

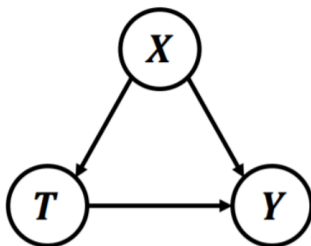
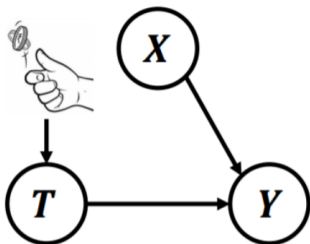
# Causal discovery

Credit : Neal

# Refresher

We have considered two settings

- RCTs [ $T \perp\!\!\!\perp (Y(0), Y(1))$ ]
- Neyman Rubin model [ $T \perp\!\!\!\perp (Y(0), Y(1))|X$ ]



# Refresher

Under these assumptions, we have addressed several questions

- Estimate the **Average Treatment Effect** (ATE)

$$ATE := \mathbb{E}[Y_i(1) - Y_i(0)]$$

- Taking into account Heterogeneous Treatment Effects and estimate the **Conditional Average Treatment Effect** (CATE)

$$CATE(x) := \mathbb{E}[Y_i(1) - Y_i(0)|X = x]$$

# Refresher

## Underlying assumptions in Course 1-Course 3

- The graph is assumed to be known
- It involves only three variables

## Some natural questions

- What about **more complex** situations?

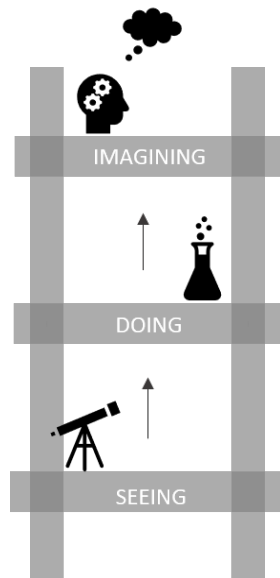
**Some answers in Course 4**

- How **can we learn** the graph?

**Some answers in Course 5**

# Challenges and principles

- Theory
  - ▶ Sometimes infeasible!
- Experts
  - ▶ Sometimes infeasible!
- Experimentations
  - ▶ Sometimes infeasible
  - ▶ Sometimes unethical
  - ▶ Costly
- Observations
  - ▶ Correlation does not imply causation!



# Challenges and principles

- In general, causal discovery from observational data is not possible.
- But it is possible under additional assumptions.
- Several approaches in the literature
  - ▶ **Constraint based methods** : run local tests of independence to create constraints on space of possible graphs.
  - ▶ **Noise based methods** : find footprints in the noise that imply causal asymmetry.
  - ▶ ...

# Principles of constrain based methods

## Main steps

- Find **skelton**
- Find **v-structures**
- Orient other edges using **basic rules**

# Principles of constrain based methods

---

## Algorithm 1 SGS

---

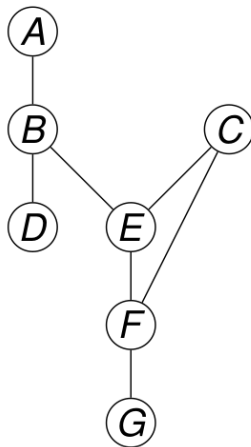
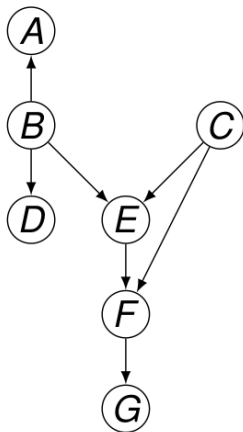
**Input:**  $P(\mathcal{V})$

**Output:** CPDAG  $\mathcal{G}^*$

- 1: Form the complete undirected graph  $\mathcal{G}^*$  on vertex set  $\mathcal{V}$
  - 2: **for** all  $X - Y$  in  $\mathcal{G}^*$   
    and subsets  $\mathcal{S} \subseteq \mathcal{V} \setminus \{X, Y\}$  **do**
  - 3:     **if**  $\exists \mathcal{S} \subseteq \mathcal{V} \setminus \{X, Y\}$  such that  $X \perp\!\!\!\perp_P Y \mid \mathcal{S}$  **then**
  - 4:         Delete edge  $X - Y$  from  $\mathcal{G}^*$
  - 5:     **end if**
  - 6: **end for**
  - 7: **for** all  $X - Z - Y$  in  $\mathcal{G}^*$  such that  $X \notin \text{Adj}(Y, \mathcal{G})$  **do**
  - 8:     **if**  $\nexists \mathcal{S} \subseteq \mathcal{V} \setminus \{X, Y\}$  such that  $Z \in \mathcal{S}$  and  $X \perp\!\!\!\perp_P Y \mid \mathcal{S}$  **then**
  - 9:         Orient  $X \rightarrow Z \leftarrow Y$  in  $\mathcal{G}^*$
  - 10:     **end if**
  - 11: **end for**
  - 12: Recursively apply rules R1-R3 until no more edges can be oriented
  - 13: **Return**  $\mathcal{G}^*$
-

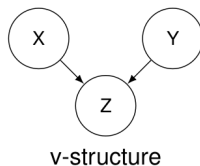
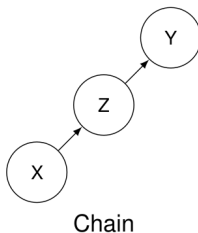
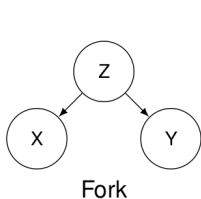


# Principles of constrain based methods



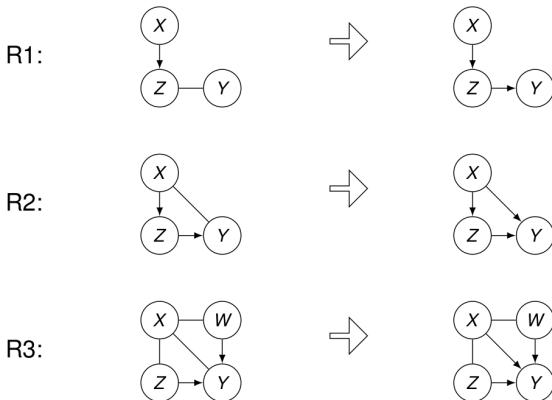
a DAG and its corresponding **skeleton**

# Principles of constrain based methods



Fork, chains and v-structures

# Principles of constrain based methods



Basic rules

# The concept of $d$ separation

## Blocked paths

A path is said to be blocked by a set of vertices  $Z$  if:

- it contains a chain  $A \rightarrow B \rightarrow C$  or a fork  $A \leftarrow B \rightarrow C$  and  $B \in Z$ , or
- it contains a collider  $A \rightarrow B \leftarrow C$  such that no descendant of  $B$  is in  $Z$

# The concept of $d$ separation

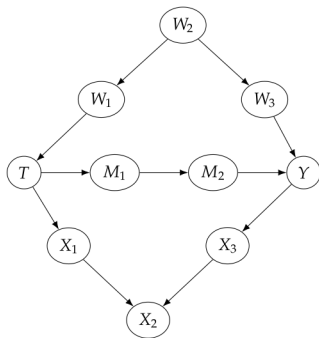
## Definition

Two (sets of) nodes  $X$  and  $Y$  are  $d$ -separated by a set of nodes  $Z$  if all of the paths between (any node in)  $X$  and (any node in)  $Y$  are blocked by  $Z$ . We denote  $X \perp\!\!\!\perp_G Y | Z$

## Theorem

Two DAGs  $G_1$  and  $G_2$  have the same  $d$ -separations iff they have the same skeleton and the same  $v$ -structures.

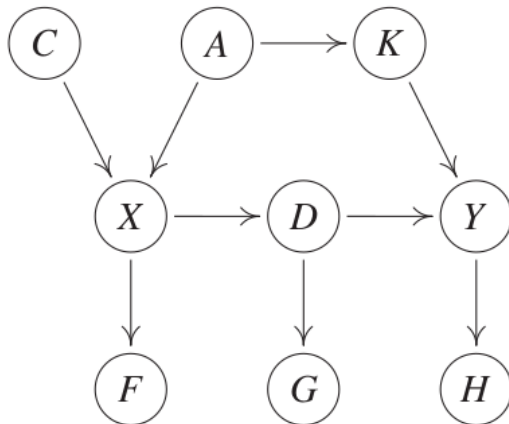
# The concept of $d$ separation



Are  $T$  and  $Y$   $d$  separated by

- the empty set?
- $\{W_2\}$ ?
- $\{W_2, M_1\}$ ?
- $\{W_1, M_2\}$ ?
- $\{W_1, M_2, X_2\}$ ?

# The concept of $d$ separation



For this DAG :  $C \perp\!\!\!\perp_G G|\{X\}$  and  $C \not\perp\!\!\!\perp_G G|\{X, H\}$

# The concept of $d$ separation

## Link with the conditional independency concept?

- Markov assumption :  $X \perp\!\!\!\perp_G Y|Z$  then  $X \perp\!\!\!\perp_{\mathbb{P}} Y|Z$
- We can assume that the converse holds, this is the **faithfulness assumption**

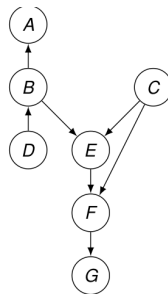
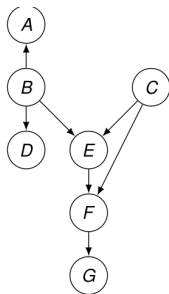
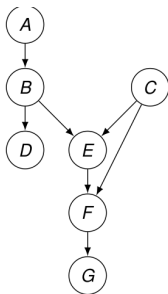


# The concept of $d$ separation

## Markov equivalence class

- Under these two assumptions, estimate the  $d$  separation in a graph consists in estimating conditional dependencies
- Graphs having the same  $d$  separation are said to be **Markov equivalent**

# The concept of $d$ separation



# The concept of $d$ separation

## Completed partially directed acyclic graph (CPDAG)

Let  $[G]$  be the Markov equivalence class of a DAG  $G$ . The CPDAG  $G^*$  of  $G$  is the graph:

- With the same skeleton as  $G$ ;
- Where an edge is directed in  $G$  iff it occurs as a directed edge with the same orientation in every graph in  $[G]$ ;
- All other edges are undirected.

# The concept of $d$ separation

- CPDAG coincide with Markov equivalence classes under two additional assumption : causal sufficiency (no latent variables) and acyclicity
- Rules R1-R3 ensure these assumptions

# The PC algorithm

- PC algorithm : optimized version of SGS
- Infer causal structure with the PC algorithm?
  - ▶ Infer mutual dependencies between variables : **skeleton of the causal graph**
  - ▶ Distinguish between causes and effects : **orientation of the v-structures of the causal graph**

# The PC algorithm

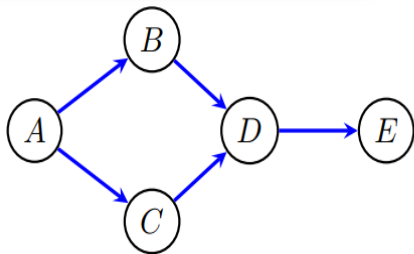
Independence tests: some examples

Type of variable	An example of independence test
Discrete	$\chi^2$ test
Gaussian	Test based on the precision matrix
Non Gaussian continuous	Non parametric tests

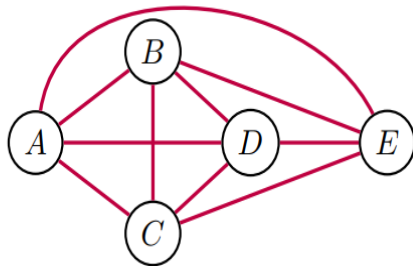
See notebook `CI.ipynb` for more details

# The PC algorithm

## An example



Unknown true graph

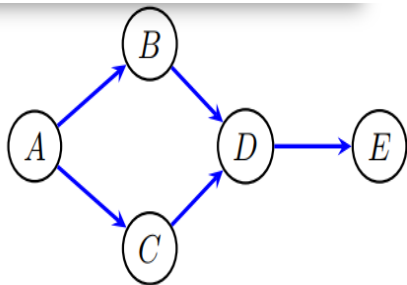


Graph after  $k = 0$

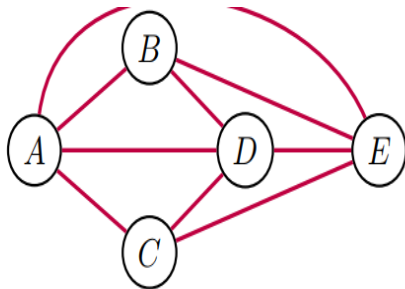
Figure: The initial graph is complete

# The PC algorithm

## An example



Unknown true graph



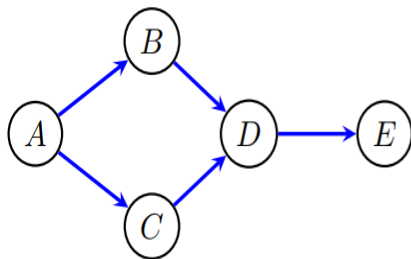
Graph

Figure: One conditions with respect to  $S = \emptyset$

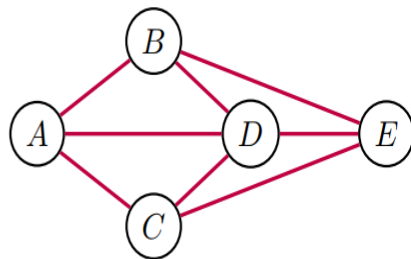


# The PC algorithm

## An example



Unknown true graph

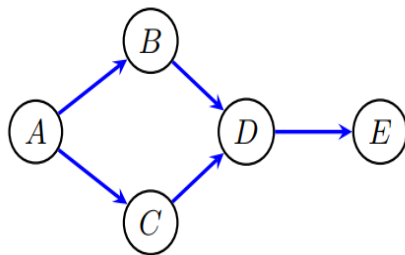


Graph

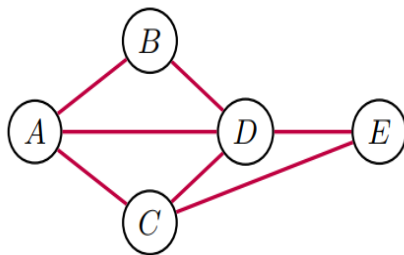
Figure:  $B$  and  $C$  are separated with respect to  $A$

# The PC algorithm

## An example



Unknown true graph

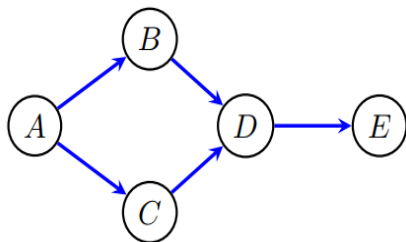


Graph

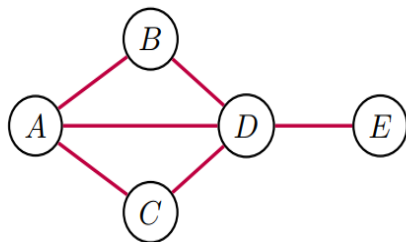
Figure:  $A$  and  $E$  are separated with respect to  $D$

# The PC algorithm

## An example



Unknown true graph

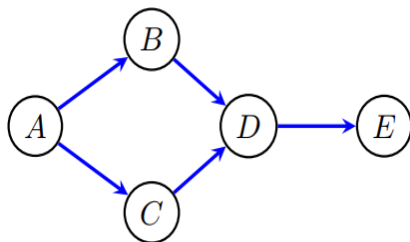


Graph after  $k = 1$

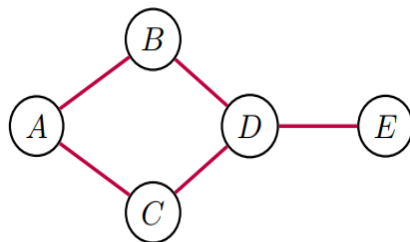
Figure:  $B$  and  $E$  are separated with respect to  $D$

# The PC algorithm

## An example



Unknown true graph

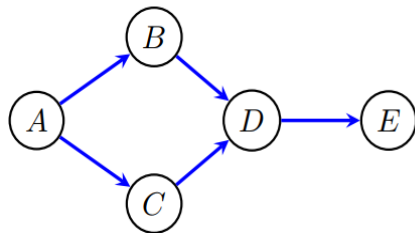


Graph after  $k = 2$

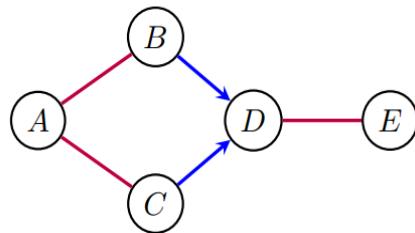
Figure:  $C$  and  $E$  are separated with respect to  $D$

# The PC algorithm

## An example



Unknown true graph

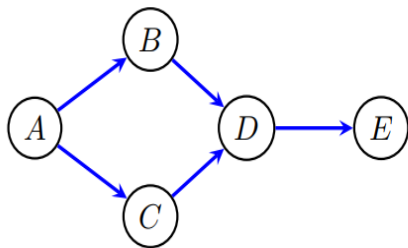


Graph after  $k = 1$

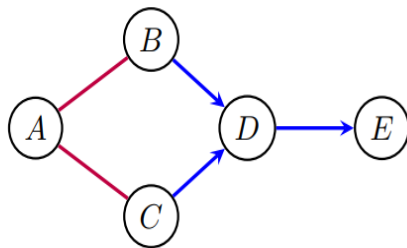
Figure:  $A$  and  $D$  are separated with respect to  $\{B, C\}$

# The PC algorithm

## An example



Unknown true graph



CPDAG

Figure:  $D$  is not in the set of nodes separating  $B$  and  $C$

# The PC algorithm

## An example

---

**Algorithm 1** The  $PC_{pop}$ -algorithm

---

- 1: **INPUT:** Vertex Set  $V$ , Conditional Independence Information
  - 2: **OUTPUT:** Estimated skeleton  $C$ , separation sets  $S$  (only needed when directing the skeleton afterwards)
  - 3: Form the complete undirected graph  $\tilde{C}$  on the vertex set  $V$ .
  - 4:  $\ell = -1$ ;  $C = \tilde{C}$
  - 5: **repeat**
  - 6:    $\ell = \ell + 1$
  - 7:   **repeat**
  - 8:     Select a (new) ordered pair of nodes  $i, j$  that are adjacent in  $C$  such that  $|adj(C, i) \setminus \{j\}| \geq \ell$
  - 9:     **repeat**
  - 10:       Choose (new)  $\mathbf{k} \subseteq adj(C, i) \setminus \{j\}$  with  $|\mathbf{k}| = \ell$ .
  - 11:       **if**  $i$  and  $j$  are conditionally independent given  $\mathbf{k}$  **then**
  - 12:         Delete edge  $i, j$
  - 13:         Denote this new graph by  $C$
  - 14:         Save  $\mathbf{k}$  in  $S(i, j)$  and  $S(j, i)$
  - 15:       **end if**
  - 16:     **until** edge  $i, j$  is deleted or all  $\mathbf{k} \subseteq adj(C, i) \setminus \{j\}$  with  $|\mathbf{k}| = \ell$  have been chosen
  - 17:   **until** all ordered pairs of adjacent variables  $i$  and  $j$  such that  $|adj(C, i) \setminus \{j\}| \geq \ell$  and  $\mathbf{k} \subseteq adj(C, i) \setminus \{j\}$  with  $|\mathbf{k}| = \ell$  have been tested for conditional independence
  - 18: **until** for each ordered pair of adjacent nodes  $i, j$ :  $|adj(C, i) \setminus \{j\}| < \ell$ .
-

# The PC algorithm

## An example

---

**Algorithm 2** Extending the skeleton to a CPDAG

---

**INPUT:** Skeleton  $G_{skel}$ , separation sets  $S$

**OUTPUT:** CPDAG  $G$

**for all** pairs of nonadjacent variables  $i, j$  with common neighbour  $k$  **do**

**if**  $k \notin S(i, j)$  **then**

    Replace  $i - k - j$  in  $G_{skel}$  by  $i \rightarrow k \leftarrow j$

**end if**

**end for**

In the resulting PDAG, try to orient as many undirected edges as possible by repeated application of the following three rules:

**R1** Orient  $j - k$  into  $j \rightarrow k$  whenever there is an arrow  $i \rightarrow j$  such that  $i$  and  $k$  are nonadjacent.

**R2** Orient  $i - j$  into  $i \rightarrow j$  whenever there is a chain  $i \rightarrow k \rightarrow j$ .

**R3** Orient  $i - j$  into  $i \rightarrow j$  whenever there are two chains  $i - k \rightarrow j$  and  $i - l \rightarrow j$  such that  $k$  and  $l$  are nonadjacent.

**R4** Orient  $i - j$  into  $i \rightarrow j$  whenever there are two chains  $i - k \rightarrow l$  and  $k \rightarrow l \rightarrow j$  such that  $k$  and  $l$  are nonadjacent.

---



# Cause or consequence?

- Can we **distinguish** cause from effect?
- That is distinguish between these two causal graphs

$$X \rightarrow Y$$

or

$$Y \rightarrow X$$

using **observational data**.

Not always possible!

# Cause or consequence?

The example of linear structural equation [f linear]

$X$  cause  $Y$  if there exists  $a \in \mathbb{R}, \varepsilon^Y$  s.t.

$$Y = aX + \varepsilon^Y, X \perp\!\!\!\perp \varepsilon^Y.$$

Distinguish cause from consequence? [Shimizu et al., 2006]

Assume that  $Y = aX + \varepsilon^Y, X \perp\!\!\!\perp \varepsilon^Y$  where all r.v. are continuous. Then

$$\exists b \in \mathbb{R}, \varepsilon^X \text{ s.t. } X = bY + \varepsilon^X, Y \perp\!\!\!\perp \varepsilon^X$$

iff  $(X, \varepsilon^X)$  are Gaussian random variables.

Existence of a non-linear extension of this result.

## More on the Bigaussian case (1)

See notebook `noise.ipynb`

### Causal model in the Bigaussian case

Let  $(X, Y) \sim \mathcal{N}((0, 0), \Sigma)$ .

$$Y = aX + \varepsilon^Y, \varepsilon^Y \perp X \text{ where } X, \varepsilon^Y \sim \mathcal{N}(0, \sigma) \text{ with } a = C_{X,Y}/V_X$$

- Caveat :  $X = \frac{1}{a}(Y - \varepsilon^Y)$  but  $\varepsilon^Y \not\perp Y : \varepsilon^Y$  and  $Y$  not **independants**. Non **causal model**.
- There exists  $(b, \epsilon^X)$  s.t.

$$X = bY + \epsilon^X, \epsilon^X \perp Y,$$

$$\text{with } b = \frac{aV_X}{a^2V_X + V_{\epsilon^Y}}$$

## More on the Bigaussian case (2)

Simulations : sample size  $n = 2000$ .

	1	2
1	10.00	3.00
2	3.00	2.00

Table:  $\Sigma$

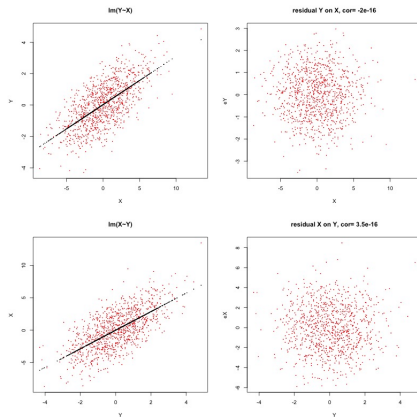
	Estimate	Std. Error	t value	Pr(> t )
(Intercept)	0.0234	0.0336	0.70	0.4868
X	0.3081	0.0110	28.03	0.0000

Table: Residual standard error: 1.063 on 998 degrees of freedom. Multiple R-squared: 0.4405, Adjusted R-squared: 0.44 ; F-statistic: 785.8 on 1 and 998 DF, p-value:  $< 2.2e - 16$

	Estimate	Std. Error	t value	Pr(> t )
(Intercept)	-0.0052	0.0725	-0.07	0.9430
Y	1.4300	0.0510	28.03	0.0000

Table: Residual standard error: 2.291 on 998 degrees of freedom; Multiple R-squared: 0.4405, Adjusted R-squared: 0.44 ; F-statistic: 785.8 on 1 and 998 DF, p-value:  $< 2.2e - 16$

## More on the Bigaussian case (3)

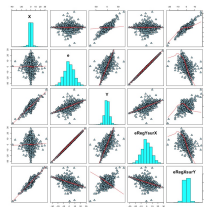


**Figure:** In Gaussian case, we cannot distinguish cause from effect

# Non Gaussian example

## Model

$$Z \sim \mathcal{N}(0, 1), X = Z^3, \epsilon^Y \sim \mathcal{N}(0, 9) \text{ and } Y = 2X + \epsilon^Y$$



# Non Gaussian example

- In the NON Gaussian setting cause can be distinguished from effect
- Independence tested and accepted for couples  $(X, \epsilon^Y)$ ,  $(X, eRegY_{surX})$  but not  $(Y, \epsilon^Y)$ , ni  $(Y, eRegX_{surY})$ .
- $X$  has an influence on  $Y$  but  $Y$  does not influence  $X$ .

# Noise based algorithm

## Theorem (LiNGAM)

Assume a linear SCM with graph  $G = (V, E)$  and a compatible distribution  $P(V)$  such that for all  $Y \in V$

$$Y = \sum_{X \in Pa(Y)} a_{xy} X + \xi_Y$$

where all  $\xi_Y$  are jointly independent and non-Gaussian distributed. Additionally, we require that for all  $Y \in V$ ,  $X \in Pa(Y)$ ,  $a_{xy} \neq 0$ . Then, the graph  $G$  is identifiable from  $P(V)$ .



# Noise based algorithm

## LINGAM

---

**Algorithm 1** LiNGAM

---

**Input:**  $P(\mathcal{V})$ **Output:**  $\mathcal{G}$ 

```
1: Form an empty graph  $\mathcal{G}$  on vertex set  $\mathcal{V} = \{X_1, \dots, X_p\}$ 
2: Let  $S = \{1, \dots, p\}$  and  $\mathcal{T} = []$ 
3: repeat
4:    $H = []$ 
5:   for  $i \in S$  do
6:     for  $j \in S \setminus \{i\}$  do
7:        $\hat{\xi}_{ij} = X_j - \frac{\text{cov}(X_i, X_j)}{\text{var}(X_i)} X_i$ 
8:     end for
9:      $h = \sum_{j \in S \setminus \{i\}} \lambda(X_i, \hat{\xi}_{ij})$ 
10:     $H = [H, h]$ 
11:   end for
12:    $i^* = \arg \min_{i \in S} H$ 
13:    $S = S \setminus \{i^*\}$ 
14:    $\mathcal{T} = [\mathcal{T}, i^*]$ 
15:    $\forall j \in S, X_j = \hat{\xi}_{i^*j}$ 
16: until  $|S| = 0$ 
17: Append( $\mathcal{T}, S_0$ )
18: Construct a strictly lower triangular matrix by following the order in  $\mathcal{T}$ , and estimate the connection strengths  $a_{i,j}$  by using some conventional covariance-based regression.
19: if  $a_{i,j} > 0$  then
20:   Add  $X_i \rightarrow X_j$  to  $\mathcal{G}$ 
21: end if
22: Return  $\mathcal{G}$ 
```

---

# Noise based algorithm

## Theorem (Anm)

Assume a linear SCM with graph  $G = (V, E)$  and a compatible distribution  $P(V)$  such that for all  $Y \in V$

$$Y = f((X \in Pa(Y)) + \xi_Y$$

where all  $\xi_Y$  are jointly independent. Then, the graph  $G$  is identifiable from  $P(V)$ .

# Noise based algorithm

## ANM

---

**Algorithm 2** ANM

---

**Input:**  $P(\mathcal{V})$ **Output:**  $\mathcal{G}$ 

```
1: Form an empty graph  $\mathcal{G}$  on vertex set  $\mathcal{V} = \{X_1, \dots, X_p\}$ 
2: Let  $S = \{1, \dots, p\}$  and  $\mathcal{T} = []$ 
3: repeat
4:    $H = []$ 
5:   for  $j \in S$  do
6:      $\hat{f}_j$ : Regress  $X_j$  on  $\{X_i\}_{i \in S \setminus \{j\}}$ 
7:      $\hat{\xi}_j = X_j - \hat{f}_j(X_i)$ 
8:      $h = \mathcal{I}(\{X_i\}_{i \in S \setminus \{j\}}, \hat{\xi}_j)$ 
9:      $H = [H, h]$ 
10:  end for
11:   $i^* = \arg \min_{i \in S} H$ 
12:   $S = S \setminus \{i^*\}$ 
13:   $\mathcal{T} = [i^*, \mathcal{T}]$ 
14: until  $|S| = 0$ 
15: for  $j \in \{2, \dots, p\}$  do
16:   for  $i \in \{\mathcal{T}_1, \dots, \mathcal{T}_{j-1}\}$  do
17:      $\hat{f}_j$ : Regress  $X_j$  on  $\{X_k\}_{k \in \{\mathcal{T}_1, \dots, \mathcal{T}_{j-1}\} \setminus \{i\}}$ 
18:      $\hat{\xi}_j = X_j - \hat{f}_j(X_i)$ 
19:     if  $\{X_k\}_{k \in \{\mathcal{T}_1, \dots, \mathcal{T}_{j-1}\} \setminus \{i\}} \not\perp_P \hat{\xi}_j$  then
20:       Add  $X_i \rightarrow X_j$  to  $\mathcal{G}$ 
21:     end if
22:   end for
23: end for
24: Return  $\mathcal{G}$ 
```

---