

## SUPPLEMENTARY MATERIAL

### A $\Sigma$ -CG UNDER MARGINALISATION AND CONDITIONING

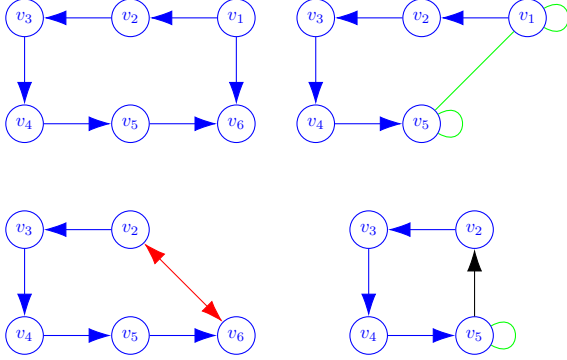


Figure 5: A directed *acyclic* graph  $G$  (top left) as the  $\sigma$ -connection graph ( $\sigma$ -CG) with the  $\sigma$ -equivalence classes  $\{\{v_1\}, \dots, \{v_6\}\}$ . We have the  $\sigma$ -separation:  $\{v_3\} \perp_G^{\sigma} \{v_5\} \mid \{v_2, v_4, v_6\}$ . The combination of marginalizing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

**Theorem A.1** ( $\sigma$ -Separation under Marginalisation). *Let  $G$  be a  $\sigma$ -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  $Z$ - $\sigma$ -open path in  $G$  then every occurrence of  $w$  in  $\pi$  is as a non-collider. If we have  $\dots v \rightarrow w \rightarrow \dots$  in  $\pi$  and  $v \notin Z$  then marginalising out  $w$  keeps  $\pi^W$   $Z$ - $\sigma$ -open in  $G^W$ . If  $v \in Z$  then  $v \in \sigma(w)$  by the  $Z$ - $\sigma$ -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  $Z$ - $\sigma$ -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  $Z$ - $\sigma$ -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  $Z$ - $\sigma$ -open in  $G$  (whether  $v$  or  $z$  are colliders or not). So we can construct a  $Z$ - $\sigma$ -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** ( $\sigma$ -Separation under Conditioning). *Let  $G$  be a  $\sigma$ -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  $(Z \cup C)$ - $\sigma$ -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  $(Z \cup C)$ - $\sigma$ -open. So the path

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

is then  $Z$ - $\sigma$ -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:

$$X \perp_{G_C}^{\sigma} Y | Z \implies X \perp_G^{\sigma} Y | Z \cup C.$$

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  $C$ - $\sigma$ -open and where the triple situation for the endnodes stays the same. So it is clearly  $(Z \cup C)$ -open in  $G$ . This shows:

$$X \perp_G^{\sigma} Y | Z \cup C \implies X \perp_{G_C}^{\sigma} Y | Z.$$

□

## B THE $\Sigma$ -SEPARATION CRITERION FOR MSCMS

The trick to prove the  $\sigma$ -separation criterion is to transform the  $\sigma$ -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  $\sigma$ -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  $\sigma$ -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  $E' := E \cup \{w \rightarrow v | v \in P\}$

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:

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

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:

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

□

**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:

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

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:

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

*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  $\pi'$  from  $W$  to  $Y$  in  $G'$  lifts to a  $Z$ - $\sigma$ -open path  $\pi$  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  $\sigma$ -separation criterion w.r.t. the induced  $\sigma$ -connection graph ( $\sigma$ -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:

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

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  $\sigma$ -CG of  $M$  conditioned on  $C$ . Then the observed variables  $(X_v)_{v \in V \setminus C}$  satisfy the  $\sigma$ -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  $\sigma$ -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:

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

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

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

[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  $\sigma$ -connection graphs. A computation graph of  $\sigma$ -connection graphs is a DAG with  $\sigma$ -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  $\sigma$ -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  $\sigma$ -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  $\sigma$ -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  $\sigma$ -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  $\sigma$ -structure is inherited in a trivial fashion during the marginalisation and conditioning operations, it only needs to be found once (namely in the original  $\sigma$ -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 %%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%% MARGINALIZATION %%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%
151
152 %% X---Y => X---Y
153 tt(X,Y,C,J,M) :- tt(X,Y,C,J,Msub),
154                 X <= Y,
155                 not ismember(C,X), not ismember(C,Y),
156                 not ismember(M,X), not ismember(M,Y),
157                 node(X),node(Y),
158                 marginalize(C,J,Msub,Z,M).
159
160 %% X-->Z---Y => X---Y
161 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:  $\sigma$ -separation (this work), d-separation (allowing for cycles, [19]) and d-separation (acyclic, [19]). Interestingly, the differences between  $\sigma$ -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  $\sigma$ -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.

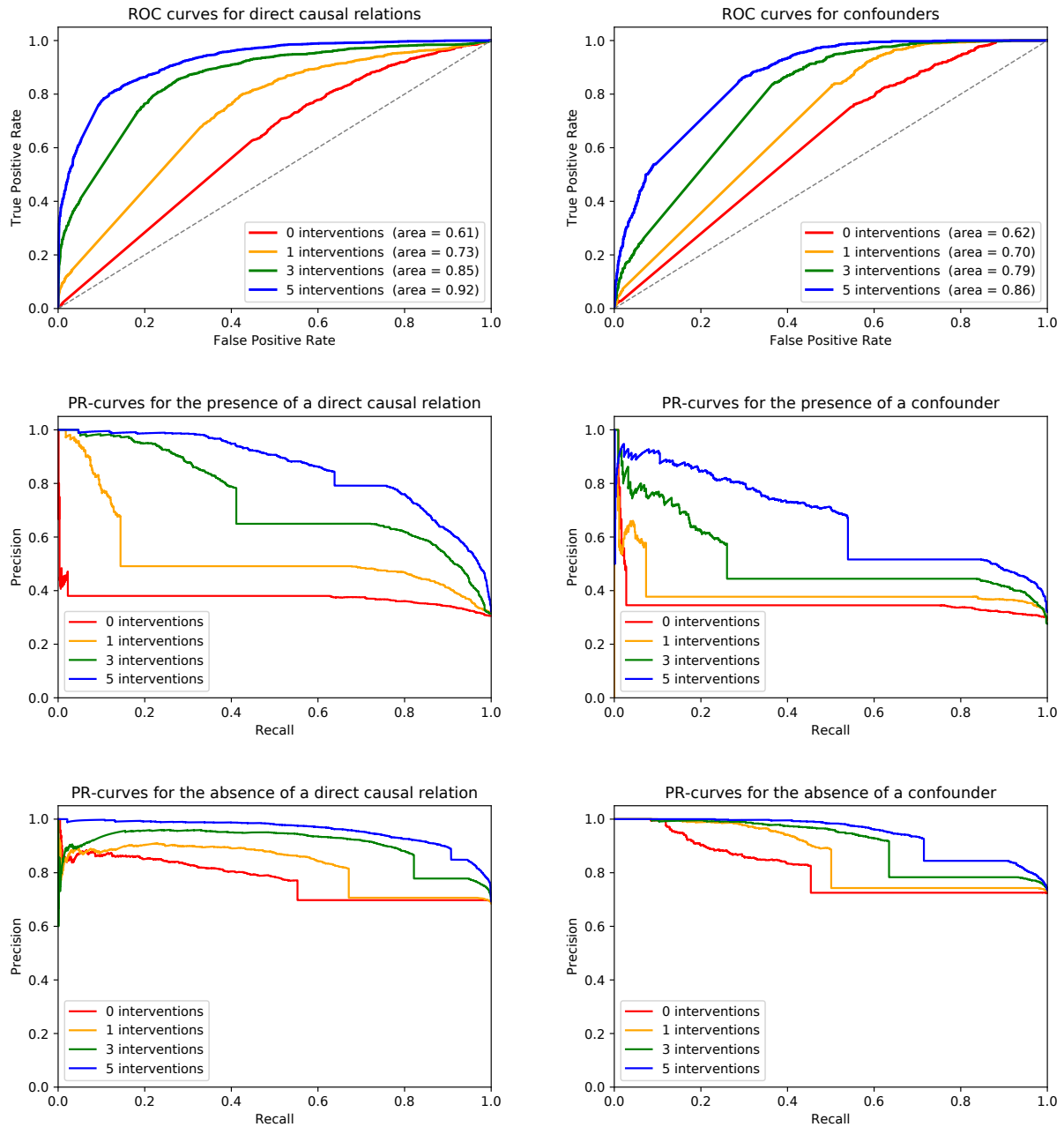


Figure 6: ROC curves (top) and PR curves (center, bottom) for directed edges (left) and confounders (right), for different numbers of single-variable interventions. All results shown here use  $\sigma$ -separation.

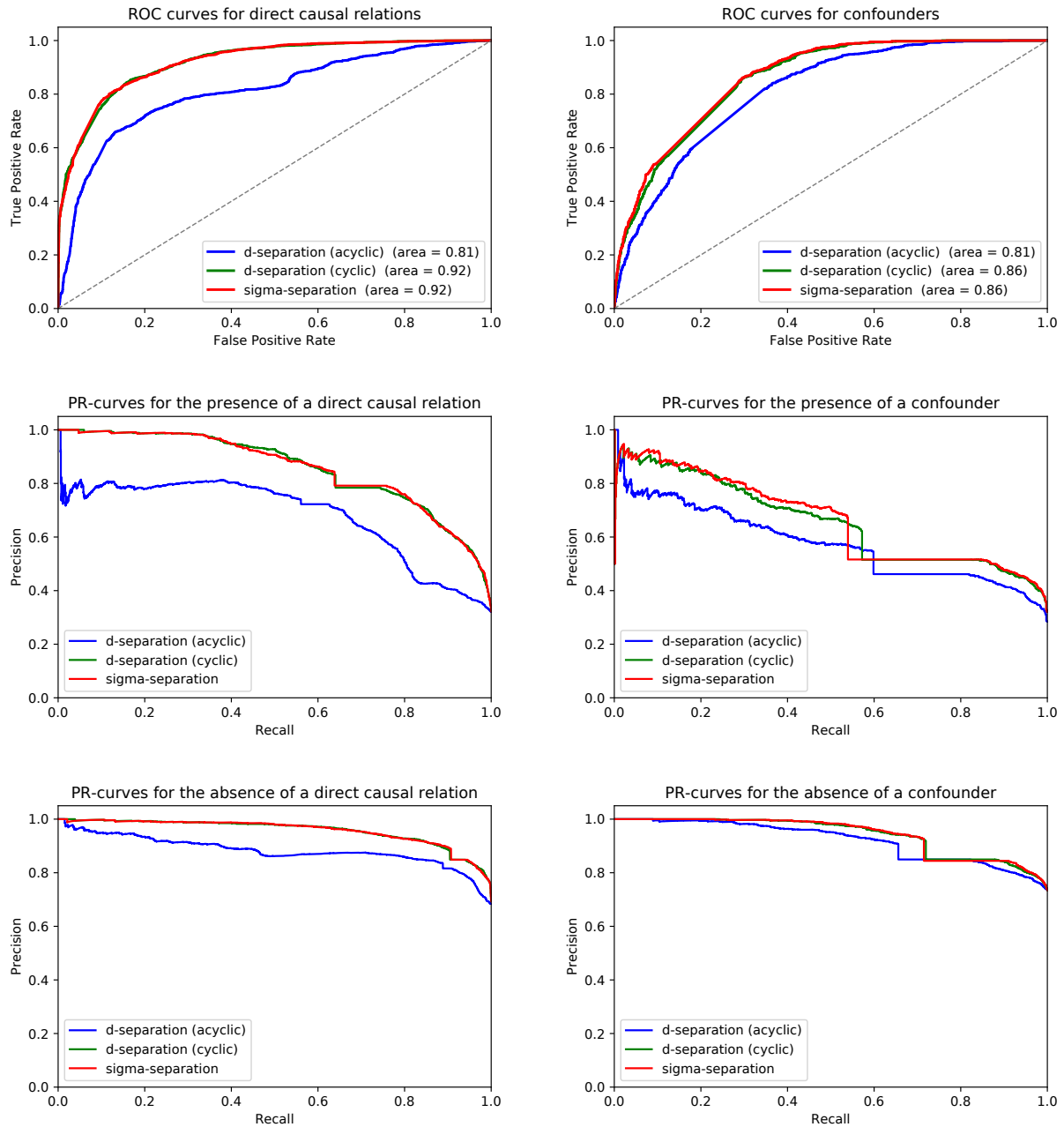


Figure 7: ROC curves (top) and PR curves (center, bottom) for directed edges (left) and confounders (right), for different encodings. All results shown here use observational and 5 interventional data sets.