

# NICTA

## Introduction

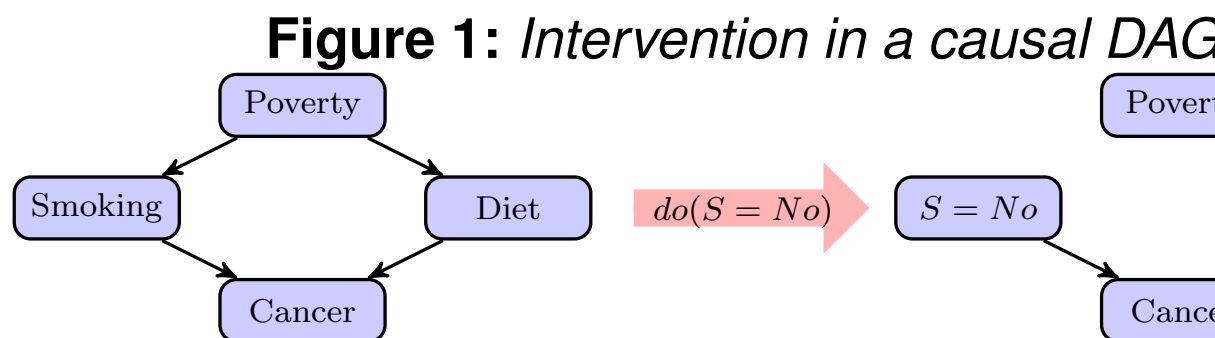
Would cutting salt intake reduce heart attack risk?  
Would increasing the minimum wage increase unemployment?  
These problems differ from the traditional machine learning problems in that they require us to predict the consequences of an intervention, which may change the properties of the population from which our data is sampled. Without explicitly testing the intervention or making assumptions to estimate its effect, such inference is impossible. There

its effect, such inference is impossible. There are many important questions where direct experimentation is expensive, unethical or impossible. Recent research on causal inference has clarified what assumptions allow causality to be estimated and has demonstrated causal inference can be performed. Every-thing is possible with some very general assumptions.

## Causal Frameworks

### Causal Directed Acyclic Graphs

- A causal DAG is a Bayesian network where the edges are defined to mean  $A$  causes  $B$ .
- Variables are independent of their non-effects given their direct causes (Causal Markov Property)
- An intervention that sets a subset of variables  $\mathbf{X}$  is denoted  $do(\mathbf{X} = \mathbf{x})$ , has a simple graphical representation in a causal DAG,  $G$ . All links entering intervention variables,  $\mathbf{X}$ , are deleted, resulting in the mutilated graph  $G_{\overline{\mathbf{X}}}$  (figure 1). Thus, a causal DAG represents all possible interventional distributions over its variables.



### (Causal) Structural Equation Models (SEMs)

- Represent each variable as a deterministic function of its direct causes and a noise term, where the noise terms are mutually independent.
- If the set of equations does not create a cycle, then the model is acyclic and can be solved for the values of the variables.

- If the set of equations does not create a cycle, the Causal Markov Property holds and the SEM is a DAG, (but not visa-versa - SEMs can encode more information).

## Counterfactuals

- Counterfactuals are statements about what would have occurred under alternate realities where some specified variables have different values. For example, consider people taking a medication.

For an individual,  $i$ , let: 
$$\begin{cases} y_i^0 = & \text{outcome if } x_i = 0 \\ y_i^1 = & \text{outcome if } x_i = 1 \end{cases}$$

- We can define a random variable  $Y^1$ , where  $P(Y^1)$  is the distribution of outcome,  $Y$ , that would occur if everyone was treated. Similarly  $P(Y^0)$  is the distribution of outcome if no-one was treated.
- If  $(X \perp\!\!\!\perp Y^0 | Z)$  &  $(X \perp\!\!\!\perp Y^1 | Z)$ :  $\leftarrow$  Ignorability. If these conditions hold, we can calculate counterfactual distributions for unserved ones:

$$P(Y^1 | Z) = P(Y | X = 1, Z) \text{ and } P(Y^0 | Z) = P(Y | X = 0, Z)$$

- Distributions over counterfactual variables that correspond to interventions can be translated directly to the observed distribution  $P(Y^1) = P(Y | do(X = 1))$ . However we cannot answer queries with counterfactual variables that are not intervened on. For example: *what is the probability that a patient who was not treated and died, would have recovered if they had been treated?* This query asks about the joint distribution of  $P(Y^0, Y^1)$ .

# Causal Inference



ask? Would  
nt? Causal  
ing setting  
es of an in-  
he distribu-  
perimentally  
o constrain  
are many

group	placebo	treatment	probability
1	die	die	$\alpha = P(Y^0 = \text{die})$
2	die	recover	$\beta = P(Y^0 = \text{recover})$
3	recover	die	$\gamma = P(Y^1 = \text{die})$
4	recover	recover	$\delta = P(Y^1 = \text{recover})$

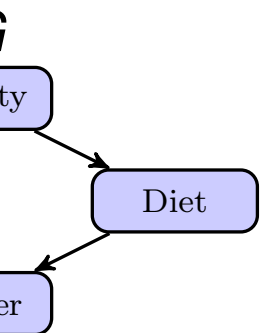
- Counterfactuals can be defined in terms of  $f(X = 0, \epsilon_Y)$
- The ignorability assumption  $(X \perp\!\!\!\perp Y^0 | \mathbf{Z})$  is weaker than that implied by mutually independent

are many  
is expen-  
causality  
to be deter-  
and discov-  
ns.

$A \rightarrow B$  is

given their

$X$  to  $x$ , de-  
representation  
ed on vari-  
ed network  
s the set of  
variables.



## SEMs)

ction of its  
e terms are

e then the

rors in a SEM,  $(X \perp\!\!\!\perp Y^0, Y^1 | Z)$ . The assumption is made equivalent by utilizing SEMs with a dependence of errors assumption [9].

## Causal Inference

Consider the problem where the causal DAG is known from theory or prior knowledge, and we wish to determine the outcome of an intervention of the form  $P(Y | do(X=x))$  from observational data.

- If there are no latent variables, we can compute the effect of any intervention by simply multiplying the joint distribution of the mutilated network (figure 1).

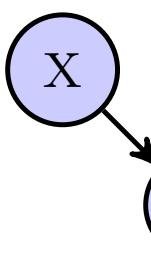
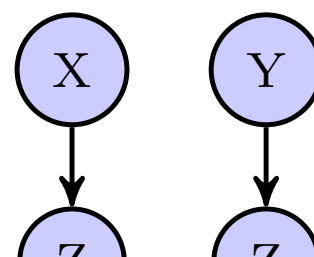
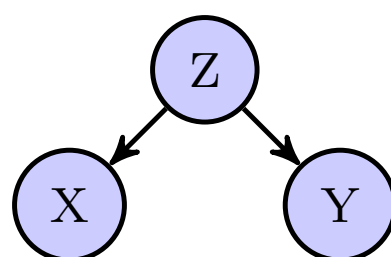
## The Do Calculus

The do-calculus consists of three rules, which are based on d-separation in a causal DAG [7]. It is complete: an effect is non-parametrically identifiable if and only if the corresponding interventional query can be reduced to an observational query via these rules [10].

**Figure 2:** *d*-separation allows us to read conditional independences off a DAG. If a set of variables  $Z$  d-separates  $X$  and  $Y$  in  $G$  then  $(X \perp\!\!\!\perp Y | Z)$  in all distributions compatible with  $G$ .

(a)  $(X \perp\!\!\!\perp Y | Z)$

(b) v-structure



then the  
is a causal  
more infor-

ould happen  
d thing dif-  
dical drug:

0 (not treated)

1 (treated)

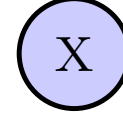
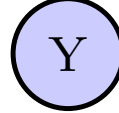
$(Y^1)$  is the  
f everyone  
of outcome

y Assump-

s from ob-

$|X = 0, Z)$

correspond  
e do nota-  
can phrase  
ot interven-  
*t Joe, who*  
*red had he*  
distribution



## The Three Rules

1. A causal DAG remains a causal DAG after  $\mathbf{X}$  so d-separation still applies.

$$\text{if } (\mathbf{Y} \perp\!\!\!\perp \mathbf{W} | \mathbf{X}) \text{ in } G_{\overline{\mathbf{X}}} \\ P(\mathbf{Y} | do(\mathbf{X} = \mathbf{x}), \mathbf{W} = \mathbf{w}) = P(\mathbf{Y} | \mathbf{W} = \mathbf{w})$$

2. If  $\mathbf{Y}$  is independent of *how* variables  $\mathbf{X}$  take their values then the effect on  $\mathbf{Y}$  of setting  $\mathbf{X}$  to some value is equivalent to observing it take that value. If this is the case then the corresponding ignoreability assumption in the counterfactual framework is satisfied.

$$\text{if } (\mathbf{Y} \perp\!\!\!\perp \hat{\mathbf{X}} | \mathbf{X}, \mathbf{Z}) \text{ in } G^{\dagger} \\ P(\mathbf{Y} | do(\mathbf{X} = \mathbf{x}), \mathbf{Z}) = P(\mathbf{Y} | \mathbf{X} = \mathbf{x}, \mathbf{Z})$$

3. If there is no direct causal path from  $\mathbf{X}$  to  $\mathbf{Y}$  then a intervention on  $\mathbf{X}$  does not change the distribution of  $\mathbf{Y}$ .

$$\text{if } (\mathbf{Y} \perp\!\!\!\perp \hat{\mathbf{X}} | \mathbf{Z}) \text{ in } G^{\dagger} \\ P(\mathbf{Y} | do(\mathbf{X} = \mathbf{x}), \mathbf{Z}) = P(\mathbf{Y} | \mathbf{Z})$$

(For readability, this is a simplified version of the do-calculus that only applies to a single variable or cases where it is sufficient for identifiability of the effect of interest on all variables together. The fully general version is only slightly more complex, see [7])

# ence in Machin

Finnian Lattimore

Australian National

finnlattimore@gm

of group

$(X^0 = 0, Y^1 = 0)$

$(X^0 = 0, Y^1 = 1)$

$(X^0 = 1, Y^1 = 0)$

$(X^0 = 1, Y^1 = 1)$

of SEMs  $Y^0 \sim$

&  $(X \perp\!\!\!\perp Y^1 | Z)$

independent er-

## Causal Discovery

Causal discovery attempts to infer causal structure based on more general assumptions.

## Independence Based Methods

### Without Latent Variables

1. We assume our distribution  $P$  was a

ptions can be  
modified inde-

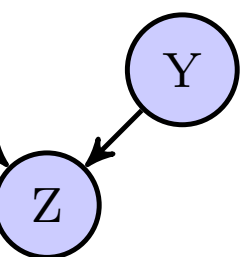
is known, due  
to infer the out-  
( $X = x$ )) using

compute outcome  
factors in the

which result from  
complete: a causal  
and only if the in-  
observational one

conditional inde-  
d-separates  $X$   
s  $P$  compatible

ture, ( $X \perp\!\!\!\perp Y$ )



1. We assume our distribution  $P$  was generated by an (unknown) causal DAG over our variables  $V$  (*causal sufficiency*)
2. We assume that all the conditional independencies in  $P$  are implied by d-separation in the true DAG (*faithfulness*)
3. Finding the causal structure equates to finding a DAG in the Markov equivalence class for the given  $P$

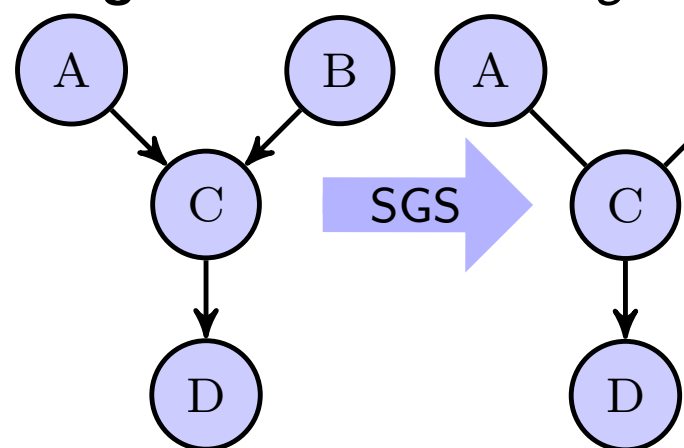
### Algorithm 1: SGS or IC Algorithm

**Input:** A Distribution  $P$  over variables  $V$

**Output:** A partially directed network (PDAG) representing the Markov equivalence class for the given  $P$

1. Create a complete undirected graph on  $V$ . For all pairs of variables  $(a, b) \in V$  search for a set  $S_{ab}$  such that  $a$  and  $b$  are d-separated by  $S_{ab}$ . If such a set exists, delete the link  $a - b$ .
2. For all pairs of unlinked-nodes  $(\alpha, \beta)$  search for a common neighbour  $c$ , if  $c \notin S_{\alpha\beta}$  direct links towards  $c$ .
3. Recursively direct any remaining links. Only one orientation that does not create additional v-structures ( $\bullet \rightarrow \bullet \leftarrow \bullet$ ).

**Figure 3: The SGS Algorithm**



- The SGS algorithm is infeasible in practice for large numbers of variables



an intervention

$$P(\mathbf{Y} | do(\mathbf{X} = \mathbf{x}))$$

make their values  
value is equiv-  
rule is satisfied  
in the coun-

$$P(\mathbf{Y} | \mathbf{X} = \mathbf{x}, \mathbf{Z})$$

to  $\mathbf{Y}$  then inter-  
tion of  $\mathbf{Y}$ .

$$P(\mathbf{Y} | \mathbf{Z})$$

covers interventions  
to consider interven-  
lightly more complex

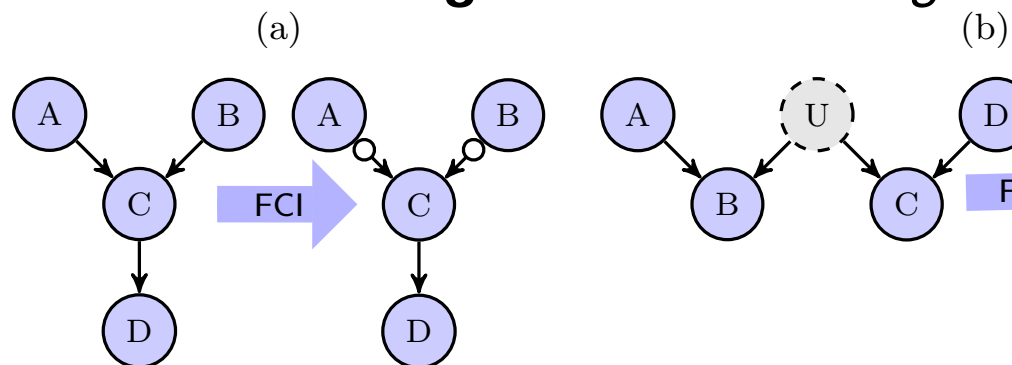
exponential number of (high order) con-  
dence tests it requires.

- The PC algorithm [11] modifies the SC algorithm to exploit any sparsity in the true network, leading to better average case performance.

## With Latent Variables

- For every latent structure there is a directed acyclic graph (DAG) in which every latent variable has exactly two children [12]. This key insight led to the FCI algorithm [11], which generalizes the PC algorithm to latent variables.
- The FCI algorithm returns an equivalence class of minimal Ancestral Graphs (a generalization of DAGs) as DAGs are not closed under marginalization.

**Figure 4: The FCI Algorithm**



- The FCI algorithm discovers all aspects of the true structure identifiable from conditional independence [13].

# ne Learning:

more

University

ail.com



structure from data

generated by some

- It can be made to require a worst case (more than exponential) number of conditional independence tests for sparse graphs [1].
- Implementations of both the PC and the ICD algorithms are available in the R package pcalg [5].

## Beyond independence

Independence based methods have been shown to require only very general, non-parametric assumptions. However they cannot distinguish between

generated by some  
observed variables

dependences in  $P$   
ue causal network

to finding perfect

nm [11, 7].

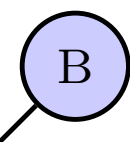
representing the  
ting causal model.

ver  $V$ . For all pairs  
 $S_{ab}$  s.t  $a \perp\!\!\!\perp b | S_{ab}$ . If

) with a common  
ards  $c$ .

s for which there is  
ate a cycle or any

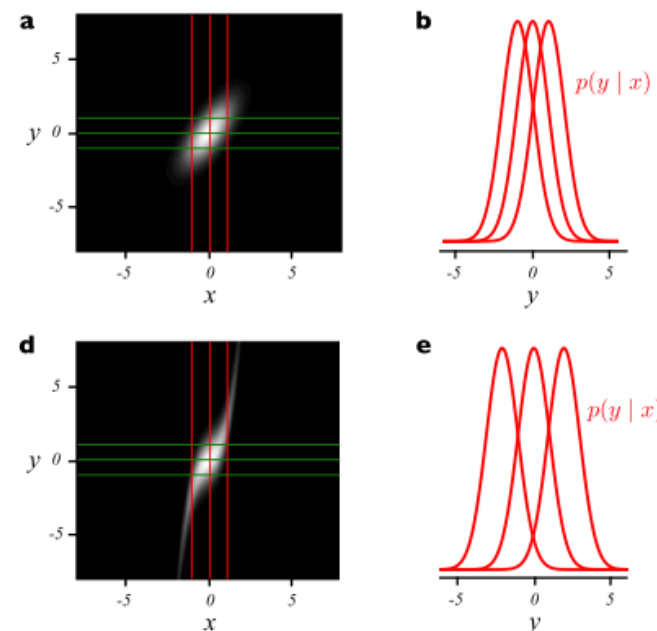
ithm



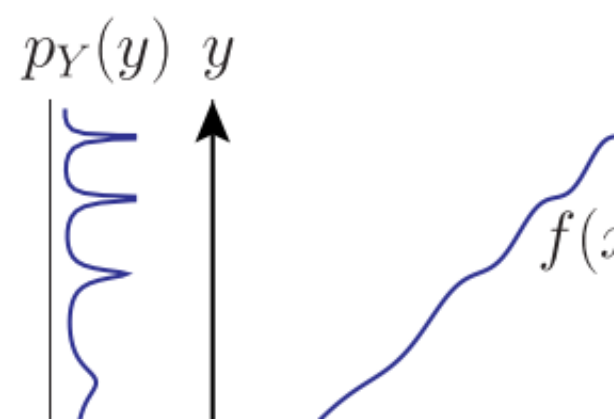
practise due to the

with equivalent dependency structure  
 $X \rightarrow Y$  and  $X \leftarrow Y$ .

**Figure 5:** Figure from [3]. Additive  
 $f(X) + \epsilon$ , are identifiable for most com  
but not in linear-gaussian case



**Figure 6:** Figure from [2]. The causal  
tifiable even where the relationship b  
terministic and invertible. Let  $Y = f$   
most input distributions,  $p_X(x)$ , the c  
higher where  $f'$  is small and a large  
similar values of  $Y$ . If  $X$  causes  $Y$  (   
we would expect  $f$  and  $p_X(x)$  to be  
should be correlated with  $f'$ .



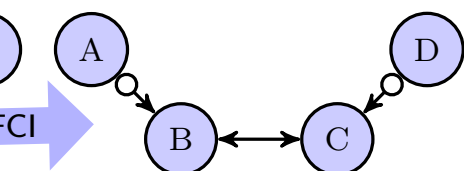
conditional indepen-

GS algorithm to ex-  
leading to much bet-

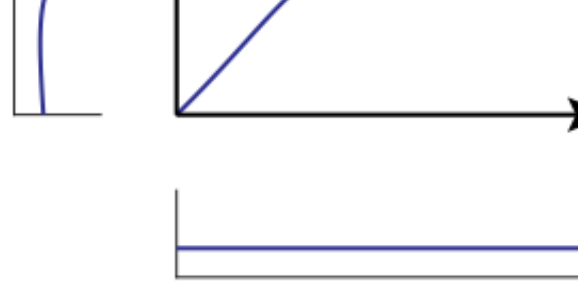
dependency equiva-  
able is a root node  
y to the FCI algo-  
algorithm to handle

ence class of Max-  
on of DAGs) since  
ation (figure 4b).

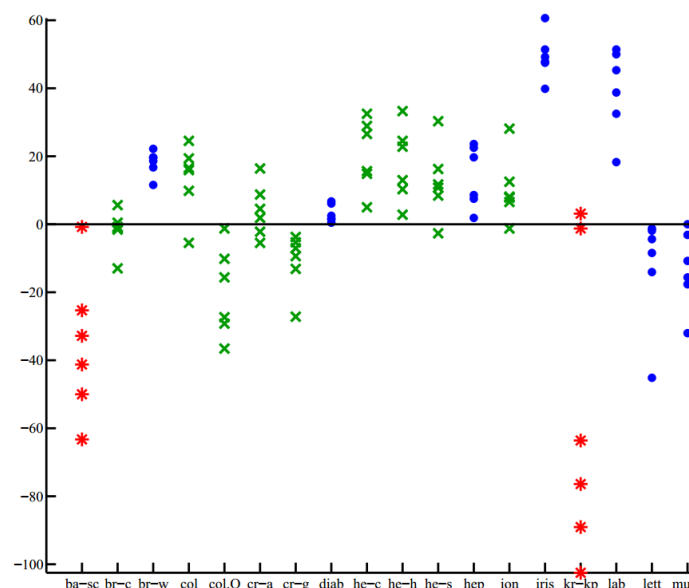
thm



cts of causal struc-  
pendence relations



**Figure 7:** Figure from [4]. The of mechanism and input can be g deterministic setting. If  $X \rightarrow Y$  th should be independent, but  $P(Y)$  Therefore, semi-supervised learning fit when trying to learn in the causal  $P(Y|X)$ ) but could help when learning rection.



**Learning what causality looks like**

Suppose we had  $M$  different causal

$$D = \{\{x_j, y_j\}_{j=1}^{N_i}, l$$

# A review



case polynomial (rather  
conditional independence

and FCI algorithm are  
5]

the advantage that they  
parametric, assumptions.  
between causal graphs

where  $l_i$  is a binary label that indicates  
in dataset  $i$ .

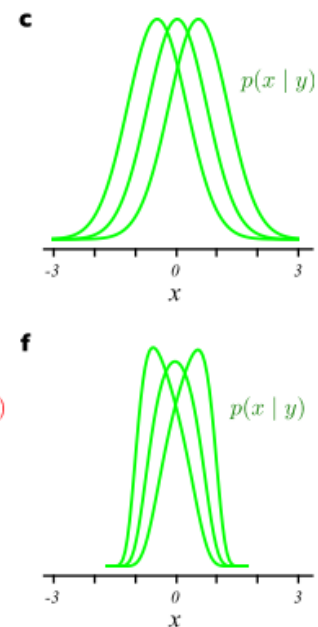
- Kernel mean embedding allows us to  
and transform it to a point in space
- We expect there to be differences  
between  $P(X)$ ,  $P(Y)$  and  $P(Y|X)$

## Algorithm

1. Let  $\mu$  be a kernel mean embedding

tion  $P$  into some Hilbert space

For example, between  
 noise models,  $Y =$   
 combinations of  $f$  and  $P(\epsilon)$



direction can be iden-  
 between  $X$  and  $Y$  is de-  
 $f(X) \Leftrightarrow X = g(Y)$ . For  
 distribution  $p_Y(y)$  will be  
 er region of  $X$  maps to  
 (but not if  $Y$  causes  $X$ )  
 independent and  $p_Y(y)$

$x)$

1. For each data set  $i = 1 \dots M$ , that approximates  $\mu(P(X)), \mu(P(Y))$
2. For each data set  $i = 1 \dots M$ , that approximates  $\mu(P(X)), \mu(P(Y))$
3. Apply a standard classification algorithm to  $Y$  or  $Y \rightarrow X$

## With more than two variables

- If you can come up with a set of distributions  $(P(Y), P(\epsilon), f)$ , that guarantee the identifiability of the structural equation model  $Y = f(X) + \epsilon$ , to get the conditions under which the model is identifiable [8].

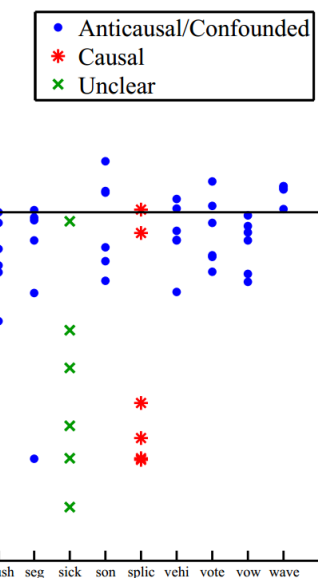
## Reference

- [1] Tom Claassen, Joris Mooij, and Peter D. van der Schaar. Sparse causal models is not identifiable. In *UAI*, 2010.
- [2] Povilas Danušis, Dominik Janzing, Bastian Steudel, and Bernhard Schölkopf. Inferring causal structure from data. In *UAI*, 2010.
- [3] Patrick Hoyer, Dominik Janzing, and Bernhard Schölkopf. Linear causal discovery with graphical models. In *NIPS*, 2009.
- [4] Dominik Janzing, Jonas Peters, Jonas Peters, Joris M. Mooij, and Bernhard Schölkopf. Causal and anticausal learning. In *UAI*, 2010.
- [5] Markus Kalisch, Martin Mächler, and Peter D. van der Schaar. Causal discovery using graphical models. In *UAI*, 2010.

→  $x$

$$= p_X(x)$$

idea of independence  
generalized to the non-  
when  $P(X)$  and  $P(Y|X)$   
and  $P(X|Y)$  are not.  
should yield no bene-  
direction, (estimating  
in the anti-causal di-



[6]

pairs data sets.

$$\{i\}_{i=1}^M$$

JSS, VV(li), 2012.

- [6] David Lopez-Paz, Krikamol Recht. The Randomized C Prepr. arXiv1409.4366, Sept
- [7] Judea Pearl. *Causality: m* ence. MIT Press, Cambridge
- [8] Jonas Peters, Joris Mooij, D hard Schölkopf. Causal disc tive noise models. *JMLR*, 15
- [9] Thomas S. Richardson and world intervention graphs (S counterfactual and graphical Center for the Statistics and versity of Washington Series
- [10] Ilya Shpitser and Judea Pea methods for the causal hiera 2008.
- [11] Peter Spirtes, Clark N Glym *Causation, prediction, and* press, 2000.
- [12] TS Verma. Graphical aspec nical report, 1993.
- [13] Jiji Zhang. On the comple for causal discovery in the founders and selection bia 172(16-17):1873–1896, Nov



indicates if  $X \rightarrow Y$  or  $Y \rightarrow X$

allows us to take a distribution  $P$   
in some Hilbert space.

addresses the relationships be-  
(1) for  $X \rightarrow Y$  and  $Y \rightarrow X$

**Item 2:**

adding that maps a distribu-



e.  
construct a feature vector  
 $(P(Y)), \mu(P(X, Y))$   
an algorithm to learn if  $X \rightarrow$

a condition, on the triple  
es identifiability for the bi-  
you can extend that result  
which the multivariate case is

## ces

and Tom Heskes. Learning  
NP-hard. In *UAI*, 2013.

Janzing, Joris Mooij, Jakob  
del, Kun Zhang, and Bern-  
terministic causal relations.

zing, and Joris Mooij. Non-  
additive noise models. In

eters, Eleni Sgouritsa, Kun  
d Bernhard Schölkopf. On  
ng. In *ICML*, 2012.

chler, Diego Colombo, Mar-  
r Bühlmann. Causal infer-  
s with the R package pcalq.

with the R package peay.  
I Muandet, and Benjamin  
causation Coefficient. *arXiv*  
September 2014.

*Models, reasoning and infer-*  
e, 2000.

Dominik Janzing, and Bern-  
covery with continuous addi-  
5:2009–2053, 2014.

James M. Robins. Single  
(SWIGs): a unification of the  
al approaches to causality.  
d the Social Sciences, Uni-  
s. Working Paper 128, 2013.

arl. Complete identification  
archy. *JMLR*, 9:1941–1979,

our, and Richard Scheines.  
*search*, volume 81. MIT

cts of causal models. Tech-

teness of orientation rules  
e presence of latent con-  
as. *Artificial Intelligence*,  
ember 2008.