
Constraint-based Causal Discovery for Non-Linear Structural Causal Models with Cycles and Latent Confounders

Patrick Forré
Informatics Institute
University of Amsterdam
The Netherlands
p.d.forre@uva.nl

Joris M. Mooij
Informatics Institute
University of Amsterdam
The Netherlands
j.m.mooij@uva.nl

Abstract

We address the problem of causal discovery from data, making use of the recently proposed causal modeling framework of *modular structural causal models (mSCM)* to handle cycles, latent confounders and non-linearities. We introduce σ -connection graphs (σ -CG), a new class of mixed graphs (containing undirected, bidirected and directed edges) with additional structure, and extend the concept of σ -separation, the appropriate generalization of the well-known notion of d-separation in this setting, to apply to σ -CGs. We prove the closedness of σ -separation under marginalisation and conditioning and exploit this to implement a test of σ -separation on a σ -CG. This then leads us to the first causal discovery algorithm that can handle non-linear functional relations, latent confounders, cyclic causal relationships, and data from different (stochastic) perfect interventions. As a proof of concept, we show on synthetic data how well the algorithm recovers features of the causal graph of modular structural causal models.

1 INTRODUCTION

Correlation does not imply causation. To go beyond spurious probabilistic associations and infer the asymmetric causal relations we need sufficiently powerful models. Structural causal models (SCMs), also known as structural equation models (SEMs), provide a popular modeling framework (see [12, 25, 26, 32]) that is up to this task. Still, the problem of causal discovery from data is notoriously hard. Theory and algorithms need to address several challenges like probabilistic settings, stability under interventions, combining observational and

interventional data, latent confounders and marginalisation, selection bias and conditioning, faithfulness violations, cyclic causation like feedback loops and pairwise interactions, and non-linear functional relations in order to go beyond artificial simulation settings and become successful on real-world data.

Several algorithms for causal discovery have been introduced over the years. For the acyclic case without latent confounders, numerous constraint-based [25,32] and score-based approaches [6, 14, 17] exist. More sophisticated constraint-based [5, 25, 32, 34] and score-based approaches [4, 7, 8, 10, 11] can deal with latent confounders in the acyclic case. For the linear cyclic case, most algorithms assume no latent confounders [18, 27, 29], though some of the more recent ones allow for those [19, 30]. To the best of our knowledge, no algorithms have yet been proposed for the general non-linear cyclic case.

In this work we present a novel conditional independence constraint-based causal discovery algorithm that—up to the knowledge of the authors—is the first causal discovery algorithm that addresses most of the previously mentioned problems at once, notably non-linearities, cycles, latent confounders, and multiple interventional data sets, only excluding selection bias and faithfulness violations.

For this to work we build upon the theory of *modular structural causal models (mSCM)* introduced in [12]. mSCMs form a general and convenient class of structural causal models that can deal with cycles and latent confounders. The measure-theoretically rigorous presentation opens the door for general non-linear measurable functions and any kind of probability distributions (e.g. mixtures of discrete or continuous ones). mSCM are provably closed under any combination of perfect interventions and marginalisations (see [12]).

Unfortunately, it is known that the direct generalization of the *d-separation criterion* (also called m- or m*-separation, see [9, 24, 25, 28, 33]), which relates the conditional independencies of the model to its underlying

graphical structure, does not apply in general if the structural equations are *non-linear* and the graph contains *cycles* (see [12, 31] or example 2.17).

Luckily, one key property of mSCMs is that the variables of the mSCM always entail the conditional independences implied by σ -separation, a non-naive generalization of the d/m/m*-separation (see [12]), which also works in the presence of cycles, non-linearities and latent confounders, and reduces to d-separation in the acyclic case.

To prove the σ -separation criterion, the authors of [12] have constructed an extensive theory for directed graphs with hyperedges. As a first contribution in this paper we give a simplified but equivalent definition of mSCMs plainly in terms of directed graphs and prove the σ -separation criterion directly under weaker assumptions.

As a second contribution we extend the definition of σ -separation to mixed graphs (including also bi- and undirected edges) by introducing additional structure. We will refer to this class of mixed graphs as σ -connection graphs (σ -CG), since they are inspired by the d-connection graphs introduced by [19]. We prove that σ -CGs and σ -separation are closed under marginalisation and conditioning, in analogy with the d-connection graphs (d-CG) from [19].

The work of [19] provides an elegant approach to causal discovery using a weighted SAT solver to find the causal graph that is most compatible with (weighted) conditional independences in the data, encoding the notion of d-separation into *answer set programming (ASP)*, a declarative programming language with an expressive syntax for implementing discrete or integer optimization problems. Our third contribution is to adapt the approach of [19] by replacing d-separation by σ -separation and d-CGs by σ -CGs. The results mentioned above will then ensure all the needed properties to make the adapted algorithm of [19] applicable to general mSCMs, i.e., to non-linear causal models with cycles and latent confounders, under the additional assumptions of *no selection bias* and of σ -faithfulness.

Finally, as a proof of concept, we will show the effectiveness of our proposed algorithm in recovering features of the causal graphs of mSCMs from simulated data.

2 THEORY

2.1 MODULAR STRUCTURAL CAUSAL MODELS

Structural causal or equation models (SCM/SEM) usually start with a set of variables $(X_v)_{v \in V}$ attached to a

graph G that satisfy (or are even defined by) equations of the form:

$$X_v = g_{\{v\}}(X_{\text{Pa}^G(v)}, E_v),$$

with a function $g_{\{v\}}$ and noise variable E_v attached to each node $v \in V$. Here $\text{Pa}^G(v)$ denotes the set of (direct causal) parents of v . In *linear* models the functions $g_{\{v\}}$ are linear functions, in *acyclic* models the nodes of V form an acyclic graph G , and under *causal sufficiency* the variables E_v are independent (i.e. “no latent confounders”). The functions $g_{\{v\}}$ are usually interpreted as *local causal mechanisms* that produce the values of X_v from the values of $X_{\text{Pa}^G(v)}$ and E_v . These local mechanisms $g_{\{v\}}$ are—in the causal setting—assumed to be stable even when one intervenes upon some of the variables, i.e. one makes a causal *local compatibility* assumption. One important observation now is that one can also consider the *global mechanism* g that maps the values of the latent variables $(E_v)_{v \in V}$ to the values of the observed variables $(X_v)_{v \in V}$. The assumption of acyclicity or invertible linearity will then guarantee the *global compatibility* of all these mechanisms $g_{\{v\}}$ and g . However, if we now abstain from assuming acyclicity or linearity, the global compatibility does not follow from the local compatibility anymore (see figure 1). So in a general consistent causal setting this needs to be *guaranteed or assumed*.

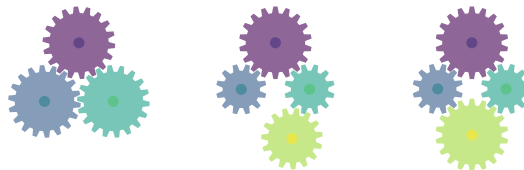


Figure 1: Gear analogy. Left: A cyclic mechanism that is locally compatible but not globally. Center: For acyclic mechanisms local compatibility implies global compatibility. Right: A cyclic mechanism that is locally and globally compatible. This shows that the assumption of global compatibility is needed when cycles are present.

The definition of modular structural causal models (mSCM) and the mentioned list of desirable properties basically follow directly from *causal* postulates:

Postulate 2.1 (Causal postulates). *The observed world appears as the projection of an extended world such that:*

1. *All latent and observed variables in this extended world are causally linked by a directed graph.*
2. *Every subsystem of this extended world can be expressed as the joint effect of its joint direct causes.*
3. *All these mechanisms are globally compatible.*

Special subsystems of interest are the loops of a graph.

Definition 2.2 (Loops). *Let $G = (V, E)$ be a directed graph (with or without cycles).*

1. A loop of G is a set of nodes $S \subseteq V$ such that for every two nodes $v_1, v_2 \in S$ there are two directed paths $v_1 \rightarrow \dots \rightarrow v_2$ and $v_2 \rightarrow \dots \rightarrow v_1$ in G with all the intermediate nodes also in S (if any). The single element sets $S = \{v\}$ are also considered as loops.
2. The strongly connected component of v in G is defined to be:

$$\text{Sc}^G(v) := \text{Anc}^G(v) \cap \text{Desc}^G(v),$$

the set of nodes that are both ancestors and descendants of v (including v itself).

3. Let $\mathcal{L}(G) := \{S \subseteq G \mid S \text{ a loop of } G\}$ be the loop set of G .

Remark 2.3. *Note that the loop set $\mathcal{L}(G)$ contains all single element loops $\{v\} \in \mathcal{L}(G)$, $v \in V$, as the smallest loops and all strongly connected components $\text{Sc}^G(v) \in \mathcal{L}(G)$, $v \in V$, as the largest loops, but also all non-trivial intermediate loops S with $\{v\} \subsetneq S \subsetneq \text{Sc}^G(v)$ inside the strongly connected components (if existent). If G is acyclic then $\mathcal{L}(G)$ only consists of the single element loops: $\mathcal{L}(G) = \{\{v\} \mid v \in V\}$.*

The definition of mSCM is made in such a way that it will automatically incorporate the causal postulates 2.1. In the following, all spaces are meant to be equipped with σ -algebras, forming standard measurable spaces, and all maps to be measurable.

Definition 2.4 (Modular Structural Causal Model, [12]). *A modular structural causal model (mSCM) by definition consists of:*

1. a set of nodes $V^+ = U \dot{\cup} V$, where elements of V correspond to observed variables and elements of U to latent variables,
2. an observation/latent space \mathcal{X}_v for every $v \in V^+$, $\mathcal{X} := \prod_{v \in V^+} \mathcal{X}_v$,
3. a product probability measure $\mathbb{P} := \mathbb{P}_U = \otimes_{u \in U} \mathbb{P}_u$ on the latent space $\prod_{u \in U} \mathcal{X}_u$,¹

¹The assumption of independence of the noise variables here is not to be confused with causal sufficiency. The noise variables here might have two or more child nodes and thus can play the role of latent confounders. The independence assumption here also does not restrict the model class. If they were dependent we would just consider them as one variable and use a different graph that encoded this.

4. a directed graph structure $G^+ = (V^+, E^+)$ with the properties:²

$$(a) V = \text{Ch}^{G^+}(U),$$

$$(b) \text{Pa}^{G^+}(U) = \emptyset,³$$

$$(c) \text{Ch}^{G^+}(u_1) \not\subseteq \text{Ch}^{G^+}(u_2) \text{ for every two distinct } u_1, u_2 \in U,³$$

where Ch^{G^+} and Pa^{G^+} stand for children and parents in G^+ , resp.,

5. a system of structural equations $g = (g_S)_{S \in \mathcal{L}(G^+): S \subseteq V}$:

$$g_S : \prod_{v \in \text{Pa}^{G^+}(S) \setminus S} \mathcal{X}_v \rightarrow \prod_{v \in S} \mathcal{X}_v,²⁴$$

that satisfy the following global compatibility conditions: For every nested pair of loops $S' \subseteq S \subseteq V$ of G^+ and every element $x_{\text{Pa}^{G^+}(S) \cup S} \in \prod_{v \in \text{Pa}^{G^+}(S) \cup S} \mathcal{X}_v$ we have the implication:

$$\begin{aligned} g_S(x_{\text{Pa}^{G^+}(S) \setminus S}) &= x_S \\ \implies g_{S'}(x_{\text{Pa}^{G^+}(S') \setminus S'}) &= x_{S'}, \end{aligned}$$

where $x_{\text{Pa}^{G^+}(S') \setminus S'}$ and $x_{S'}$ denote the corresponding components of $x_{\text{Pa}^{G^+}(S) \cup S}$.

The mSCM can be summarized by the tuple $M = (G^+, \mathcal{X}, \mathbb{P}, g)$.

Remark 2.5. *Given the mechanisms attached to the nodes $g_{\{v\}}$ the existence (and compatibility) of the other mechanisms g_S for non-trivial loops S can be guaranteed under certain conditions, e.g. trivially in the acyclic case, or if every cycle is contractive (see subsection 4.1), or more generally if the cycles are “uniquely solvable” (see [3, 12]).*

We are now going to define the actual random variables $(X_u)_{u \in V^+}$ attached to any mSCM.

Remark 2.6. *Let $M = (G^+, \mathcal{X}, \mathbb{P}, g)$ be a mSCM with $G^+ = (U \dot{\cup} V, E^+)$.*

²Even though we allow for selfloops in the directed graph G^+ we note that the causal mechanisms g_S will depend only on $\text{Pa}^{G^+}(S) \setminus S$, removing the self-dependence on the functional level. Otherwise, the functions g_S would not hold up to a direct interventional interpretation and one would want to replace them with functions that do.

³This assumption is only necessary to give the mSCM a “reduced/summarized” form. In practice one could allow for more latent variables and more complex latent structure.

⁴Note that the index set runs over all “observable loops” $S \subseteq V$, $S \in \mathcal{L}(G^+)$, which contains the usual single element sets $S = \{v\}$, which relate to the usual mechanisms $g_{\{v\}}$.

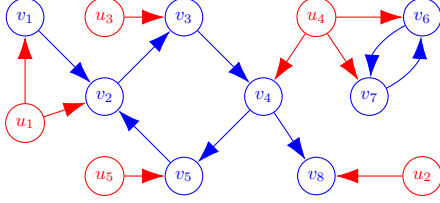


Figure 2: The graph G^+ of a modular structural causal model (mSCM). The observed variables $v_i \in V$ are in blue and the latent variables $u_j \in U$ are in red. We have the four observed strongly connected components: $\{v_1\}$, $\{v_2, v_3, v_4, v_5\}$, $\{v_6, v_7\}$, $\{v_8\}$.

1. The latent variables $(X_u)_{u \in U}$ are given by the canonical projections $E_u : \prod_{u' \in U} \mathcal{X}_{u'} \rightarrow \mathcal{X}_u$ and are jointly \mathbb{P} -independent (by 3). Sometimes we will write $(E_u)_{u \in U}$ instead of $(X_u)_{u \in U}$ to stress their interpretation as error/noise variables.
2. The observed variables $(X_v)_{v \in V}$ are inductively defined by:

$$X_v := g_{S,v}((X_w)_{w \in \text{Pa}^{G^+}(S) \setminus S}),$$

where $S := \text{Sc}^{G^+}(v) := \text{Anc}^{G^+}(v) \cap \text{Desc}^{G^+}(v)$ and where the second index v refers to the v -component of g_S . Note that the inductive definition is possible because when “aggregating” each of the biggest cycles $\text{Sc}^{G^+}(v)$ into one node then only an acyclic graph is left, which can be totally ordered.

3. By the compatibility condition for g we then have that for every $S \in \mathcal{L}(G^+)$ with $S \subseteq V$ the following equality holds:

$$X_S = g_S(X_{\text{Pa}^{G^+}(S) \setminus S}),$$

where we put $X_A := (X_v)_{v \in A}$ for subsets A .

As a consequence of the convenient definition 2.4 all the following desirable constructions (like marginalisations and interventions) are easily seen to be well-defined (for proofs see [12]). Note that already defining these constructions was a key challenge in the theory of causal models in the presence of cycles (see [3, 12], a.o.).

Definition 2.7. Let $M = (G^+, \mathcal{X}, \mathbb{P}, g)$ be a mSCM with $G^+ = (U \dot{\cup} V, E^+)$.

1. By plugging the functions g_S into each other we can define the marginalised mSCM M' w.r.t. a subset $W \subseteq V$. For example, when marginalizing out $W = \{w\}$ we can define (for the non-trivial case $w \in \text{Pa}^{G^+}(S) \setminus S$):

$$g_{S',v}(x_{\text{Pa}^{G'}(S') \setminus S'}) := g_{S,v}(x_{\text{Pa}^{G^+}(S) \setminus (S \cup \{w\})}, g_{\{w\}}(x_{\text{Pa}^{G^+}(w) \setminus \{w\}})),$$

where G' is the marginalised graph of G^+ , $S' \subseteq V'$ is any loop of G' and S the corresponding induced loop in G^+ .

2. For a subset $I \subseteq V$ and a value $x_I \in \prod_{v \in I} \mathcal{X}_v$ we define the intervened graph G' by removing all the edges from parents of I to I . We put $X_u^{\text{do}(x_I)} := X_u$ for $u \in U$ and $X_I^{\text{do}(x_I)} := x_I$ and inductively ($S := \text{Sc}^{G'}(v)$):

$$X_v^{\text{do}(x_I)} := g_{S,v}(X_{\text{Pa}^{G'}(S) \setminus S}^{\text{do}(x_I)}).$$

By selecting all functions g_S where S is still a loop in the intervened graph G' we get the post-interventional mSCM M' . These constructions give us all interventional distributions, e.g. (cf. [25]):

$$\mathbb{P}(X_A | \text{do}(x_I), X_B) := \mathbb{P}(X_A^{\text{do}(x_I)} | X_B^{\text{do}(x_I)}).$$

Instead of fixing $X_I^{\text{do}(x_I)}$ to a value x_I we could also specify a distribution \mathbb{P}'_I for it (“randomization”). In this way we define stochastic interventions $\text{do}(\xi_I)$ with an independent random variable ξ_I taking values in \mathcal{X}_I and get a $\mathbb{P}^{\text{do}(I)}$ similarly.

2.2 Σ -SEPARATION IN MSCMS AND Σ -CONNECTION GRAPHS

We now introduce σ -separation as a generalization of d-separation directly on the level of mixed graphs. To make the definition stable under marginalisation and conditioning we need to carry extra structure. The resulting graphs will be called σ -connection graphs (σ -CG), where the name is inspired by [19]. An example that shall clarify the difference between d- and σ -separation is given later in figure 2 and table 1.

Definition 2.8 (σ -Connection Graphs (σ -CG)). A σ -connection graph (σ -CG) is a mixed graph G with a set of nodes V and directed (\rightarrow), undirected (---) and bi-directed (\leftrightarrow) edges, together with an equivalence relation \sim_σ on V that has the property that every equivalence class $\sigma(v)$, $v \in V$, is a loop in the underlying directed graph structure: $\sigma(v) \in \mathcal{L}(G)$. Undirected self-loops ($v \text{---} v$) are allowed, (bi)-directed self-loops ($v \rightarrow v$, $v \leftrightarrow v$) are not.

In particular, every node is assigned to a unique fixed loop $\sigma(v)$ in G with $v \in \sigma(v)$ and two of such loops $\sigma(v_1)$, $\sigma(v_2)$ are either identical or disjoint. The reason for why we need such structure is illustrated in figure 5.

Definition 2.9 (σ -Open Path in a σ -CG). Let G be a σ -CG with set of nodes V and $Z \subseteq V$ a subset. Consider

a path π in G with $n \geq 1$ nodes:

$$v_1 \rightleftarrows \dots \rightleftarrows v_n.^5$$

The path will be called Z - σ -open if:

1. the endnodes $v_1, v_n \notin Z$, and
2. every triple of adjacent nodes in π that is of the form:

(a) collider:

$$v_{i-1} \rightleftarrows v_i \rightleftarrows v_{i+1},$$

satisfies $v_i \in Z$,

(b) (non-collider) left chain:

$$v_{i-1} \leftarrow v_i \rightleftarrows v_{i+1},$$

satisfies $v_i \notin Z$ or $v_i \in Z \cap \sigma(v_{i-1})$,

(c) (non-collider) right chain:

$$v_{i-1} \rightleftarrows v_i \rightarrow v_{i+1},$$

satisfies $v_i \notin Z$ or $v_i \in Z \cap \sigma(v_{i+1})$,

(d) (non-collider) fork:

$$v_{i-1} \leftarrow v_i \rightarrow v_{i+1},$$

satisfies $v_i \notin Z$ or $v_i \in Z \cap \sigma(v_{i-1}) \cap \sigma(v_{i+1})$,

(e) (non-collider) with undirected edge:

$$v_{i-1} \text{ --- } v_i \rightleftarrows v_{i+1},$$

$$v_{i-1} \rightleftarrows v_i \text{ --- } v_{i+1},$$

satisfies $v_i \notin Z$.

The difference between σ - and d -separation lies in the additional conditions involving $Z \cap \sigma(v_{i\pm 1})$. The intuition behind them is that the dependence structure *inside* a loop $\sigma(v_i)$ is so strong that non-colliders can only be blocked by conditioning if an edge is pointing out of the loop (see example 2.17 and table 1).

Similar to d -separation we can now define σ -separation in a σ -CG.

Definition 2.10 (σ -Separation in a σ -CG). *Let G be a σ -CG with set of nodes V . Let $X, Y, Z \subseteq V$ be subsets.*

1. We say that X and Y are σ -connected by Z or not σ -separated by Z if there exists a path π (with some $n \geq 1$ nodes) in G with one endnode in X and one endnode in Y that is Z - σ -open. In symbols this statement will be written as follows:

$$X \underset{G}{\overset{\sigma}{\not\perp}} Y | Z.$$

⁵The stacked edges are meant to be read as an ‘‘OR’’ at each place independently. We also allow for repeated nodes in the paths.

2. Otherwise, we will say that X and Y are σ -separated by Z and write:

$$X \underset{G}{\overset{\sigma}{\perp}} Y | Z.$$

Remark 2.11. 1. The finest/trivial σ -CG structure of a mixed graph G is given by $\sigma(v) := \{v\}$ for all $v \in V$. In this way σ -separation in G coincides with the usual notion of d -separation in a d -connection graph (d -CG) G (see [19]). We will take this as the definition of d -separation and d -CG in the following.

2. The coarsest σ -CG structure of a mixed graph G is given by $\sigma(v) := \text{Sc}^G(v) := \text{Anc}^G(v) \cap \text{Desc}^G(v)$ w.r.t. the underlying directed graph. Note that the definition of strongly connected component totally ignores the bi- and undirected edges of the σ -CG.
3. In any σ -CG we will always have that σ -separation implies d -separation, since every Z - d -open path is also Z - σ -open because $\{v\} \subseteq \sigma(v)$.
4. If a σ -CG G is acyclic (implying $\text{Sc}^G(v) = \{v\}$) then σ -separation coincides with d -separation.

We now want to ‘‘hide’’ or marginalise out the latent nodes $u \in U$ from the graph of any mSCM and represent their induced dependence structure with bidirected edges.

Definition 2.12 (Induced σ -CG of a mSCM). *Let $M = (G^+, \mathcal{X}, \mathbb{P}, g)$ be a mSCM with $G^+ = (U \dot{\cup} V, E^+)$. The induced σ -CG G of M , also referred to as the causal graph G of M is defined as follows:*

1. The nodes of G are all $v \in V$, i.e. all observed nodes of G^+ .
2. G contains all the directed edges of G^+ whose endnodes are both in V , i.e. observed.
3. G contains the bidirected edge $v \leftrightarrow w$ with $v, w \in V$ if and only if $v \neq w$ and there exists a $u \in U$ with $v, w \in \text{Ch}^{G^+}(u)$, i.e. v and w have a common latent confounder.
4. G contains no undirected edges.
5. We put $\sigma(v) := \text{Sc}^G(v) = \text{Anc}^G(v) \cap \text{Desc}^G(v)$.

Remark 2.13. Caution must be applied when going from G^+ to G : It is possible that three observed nodes v_1, v_2, v_3 have one joint latent common cause u_1 , which can be read off G^+ . This information will get lost when going from G^+ to G , as we will represent this with three bidirected edges. G will nonetheless capture the conditional independence relations (see Theorem 2.14).

We now present the most important ingredient for our constraint-based causal discovery algorithm, namely a generalized directed global Markov property that relates the underlying causal graph (σ -CG) G of any mSCM M to the conditional independencies of the observed random variables $(X_v)_{v \in V}$ via a σ -separation criterion.

Theorem 2.14 (σ -Separation Criterion, see Corollary B.3). *The observed variables $(X_v)_{v \in V}$ of any mSCM M satisfy the σ -separation criterion w.r.t. the induced σ -CG G . In other words, for all subsets $W, Y, Z \subseteq V$ we have the implication:*

$$W \perp\!\!\!\perp_G^\sigma Y \mid Z \implies X_W \perp\!\!\!\perp_{\mathbb{P}} X_Y \mid X_Z.$$

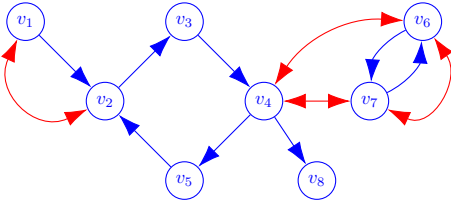


Figure 3: The induced σ -connection graph (causal graph) of the modular structural causal model (mSCM) of figure 2, with σ -equivalence classes $\{\{v_1\}, \{v_2, v_3, v_4, v_5\}, \{v_6, v_7\}, \{v_8\}\}$.

Table 1: d- and σ -separation in the σ -connection graph G from figure 3. σ -separation implies d-separation, but i.g. d-separation encodes more conditional independencies than σ -separation:

| d-separation | σ -separation |
|--|---|
| $\{v_2\} \perp\!\!\!\perp_G^d \{v_4\} \mid \{v_3, v_5\}$ | $\{v_2\} \not\perp\!\!\!\perp_G^\sigma \{v_4\} \mid \{v_3, v_5\}$ |
| $\{v_1\} \perp\!\!\!\perp_G^d \{v_6\}$ | $\{v_1\} \perp\!\!\!\perp_G^\sigma \{v_6\}$ |
| $\{v_1\} \perp\!\!\!\perp_G^d \{v_6\} \mid \{v_3, v_5\}$ | $\{v_1\} \not\perp\!\!\!\perp_G^\sigma \{v_6\} \mid \{v_3, v_5\}$ |
| $\{v_1\} \not\perp\!\!\!\perp_G^d \{v_8\}$ | $\{v_1\} \not\perp\!\!\!\perp_G^\sigma \{v_8\}$ |
| $\{v_1\} \perp\!\!\!\perp_G^d \{v_8\} \mid \{v_3, v_5\}$ | $\{v_1\} \not\perp\!\!\!\perp_G^\sigma \{v_8\} \mid \{v_3, v_5\}$ |
| $\{v_1\} \perp\!\!\!\perp_G^d \{v_8\} \mid \{v_4\}$ | $\{v_1\} \perp\!\!\!\perp_G^\sigma \{v_8\} \mid \{v_4\}$ |

If we want to infer the causal graph (σ -CG) G of a mSCM from data with help of conditional independence tests in practice we usually need to assume also the reverse implication of the σ -separation criterion from 2.14 for the observational distribution $\mathbb{P}(X_V)$ and the relevant interventional distributions $\mathbb{P}(X_V \mid \text{do}(\xi_I))$ (see 2.7). This will be called σ -faithfulness.

Definition 2.15 (σ -faithfulness). *We will say that the tuple (G, \mathbb{P}) is σ -faithful if for every three subsets $W, Y, Z \subseteq V$ we have the equivalence:*

$$W \perp\!\!\!\perp_G^\sigma Y \mid Z \iff X_W \perp\!\!\!\perp_{\mathbb{P}} X_Y \mid X_Z.$$

Remark 2.16 (Strong σ -completeness, cf. [22]). *We do believe that the generic non-linear mSCM is σ -faithful to G , as the needed conditional dependence structure in all our simulated (sufficiently non-linear) cases was observed (cf. example 2.17). But proving such a strong σ -completeness result is difficult (and to our knowledge only done for multinomial and linear Gaussian DAG models, see [22]) and left for future research. Further note that the class of linear models, which even follow d-separation, would i.g. not be considered σ -faithful. Since linear models are of measure zero in the bigger class of general mSCMs this would not contradict our conjecture.*

Example 2.17. *Consider a directed four-cycle, e.g. the subgraph $\{v_2, v_3, v_4, v_5\}$ from figure 2, where all other observed nodes are assumed to be absent. Consider only the non-linear causal mechanisms given by ($i = 3, 4, 5$):*

$$\begin{aligned} g_{\{v_2\}}(X_5, E_1) &:= \tanh(0.9 \cdot X_5 + 0.5) + E_1, \\ g_{\{v_i\}}(X_{i-1}, E_i) &:= \tanh(0.9 \cdot X_{i-1} + 0.5) + E_i, \end{aligned}$$

where E_1, E_3, E_4, E_5 are assumed to be independent. The equations $X_i = g_{\{v_i\}}(X_{i-1}, E_i)$ and the one for X_2 will imply the conditional dependence

$$X_2 \not\perp\!\!\!\perp_{\mathbb{P}} X_4 \mid (X_3, X_5),$$

which can be checked by computations and/or simulations. As one can read off table 1, d-separation fails to express the dependence, in contrast to σ -separation, which captures it correctly.

2.3 MARGINALISATION AND CONDITIONING IN Σ -CONNECTION GRAPHS

Inspired by [19] we will define marginalisation and conditioning operations on σ -connection graphs (σ -CG) and prove the closedness of σ -separation (and thus its criterion) under these operations. These are key results to extend the algorithm of [19] to the setting of mSCMs.

Definition 2.18 (Marginalisation of a σ -CG). *Let G be a σ -CG with set of nodes V and $w \in V$, $W := \{w\}$. We define the marginalised σ -CG G^W with set of nodes $V \setminus W$ via the rules for $v_1, v_2 \in V \setminus W$:*

$v_1 \overset{a}{\circ} \overset{b}{\circ} v_2 \in G^W$ with arrow heads/tails a and b if and only if there exists:

1. $v_1 \overset{a}{\circ} \overset{b}{\circ} v_2$ in G , or
2. $v_1 \overset{a}{\circ} w \overset{q}{\circ} \overset{b}{\circ} v_2$ in G , or
3. $v_1 \overset{a}{\circ} \overset{q}{\circ} w \overset{b}{\circ} v_2$ in G , or
4. $v_1 \overset{a}{\circ} \rightarrow w \leftarrow \overset{b}{\circ} v_2$ in G .

Note that directed paths in G have no colliders, so loops in G map to loops in G^W (if not empty). Thus we have the induced σ -CG structure \sim_σ on G^W .

Definition 2.19 (Conditioning of a σ -CG). *Let G be a σ -CG with set of nodes V and $c \in V$, $C := \{c\}$. We define the conditioned σ -CG G_C with set of nodes $V \setminus C$ via the rules for $v_1, v_2 \in V \setminus C$:*

$v_1 \overset{\sigma}{\leftarrow} v_2 \in G_C$ if and only if there exists:

1. $v_1 \overset{\sigma}{\leftarrow} v_2$ in G , or
2. $v_1 \overset{\sigma}{\rightarrow} c \overset{\sigma}{\leftarrow} v_2$ in G , or
3. $v_1 \overset{\sigma}{\leftarrow} c \overset{\sigma}{\leftarrow} v_2$ in G , $\sigma(v_1) = \sigma(c)$, or
4. $v_1 \overset{\sigma}{\rightarrow} c \overset{\sigma}{\rightarrow} v_2$ in G , $\sigma(c) = \sigma(v_2)$, or
5. $v_1 \overset{\sigma}{\leftarrow} c \overset{\sigma}{\rightarrow} v_2$ in G , $\sigma(v_1) = \sigma(c) = \sigma(v_2)$.

Note that directed paths in $\sigma(v)$ in G condition to directed paths, so loops in $\sigma(v)$ in G map to loops in G_C (if not empty). Thus we have a well-defined induced σ -CG structure \sim_{σ} on G_C .

The proofs of the following theorem, stating the closedness of σ -separation under marginalisation and conditioning, can be found in the supplementary material (Theorem A.1 and Theorem A.2). See also figure 5.

Theorem 2.20. *Let G be a σ -CG with set of nodes V and $X, Y, Z \subseteq V$ any subsets. For any nodes $w, c \in V \setminus (X \cup Y \cup Z)$, $W := \{w\}$, $C := \{c\}$, we then have the equivalences:*

$$\begin{aligned} X \underset{G^W}{\perp\!\!\!\perp} Y \mid Z &\iff X \underset{G}{\perp\!\!\!\perp} Y \mid Z, \\ X \underset{G_C}{\perp\!\!\!\perp} Y \mid Z &\iff X \underset{G}{\perp\!\!\!\perp} Y \mid Z \cup C. \end{aligned}$$

Corollary 2.21. *Let G be a σ -CG with set of nodes V and $X, Y, Z \subseteq V$ pairwise disjoint subsets and $W := V \setminus (X \cup Y \cup Z)$. Then we have the equivalence:*

$$X \underset{G}{\perp\!\!\!\perp} Y \mid Z \iff X \underset{G_Z^W}{\perp\!\!\!\perp} Y,$$

where G_Z^W is any σ -CG with set of nodes $X \cup Y$ obtained by marginalising out all the nodes from W and conditioning on all the nodes from Z in any order. This means that if $X = \{x\}$ and $Y = \{y\}$ then x and y are σ -separated by Z in G if and only if x and y are not connected by any edge in the σ -CG G_Z^W .

It is also tempting to introduce an intervention operator directly on the level of σ -CGs. However, since the interplay between conditioning and intervention is complicated (e.g. they do not commute i.g.) we do not investigate this further in this paper. The intervention operator on the level of mSCMs will be enough for our purposes as we assume no pre-interventional selection bias and then only encounter observational or post-interventional conditioning, which is covered by our framework.

3 ALGORITHM

In this section, we propose an algorithm for causal discovery that is based on the theory in the previous section. Given that theory, our proposed algorithm is a straightforward modification of the algorithm by [19]. The main idea is to formulate the causal discovery problem as an optimization problem that aims at finding the causal graph that best matches the data at hand. This is done by encoding the rules for conditioning, marginalisation, and intervention (see below) on a σ -CG into Answer Set Programming (ASP), an expressive declarative programming language based on stable model semantics that supports optimization [15, 20]. The optimization problem can then be solved by employing an off-the-shelf ASP solver.

3.1 CAUSAL DISCOVERY WITH Σ -CONNECTION GRAPHS

Let $M = (G^+, \mathcal{X}, \mathbb{P}, g)$ be a mSCM with $G^+ = (U \dot{\cup} V, E^+)$ and $I \subseteq V$ a subset. Consider a (stochastic) perfect intervention $\text{do}(\xi_I)$ that enforces $X_I = \xi_I$ for an independent random variable ξ_I taking values in \mathcal{X}_I . Denote the (unique) induced distribution of the intervened mSCM $M_{\text{do}(\xi_I)}$ by $\mathbb{P}^{\text{do}(I)}$, and the causal graph (i.e., induced σ -CG of the intervened mSCM on the observed variables) by $G_{\text{do}(I)} = (G_{\text{do}(I)}^+)^U$.

Under $\mathbb{P}^{\text{do}(I)}$, the observed variables $(X_v)_{v \in V}$ satisfy the σ -separation criterion w.r.t. $G_{\text{do}(I)}$ by Theorem 2.14. For the purpose of causal discovery, we will in addition assume σ -faithfulness (Definition 2.15), i.e., that each conditional independence between observed variables is due to a σ -separation in the causal graph. Taken together, and by applying Corollary 2.21, we get for all subsets $W, Y, Z \subseteq V$ the equivalences:

$$\begin{aligned} X_W \underset{\mathbb{P}^{\text{do}(I)}}{\perp\!\!\!\perp} X_Y \mid X_Z &\iff W \underset{G_{\text{do}(I)}}{\perp\!\!\!\perp} Y \mid Z \\ &\iff W \underset{(G_{\text{do}(I)})_Z}{\perp\!\!\!\perp} Y. \end{aligned} \quad (1)$$

If $W = \{w\}$ and $Y = \{y\}$ consist of a single node each, the latter can be easily checked by testing whether w is non-adjacent to y in $(G_{\text{do}(I)})_Z^{V \setminus W \cup Y \cup Z}$.

3.2 CAUSAL DISCOVERY AS AN OPTIMIZATION PROBLEM

Following [19], we formulate causal discovery as an optimization problem where a certain loss function is optimized over possible causal graphs. This loss function sums the weights of all the inputs that are violated assuming a certain underlying causal graph.

The input for the algorithm is a list $S = ((w_j, y_j, Z_j, I_j, \lambda_j))_{j=1}^n$ of weighted conditional independence statements. Here, the weighted statement $(w_j, y_j, Z_j, I_j, \lambda_j)$ with $w_j, y_j \in V$, $Z_j, I_j \subseteq V$, and $\lambda_j \in \mathbb{R} := \mathbb{R} \cup \{-\infty, +\infty\}$ encodes that $X_{w_j} \perp\!\!\!\perp_{\mathbb{P}^{\text{do}(X_{I_j})}} X_{y_j} \mid X_{Z_j}$ holds with confidence λ_j , where a finite value of λ_j gives a ‘‘soft constraint’’ and a value of $\lambda_j = \pm\infty$ imposes a ‘‘hard constraint’’. Positive weights encode that we have empirical support *in favor* of the independence, whereas negative weights encode empirical support *against* the independence (in other words, in favor of *dependence*).

As in [19], we define a loss function that measures the amount of evidence *against* the hypothesis that the data was generated by an mSCM with causal graph G , by simply summing the absolute weights of the input statements that conflict with G under the σ -Markov and σ -faithfulness assumptions:

$$\begin{aligned} \mathcal{L}(G, S) \\ := \sum_{(w_j, y_j, Z_j, I_j, \lambda_j) \in S} \lambda_j (\mathbb{1}_{\lambda_j > 0} - \mathbb{1}_{w_j \perp\!\!\!\perp_{\sigma_{\text{do}(I_j)}} y_j \mid Z_j}) \end{aligned} \quad (2)$$

where $\mathbb{1}$ is the indicator function. This loss function differs from the one used in [19] in that we use σ -separation instead of d-separation. Causal discovery can now be formulated as the optimization problem:

$$G^* = \arg \min_{G \in \mathbb{G}(V)} \mathcal{L}_R(G, S) \quad (3)$$

where $\mathbb{G}(V)$ denotes the set of all possible causal graphs with variables V .

The optimization problem (3) may have multiple optimal solutions, because the underlying causal graph may not be identifiable from the inputs. Nonetheless, some of the features of the causal graph (e.g., the presence or absence of a certain directed edge) may still be identifiable. We employ the method proposed by [21] for scoring the confidence that a certain feature f is present by calculating the difference between the optimal losses under the additional hard constraints that the feature f is present vs. that the feature f is absent in G .

In our experiments, we will use the weights proposed in [21]: $\lambda_j = \log p_j - \log \alpha$, where p_j is the p-value of a statistical test with independence as null hypothesis, and α is a significance level (e.g., 1%). This test is performed on the data measured in the context of the (stochastic) perfect intervention I_j . These weights have the desirable property that independences typically get a lower absolute weight than strong dependences. For the conditional independence test, we use a standard partial

correlation test after marginal rank-transformation of the data so as to obtain marginals with standard-normal distributions.

3.3 FORMULATING THE OPTIMIZATION PROBLEM IN ASP

In order to calculate the loss function (2), we make use of Corollary 2.21 to reduce the σ -separation test to a simple non-adjacency test in a conditioned and marginalised σ -CG, as in (1). We do this by encoding σ -CGs, Theorem 2.20 and the marginalisation and conditioning operations on σ -CGs (Definitions 2.18 and 2.19) in ASP. The details of the encoding are provided in the Supplementary Material.⁶ The optimization problem in (3) can then be solved straightforwardly by running an off-the-shelf ASP solver with as input the encoding and the weighted independence statements.

A more precise statement of the following result is provided in the Supplementary Material. The proof is basically the same as the one given in [21].

Theorem 3.1. *The algorithm for scoring features f is sound for oracle inputs and asymptotically consistent under mild assumptions.*

4 EXPERIMENTS

4.1 CONSTRUCTING MSCMS AND SAMPLING FROM MSCMS

To construct a modular structural causal model (mSCM) in practice we need to specify the compatible system of functions $(g_S)_{S \in \mathcal{L}(G)}$. The following Theorem is helpful (and a direct consequence of Banach’s fixed point theorem).

Theorem 4.1. *Consider the functions $g_{\{v\}}$ for the trivial loops $\{v\} \in \mathcal{L}(G)$, $v \in V$ and assume the following contractivity condition:*

For every non-trivial loop $S \in \mathcal{L}(G)$ and for every value $x_{\text{Pa}^{G^+}(S) \setminus S}$ the multi-dimensional function:

$$x'_S \mapsto \left(g_{\{v\}}(x'_{S \cap \text{Pa}^{G^+}(v) \setminus \{v\}}, x_{\text{Pa}^{G^+}(v) \setminus S}) \right)_{v \in S}$$

is a contraction, i.e. Lipschitz continuous with Lipschitz constant $L(x_{\text{Pa}^{G^+}(S) \setminus S}) < 1$ w.r.t. a suitable norm $\|\cdot\|$.

Then all the functions g_S for the non-trivial loops $S \in \mathcal{L}(G)$ exist, are unique and $g = (g_S)_{S \in \mathcal{L}(G)}$ forms a

⁶The full source code for the algorithm and to reproduce our experiments is available under an open source license from <https://github.com/caus-am/sigmasep>.

globally compatible system.

More constructively, for every value $x_{\text{Pa}G^+(S)\setminus S}$ and initialization $x_S^{(0)}$ the iteration scheme (using vector notations):

$$x_S^{(t+1)} := (g_{\{v\}})_{v \in S}(x_S^{(t)}, x_{\text{Pa}G^+(S)\setminus S})$$

converges to a unique limit vector x_S (for $t \rightarrow \infty$ and independent of $x_S^{(0)}$). g_S is then given by putting:

$$g_S(x_{\text{Pa}G^+(S)\setminus S}) := x_S.$$

This provides us with a method for constructing very general non-linear mSCMs (e.g. neural networks, see Section C in Supplementary Material) and to sample from them: by sampling x_U from the external distribution and then apply the above iteration scheme until convergence for all loops, yielding the limit x_V as one data point.

4.2 RESULTS ON SYNTHETIC DATA

In our experiments we will—due to computational restrictions—only allow for $d = 5$ observed nodes and $k = 2$ additional latent confounders. We sample edges independently with a probability of $p = 0.3$. We model the non-linear function $g_{\{v\}}$ as a neural network with tanh activation, bias terms that have a normal distribution with mean -0.5 and standard deviation 0.2 , and weights sampled uniformly from the L1-unit ball to satisfy the contraction condition of Theorem 4.1 (see also Supplementary Material, Section C). We simulate 0–5 single-variable interventions with random (unique) targets. For each intervened model we sample from standard-normal noise terms and compute the observations. To also detect weak dependencies in cyclic models we allow for $n = 10^4$ samples in each such model for each allowed intervention. We then run all possible conditional independence tests between every pair of single nodes and calculate their p -values. We used $\alpha = 10^{-3}$ as the threshold between dependence and independence. For computational reasons we restrict to partial correlation tests of marginal Gaussian rank-transforms of the data. These tests are then fed into the ASP solver together with our encoding of the optimization problem (3). We query the ASP solver for the confidence for the absence or presence of each possible directed and bidirected edge. We simulate 300 models and aggregate results, using the confidence scores to compute ROC- and PR-curves for features. Figure 4 shows that, as expected, our algorithm recovers more directed edges of the underlying causal graph in the simulation setting as the number of single-variable interventions increases. More results (ROC- and PR-curves for directed edges and con-

founders for different numbers of single-variable interventions and for different encodings) are provided in the Supplementary Material.

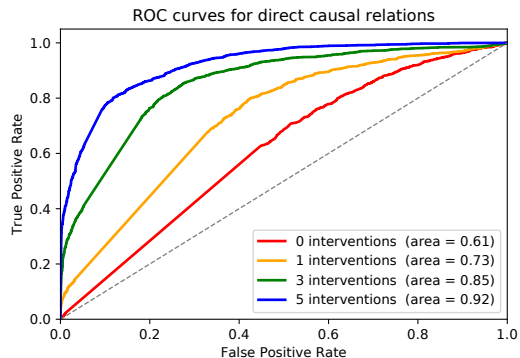


Figure 4: The ROC curves for identifying directed edges. See also figures 6 and 7 in the Supplementary Material.

5 CONCLUSION

We introduced σ -connection graphs (σ -CG) as a generalization of the d-connection graphs (d-CG) of [19] and extended the notion of σ -separation that was introduced in [12] to σ -CGs. We showed how σ -CGs behave under marginalisation and conditioning. This provides a graphical representation of how conditional independencies of modular structural causal models (mSCMs) behave under these operations. We provided a sufficient condition that allows constructing mSCMs and sampling from them. We extended the algorithm of [19] to deal with the more generally applicable notion of σ -separation instead of d-separation, thereby obtaining the first algorithm for causal discovery that can deal with cycles, nonlinearities, latent confounders and a combination of data sets corresponding to observational and different interventional settings. We illustrated the effectiveness of the algorithm on simulated data. In this work, we restricted attention to (stochastic) perfect (“surgical”) interventions, but a straightforward extension to deal with other types of interventions and to generalize the idea of randomized controlled trials can be obtained by applying the JCI framework [23]. In future work we wish to improve our algorithm to also handle selection bias, become more scalable and apply it to real world data sets.

Acknowledgements

This work was supported by the European Research Council (ERC) under the European Union’s Horizon 2020 research and innovation programme (grant agreement 639466).

References

- [1] F. Barthe, O. Guédon, S. Mendelson, and A. Naor. A probabilistic approach to the geometry of the l_p^p -ball. *Ann. Probab.*, 33(2):480–513, 2005.
- [2] V.I. Bogachev. *Measure Theory. Vol. I, II.* Springer, 2007.
- [3] S. Bongers, J. Peters, B. Schölkopf, and J. M. Mooij. Theoretical aspects of cyclic structural causal models. *arXiv.org preprint*, arXiv:1611.06221v2 [stat.ME], 2018.
- [4] T. Claassen and T. Heskes. Bayesian probabilities for constraint-based causal discovery. In *IJCAI-13*, pages 2992–2996, 2013.
- [5] T. Claassen, J.M. Mooij, and T. Heskes. Learning Sparse Causal Models is not NP-hard. In *UAI-13*, pages 172–181, 2013.
- [6] G. F. Cooper and E. Herskovits. A Bayesian method for the induction of probabilistic networks from data. *Machine Learning*, 9:309–347, 1992.
- [7] M. Drton, M. Eichler, and T. Richardson. Computing maximum likelihood estimates in recursive linear models with correlated errors. *Journal of Machine Learning Research*, 10:2329–2348, 2009.
- [8] R. Evans and T. Richardson. Maximum likelihood fitting of acyclic directed mixed graphs to binary data. In *UAI-10*, 2010.
- [9] R.J. Evans. Graphs for margins of Bayesian networks. *Scand. J. Stat.*, 43(3):625–648, 2016.
- [10] R.J. Evans. Margins of discrete Bayesian networks. *arXiv:1501.02103*, pages 1–41, 2017. Submitted to *Annals of Statistics*.
- [11] R.J. Evans and T.S. Richardson. Markovian acyclic directed mixed graphs for discrete data. *Ann. Statist.*, 42(4):1452–1482, 08 2014.
- [12] P. Forré and J. M. Mooij. Markov properties for graphical models with cycles and latent variables. *arXiv:1710.08775*, 2017.
- [13] M. Gebser, R. Kaminski, B. Kaufmann, and T. Schaub. *Clingo = ASP + control: Extended report.* Technical report, University of Potsdam, 2014. <http://www.cs.uni-potsdam.de/wv/pdfformat/gekakasc14a.pdf>.
- [14] D. Geiger and D. Heckerman. Learning Gaussian networks. In *UAI-94*, pages 235–243, 1994.
- [15] M. Gelfond. Answer sets. In *Handbook of Knowledge Representation*, pages 285–316. 2008.
- [16] I. Goodfellow, Y. Bengio, and A. Courville. *Deep Learning.* MIT Press, 2016.
- [17] D. Heckerman, C. Meek, and G. Cooper. A Bayesian approach to causal discovery. In C. Glymour and G. F. Cooper, editors, *Computation, Causation, and Discovery*, pages 141–166. MIT Press, 1999.
- [18] A. Hyttinen, F. Eberhardt, and P.O. Hoyer. Causal discovery for linear cyclic models. In *Proceedings of the Fifth European Workshop on Probabilistic Graphical Models*, 2010.
- [19] A. Hyttinen, F. Eberhardt, and M. Järvisalo. Constraint-based causal discovery: Conflict resolution with answer set programming. In *UAI-14*, pages 340–349, 2014.
- [20] V. Lifschitz. What is Answer Set Programming? In *AAAI Conference on Artificial Intelligence*, pages 1594–1597, 2008.
- [21] S. Magliacane, T. Claassen, and J.M. Mooij. Ancestral causal inference. In *NIPS-16*, pages 4466–4474. 2016.
- [22] C. Meek. Strong completeness and faithfulness in Bayesian networks. In *UAI-95*, pages 411–418, 1995.
- [23] J. M. Mooij, S. Magliacane, and T. Claassen. Joint causal inference from multiple contexts. *arXiv.org preprint*, <https://arxiv.org/abs/1611.10351v3> [cs.LG], March 2018.
- [24] J. Pearl. Fusion, propagation and structuring in belief networks. Technical Report 3, UCLA Computer Science Department, 1986. Technical Report 850022 (R-42).
- [25] J. Pearl. *Causality: Models, reasoning, and inference.* Cambridge University Press, 2nd edition, 2009.
- [26] J. Peters, D. Janzing, and B. Schölkopf. *Elements of Causal Inference: Foundation and Learning Algorithms.* MIT press, 2017.
- [27] T. Richardson. A discovery algorithm for directed cyclic graphs. In *UAI-96*, pages 454–461. 1996.
- [28] T. Richardson. Markov properties for acyclic directed mixed graphs. *Scand. J. Stat.*, 30(1):145–157, 2003.
- [29] T. Richardson and P. Spirtes. Automated discovery of linear feedback models. In C. Glymour and G. F. Cooper, editors, *Computation, Causation, and Discovery*, pages 253–304. MIT Press, 1999.
- [30] D. Rothenhäusler, C. Heinze, J. Peters, and N. Meinshausen. BACKSHIFT: Learning causal cyclic graphs from unknown shift interventions. In *NIPS-15*, pages 1513–1521. 2015.
- [31] P. Spirtes. Directed cyclic graphical representations of feedback models. In *UAI-95*, pages 491–499, 1995.
- [32] P. Spirtes, C. Glymour, and R. Scheines. *Causation, Prediction, and Search.* MIT press, 2000.
- [33] T.S. Verma and J. Pearl. Causal Networks: Semantics and Expressiveness. *UAI-90*, 4:69–76, 1990.
- [34] J. Zhang. On the completeness of orientation rules for causal discovery in the presence of latent confounders and selection bias. *Artificial Intelligence*, 172(16-17):1873–1896, 2008.

SUPPLEMENTARY MATERIAL

A Σ -CG UNDER MARGINALISATION AND CONDITIONING

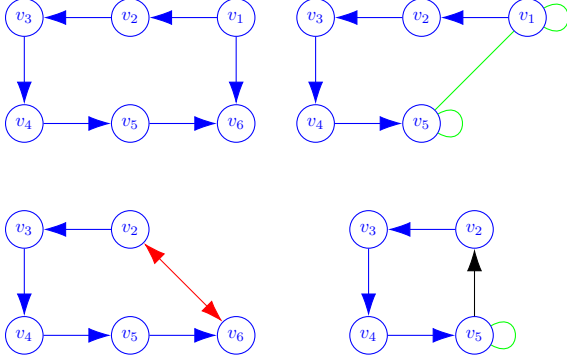


Figure 5: A directed *acyclic* graph G (top left) as the σ -connection graph (σ -CG) with the σ -equivalence classes $\{\{v_1\}, \dots, \{v_6\}\}$. We have the σ -separation: $\{v_3\} \perp_G^{\sigma} \{v_5\} \mid \{v_2, v_4, v_6\}$. The combination of marginalizing out v_1 and conditioning on v_6 introduces cycles. Without keeping track of above σ -equivalence classes we would not get the corresponding σ -separation: $\{v_3\} \perp_{G_{\{v_1\}}^{\sigma}}^{\sigma} \{v_5\} \mid \{v_2, v_4\}$ in the bottom right σ -CG.

Theorem A.1 (σ -Separation under Marginalisation). *Let G be a σ -CG with set of nodes V and $W, X, Y, Z \subseteq V$ subsets with $W = \{w\}$ and $w \notin X \cup Y \cup Z$. Then we have the equivalence:*

$$X \perp_G^{\sigma} Y \mid Z \iff X \perp_{G^W}^{\sigma} Y \mid Z.$$

Proof. If $\pi = x \dots y$ is a Z - σ -open path in G then every occurrence of w in π is as a non-collider. If we have $\dots v \rightarrow w \rightarrow \dots$ in π and $v \notin Z$ then marginalising out w keeps π^W Z - σ -open in G^W . If $v \in Z$ then $v \in \sigma(w)$ by the Z - σ -openness. Since $\sigma(w) \in \mathcal{L}(G)$ is a loop we find elements $v_i \in \sigma(w)$, $r \geq 0$, $i = 1, \dots, r$, and a path

$$\dots v \leftarrow v_1 \leftarrow \dots \leftarrow v_r \leftarrow w \rightarrow \dots$$

We do the same replacement on the right hand side of w if necessary. Then marginalising this path w.r.t. W gives a Z - σ -open path in G^W . This shows:

$$X \perp_G^{\sigma} Y \mid Z \implies X \perp_{G^W}^{\sigma} Y \mid Z.$$

Now let $x \dots y$ be a Z - σ -open path in G^W . Then every edge lifts to a subpath in G where w only occurs as a

non-collider. If a path $\dots v \leftarrow z \dots$ in G^W with $z \in Z \cap \sigma(v)$ in G comes from $\dots v \leftarrow w \leftarrow z \dots$ or $v \leftrightarrow w \rightarrow w \leftarrow z$ then we can, again since $\sigma(v) \in \mathcal{L}(G)$ is a loop, find nodes $v_i \in \sigma(v)$, $r \geq 0$, $i = 1, \dots, r$, and a path in G of the form:

$$\dots v \rightarrow v_1 \rightarrow \dots \rightarrow v_r \rightarrow z \dots$$

which is in any case Z - σ -open in G (whether v or z are colliders or not). So we can construct a Z - σ -open path in G and we get:

$$X \perp_G^{\sigma} Y \mid Z \implies X \perp_{G^W}^{\sigma} Y \mid Z.$$

□

Theorem A.2 (σ -Separation under Conditioning). *Let G be a σ -CG with set of nodes V and $C, X, Y, Z \subseteq V$ subsets with $C = \{c\}$ and $c \notin X \cup Y \cup Z$. Then we have the equivalence:*

$$X \perp_G^{\sigma} Y \mid Z \cup C \iff X \perp_{G_C}^{\sigma} Y \mid Z.$$

Proof. Let $C = \{c\}$. Let $\pi = x \dots y$ be a $(Z \cup C)$ - σ -open path in G with a minimal number of arrowheads pointing to nodes in C . Then at c (if c occurs) there are no undirected edges. So we have the cases:

1. fork: $\dots v_1 \leftarrow c \rightarrow v_2 \dots$ in G with $\sigma(v_1) = \sigma(c) = \sigma(v_2)$. Then $\dots v_1 \leftrightarrow v_2 \dots$ is in G_C with $\sigma(v_1) = \sigma(v_2) = \sigma(c)$. So the triple situation for v_1, v_2 stays the same in G_C .
2. right chain: $\dots v_1 \leftrightarrow c \rightarrow v_2 \dots$ in G with $\sigma(c) = \sigma(v_2)$. Then $\dots v_1 \leftrightarrow v_2 \dots$ is in G_C with $\sigma(v_2) = \sigma(c)$. So the triple situation for v_1, v_2 stays the same in G_C .
3. left chain: similar to right chain.
4. collider: $\dots v_1 \leftrightarrow c \leftrightarrow v_2 \dots$ in G . Then $\dots v_1 \leftrightarrow v_2 \dots$ is in G_C . So the triple situation for v_1, v_2 stays the same in G_C .
5. collider: $\dots v_1 \rightarrow c \leftarrow v_2 \dots$ in G with $v_1, v_2 \notin Z$. Then $\dots v_1 \rightarrow v_2 \dots$ is in G_C with v_1, v_2 non-collider. So it is Z -open at v_1, v_2 .
6. collider: $\dots \leftrightarrow v_1 \rightarrow c \leftrightarrow v_2 \dots$ in G with $v_1 \in Z \cap \sigma(c)$. Since $\sigma(c) \in \mathcal{L}(G)$ is a loop there is a path in G with $w_i \in \sigma(c)$, $r \geq 0$, $i = 1, \dots, r$, of the form:

$$\dots \leftrightarrow v_1 \leftarrow w_1 \leftarrow \dots \leftarrow w_r \leftarrow c \leftrightarrow v_2 \dots,$$

which then is $(Z \cup C)$ - σ -open. So the path

$$\dots \rightleftarrows v_1 \leftarrow w_1 \leftarrow \dots \leftarrow w_r \rightleftarrows v_2 \dots$$

is then Z - σ -open in G_C .

7. collider: $\dots v_1 \rightarrow c \leftrightarrow v_2 \dots$ in G with $v_1 \notin Z$. Then $\dots v_1 \rightarrow v_2 \dots$ is in G_C and Z -open.
8. collider: as before with v_1 and v_2 swapped. Same arguments.

These cover all cases and we have shown:

$$X \underset{G_C}{\perp\!\!\!\perp} Y | Z \implies X \underset{G}{\perp\!\!\!\perp} Y | Z \cup C.$$

Now let $x \dots y$ be a Z -open path in G_C . Then the rules for conditioning lift every edge in G_C to an edge or triple in G , where the triple situation for c is C - σ -open and where the triple situation for the endnodes stays the same. So it is clearly $(Z \cup C)$ -open in G . This shows:

$$X \underset{G}{\perp\!\!\!\perp} Y | Z \cup C \implies X \underset{G_C}{\perp\!\!\!\perp} Y | Z.$$

□

B THE Σ -SEPARATION CRITERION FOR MSCMS

The trick to prove the σ -separation criterion is to transform the σ -connection graph G of an mSCM, which has no undirected edges and can be seen as a directed mixed graph (DMG), into an acyclic directed mixed graph (ADMG) that encodes the same conditional independencies in terms of the well known d-separation. This also shows that every σ -separation-equivalence-class contains an acyclic graph (if one only looks at the observational distributions). Caution: the constructed ADMG is not well-behaved under marginalisation or interventions. We will refer to the d-separation criterion as the *directed global Markov property (dGMP)* and to the σ -separation criterion as the *generalized directed global Markov property (gdGMP)* in the following.

Lemma B.1. *Let $G = (V, E, H)$ be an acyclic directed mixed graph (ADMG) and $(X_v)_{v \in V}$ be random variables that satisfy the dGMP w.r.t. G . Let E_w be a random variable independent of $(X_v)_{v \in V}$ and X_w be another random variable, $w \notin V$, given by a functional relation:*

$$X_w = f((X_v)_{v \in P}, E_w),$$

where $P \subseteq V$ is a subset of nodes. Let $G' = (V', E', H')$ be the ADMG with set of nodes $V' := V \cup \{w\}$, set of edges $E' := E \cup \{v \rightarrow w | v \in P\}$

and set of bidirected edges $H' := H$. Then $G \subseteq G'$ is an ancestral sub-ADMG and $(X_v)_{v \in V'}$ satisfies the dGMP w.r.t. G' .

Proof. Since w is a childless node in G' clearly G' is acyclic, $\text{Pa}^{G'}(w) = P$ and $G \subseteq G'$ is an ancestral sub-ADMG. So there exists a topological order $<$ for G' such that w is the last element. Since for an ADMG the directed global Markov property (dGMP) is equivalent to the ordered local Markov property (oLMP) w.r.t. any topological order (see [12,28]) we only need to check the local independence:

$$\{w\} \underset{\mathbb{P}}{\perp\!\!\!\perp} A \setminus \{w\} | \partial_{A^{\text{mor}}}(w)$$

for every ancestral $A \subseteq G'$ with $w \in A$. Since $\partial_{A^{\text{mor}}}(w) = \text{Pa}^{G'}(w) = P$ and $A \setminus \{w\} \subseteq V$ the statement follows directly from the implication:

$$E_w \underset{\mathbb{P}}{\perp\!\!\!\perp} (X_v)_{v \in V} \implies f((X_v)_{v \in P}, E_w) \underset{\mathbb{P}}{\perp\!\!\!\perp} (X_v)_{v \in A \setminus \{w\}} | (X_v)_{v \in P}.$$

□

Theorem B.2. *Let $G = (V, E, H)$ be a directed mixed graph (DMG) and $\mathcal{S}(G)$ the set of its strongly connected components. Assume that we have:*

1. random variables $(X_v)_{v \in V}$,
2. random variables $(E_v)_{v \in V}$ that jointly satisfy the dGMP w.r.t. the bidirected graph (V, \emptyset, H) , i.e. for every $W, Y \subseteq V$ we have the implication:

$$W \underset{(V, \emptyset, H)}{d}{\perp\!\!\!\perp} Y \implies (E_v)_{v \in W} \underset{\mathbb{P}}{\perp\!\!\!\perp} (E_v)_{v \in Y},$$

3. a tuple of functions $(g_S)_{S \in \mathcal{S}(G)}$ indexed by the strongly connected components S of G ,

such that we have the following equations for $S \in \mathcal{S}(G)$:

$$(X_v)_{v \in S} = g_S((X_w)_{w \in \text{Pa}^G(S) \setminus S}, (E_w)_{w \in S}).$$

Then $(X_v)_{v \in V}$ satisfy the general directed global Markov property (gdGMP) w.r.t. the DMG G , i.e. for every three subsets $W, Y, Z \subseteq V$ we have the implication:

$$W \underset{G}{\perp\!\!\!\perp} Y | Z \implies (X_v)_{v \in W} \underset{\mathbb{P}}{\perp\!\!\!\perp} (X_v)_{v \in Y} | (X_v)_{v \in Z}.$$

Proof. By assumption we have that $(E_v)_{v \in V}$ satisfies the dGMP w.r.t. the ADMG (V, \emptyset, H) . By lemma B.1 we can inductively add:

$$X_v = g_{S,v}((X_w)_{w \in \text{Pa}^G(S) \setminus S}, (E_w)_{w \in S})$$

for $v \in V$ where $S = \text{Sc}^G(v)$. We then finally get an ADMG G' with nodes $(E_v)_{v \in V}$ and $(X_v)_{v \in V}$ that satisfy the *dGMP* w.r.t. this G' . This implies that for $W, Y, Z \subseteq V$ we have:

$$W \perp_{G'}^d Y | Z \implies (X_v)_{v \in W} \perp_{\mathbb{P}} (X_v)_{v \in Y} | (X_v)_{v \in Z}.$$

It is thus left to show that we also have the implication:

$$W \perp_G^\sigma Y | Z \implies W \perp_{G'}^d Y | Z.$$

For this it is enough to show that every Z -d-open path π' from W to Y in G' lifts to a Z - σ -open path π from W to Y in G . The construction is straightforward. For details see [12]. \square

Corollary B.3. *The observed variables $(X_v)_{v \in V}$ of any mSCM $M = (G^+, \mathcal{X}, \mathbb{P}, g)$, $G^+ = (U \dot{\cup} V, E^+)$, satisfy the σ -separation criterion w.r.t. the induced σ -connection graph (σ -CG) G .*

Proof. For $v \in V$ we put $E_v := (E_u)_{\substack{u \in U \\ v \in \text{Ch}^{G^+}(u)}}$. The $(E_v)_{v \in V}$ then entail the conditional independence relations implied by d-separation of the bidirected graph (V, \emptyset, H) . Furthermore, for $S \in \mathcal{S}(G)$ we have equations:

$$X_S = g_S(X_{\text{Pa}^\sigma(G) \setminus S}, E_S).$$

The claim then directly follows from B.2. \square

As a motivation for future work on selection bias we state the following direct corollary.

Corollary B.4 (mSCM with context). *Let $M = (G^+, \mathcal{X}, \mathbb{P}, g)$ be a mSCM with $G^+ = (U \dot{\cup} V, E^+)$ and $C \subseteq V$ a subset. Let $G_C = (G^+)_C^U$ be the induced σ -CG of M conditioned on C . Then the observed variables $(X_v)_{v \in V \setminus C}$ satisfy the σ -separation criterion w.r.t. G_C and w.r.t. the regular conditional probability distribution $\mathbb{P}^{X_C = x_C}$ given $X_C = x_C$ (for \mathbb{P}^{X_C} -almost-all values $x_C \in \mathcal{X}_C$): For all subsets $W, Y, Z \subseteq V \setminus C$ we have the implication:*

$$W \perp_{G_C}^\sigma Y | Z \implies X_W \perp_{\mathbb{P}^{X_C = x_C}} X_Y | X_Z.$$

Proof. The lhs is equivalent to $W \perp_G^\sigma Y | Z \cup C$ (see Theorem 2.20, Theorem A.2, resp.) and this implies $X_W \perp_{\mathbb{P}} X_Y | X_Z, X_C$ (see Theorem 2.14, Corollary B.3, resp.), which implies the claim on the rhs (for \mathbb{P}^{X_C} -almost-all values $x_C \in \mathcal{X}_C$). \square

The last corollary can be used as a starting point for conditional independence constraint-based causal discovery

in the presence of (unknown) *selection bias* given by the unknown context C and $x_C \in \mathcal{X}_C$ (in addition to non-linear functional relations, cycles and latent confounders etc.).

C NEURAL NETWORKS AS MSCMS

For constructing causal mechanisms we could use any parametric or non-parametric family of functions. Since we want to stay as general as possible and also make use of the practical advantages of parametric models we represent/approximate the structural functions $g_{\{v\}}$, $v \in V$ by *universal approximators*. A well known class of universal approximators are *neural networks* (see e.g. [16]). A neural network is a function that is constructed from several compositions of linear maps and a fixed one-dimensional *activation* function h . A sufficient condition to have the universal approximation property is if one assumes h be continuous, non-polynomial and piecewise differentiable. A further advantage of neural networks is that the *hidden units* (given by composition of functions $z \mapsto h(w^T z + b)$) can be interpreted as intermediate variables of an extended structural causal model. This means that by modelling the hidden units of every $g_{\{v\}}$ explicitly as a node in an extended graph we can restrict—for the analysis purposes here—to this extended setting, where now the functions $g_{\{v\}}$ (the index $\{v\}$ refers to the trivial loop) are of the form:

$$\begin{aligned} & g_{\{v\}}(x_{\text{Pa}^{G^+}(v) \setminus \{v\}}) \\ &= h\left(\sum_{k \in \text{Pa}^{G^+}(v) \setminus \{v\}} A_{v,k} \cdot x_k + b_v\right), \end{aligned}$$

with weights $A_{v,k}$ and biases b_v .

Further note that introducing or marginalizing intermediate variables will not change the outcome of the σ -separation criterion defined in Definition 2.10, Theorem 2.14, and Theorem 2.20 (also see [12]). So also this part is compatible with our theory.

Theorem C.1. *The conditions for the contractiveness of the iterations scheme from subsection 4.1 are satisfied if the following three points hold:*

1. $\sup_z |h'(z)| \leq C$ with $0 < C < \infty$, and
2. $A_{v,k} := 0$ for $k \notin \text{Pa}^{G^+}(v) \setminus \{v\}$, and
3. $\|(A_{v,k})_{v,k \in S}\| < \frac{1}{C}$ for every non-trivial loop $S \subseteq G$, where $\|\cdot\|$ can be one of the matrix norms: $\|\cdot\|_p$, $p \geq 1$, or $\|\cdot\|_\infty$.

In this case the functions $(g_{\{v\}})_{v \in V}$ will constitute a well-defined mSCM.

Note that we can put $C = 1$ for popular activation functions $h(z)$ like $\tanh(z)$, $\text{ReLU}(z) = \max(0, z)$, $\sigma(z) =$

$\frac{1}{1+\exp(-z)}$, LeakyRelu, SoftPlus($z = \ln(1 + e^z)$), etc.. Further note that by using one of these activation functions $h(z)$ and $\|\cdot\| = \|\cdot\|_\infty$ all the conditions are satisfied if we choose the $A_{v,k}$ such that for all $v \in V$:

$$\sum_{k \in \text{Pa}^{G^+}(v) \setminus \{v\}} |A_{v,k}| < 1$$

and $A_{v,k} := 0$ for $k \notin \text{Pa}^{G^+}(v) \setminus \{v\}$.

Furthermore, we can then iterate the whole system for given error value x_U and initialization $x_V^{(0)}$:

$$\begin{aligned} x_V^{(t+1)} &:= (g_{\{v\}})_{v \in V}(x_V^{(t)}, x_U) \\ &= h\left(A_{V^+} \cdot \begin{pmatrix} x_V^{(t)} \\ x_U \end{pmatrix} + b_V\right) \end{aligned}$$

and reach a unique fixed point x_V . This analysis also holds if we have the error variables outside of the activation function as additive noise.

Proof. For a non-trivial loop $S \subseteq G$ we want to show that for every value $x_{\text{Pa}^{G^+}(S) \setminus S}$ and initialization $x_S^{(0)}$ the iteration (using vector and matrix notations):

$$\begin{aligned} x_S^{(t+1)} &:= \\ &h\left(A_{\text{Pa}^{G^+}(S) \cup S} \cdot \begin{pmatrix} x_S^{(t)} \\ x_{\text{Pa}^{G^+}(S) \setminus S} \end{pmatrix} + b_S\right) \end{aligned}$$

converges to a unique point x_S (for $t \rightarrow \infty$) under the three stated assumptions in the text.

For applying Banach's fixed point theorem we need to show that for every value $x_{\text{Pa}^{G^+}(S) \setminus S}$ we have a bounded partial Jacobian (S a non-trivial loop):

$$\sup_{x_S} \|J_S(x_{\text{Pa}^{G^+}(S) \cup S})\| \leq L(x_{\text{Pa}^{G^+}(S) \setminus S}) < 1$$

where $L(x_{\text{Pa}^{G^+}(S) \setminus S})$ is a constant smaller than 1 and $\|\cdot\|$ is a suitable matrix norm. In our case we have:

$$\begin{aligned} &J_S(x_{\text{Pa}^{G^+}(S) \cup S}) \\ &:= \left(\frac{\partial g_{\{v\}}}{\partial x_k}\right)_{v,k \in S}(x_{\text{Pa}^{G^+}(S) \cup S}) \\ &= \nabla_{x_S} h\left(A_{\text{Pa}^{G^+}(S) \cup S} \cdot \begin{pmatrix} x_S \\ x_{\text{Pa}^{G^+}(S) \setminus S} \end{pmatrix} + b_S\right) \\ &= \left(h' \left(\sum_{j \in \text{Pa}^{G^+}(v) \setminus \{v\}} A_{v,j} \cdot x_j + b_v\right) \cdot A_{v,k} \cdot \mathbf{1}_{k \in \text{Pa}^{G^+}(v) \setminus \{v\}}\right)_{v,k \in S} \\ &= \text{diag}(h') \cdot (A_S \odot \mathbf{1}_S). \end{aligned}$$

Here $\text{diag}(h')$ refers to the diagonal matrix with the corresponding values of h' and $\mathbf{1}_S$ is the adjacency matrix as indicated on the line above.

If $|h'(z)| \leq C < \infty$ and $\|\cdot\|$ is either $\|\cdot\|_p$, $p \geq 1$, or $\|\cdot\|_\infty$ then $\|\text{diag}(h')\| \leq C$. If, furthermore, $\|A_S \odot \mathbf{1}_S\| < \frac{1}{C}$ then we get:

$$\begin{aligned} \|J_S\| &\leq \|\text{diag}(h')\| \cdot \|A_S \odot \mathbf{1}_S\| \\ &\leq C \cdot \|A_S \odot \mathbf{1}_S\| =: L \\ &< C \cdot \frac{1}{C} \\ &= 1. \end{aligned}$$

Note that we can represent $A \odot \mathbf{1}$ in a single matrix A if we put $A_{v,k} := 0$ whenever $k \notin \text{Pa}^{G^+}(v) \setminus \{v\}$. From the above then follows that the map of the iteration scheme becomes contractive and the series thus converges to a unique fixed point x_S . g_S can then be defined via:

$$g_S(x_{\text{Pa}^{G^+}(S) \setminus S}) := x_S.$$

The system $(g_S)_{S \in \mathcal{L}(G)}$ is also compatible. Indeed, the convergence shows that the above element $(x_{\text{Pa}^{G^+}(S) \setminus S}, x_S)$ simultaneously solves the system $x_v = g_{\{v\}}(x_{\text{Pa}^{G^+}(v) \setminus \{v\}})$, $v \in S$. So for a loop $S' \subseteq S$ the corresponding components $(x_{\text{Pa}^{G^+}(S') \setminus S'}, x_{S'})$ simultaneously solves the system $x_v = g_{\{v\}}(x_{\text{Pa}^{G^+}(v) \setminus \{v\}})$, $v \in S'$. Since also the solution for the loop S' is unique we get:

$$g_{S'}(x_{\text{Pa}^{G^+}(S') \setminus S'}) = x_{S'},$$

which shows the compatibility. The measurability of this map follows from a measurable choice theorem (see [2]) as explained in [3]. \square

If we want to uniformly sample weights for the parent nodes one can use the following:

Remark C.2 (See [1]). *To uniformly sample from the d -dimensional L_p -ball $B_p^d := \{x \in \mathbb{R}^d : \|x\|_p \leq 1\}$ we can sample i.i.d. $y_1, \dots, y_d \sim p(t) = \frac{1}{2\Gamma(1+\frac{1}{p})} e^{-|t|^p}$, $t \in \mathbb{R}$ and $z \sim p(s) = e^{-s}$, $s \geq 0$. Then $x = \frac{(y_1, \dots, y_d)^T}{(\sum_{j=1}^d |y_j|^p + z)^{1/p}}$ is uniformly sampled from B_p^d .*

D MORE DETAILS ON THE ALGORITHM

D.1 SCORING FEATURES

In order to score features, which can be defined as Boolean functions of the causal graph G , we define a modified loss function

$$\begin{aligned} \mathcal{L}(G, S, f) &:= \\ &\sum_{(w_j, y_j, Z_j, I_j, \lambda_j) \in S} \lambda_j (\mathbf{1}_{\lambda_j > 0} - \mathbf{1}_{w_j} \perp_{G_{\text{do}(I_j)}}^{\sigma} y_j \mid Z_j) \mathbf{1}_{f(G)} \end{aligned} \quad (4)$$

[21] proposed to score the confidence of a feature with

$$\begin{aligned} C(S, f) &:= \min_{G \in \mathbb{G}(V)} \mathcal{L}(G, S, \neg f) \\ &\quad - \min_{G \in \mathbb{G}(V)} \mathcal{L}(G, S, f). \end{aligned} \quad (5)$$

They showed that this scoring method is sound for oracle inputs.

Theorem D.1. For any feature f , the confidence score $C(S, f)$ of (5) is sound for oracle inputs with infinite weights. In other words, $C(S, f) = \infty$ if f is identifiable from the inputs, $C(S, f) = -\infty$ if $\neg f$ is identifiable from the inputs, and $C(S, f) = 0$ if f is unidentifiable from the inputs.

Furthermore, they showed that the scoring method is asymptotically consistent under a consistency condition on the statistical independence test.

Theorem D.2. Assume that the weights are asymptotically consistent, meaning that

$$\log p_N - \log \alpha_N \xrightarrow{P} \begin{cases} -\infty & H_1 \\ +\infty & H_0, \end{cases} \quad (6)$$

as the number of samples $N \rightarrow \infty$, where the null hypothesis H_0 is independence and the alternative hypothesis H_1 is dependence. Then for any feature f , the confidence score $C(S, f)$ of (5) is asymptotically consistent, i.e., $C(S, f) \rightarrow \infty$ in probability if f is identifiably true, $C(S, f) \rightarrow -\infty$ in probability if f is identifiably false, and $C(S, f) \rightarrow 0$ in probability otherwise.

By using the scoring method of [21] as explained above, our algorithm inherits these desirable properties.

D.2 ENCODING IN ANSWER SET PROGRAMMING

In order to test whether a causal graph G entails a certain independence, we create a computation graph of σ -connection graphs. A computation graph of σ -connection graphs is a DAG with σ -connection graphs as nodes, and directed edges that correspond with the operations of conditioning and marginalisation. The “source node” of an encoding DAG is an (intervened) causal graph. The “sink” nodes are σ -connection graphs that consist of only two variables (because all other variables have been conditioned or marginalised out) that can be reached from the source node by applying a sequence of conditioning and marginalisation operations. Testing a σ -separation statement in the intervened causal graph reduces to testing for adjacency in the corresponding sink node.

Since interventions and conditioning do not commute, one has to take care to employ these operations in the right ordering. We define the computation graph in such a way that intervention operations are performed first, followed by marginalisations, and finally conditioning operations. At each stage, we always remove the node with the highest possible label first, which means that our computation graph is actually a computation tree.

Below we provide the source code of the essential part of the algorithm, using the ASP syntax for `clingo 4`. It is based upon the source code provided by [19]. The differences to [19], i.e. of σ -separation vs. d-separation, are indicated with “(*sigma*)” in the comments, i.e. at lines 100, 128, 138. Note that the main difference between the encoding of d-separation and σ -separation is that in the non-collider case (see definition 2.9) we need to check in which strongly connected component $\sigma(v)$ the non-collider node lies in comparison to its adjacent nodes. This boils down to checking ancestral relations. Since the σ -structure is inherited in a trivial fashion during the marginalisation and conditioning operations, it only needs to be found once (namely in the original σ -CG induced by the mSCM).

We used the state-of-the-art ASP solver `clingo 4` [13] in our experiments to run the ASP program.


```

82
83 %% X-->Y => X-->Y
84 th(X,Y,C,J,M) :- th(X,Y,Csub,J,M),
85     X != Y,
86     not ismember(C,X), not ismember(C,Y),
87     not ismember(M,X), not ismember(M,Y),
88     node(X),node(Y),
89     condition(Csub,Z,C,J,M).
90
91 %% X-->Z<->Y => X-->Y
92 th(X,Y,C,J,M) :- th(X,Z,Csub,J,M),
93     { hh(Z,Y,Csub,J,M); hh(Y,Z,Csub,J,M) } >= 1,
94     X != Y,
95     not ismember(C,X), not ismember(C,Y),
96     not ismember(M,X), not ismember(M,Y),
97     node(X),node(Y),
98     condition(Csub,Z,C,J,M).
99
100 %% X-->Z-->Y (anc of Z) => X-->Y (sigma)
101 th(X,Y,C,J,M) :- th(X,Z,Csub,J,M), th(Z,Y,Csub,J,M),
102     ancestor(Y,Z,J),
103     X != Y,
104     not ismember(C,X), not ismember(C,Y),
105     not ismember(M,X), not ismember(M,Y),
106     node(X),node(Y),node(Z),
107     condition(Csub,Z,C,J,M).
108
109 %%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%
110
111 %% X<->Y => X<->Y
112 hh(X,Y,C,J,M) :- hh(X,Y,Csub,J,M),
113     X < Y,
114     not ismember(C,X), not ismember(C,Y),
115     not ismember(M,X), not ismember(M,Y),
116     node(X),node(Y),
117     condition(Csub,Z,C,J,M).
118
119 %% X<->Z<->Y => X<->Y
120 hh(X,Y,C,J,M) :- { hh(Z,X,Csub,J,M); hh(X,Z,Csub,J,M) } >= 1,
121     { hh(Z,Y,Csub,J,M); hh(Y,Z,Csub,J,M) } >= 1,
122     X < Y,
123     not ismember(C,X), not ismember(C,Y),
124     not ismember(M,X), not ismember(M,Y),
125     node(X),node(Y),
126     condition(Csub,Z,C,J,M).
127
128 %% X<->Z-->Y (anc of Z) => X<->Y (sigma)
129 hh(X,Y,C,J,M) :- { hh(Z,X,Csub,J,M); hh(X,Z,Csub,J,M) } >= 1,
130     th(Z,Y,Csub,J,M),
131     ancestor(Y,Z,J),
132     X < Y,
133     not ismember(C,X), not ismember(C,Y),
134     not ismember(M,X), not ismember(M,Y),
135     node(X),node(Y),node(Z),
136     condition(Csub,Z,C,J,M).
137
138 %% (anc of Z) X<-Z-->Y (anc of Z) => X<->Y (sigma)
139 hh(X,Y,C,J,M) :- th(Z,X,Csub,J,M), th(Z,Y,Csub,J,M),
140     ancestor(X,Z,J),
141     ancestor(Y,Z,J),
142     X < Y,
143     not ismember(C,X), not ismember(C,Y),
144     not ismember(M,X), not ismember(M,Y),
145     node(X),node(Y),node(Z),
146     condition(Csub,Z,C,J,M).
147
148 %%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%
149
150
151 %%%%%%%%% MARGINALIZATION %%%%%%%%%
152
153 %% X---Y => X---Y
154 tt(X,Y,C,J,M) :- tt(X,Y,C,J,Msub),
155     X <= Y,
156     not ismember(C,X), not ismember(C,Y),
157     not ismember(M,X), not ismember(M,Y),
158     node(X),node(Y),
159     marginalize(C,J,Msub,Z,M).
160
161 %% X-->Z---Y => X---Y
162 tt(X,Y,C,J,M) :- th(X,Z,C,J,Msub),

```



```

244 %% X<->Y => X<->Y
245 hh(X,Y,C,J,M) :- hh(X,Y,C,J,Msub),
246 X < Y,
247 not ismember(C,X), not ismember(C,Y),
248 not ismember(M,X), not ismember(M,Y),
249 node(X),node(Y),
250 marginalize(C,J,Msub,Z,M).
251
252 %% X<->Z-->Y => X<->Y
253 hh(X,Y,C,J,M) :- { hh(X,Z,C,J,Msub); hh(Z,X,C,J,Msub) } >= 1,
254 th(Z,Y,C,J,Msub),
255 X < Y,
256 not ismember(C,X), not ismember(C,Y),
257 not ismember(M,X), not ismember(M,Y),
258 node(X),node(Y),
259 marginalize(C,J,Msub,Z,M).
260
261 %% X<->Z-->Y => X<->Y
262 hh(X,Y,C,J,M) :- th(Z,X,C,J,Msub),
263 th(Z,Y,C,J,Msub),
264 X < Y,
265 not ismember(C,X), not ismember(C,Y),
266 not ismember(M,X), not ismember(M,Y),
267 node(X),node(Y),
268 marginalize(C,J,Msub,Z,M).
269
270 %% X<->Z<->Y => X<->Y
271 hh(X,Y,C,J,M) :- th(Z,X,C,J,Msub),
272 { hh(Y,Z,C,J,Msub); hh(Z,Y,C,J,Msub) } >= 1,
273 X < Y,
274 not ismember(C,X), not ismember(C,Y),
275 not ismember(M,X), not ismember(M,Y),
276 node(X),node(Y),
277 marginalize(C,J,Msub,Z,M).
278
279 %% X<->Z---Z<->Y => X<->Y
280 hh(X,Y,C,J,M) :- { hh(X,Z,C,J,Msub); hh(Z,X,C,J,Msub) } >= 1,
281 { hh(Y,Z,C,J,Msub); hh(Z,Y,C,J,Msub) } >= 1,
282 tt(Z,Z,C,J,Msub),
283 X < Y,
284 not ismember(C,X), not ismember(C,Y),
285 not ismember(M,X), not ismember(M,Y),
286 node(X),node(Y),
287 marginalize(C,J,Msub,Z,M).
288
289 %%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%
290
291 %%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%
292 %%%%%%%%%%% LOSS FUNCTION %%%%%%%%%%%
293
294 fail(X,Y,C,J,M,W) :- th(X,Y,C,J,M), indep(X,Y,C,J,M,W), X<Y.
295 fail(X,Y,C,J,M,W) :- th(Y,X,C,J,M), indep(X,Y,C,J,M,W), X<Y.
296 fail(X,Y,C,J,M,W) :- hh(X,Y,C,J,M), indep(X,Y,C,J,M,W), X<Y.
297 fail(X,Y,C,J,M,W) :- tt(X,Y,C,J,M), indep(X,Y,C,J,M,W), X<Y.
298
299 fail(X,Y,C,J,M,W) :- not th(X,Y,C,J,M),
300 not th(Y,X,C,J,M),
301 not hh(X,Y,C,J,M),
302 not tt(X,Y,C,J,M),
303 dep(X,Y,C,J,M,W), X<Y.
304
305 %%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%
306
307 %%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%
308 %%%%%%%%%%% OPTIMIZATION PROBLEM %%%%%%%%%%%
309
310 #minimize{W,X,Y,C,J,M:fail(X,Y,C,J,M,W) }.

```

E EXPERIMENTAL RESULTS

Here we provide additional visualisations of the results of our experiments, for which no space was left in the main paper.

Figure 6 shows ROC curves and PR curves for detecting directed edges (i.e., direct causal relations) and for detecting latent confounders in the causal graph. Results are shown for the purely observational setting (“0 interventions”) and for a combination of observational and interventional data (“1–5 interventions”) where the targets of the stochastic surgical interventions are single variables chosen randomly, without replacement. Clearly, making use of interventional data is beneficial for causal discovery.

Figure 7 shows similar curves, now for 5 interventions only, but for different encodings: σ -separation (this work), d-separation (allowing for cycles, [19]) and d-separation (acyclic, [19]). Interestingly, the differences between σ -separation and d-separation turn out to be quite small in our simulation setting. The difference is largest for the detection of confounders. On the other hand, the difference between assuming acyclicity and allowing for cycles is much more pronounced, and is also significant for the detection of direct causal relations.

We expect that when going to larger graphs with more variables and with nested loops, the differences between σ -separation and d-separation should increase. However, due to computational restrictions we were not able to perform sufficiently many experiments in this regime to gather enough empirical support for that hypothesis and leave this for future research.

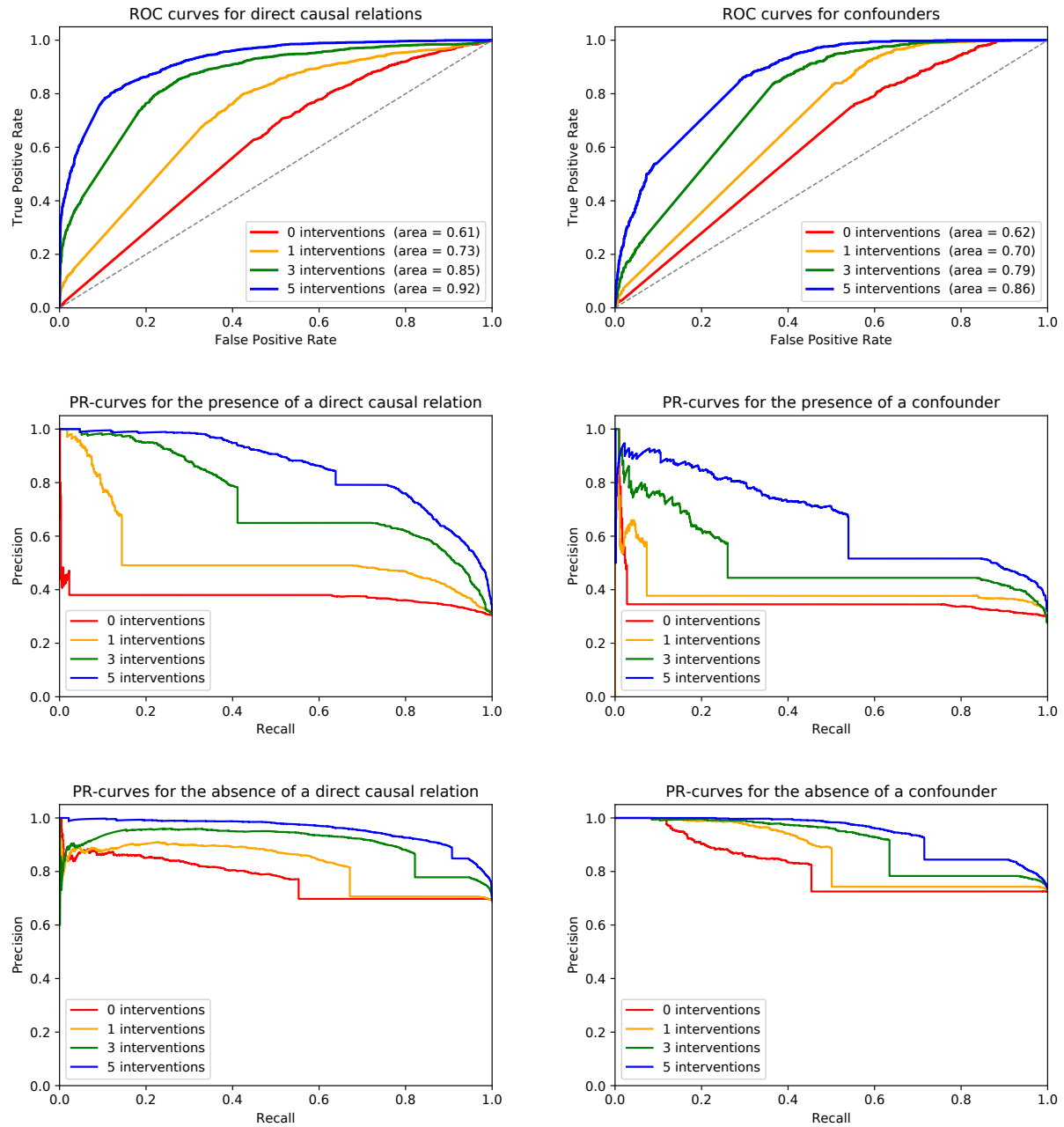


Figure 6: ROC curves (top) and PR curves (center, bottom) for directed edges (left) and confounders (right), for different numbers of single-variable interventions. All results shown here use σ -separation.

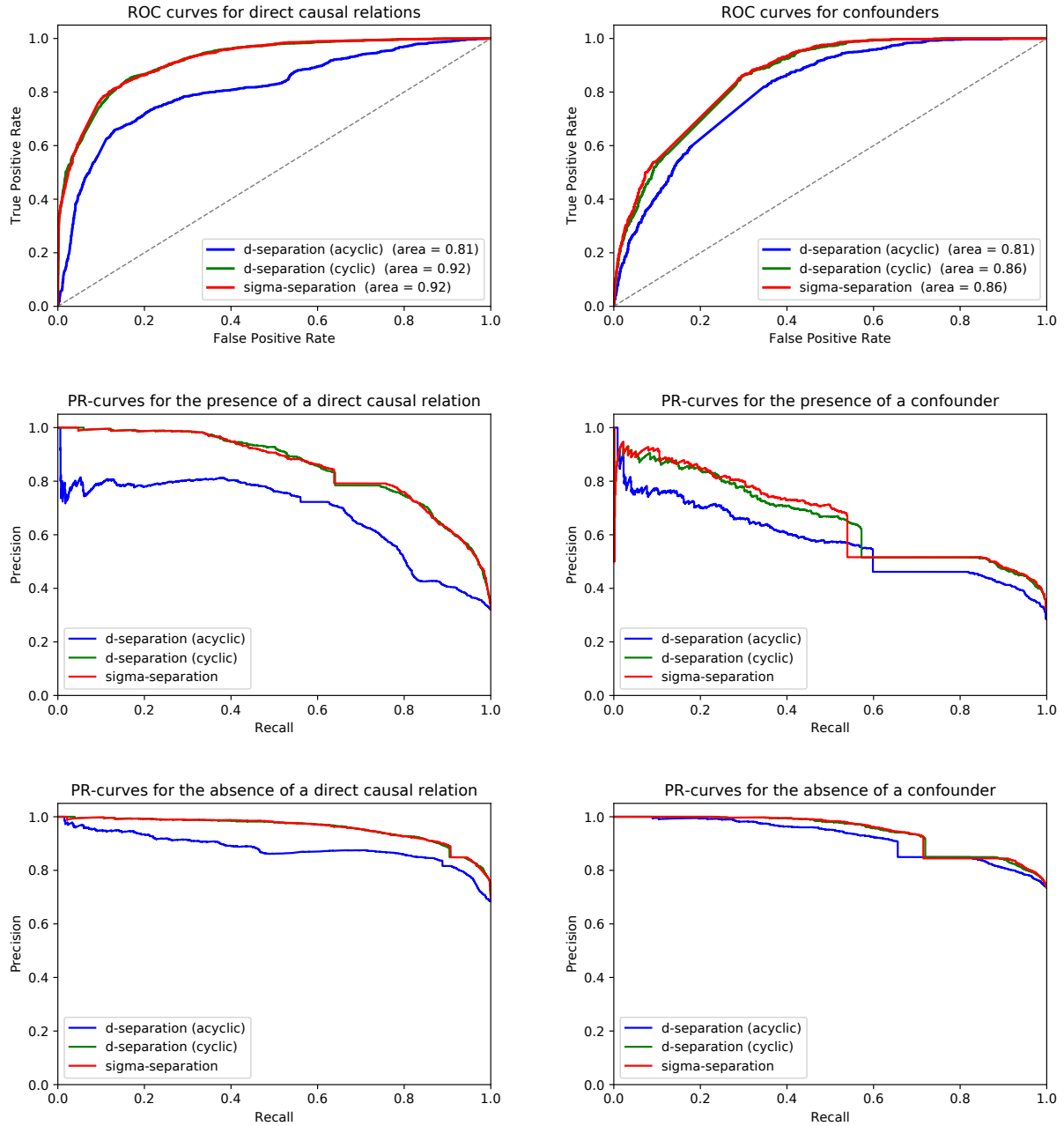


Figure 7: ROC curves (top) and PR curves (center, bottom) for directed edges (left) and confounders (right), for different encodings. All results shown here use observational and 5 interventional data sets.