x_1, \dots, x_d) = \prod_{i=1}^d p(x_i | x_{pa(i)})$$

Modularity suggests:

$p(x_1 | x_{pa(1)}), \dots, p(x_d | x_{pa(d)})$ are “independent”

Special case:

$p(\text{cause}), p(\text{effect} | \text{cause})$ are “independent”

Idea 1: Semi-supervised learning

Consider a Markov factorization w.r.t. causal DAG:

$$p(x_1, \dots, x_d) = \prod_{i=1}^d p(x_i | x_{pa(i)})$$

Modularity suggests:

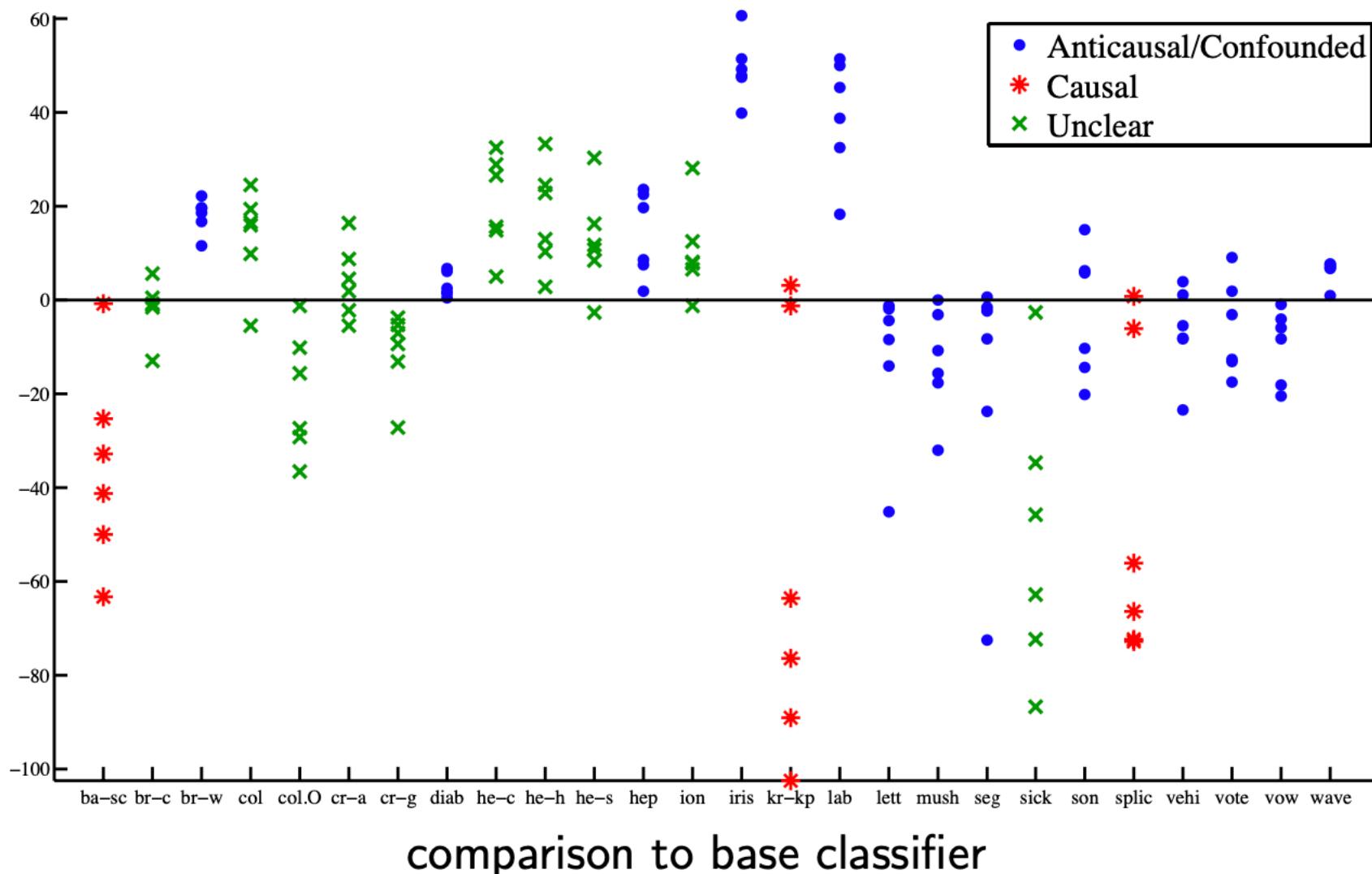
$p(x_1 | x_{pa(1)}), \dots, p(x_d | x_{pa(d)})$ are “independent”

Special case:

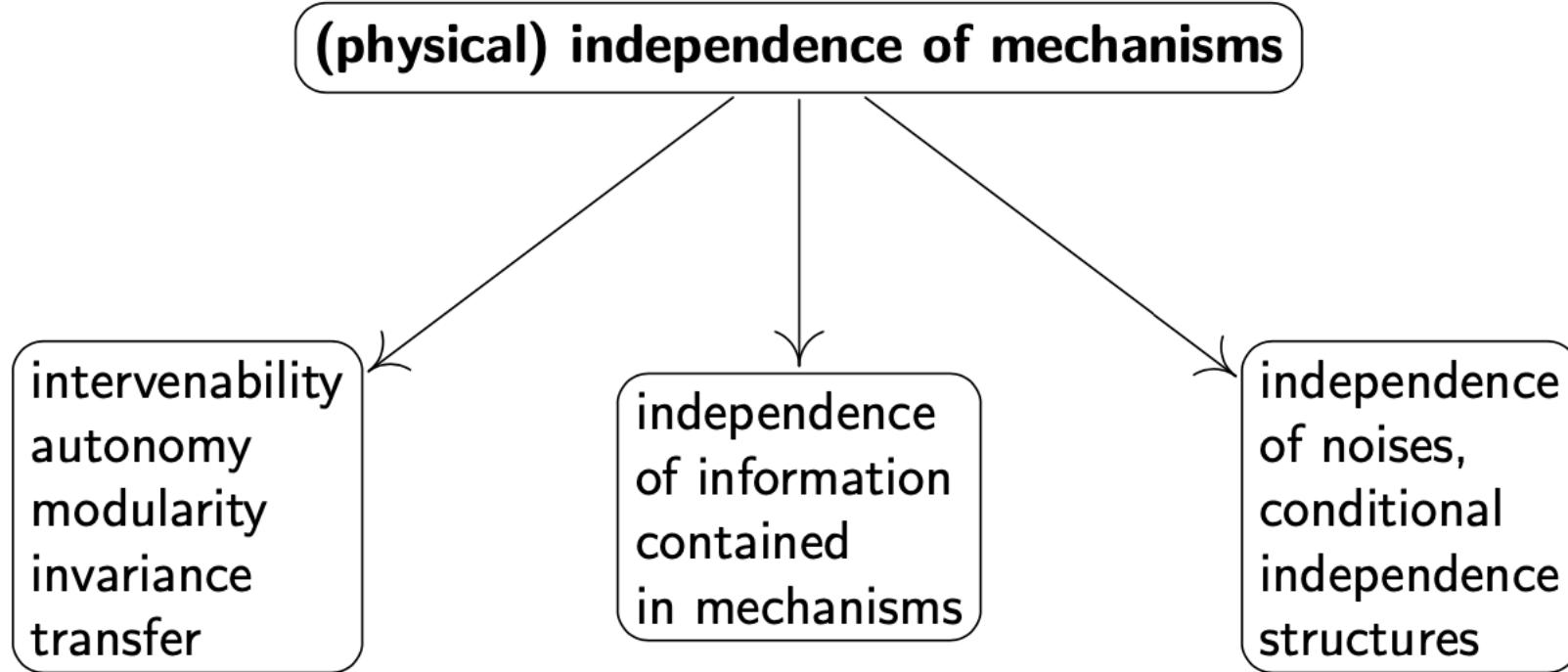
$p(\text{cause}), p(\text{effect} | \text{cause})$ are “independent”

But then: Semi-supervised Learning does not work from cause to effect.

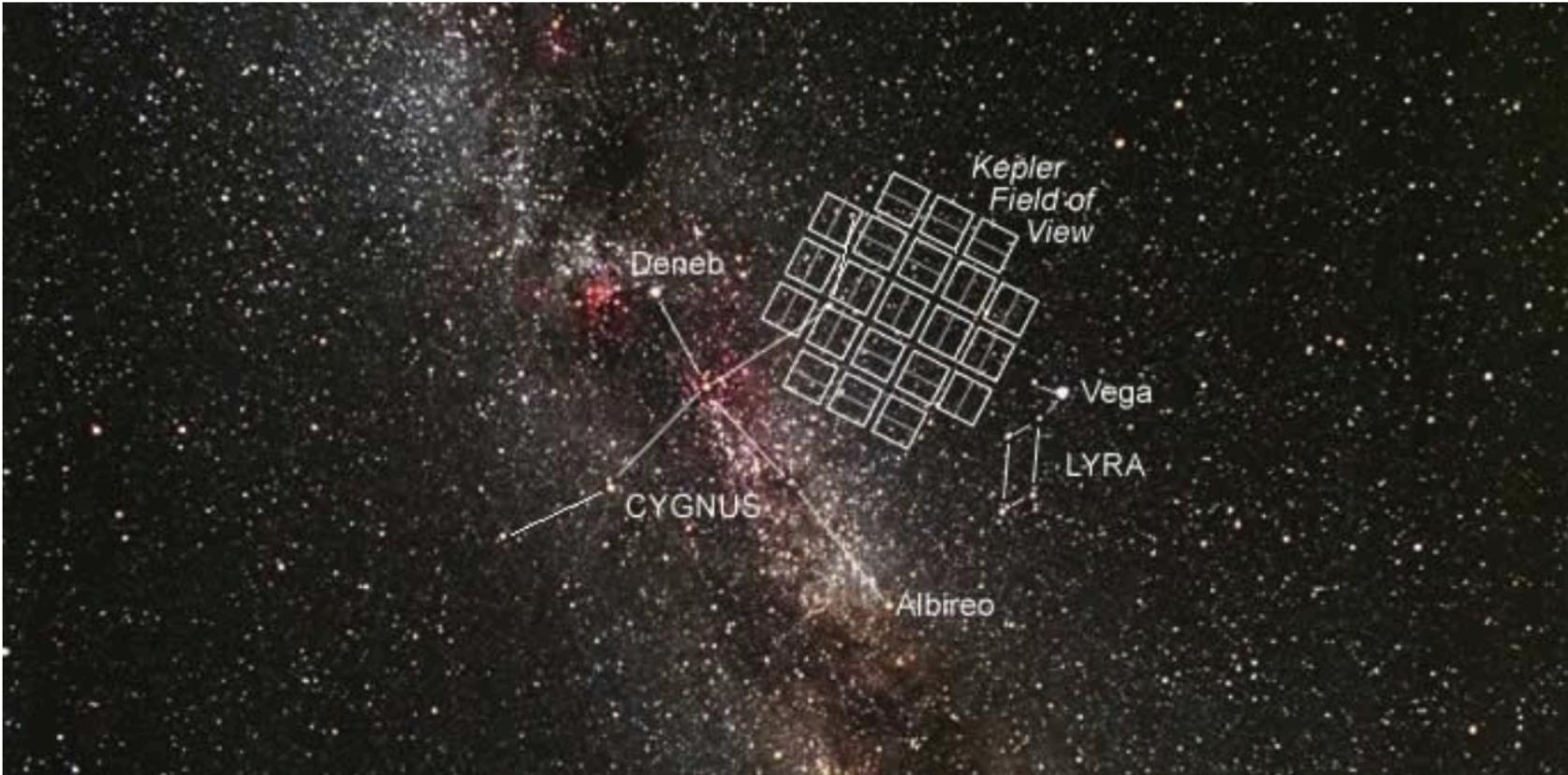
Idea 1: Semi-supervised learning



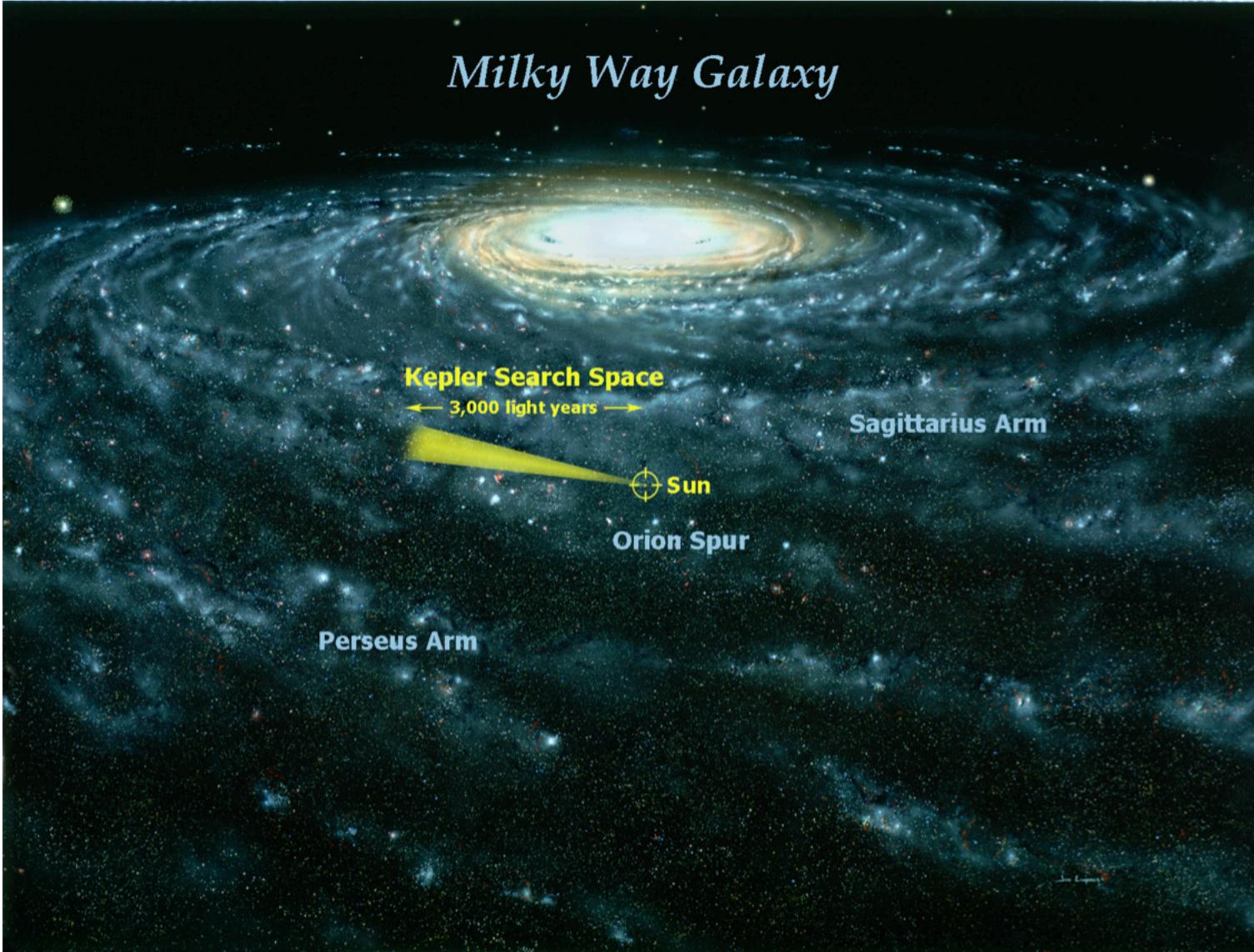
Idea 1: Semi-supervised learning



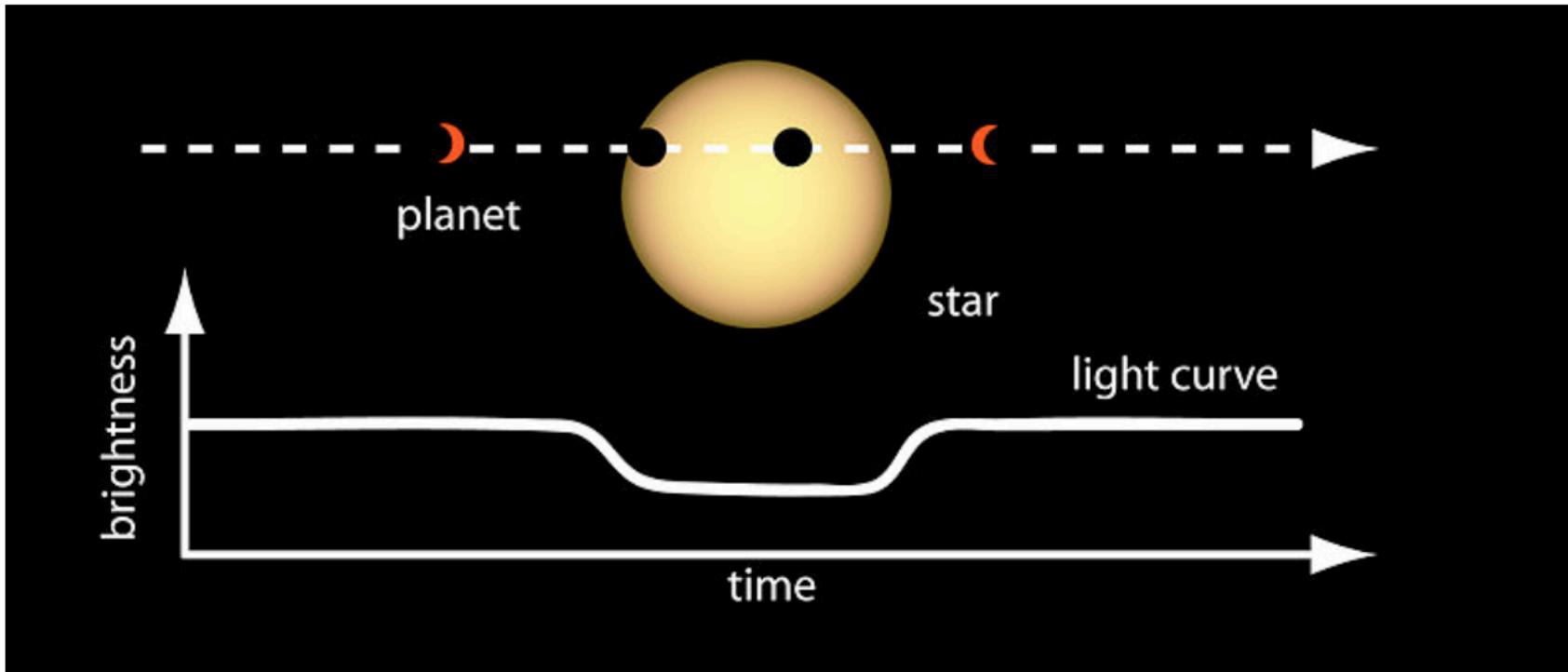
Idea 2: half-sibling regression



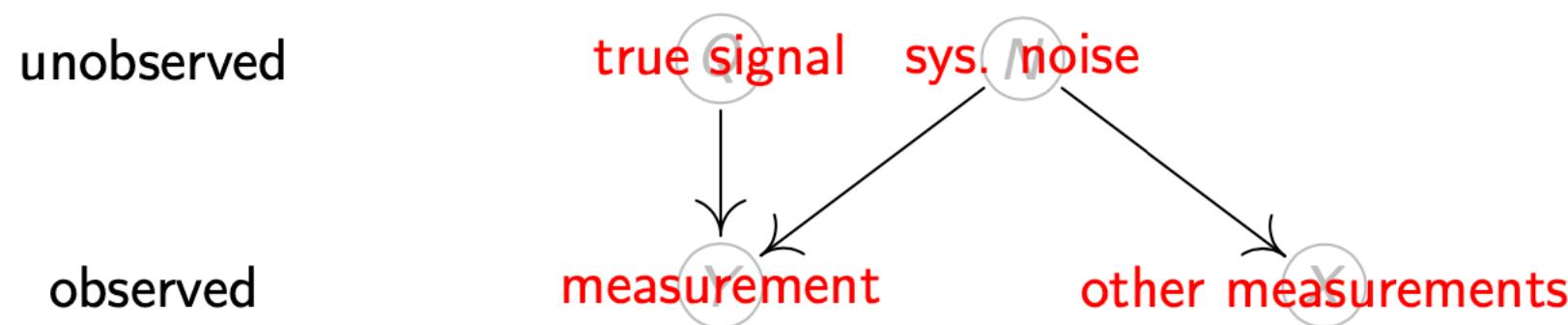
Idea 2: half-sihlinσ regression



Idea 2: half-sibling regression



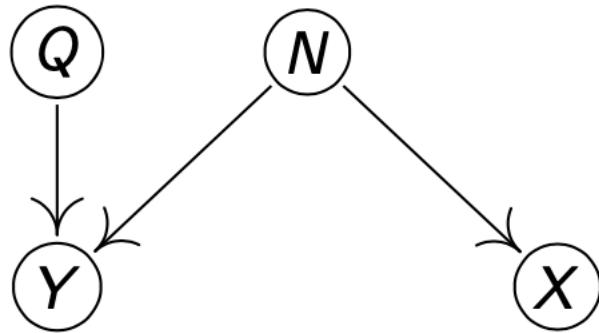
Idea 2: half-sibling regression



Idea 2: half-sibling regression

unobserved

observed

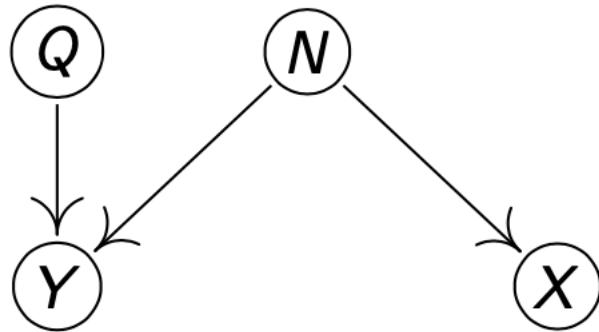


Assume $Y = f(N) + Q$.

Idea 2: half-sibling regression

unobserved

observed



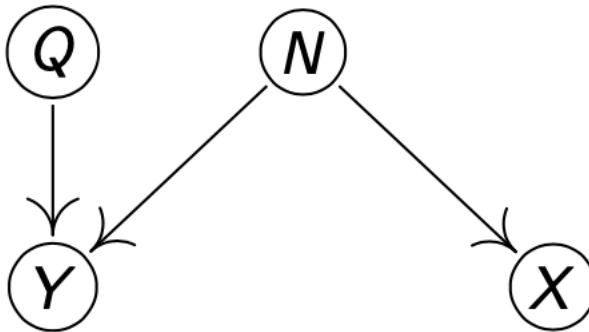
Assume $Y = f(N) + Q$.

Proposed idea:

Remove everything from Y explained by X .

Idea 2: half-sibling regression

unobserved



observed

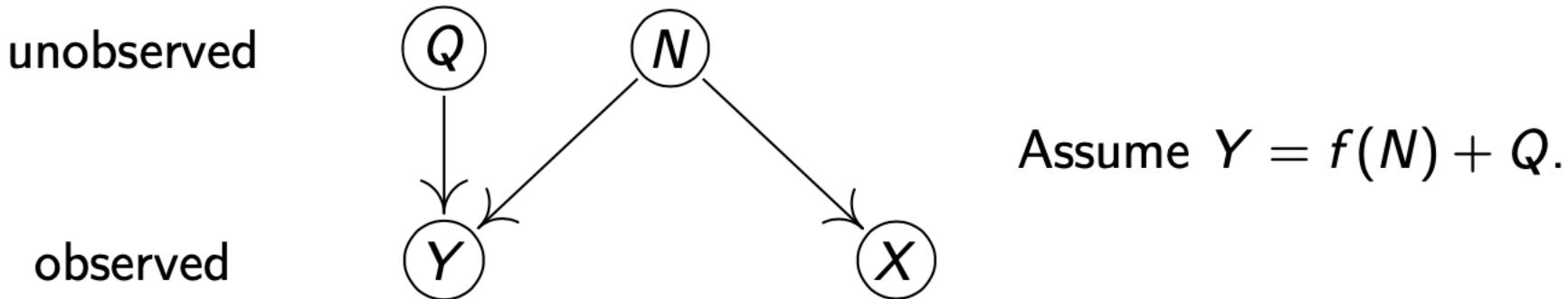
Assume $Y = f(N) + Q$.

Proposed idea:

Remove everything from Y explained by X .

Or: $\hat{Q} := Y - \mathbf{E}[Y | X]$.

Idea 2: half-sibling regression



Proposed idea:

Remove everything from Y explained by X .

Or: $\hat{Q} := Y - \mathbf{E}[Y | X]$.

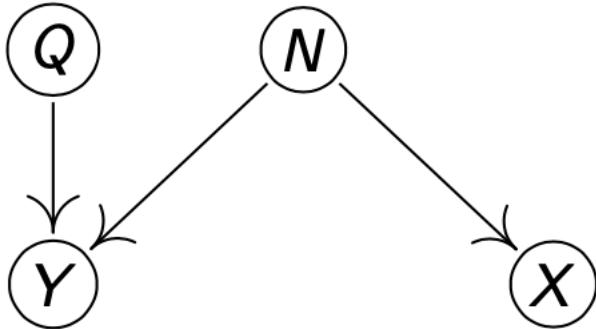
Proposition

Convergence against “correct” signal Q (up to reparameterization) if

- *perfect reconstruction:* $\exists \psi$ such that $f(N) = \psi(X)$

Idea 2: half-sibling regression

unobserved



observed

Assume $Y = f(N) + Q$.

Proposed idea:

Remove everything from Y explained by X . Or: $\hat{Q} := Y - \mathbb{E}[Y | X]$.

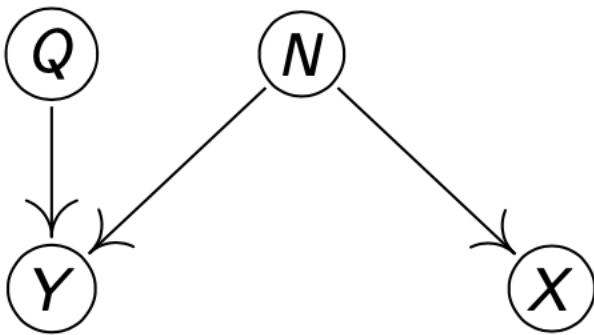
Proposition

Convergence against “correct” signal Q (up to reparameterization) if

- *perfect reconstruction: $\exists \psi$ such that $f(N) = \psi(X)$*
- *low noise: $X = g(N) + s \cdot R$ and $s \rightarrow 0$*

Idea 2: half-sibling regression

unobserved



Assume $Y = f(N) + Q$.

Proposed idea:

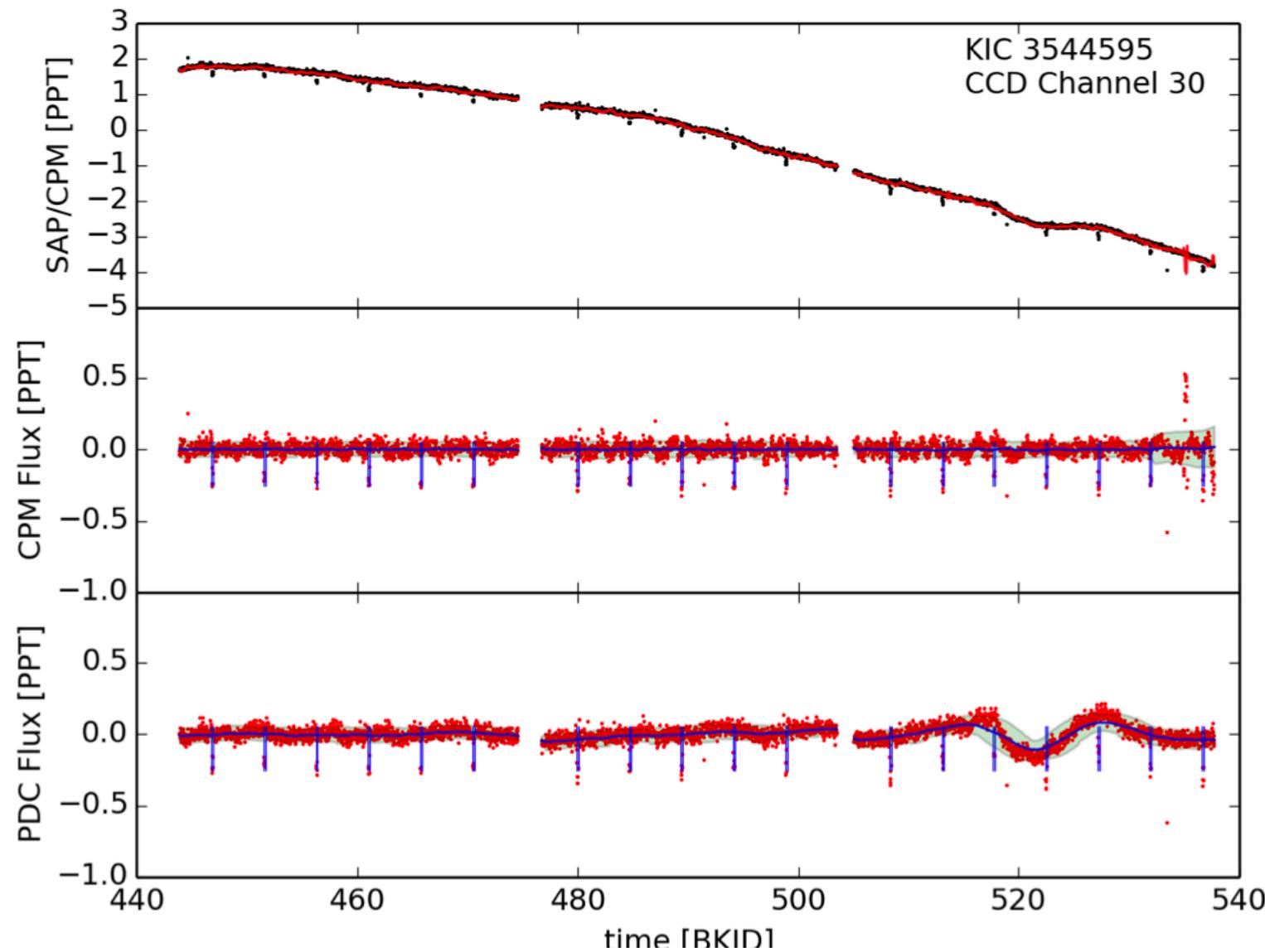
Remove everything from Y explained by X . Or: $\hat{Q} := Y - \mathbf{E}[Y | X]$.

Proposition

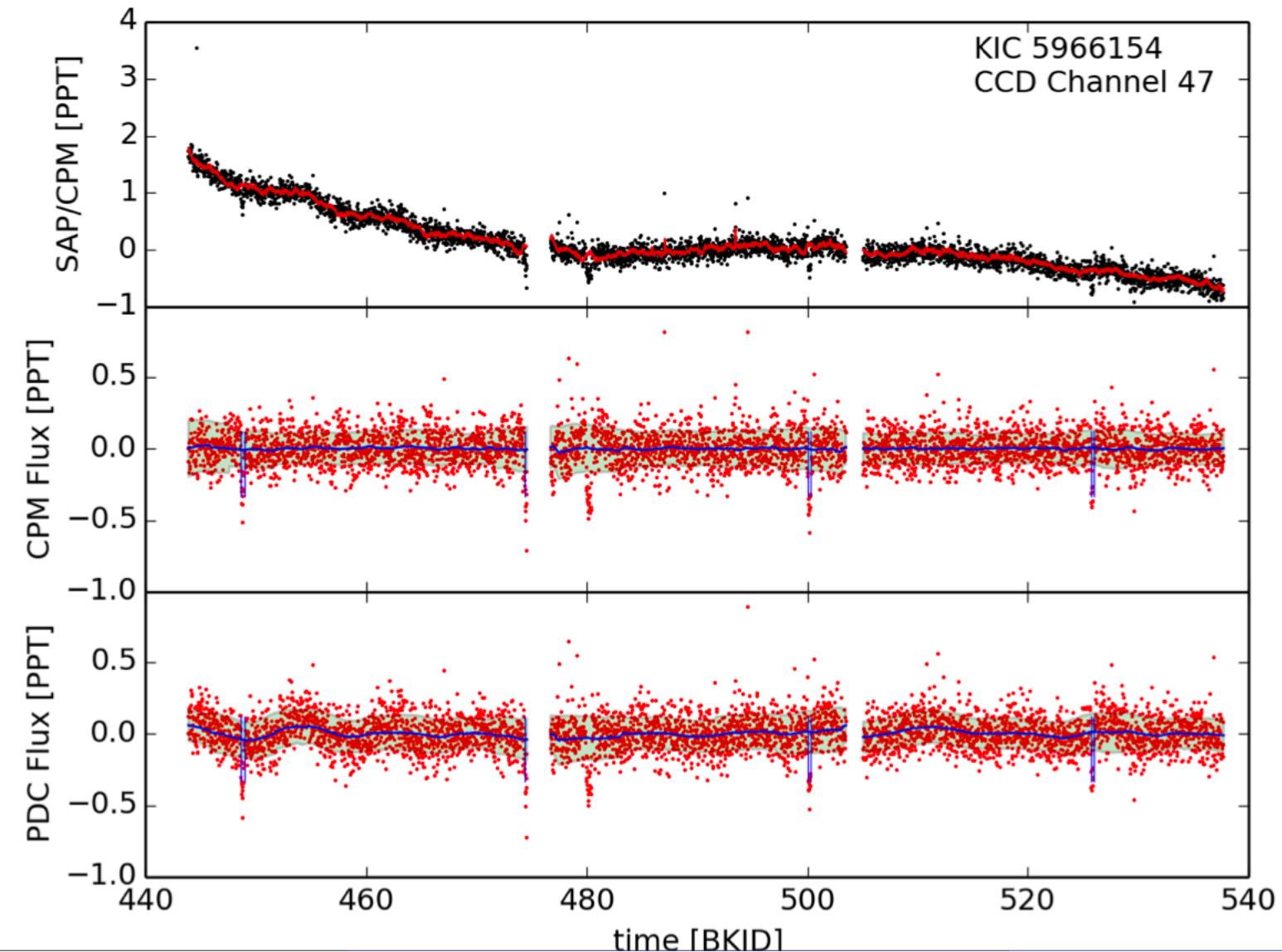
Convergence against “correct” signal Q (up to reparameterization) if

- perfect reconstruction: $\exists \psi$ such that $f(N) = \psi(X)$
- low noise: $X = g(N) + s \cdot R$ and $s \rightarrow 0$
- limit of infinitely many X ’s: $X_i = g_i(N) + R_i, i = 1, \dots$

Idea 2: half-sibling regression



Idea 2: half-sibling regression



Idea 2: half-sibling regression

