
Disentangling Mixtures of Unknown Causal Interventions

Supplementary Material

Abhinav Kumar¹

Gaurav Sinha²

¹Paypal, Hyderabad, Telangana, India

²Adobe Research, Bangalore, Karnataka, India

A PROOFS FROM SECTION 4

In this section, we provide missing proofs of the lemmas stated in Section 4.

A.1 PROOF OF LEMMA 4.1

Note that we have assumed $a_i > 0$ for all $i \in [k]$. We iterate from $i = 1$ to $k - 1$ and apply the following row transformations on our matrix.

$$R_i \mapsto R_i - \left(\frac{a_i}{a_k}\right)R_k$$

This results in the following linear system.

$$\begin{bmatrix} c & 0 & \cdot & -\frac{a_1}{a_k}c \\ 0 & c & \cdot & -\frac{a_2}{a_k}c \\ \cdot & \cdot & \cdot & \cdot \\ \cdot & \cdot & \cdot & \cdot \\ -a_k & -a_k & \cdot & c - a_k \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \\ \cdot \\ \cdot \\ x_k \end{bmatrix} = \begin{bmatrix} \tilde{b}_1 \\ \tilde{b}_2 \\ \cdot \\ \cdot \\ \tilde{b}_k \end{bmatrix}$$

where $\tilde{b}_i = (b_i - \frac{a_i}{a_k}b_k)$ for all $i \in [k - 1]$ and $\tilde{b}_k = b_k$. Since $c > 0$, this matrix is easily seen to have rank $\geq k - 1$. Using $c = \sum_{i=1}^k a_i$ we can easily check that the last row R_k is $\frac{-a_k}{c}(R_1 + \dots + R_{k-1})$, implying that it has rank $k - 1$. Since the system is assumed to have at least one solution, it actually has infinitely many solutions. The null space of this matrix is the one dimensional space spanned by $\mathbf{w} = (\frac{a_1}{a_k}, \dots, \frac{a_k}{a_k})^T$ which has all positive entries since $a_i > 0, i \in [k]$. Assume there are two distinct solutions $\mathbf{u} = (u_1, \dots, u_k)^T$ and $\mathbf{v} = (v_1, \dots, v_k)^T$ in $\mathbb{R}_{\geq 0}^k$ such that both have at least one of their co-ordinates 0, then $\mathbf{u} - \mathbf{v}$ belongs to the null space i.e. $\mathbf{u} - \mathbf{v} = \lambda \mathbf{w}$ for some non-zero scalar λ . If the same co-ordinate of \mathbf{u}, \mathbf{v} are 0 i.e. for some $i \in [k], u_i = v_i = 0$, then since $\frac{a_i}{a_k}$ is non-zero $\lambda = 0 \Rightarrow \mathbf{u} = \mathbf{v}$, a contradiction. If different co-ordinates of \mathbf{u}, \mathbf{v} are 0, say $u_1 = v_2 = 0$, then since $\mathbf{u}, \mathbf{v} \in \mathbb{R}_{\geq 0}^k$, $u_1 - v_1$ is negative and $u_2 - v_2$ is positive. This is not possible since both these quantities should have the same sign as λ , as all co-ordinates of \mathbf{w} are strictly positive. Therefore we arrive at a contradiction and there is a unique solution.

Having proved this uniqueness, finding the solution is easy. We perform the above-mentioned row transformations and obtain the general solution. Then for each $i \in [k]$, we set $x_i = 0$ and try to solve for the other variables $x_1, \dots, x_{i-1}, x_{i+1}, \dots, x_k \in \mathbb{R}_{\geq 0}$. By the above argument, we will get a valid solution for only one such i , which we return as the unique solution. Clearly it takes $k^{O(1)}$ time.

A.2 PROOF OF LEMMA 4.2

1. Since $V_1 \prec \dots \prec V_{N+1}$ is a topological order, marginalizing over V_{N+1} in the factorization $\mathbb{P}(\mathbf{V}) = \prod_{i=1}^{N+1} \mathbb{P}(V_i | \mathbf{pa}(V_i))$ will give the factorization $\mathbb{P}(\mathbf{V}_N) = \prod_{i=1}^N \mathbb{P}(V_i | \mathbf{pa}(V_i))$ which is the factorization over \mathcal{G}_N .
2. Let $\mathcal{T} = \{(\pi_1, \mathbf{t}_1), \dots, (\pi_m, \mathbf{t}_m)\}$ be a set of intervention tuples generating $\mathbb{P}_{mix}(\mathbf{V})$ and $C_{V_{N+1}} = \{v^1, \dots, v^k\}$. Marginalizing with respect to V_{N+1} for a single target $\mathbf{t}_i, i \in [m]$ looks like,

$$\sum_{l=1}^k \mathbb{P}_{\mathbf{t}_i}(\mathbf{V}_N, v^l) = \begin{cases} \mathbb{P}_{\mathbf{t}_i}(\mathbf{V}_N) & : V_{N+1} \notin T_i \\ \mathbb{P}_{\mathbf{t}_i \setminus \{v^j\}}(\mathbf{V}_N) & : v^j \in \mathbf{t}_i \end{cases} \quad (1)$$

Applying this marginalization to Equation 2, i.e., $\mathbb{P}_{mix}(\mathbf{V}) = \sum_{i=1}^m \pi_i \mathbb{P}_{\mathbf{t}_i}(\mathbf{V})$, gives a convex linear combination of different $\mathbb{P}_{\mathbf{s}}(\mathbf{V}_N)$, where \mathbf{s} are values of some set of variables $\mathcal{S} \subset \mathbf{V}_N$, implying that $\mathbb{P}_{mix}(\mathbf{V}_N)$ is a mixture of interventions on \mathcal{G}_N .

3. This is straight-forward. To query $\mathbb{P}(v_1, \dots, v_N)$, we query $\mathbb{P}(v_1, \dots, v_N, v_{N+1})$ for all $v_{N+1} \in C_{V_{N+1}}$ and sum them up. The same can be done to create access for $\mathbb{P}_{mix}(\mathbf{V}_N)$. Since we are summing at most k_{max} terms, in $O(k_{max})$ time we can simulate access to both $\mathbb{P}(\mathbf{V}_N)$, $\mathbb{P}_{mix}(\mathbf{V}_N)$.

A.3 PROOF OF LEMMA 4.3

This follows by the marginalization equation (Equation 1 in Appendix A.2).

A.4 PROOF OF LEMMA 4.5

Let $\mathbf{s} \in \mathcal{S}_r$ where $r > i$. This means that \mathbf{s} is either \mathbf{s}_r or $\mathbf{s}_r \cup \{v\}$ for some $v \in C_{V_{N+1}}$. We show that $\mathbb{P}_{\mathbf{s}_r}(v_{i,l}) = \mathbb{P}_{\mathbf{s}_r}(\mathbf{s}_i \cup \mathbf{s}_{-i} \cup \{v^l\}) = 0$ for all $v^l \in C_{V_{N+1}}$. The proof for $\mathbf{s} = \mathbf{s}_r \cup \{v\}$ is identical. Since $i < r$, we get that $\mathbf{s}_r \not\subseteq \mathbf{s}_i$. Now there are two cases, either the set of variables $\mathcal{S}_r \subseteq \mathcal{S}_i$ or $\mathcal{S}_r \not\subseteq \mathcal{S}_i$. In the first case, since $\mathbf{s}_r \not\subseteq \mathbf{s}_i$ we get that there is some variable $V_j \in \mathcal{S}_r, \mathcal{S}_i$ such that different values v_j^r and v_j^i belong to \mathbf{s}_r and \mathbf{s}_i respectively implying that $\mathbb{P}_{\mathbf{s}_r}(\mathbf{s}_i \cup \mathbf{s}_{-i} \cup \{v\}) = 0$. In the second case, there is some variable $V_j \in \mathcal{S}_r$ ($j \in [N]$) that is not in \mathcal{S}_i . Note that the ‘‘missing value’’ $\bar{v}_j \in C_{V_j}$ (i.e. the one that is missing from all targets $\mathbf{s}_j, j \in [q]$) belongs to \mathbf{s}_{-i} (since $V_j \notin \mathcal{S}_i$) but it cannot belong to \mathbf{s}_r (since it is missing from all $\mathbf{s}_1, \dots, \mathbf{s}_q$) $\Rightarrow \mathbb{P}_{\mathbf{s}_r}(\mathbf{s}_i \cup \mathbf{s}_{-i} \cup \{v\}) = 0$.

A.5 PROOF OF LEMMA 4.6

Let $\mathbf{s} \in \mathcal{S}_i$. Note that \mathbf{s} is either \mathbf{s}_i or $\mathbf{s}_i \cup \{v\}$ for some $v \in C_{V_{N+1}}$. Using Equation 1 we get that for all $\mathbf{s} \in \mathcal{S}_i$, marginalization of V_{N+1} in $\pi_{\mathbf{s}} \mathbb{P}_{\mathbf{s}}(\mathbf{V})$ gives $\pi_{\mathbf{s}} \mathbb{P}_{\mathbf{s}_i}(\mathbf{V}_N)$. Marginalizing V_{N+1} in Equation 3, converts the left hand side to $\mathbb{P}_{mix}(\mathbf{V}_N)$ and right-hand side to $\sum_{i=1}^q (\sum_{\mathbf{s} \in \mathcal{S}_i} \pi_{\mathbf{s}}) \mathbb{P}_{\mathbf{s}_i}(\mathbf{V}_N)$ giving a set of intervention tuples that generates $\mathbb{P}_{mix}(\mathbf{V}_N)$. By the inductive hypothesis, $\mathbb{P}_{mix}(\mathbf{V}_N)$ is generated by the unique set of intervention tuples \mathcal{S} satisfying Assumption 3.1. Since the intervention targets in \mathcal{S} and the ones in the set of intervention tuples we just obtained are the same, using uniqueness of \mathcal{S} we get that $\mu_i = \sum_{\mathbf{s} \in \mathcal{S}_i} \pi_{\mathbf{s}}$.

A.6 PROOF OF LEMMA 4.7

Note that since V_{N+1} is the last node in the topological order, using the definition of interventions, we can conclude that,

$$\mathbb{P}_{\mathbf{s}_i \cup \{v^j\}}(\mathbf{s}_i \cup \mathbf{s}_{-i} \cup \{v^l\}) = \mathbb{P}_{\mathbf{s}_i}(\mathbf{s}_i \cup \mathbf{s}_{-i}) \delta_{v^j, v^l}$$

where $v^l, v^j \in C_{V_{N+1}} = \{v^1, \dots, v^k\}$. Recall that $\mathbf{v}_{i,l} = \mathbf{s}_i \cup \mathbf{s}_{-i} \cup \{v^l\}$. Now, on substituting for Δ using Equation 4 into Equation 5, we obtain,

$$\begin{aligned} \mathbb{P}_{mix}(\mathbf{v}_{i,l}) - \mu_i \mathbb{P}_{\mathbf{s}_i}(\mathbf{v}_{i,l}) - \sum_{j=1}^{i-1} \sum_{\mathbf{s} \in \mathcal{S}_j} \pi_{\mathbf{s}} \mathbb{P}_{\mathbf{s}}(\mathbf{v}_{i,l}) = \\ \sum_{j \in [k]} \pi_{\mathbf{s}_i \cup \{v^j\}} \left(\mathbb{P}_{\mathbf{s}_i}(\mathbf{s}_i \cup \mathbf{s}_{-i}) \delta_{v^j, v^l} - \mathbb{P}_{\mathbf{s}_i}(\mathbf{v}_{i,l}) \right) \end{aligned}$$

Note that all the unknown variables are on the right-hand side of this equation. Varying $l \in [k]$, gives us a linear system of equations satisfied by scalars $\pi_{\mathbf{s}_i \cup \{v^l\}}$.

$$\begin{bmatrix} c - a_1 & -a_1 & \cdot & \cdot & -a_1 \\ -a_2 & c - a_2 & \cdot & \cdot & -a_2 \\ \cdot & \cdot & \cdot & \cdot & \cdot \\ -a_k & -a_k & \cdot & \cdot & c - a_k \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \\ \cdot \\ x_k \end{bmatrix} = \begin{bmatrix} b_1 \\ b_2 \\ \cdot \\ b_k \end{bmatrix}$$

In the above system, we have renamed the known values as follows. For $l \in [k]$, denote

$$\begin{aligned} a_l &= \mathbb{P}_{\mathbf{s}_i}(\mathbf{v}_{i,l}), \\ b_l &= \mathbb{P}_{mix}(\mathbf{v}_{i,l}) - \mu_i \mathbb{P}_{\mathbf{s}_i}(\mathbf{v}_{i,l}) - \sum_{j=1}^{i-1} \sum_{\mathbf{s} \in \mathcal{S}_j} \pi_{\mathbf{s}} \mathbb{P}_{\mathbf{s}}(\mathbf{v}_{i,l}), \\ c &= \mathbb{P}_{\mathbf{s}_i}(\mathbf{s}_i \cup \mathbf{s}_{-i}) \end{aligned}$$

All a_l 's are probabilities from interventional distributions and can be computed as product of conditional probabilities. Thus, by Assumption 3.2, $a_l > 0$ for all $l \in [k]$. It's easy to see that $c = \sum_{l \in [k]} a_l$, by the sum rule of probability. By statement of Lemma 4.3, for each $i \in [q]$ and $l \in [k]$, $\pi_{\mathbf{s}_i \cup \{v^l\}} \geq 0$. Since we are only considering set of intervention tuples which satisfy Assumption 3.1, there is some $l \in [k]$ such that $\pi_{\mathbf{s}_i \cup \{v^l\}} = 0$. On constraining the variables x_1, \dots, x_k in the above system to these conditions (i.e. $x_i \geq 0$ for all $i \in [q]$ and $x_i = 0$ for some $i \in [q]$), by Lemma 4.1 we are guaranteed a unique solution. Therefore there is a unique tuple $(\pi_{\mathbf{s}_i \cup \{v^1\}}, \dots, \pi_{\mathbf{s}_i \cup \{v^k\}})$ satisfying these requirements \Rightarrow Equation 5 has a unique solution which is easily computed in $k^{O(1)}$ time using the technique described in proof of Lemma 4.1.

B NON-NECESSITY OF ASSUMPTION 3.2

With the help of an example we argue that Assumption 3.2 is not necessary in Theorem 3.1.

Example B.1. Consider a causal Bayesian Network

$$V_1 \rightarrow V_2$$

defined over two binary variable $\mathbf{V} = \{V_1, V_2\}$ with $C_{V_i} = \{0, 1\}$, $i \in [2]$. Further, define CPDs, $\mathbb{P}(V_1 = 1) = 0.5$, $\mathbb{P}(V_2 = 1|V_1 = 0) = 0.5$, and $\mathbb{P}(V_2 = 1|V_1 = 1) = 0$. Clearly $\mathbb{P}(V_2 = 1, V_1 = 1) = 0$ implying that this CBN doesn't satisfy Assumption 3.2.

Let $\mathbb{P}_{mix}(\mathbf{V})$ be a mixture distribution defined as

$$\frac{1}{2} \mathbb{P}(\mathbf{V} | do(V_1 = 0)) + \frac{1}{2} \mathbb{P}(\mathbf{V} | do(V_1 = 0, V_2 = 0))$$

This mixture satisfies Assumption 3.1. Our algorithm first marginalizes on V_2 and tries to find the unique set of intervention targets for $\mathbb{P}_{mix}(V_1)$. For this sub-problem, all steps of Algorithm 1 go through (since the distribution $\mathbb{P}(V_1)$ satisfies positivity), and the correct components get identified. Note that, for this sub-problem the algorithm identifies that $\mathbb{P}_{mix}(V_1) = \mathbb{P}(V_1 | do(V_1 = 0))$. Now it tries to lift this computed target ($V_1 = 0$) to targets for the full mixture $\mathbb{P}_{mix}(\mathbf{V})$.

Since the algorithm does not try to lift the target ($V_1 = 1$) (as it was not found as a target for $\mathbb{P}_{mix}(V_1)$), it does not require $\mathbb{P}(V_2|V_1 = 1)$ to be non-zero. This can be easily checked in the lifting process described in Section 4.2.1. We do not repeat the steps of our algorithm here and encourage the reader to work through the lifting steps outlined in Section 4.2.1 and obtain unique solutions proving our point that Assumption 3.2 is not necessary and can be weakened.

C WORKED-OUT EXAMPLES

In this section, we illustrate the workings of Algorithm 1 (in the main paper) using two worked-out examples. Example C.1 is simpler and uses a mixture distribution on a CBN with just two nodes. It does not really require all of the crucial ideas from the lifting procedure described in Section 4.2.1. However, we believe it is important since it gives a good broad understanding of the entire algorithm. Example C.2 is complicated enough (using a mixture distribution on a CBN with three nodes) to highlight some of the key novelties of our lifting procedure in Section 4.2.1. We urge the reader to first work through Example C.1 and then through Example C.2 to get a full understanding of the critical ideas that make our proof of Theorem 3.1 work.

Example C.1. Consider a causal Bayesian Network

$$V_1 \rightarrow V_2$$

defined over two binary variable $\mathbf{V} = \{V_1, V_2\}$ with $C_{V_i} = \{0, 1\}$, $i \in [2]$. Further, define CPDs, $\mathbb{P}(V_1 = 1) = 0.5$, $\mathbb{P}(V_2 = 1|V_1) = 0.5$ for both $V_1 = 0$ and $V_1 = 1$. Clearly, this CBN satisfies Assumption 3.2. Let $\mathbb{P}_{mix}(\mathbf{V})$ be a mixture distribution defined as

$$\frac{1}{2}\mathbb{P}(\mathbf{V}|do(V_1 = 0)) + \frac{1}{2}\mathbb{P}(\mathbf{V}|do(V_1 = 0, V_2 = 0))$$

This mixture satisfies Assumption 3.1. On marginalizing variable V_2 , the mixture becomes $\mathbb{P}_{mix}(V_1) = \mathbb{P}(V_1|do(V_1 = 0))$ which also satisfies Assumption 3.1, as a mixture of interventions on the CBN comprising of the single variable V_1 . A general mixture distribution on variable V_1 looks like,

$$\pi_0\mathbb{P}(V_1|do(V_1 = 0)) + \pi_1\mathbb{P}(V_1|do(V_1 = 1)) + (1 - \pi_0 - \pi_1)\mathbb{P}(V_1)$$

Varying V_1 in $\{0, 1\}$ gives the system of equations,

$$\begin{bmatrix} 1 - 0.5 & -0.5 \\ -0.5 & 1 - 0.5 \end{bmatrix} \begin{bmatrix} \pi_0 \\ \pi_1 \end{bmatrix} = \begin{bmatrix} \mathbb{P}_{mix}(V_1 = 0) - 0.5 \\ \mathbb{P}_{mix}(V_1 = 1) - 0.5 \end{bmatrix}$$

Lemma 4.1 highlights why such a system has a unique non-negative solution if we assume that at least one of π_0, π_1 is 0, i.e. the set of intervention tuples we are trying to construct satisfy Assumption 3.1. Lemma 4.1 also provides an efficient algorithm to find the unique solution giving $\mathbb{P}_{mix}(V_1) = \mathbb{P}(V_1|do(V_1 = 0))$. Now we have one intervention target ($V_1 = 0$) at hand, which we will try to lift. Note that, the possible lifts of such a target are $(V_1 = 0)$, $(V_1 = 0, V_2 = 0)$, $(V_1 = 0, V_2 = 1)$. Our next step will search within this space of targets and try to complete the construction. A general solution for mixtures within this space would look like

$$\mu_0\mathbb{P}(\mathbf{V}|do(V_1 = 0)) + \mu_1\mathbb{P}(\mathbf{V}|do(V_1 = 0, V_2 = 0)) + \mu_2\mathbb{P}(\mathbf{V}|do(V_1 = 0, V_2 = 1))$$

where $\mu_i \geq 0, i \in \{0, 1, 2\}$. So we want to find μ_0, μ_1, μ_2 such that the above general solution becomes equal to $\mathbb{P}_{mix}(\mathbf{V})$. Since we already know that $\mathbb{P}_{mix}(V_1) = \mathbb{P}(V_1|do(V_1 = 0))$, marginalizing on V_2 implies that $\mu_0 + \mu_1 + \mu_2 = 1$. We substitute $\mu_0 = 1 - \mu_1 - \mu_2$, equate the above general mixture to $\mathbb{P}_{mix}(\mathbf{V})$, and re-arrange terms to get,

$$\begin{aligned} \mathbb{P}_{mix}(\mathbf{V}) - \mathbb{P}(\mathbf{V}|do(V_1 = 0)) &= \mu_1(\mathbb{P}(\mathbf{V}|do(V_1 = 0, V_2 = 0)) - \mathbb{P}(\mathbf{V}|do(V_1 = 0))) \\ &\quad + \mu_2(\mathbb{P}(\mathbf{V}|do(V_1 = 0, V_2 = 1)) - \mathbb{P}(\mathbf{V}|do(V_1 = 0))) \end{aligned}$$

We evaluate this mixture at settings $V_1 = 0, V_2 = 0$, and $V_1 = 0, V_2 = 1$, to get a system of linear equations in (μ_1, μ_2) ,

$$\begin{bmatrix} 1 - 0.5 & -0.5 \\ -0.5 & 1 - 0.5 \end{bmatrix} \begin{bmatrix} \mu_1 \\ \mu_2 \end{bmatrix} = \begin{bmatrix} \mathbb{P}_{mix}(0, 0) - 0.5 \\ \mathbb{P}_{mix}(0, 1) - 0.5 \end{bmatrix}$$

Again, since we are solving for a set of intervention tuples satisfying Assumption 3.1, one of μ_1, μ_2 would be 0. This used with Lemma 4.1 gives the unique solution $\mu_1 = 0.5, \pi_2 = 0 \Rightarrow \mu_0 = 0.5$, thereby identifying the correct set of intervention tuples.

We would like to note that the example we illustrated above is rather simple and does not capture some main non-trivial aspects of our Algorithm. However, we think it is important as a warm up exercise. A more involved example that would bring out the crucial ideas of our proof is presented below in Example C.2.

Example C.2. Consider a causal Bayesian Network defined over three binary variables $\mathbf{V} = \{V_1, V_2, V_3\}$ taking values in $C_{V_i} = \{0, 1\}$, $i \in [3]$, with $\mathbf{pa}(V_1) = \emptyset$, $\mathbf{pa}(V_2) = \{V_1\}$ and $\mathbf{pa}(V_3) = \{V_1, V_2\}$. Let CPDs be such that Assumption 3.2 is satisfied. Consider a mixture distribution

$$\mathbb{P}_{mix}(\mathbf{V}) = \mu_0 \mathbb{P}(\mathbf{V} | do(V_1 = 0)) + \mu_1 \mathbb{P}(\mathbf{V} | do(V_1 = 0, V_2 = 0)) + \mu_2 \mathbb{P}(\mathbf{V} | do(V_1 = 0, V_2 = 0, V_3 = 0))$$

with positive scalars μ_0, μ_1, μ_2 satisfying $\mu_0 + \mu_1 + \mu_2 = 1$. Marginalizing on V_3 , gives,

$$\mathbb{P}_{mix}(V_1, V_2) = \pi_0 \mathbb{P}(V_1, V_2 | do(V_1 = 0)) + \pi_1 \mathbb{P}(V_1, V_2 | do(V_1 = 0, V_2 = 0))$$

Our inductive hypothesis assumes that this mixture on smaller number of nodes is generated by a unique set of intervention tuples that satisfies Assumption 3.1 and also that this set can be efficiently computed. This will give us access to the two scalars π_0, π_1 and to the two targets $(V_1 = 0)$ and $(V_1 = 0, V_2 = 0)$ (we call them the currently computed targets). We need to lift these targets to targets for the original mixture distribution like we did in Example C.1. However, the situation is not as simple here. Note that $(V_1 = 0)$ can be lifted to one of the three targets $(V_1 = 0), (V_1 = 0, V_3 = 0), (V_1 = 0, V_3 = 1)$. Similarly $(V_1 = 0, V_2 = 0)$ can be lifted to one of the three targets $(V_1 = 0, V_2 = 0), (V_1 = 0, V_2 = 0, V_3 = 0), (V_1 = 0, V_2 = 0, V_3 = 1)$. So there are 6 possible targets in the original mixture and therefore a general solution for our mixture can be written using 6 new variables (say $\delta_0, \delta_1, \delta_2, \delta_3, \delta_4, \delta_5$) such that,

$$\begin{aligned} \mathbb{P}_{mix}(\mathbf{V}) = & \delta_0 \mathbb{P}(\mathbf{V} | do(V_1 = 0)) + \delta_1 \mathbb{P}(\mathbf{V} | do(V_1 = 0, V_3 = 0)) + \delta_2 \mathbb{P}(\mathbf{V} | do(V_1 = 0, V_3 = 1)) \\ & + \delta_3 \mathbb{P}(\mathbf{V} | do(V_1 = 0, V_2 = 0)) + \delta_4 \mathbb{P}(\mathbf{V} | do(V_1 = 0, V_2 = 0, V_3 = 0)) + \delta_5 \mathbb{P}(\mathbf{V} | do(V_1 = 0, V_2 = 0, V_3 = 1)) \end{aligned} \quad (2)$$

where all δ_i are non-negative. By marginalizing on V_3 , and using the solution we got from the inductive hypothesis (like we did in Example C.1), we can show that,

$$\pi_0 = \delta_0 + \delta_1 + \delta_2, \quad \pi_1 = \delta_3 + \delta_4 + \delta_5$$

Now the two main non trivial ingredients needed from here are:

- Deciding the order in which the currently computed targets should be lifted, and
- Deciding the settings for \mathbf{V} that would give linear systems where we can argue about unique solutions like in Example C.1.

For the first one, we lift the currently computed targets in an order which does not violate set inclusion for these targets, i.e. we first lift $(V_1 = 0)$ and then lift $(V_1 = 0, V_2 = 0)$. This can be done by considering any extension of the set inclusion partial order on these targets. Then for lifting the target $(V_1 = 0)$, we choose to evaluate on the settings $\mathbf{v}_1 = (V_1 = 0, V_2 = 1, V_3 = 0), \mathbf{v}_2 = (V_1 = 0, V_2 = 1, V_3 = 1)$. Here, we pick the value of V_2 that is missing from the currently computed target under consideration i.e. $(V_1 = 0)$. There will always be one such missing value (it follows from Assumption 3.1). Evaluating on these settings simplifies our equation drastically. For $l \in [2]$, we get,

$$\mathbb{P}_{mix}(\mathbf{v}_l) = \delta_0 \mathbb{P}(\mathbf{v}_l | do(V_1 = 0)) + \delta_1 \mathbb{P}(\mathbf{v}_l | do(V_1 = 0, V_3 = 0)) + \delta_2 \mathbb{P}(\mathbf{v}_l | do(V_1 = 0, V_3 = 1))$$

Basically all possible lifts of the other currently computed target i.e. $(V_1 = 0, V_2 = 0)$ vanish and we have a much simpler system of equations at hand. From here the solution follows exactly like the previous example. We substitute $\delta_0 = \pi_0 - \delta_1 - \delta_2$ and rearrange to get a linear system in 2 equations and 2 variables δ_1, δ_2 . Similar to the argument made in Example C.1, at least one of δ_1, δ_2 will be 0 and therefore this system has a unique

solution (using Lemma 4.1) giving values of $\delta_0, \delta_1, \delta_2$. These can then be substituted back in Equation 2 reducing the number of variables to 3 (i.e. $\delta_3, \delta_4, \delta_5$). Again we substitute $\delta_3 = \pi_1 - \delta_4 - \delta_5$ and reduce the equation to just two unknowns. Finally by using settings $\mathbf{v}_1 = (V_1 = 0, V_2 = 0, V_3 = 0)$, $\mathbf{v}_2 = (V_1 = 0, V_2 = 0, V_3 = 1)$ in Equation 2 we will be left with 2 equations in 2 variables. Exactly like our argument for lifting of $(V_1 = 0)$ (i.e. using Lemma 4.1), we can show that, when at least one of δ_4, δ_5 is 0, this system has a unique solution as well. Lemma 4.1 would efficiently give us the values of δ_4, δ_5 and therefore of δ_3 . Thus, we uniquely identify a set of intervention tuples that satisfies Assumption 3.1 and generates $\mathbb{P}_{mix}(\mathbf{V})$.

D FINITE SAMPLE ALGORITHM

In a real world scenario, we will only have finitely many samples from the distributions $\mathbb{P}(\mathbf{V})$ and $\mathbb{P}_{mix}(\mathbf{V})$. We still assume access to the underlying causal graph G (the CPDs might not be known though). In this situation, we modify Algorithm 1 (from the main paper) slightly to make it work with finitely many samples. The resulting algorithm is presented in Algorithm 1 (Appendix D). Let the sets containing the samples be $\mathcal{B} = \{\mathbf{b}_1, \dots, \mathbf{b}_M\}$ where $\mathbf{b}_i \sim \mathbb{P}(\mathbf{V})$ and $\mathcal{B}_{mix} = \{\mathbf{b}_1^{mix}, \dots, \mathbf{b}_M^{mix}\}$ where $\mathbf{b}_j^{mix} \sim \mathbb{P}_{mix}(\mathbf{V})$. As a preprocessing step, we estimate the distributions $\mathbb{P}(\mathbf{V})$ and $\mathbb{P}_{mix}(\mathbf{V})$ as $\hat{\mathbb{P}}(\mathbf{V})$ and $\hat{\mathbb{P}}_{mix}(\mathbf{V})$ (respectively) using samples in \mathcal{B} and \mathcal{B}_{mix} respectively. $\hat{\mathbb{P}}(\mathbf{V})$ is estimated by estimating all the CPDs using maximum likelihood estimation (MLE). In our implementation, we use a function from the pgmpy library (Ankan and Panda [2015]), to compute these MLE estimates. We enforce Assumption 3.2 on $\hat{\mathbb{P}}(\mathbf{V})$, by enforcing it on all it's CPDs, using a small positive parameter δ (chosen by us). When $\hat{\mathbb{P}}(v_i|\mathbf{pa}(v_i)) = 0$ for some v_i and some setting of parents $\mathbf{pa}(v_i)$, we update,

$$\hat{\mathbb{P}}(V_i|\mathbf{pa}(v_i)) \leftarrow \hat{\mathbb{P}}(V_i|\mathbf{pa}(v_i)) + \delta$$

for all values of V_i , and then re-normalize to make it a probability distribution again. Marginal $\hat{\mathbb{P}}_{mix}(\mathbf{V} = \mathbf{v})$ is calculated using relative frequency of the occurrence of $\mathbf{V} = \mathbf{v}$ in the samples inside \mathcal{B}_{mix} . These estimated distributions are then used as inputs in Algorithm 1 (Appendix D). We use another small positive parameter ϵ as input to Algorithm 1 (Appendix D) which prunes each recovered set of intervention tuples computed during the algorithm, by only keeping mixing coefficients greater than ϵ . It's easy to see that the time complexity of Algorithm 1 (Appendix D), including the estimation of probabilities from samples is:

$$\left(\frac{N k_{max}^d M}{\epsilon} \right)^{O(1)}$$

Here N is number of nodes in the causal graph G , d is the maximum in-degree of any node in G , k_{max} is maximum number of values that any node in G can take and M is the number of samples present in \mathcal{B} and \mathcal{B}_{mix} .

Remark. *Since our algorithm's run-time depends on ϵ , we need to carefully select it's value. Setting it too small could increase the run time whereas setting it too big could lead to wrongfully pruning intervention targets (with significant mixing proportions) present in the mixture.*

E ADDITIONAL SIMULATIONS

E.1 EFFECT OF GRAPH SIZE

Figure 1 (Appendix E) shows the variation of performance of Algorithm 1 (Appendix D) keeping the number of samples fixed at $\sim 10^6$. We observe that recall decreases and root-mean-squared error in mixing coefficient increases very quickly as the number of nodes increases in the graph. Even though this is expected since error is accumulated as we successively add nodes and find new intervention targets, such performance for a very large sample size indicates bad dependence of sample complexity on the number of nodes. Improving this needs more exploration and is left for future work.

E.2 EFFECT OF GRAPH TYPE

In Figure 2 (Appendix E), we demonstrate performance of Algorithm 1 (Appendix D) for CBNs generated from two different family of random graphs (Erdős-Rényi (ER) and Scale-Free (SF)). We observe no significant

Algorithm 1: DISENTANGLE-FINITE

input : \mathbf{V} , Causal Graph G , $\hat{\mathbb{P}}(\mathbf{V})$, $\hat{\mathbb{P}}_{mix}(\mathbf{V})$, ϵ
output : Set of intervention tuples \mathcal{T}

1. When $|\mathbf{V}| = 1$, setup the linear system in Equation 1 (say $\mathbf{A}\mathbf{x} = \mathbf{b}$) using the estimated distributions. Similar to the technique described in Lemma 4.1, set one variable to 0 at a time giving solution (π_1, \dots, π_k) corresponding to targets $(\mathbf{t}_1, \dots, \mathbf{t}_k)$ as described in Section 4.1. For every variable that is set to 0, create a set $\mathcal{T} = \{(\mathbf{t}_i, \pi_i) : i \in [k]\}$ containing the solution. For every such \mathcal{T} , iterate through the tuples (\mathbf{t}_i, π_i) in it. If some $\pi_i < 0$, set $\pi_i \leftarrow 0$. Compute the score $r(\mathcal{T}) = \|\mathbf{A}\boldsymbol{\pi} - \mathbf{b}\|^2$, where $\boldsymbol{\pi} = (\pi_1, \dots, \pi_k)$. Next, select \mathcal{T} with the smallest value of $r(\mathcal{T})$. For this selected \mathcal{T} , check if $1 - \sum_{i=1}^k \pi_i < \epsilon$. If yes, renormalize $\pi_i \leftarrow \pi_i / (\sum_{i=1}^k \pi_i)$. If no, add the tuple $(\mathbf{t}_0, 1 - \sum_{i=1}^k \pi_i)$ to \mathcal{T} . Only keep the tuples with strictly positive mixing coefficients i.e. $\mathcal{T} \leftarrow \{(\mathbf{t}_i, \pi_i) \in \mathcal{T} : \pi_i > 0\}$. **return** \mathcal{T} .
 2. Let $V_1 \prec \dots \prec V_{N+1}$ denote a topological order in G . Marginalize on V_{N+1} to create access to $\hat{\mathbb{P}}_{mix}(\mathbf{V}_N)$ and $\hat{\mathbb{P}}(\mathbf{V}_N)$ where $\mathbf{V}_N = (V_1, \dots, V_N)$. Construct $G_N = G \setminus \{V_{N+1}\}$. Recursively call this algorithm with inputs G_N , $\hat{\mathbb{P}}(\mathbf{V}_N)$, $\hat{\mathbb{P}}_{mix}(\mathbf{V}_N)$, and obtain a set of intervention tuples $\mathcal{S} = \{(\mathbf{s}_1, \mu_1), \dots, (\mathbf{s}_q, \mu_q)\}$. Let $\mathbf{s}_1, \dots, \mathbf{s}_q$ be ordered such that $i \leq j$ implies that $\mathbf{s}_j \not\subseteq \mathbf{s}_i$. For all $i \in [N]$, by inspecting \mathbf{s}_j , identify $\bar{v}_i \in C_{V_i}$ such that $\bar{v}_i \notin \mathbf{s}_j$ for any $j \in [q]$. Define $\mathbf{s}_{-j} = \{\bar{v}_i : V_i \notin \mathbf{s}_j\}$. Let $C_{V_{N+1}} = \{v^1, \dots, v^k\}$. For each $i \in [q]$ and $l \in [k]$, create setting $\mathbf{v}_{i,l} = \mathbf{s}_i \cup \mathbf{s}_{-i} \cup \{v^l\}$.
 3. For each fixed $i \in [q]$, evaluate distributions for different $\mathbf{v}_{i,l}$, $l \in [k]$, to setup the system of equations (say $\mathbf{A}\mathbf{x} = \mathbf{b}$) described in Equation 5. Similar to the technique described in Lemma 4.7 (which in turn uses Lemma 4.1), set one variable to 0 at a time giving solution $(\pi_{\mathbf{s}_i \cup \{v^1\}}, \dots, \pi_{\mathbf{s}_i \cup \{v^k\}})$ corresponding to targets $(\mathbf{s}_i \cup \{v^1\}, \dots, \mathbf{s}_i \cup \{v^k\})$ as described in Section 4.2. For every variable that is set to 0, create a set $\mathcal{T} = \{(\mathbf{s}_i \cup \{v^l\}, \pi_{\mathbf{s}_i \cup \{v^l\}}) : l \in [k]\}$ containing the solution. For every such \mathcal{T} , iterate through the tuples $(\mathbf{s}_i \cup \{v^l\}, \pi_{\mathbf{s}_i \cup \{v^l\}})$ in it. If some $\pi_{\mathbf{s}_i \cup \{v^l\}} < 0$, set $\pi_{\mathbf{s}_i \cup \{v^l\}} \leftarrow 0$. Compute the score $r(\mathcal{T}) = \|\mathbf{A}\boldsymbol{\pi} - \mathbf{b}\|^2$, where $\boldsymbol{\pi} = (\pi_{\mathbf{s}_i \cup \{v^1\}}, \dots, \pi_{\mathbf{s}_i \cup \{v^k\}})$. Next, select \mathcal{T} with the smallest value of $r(\mathcal{T})$. For this selected \mathcal{T} , check if $\mu_i - \sum_{l=1}^k \pi_{\mathbf{s}_i \cup \{v^l\}} < \epsilon$. If yes, renormalize $\pi_{\mathbf{s}_i \cup \{v^l\}} \leftarrow (\mu_i \times \pi_{\mathbf{s}_i \cup \{v^l\}}) / (\sum_{l=1}^k \pi_{\mathbf{s}_i \cup \{v^l\}})$. If no, add the tuple $(\mathbf{s}_i, \mu_i - \sum_{l=1}^k \pi_{\mathbf{s}_i \cup \{v^l\}})$ to \mathcal{T} . At the end of this process, collect all the intervention tuples thus obtained (for all $i \in [q]$), in the set \mathcal{T} .
 4. Find the excluded value of node V_{N+1} , i.e. the value which is not present in any target in \mathcal{T} . If no such value exists, find $v \in C_{V_{N+1}}$ which minimizes $\sum_{i=1}^q \pi_{\mathbf{s}_i \cup \{v\}}$. For each $i \in [q]$, set $\pi_{\mathbf{s}_i \cup \{v\}} \leftarrow 0$. For each $i \in [q]$, renormalize the mixing coefficients $\pi_{\mathbf{s}_i \cup \{v^l\}} \leftarrow (\pi_{\mathbf{s}_i \cup \{v^l\}} \times \mu_i) / (\sum_{l=1}^k \pi_{\mathbf{s}_i \cup \{v^l\}})$. Only keep the tuples with strictly positive mixing coefficients in \mathcal{T} i.e. $\mathcal{T} \leftarrow \{(\mathbf{s}, \pi_{\mathbf{s}}) \in \mathcal{T} : \pi_{\mathbf{s}} > 0\}$. **return** \mathcal{T}
-

difference in performance for these models and make a conjecture that only high level graph parameters (such as number of nodes, edges, in-degree etc.) might be having an impact on performance and the topology (given these parameters) might not be that crucial.

F EVALUATION METRICS

Let \mathcal{T} denote the actual set of intervention targets and $\hat{\mathcal{T}}$ denote the set of intervention targets computed by our algorithm. Let $\pi_{\mathbf{t}}$, $\hat{\pi}_{\mathbf{s}}$ denote mixing coefficients of target \mathbf{t} , \mathbf{s} in \mathcal{T} and $\hat{\mathcal{T}}$ respectively. We use the following evaluation metrics to evaluate the performance of our algorithm.

1. **Recall:** Proportion of number of targets in \mathcal{T} that were correctly identified in $\hat{\mathcal{T}}$

$$\text{Recall} = \frac{|\mathcal{T} \cap \hat{\mathcal{T}}|}{|\hat{\mathcal{T}}|}.$$

2. **Root Mean Squared Error:** Root-mean-squared error (RMSE) in the mixing coefficients.

$$\text{RMSE} = \sqrt{\frac{\sum_{t \in \mathcal{T} \cap \hat{\mathcal{T}}} (\pi_t - \hat{\pi}_t)^2 + \sum_{t \in (\mathcal{T} \setminus \hat{\mathcal{T}})} (\pi_t)^2 + \sum_{t \in (\hat{\mathcal{T}} \setminus \mathcal{T})} (\hat{\pi}_t)^2}{|\mathcal{T} \cup \hat{\mathcal{T}}|}}$$

3. **False-Positive RMSE:** RMSE in the mixing coefficients of the incorrectly identified targets.

$$\text{FP-RMSE} = \sqrt{\frac{\sum_{t \in (\hat{\mathcal{T}} \setminus \mathcal{T})} (\hat{\pi}_t)^2}{|\hat{\mathcal{T}} \setminus \mathcal{T}|}}$$

4. **False-Negative RMSE:** RMSE in the mixing coefficients of targets not identified.

$$\text{FN-RMSE} = \sqrt{\frac{\sum_{t \in (\mathcal{T} \setminus \hat{\mathcal{T}})} (\pi_t)^2}{|\mathcal{T} \setminus \hat{\mathcal{T}}|}}$$

References

Ankur Ankan and Abinash Panda. pgmpy: Probabilistic graphical models using python. In *Proceedings of the 14th Python in Science Conference (SCIPY 2015)*. Citeseer, 2015.

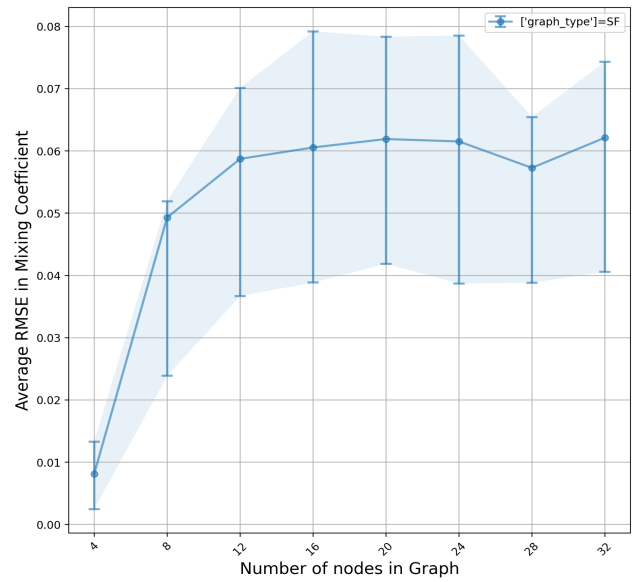
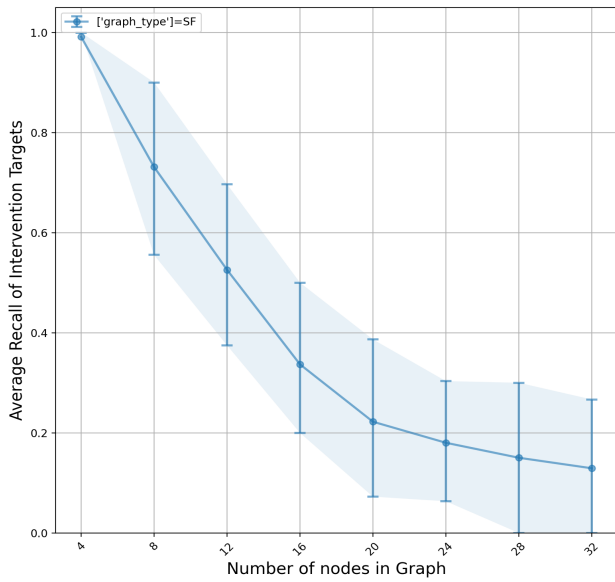


Figure 1: Performance of Algorithm 1 (Appendix D) as a function of number of nodes in the graph. The error bars show the 20% and 80% quantile respectively.

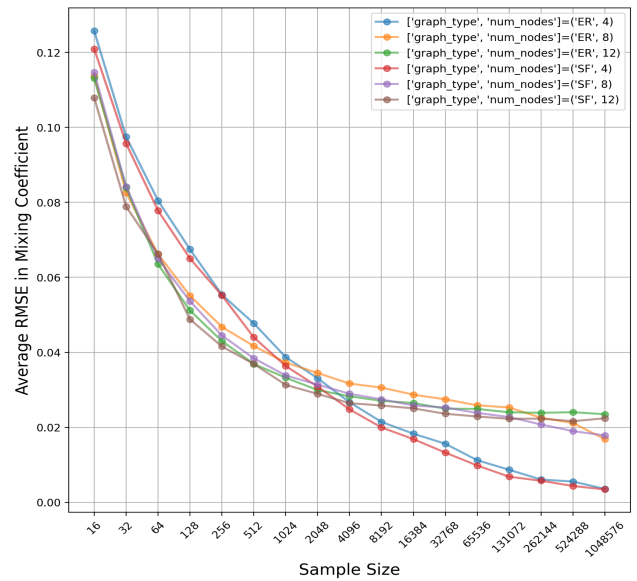
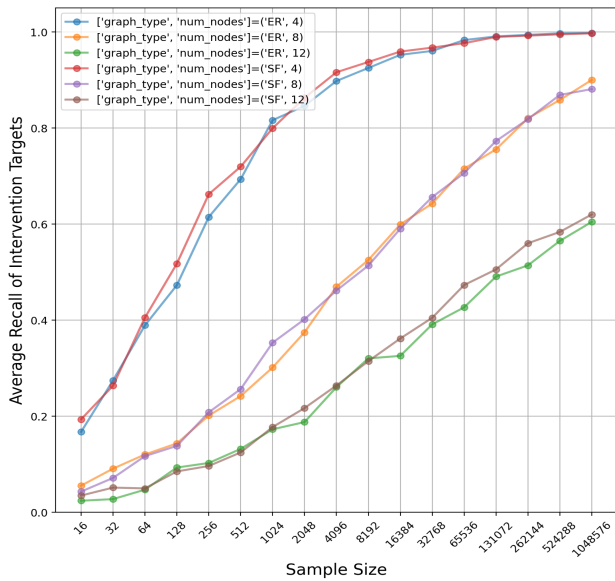


Figure 2: Comparison of performance of Algorithm 1 (Appendix D) for CBNs generated from Erdős-Rényi (ER) and Scale Free (SF) models.