

## I. CAUSAL GRAPHICAL MODEL WITH SOFT INTERVENTIONS

To represent causal effects, consider a DAG with structure  $(\mathcal{V}, \mathcal{B})$ , where  $\mathcal{V} = [N] \doteq \{1, \dots, N\}$  is the set of  $N$  nodes and  $\mathcal{B}$  is the set of directed edges. The observational (without intervention) edge-weight matrix  $\mathbf{B} \in \mathbb{R}^{N \times N}$  captures the strength of causal effects, where the  $(i, j)$ -th entry represents the weight of the edge  $i \rightarrow j$ .

To model causal effects under intervention, consider node-wise intervention, defined as

$$\mathbf{a} = (a_1, \dots, a_N)^\top \in \{0, 1\}^N, \quad (1)$$

where  $a_i$  represents whether node  $i$  is intervened (1) or not (0). Specifically, instead of hard interventions, we consider soft interventions, which do not cut off causal relationships between the intervened node and its parents, but change the incoming edges to the node. We denote the sets of parents of node  $i$  by  $\mathcal{P}_i(a_i)$ , the estimated set of parents by  $\hat{\mathcal{P}}_i(a_i)$ . The set difference of the estimated and true parent sets is denoted by  $\hat{\mathcal{P}}_i \setminus \mathcal{P}_i(a_i)$ .

Further, we denote the interventional edge-weight matrix by  $\mathbf{B}' \in \mathbb{R}^{N \times N}$ , such that the post-intervention weight matrix  $\mathbf{B}_a$  can be constructed as

$$[\mathbf{B}_a]_i = \mathbb{I}(a_i = 1) \mathbf{B}'_i + \mathbb{I}(a_i = 0) \mathbf{B}_i, \quad (2)$$

where  $\mathbb{I}(\cdot)$  is the indicator function and  $[\cdot]_i$  represents the  $i$ -th column of a matrix. The  $i$ -th column of the post-intervention weight matrix determines the set of parents of node  $i$  and how these parents causally influence node  $i$ .

As a result of the intervention, the vector of stochastic values associated with the nodes is represented by  $\mathbf{x} \in \mathbb{R}^N$ . The causal relationship among nodes is described by a linear structure equation model (LinSEM),

$$\mathbf{x} = (\mathbf{B}_a)^\top \mathbf{x} + \mathbf{e}, \quad (3)$$

where  $\mathbf{e}$  is a vector of exogenous variables, independent of  $\mathbf{x}$ . We assume that  $\mathbf{e}$  contains independent elements, with known means and unknown variances represented by  $\boldsymbol{\nu}$  and  $\boldsymbol{\epsilon}$ . The causal relationship described in (3) can be further manipulated, resulting in

$$\mathbf{x} = (\mathbf{I} - \mathbf{B}_a)^{-\top} \mathbf{e} \doteq (\mathbf{C}_a)^\top \mathbf{e}, \quad (4)$$

where  $\mathbf{C}_a$  denotes the post-intervention flow-weight matrix, whose  $(i, j)$ -th entry represents the net flow weight from node  $i$  to  $j$ . In this way, each variable  $x_i$  can be considered as a linear combination of exogenous variables in  $\mathbf{e}$ , weighted by the strength of causal flow. Thus, under a specific interven-

tion  $\mathbf{a}$ ,  $\mathbf{x}$  follows a multivariate distribution with mean and covariance defined as

$$\boldsymbol{\mu}(\mathbf{a}) \doteq \mathbb{E}[\mathbf{x}|\mathbf{a}] = (\mathbf{C}_a)^\top \boldsymbol{\nu}, \quad (5)$$

$$\boldsymbol{\Sigma}(\mathbf{a}) \doteq \mathbb{E}[(\mathbf{x} - \boldsymbol{\mu}(\mathbf{a}))(\mathbf{x} - \boldsymbol{\mu}(\mathbf{a}))^\top | \mathbf{a}] \quad (6)$$

$$= (\mathbf{C}_a)^\top \text{diag}(\boldsymbol{\epsilon}) \mathbf{C}_a. \quad (7)$$

Further, given an intervention selection strategy  $\pi$ , the mean and variance can be defined as

$$\boldsymbol{\mu}^\pi \doteq \mathbb{E}_\pi[\mathbf{x}] = \sum_{\mathbf{a}} \pi(\mathbf{a}) \mathbb{E}[\mathbf{x}|\mathbf{a}] = \sum_{\mathbf{a}} \pi(\mathbf{a}) \boldsymbol{\mu}(\mathbf{a}), \quad (8)$$

$$\boldsymbol{\Sigma}^\pi \doteq \mathbb{E}_\pi[(\mathbf{x} - \boldsymbol{\mu}(\mathbf{a}))(\mathbf{x} - \boldsymbol{\mu}(\mathbf{a}))^\top] = \sum_{\mathbf{a}} \pi(\mathbf{a}) \boldsymbol{\Sigma}(\mathbf{a}), \quad (9)$$

where  $\pi(\mathbf{a})$  stands for the intervention selection probability.

## II. PROOFS

With the whole set of potential parents, the estimated weights,  $\hat{\mathbf{B}}_{j, \hat{\mathcal{P}}_j}$  and residuals,  $\mathbf{r}_j(\hat{\mathcal{P}}_j)$  are given by linear minimum mean square error (MMSE) estimation as

$$\hat{\mathbf{B}}_{j, \hat{\mathcal{P}}_j} = (\mathbf{X}_{\hat{\mathcal{P}}_j}^\top \mathbf{X}_{\hat{\mathcal{P}}_j})^{-1} \mathbf{X}_{\hat{\mathcal{P}}_j}^\top (\mathbf{X}_j - \nu_j \mathbf{1}), \quad (10)$$

$$\mathbf{r}_j(\hat{\mathcal{P}}_j) = [\mathbf{I} - \mathbf{X}_{\hat{\mathcal{P}}_j} (\mathbf{X}_{\hat{\mathcal{P}}_j}^\top \mathbf{X}_{\hat{\mathcal{P}}_j})^{-1} \mathbf{X}_{\hat{\mathcal{P}}_j}^\top] (\mathbf{X}_j - \nu_j \mathbf{1}). \quad (11)$$

Denote the projection matrix onto the column space of  $\mathbf{X}_{\mathcal{P}}$  by

$$\Phi_{\mathbf{X}}(\mathcal{P}) \doteq \mathbf{X}_{\mathcal{P}} (\mathbf{X}_{\mathcal{P}}^\top \mathbf{X}_{\mathcal{P}})^{-1} \mathbf{X}_{\mathcal{P}}^\top, \quad (12)$$

which allows us to rewrite the residual vector as

$$\mathbf{r}_j(\hat{\mathcal{P}}_j) = [\mathbf{I} - \Phi_{\mathbf{X}}(\hat{\mathcal{P}}_j)] (\mathbf{X}_j - \nu_j \mathbf{1}) \quad (13)$$

$$= [\mathbf{I} - \Phi_{\mathbf{X}}(i)] [\mathbf{I} - \Phi_{\mathbf{X}}(\hat{\mathcal{P}}_j \setminus i)] (\mathbf{X}_j - \nu_j \mathbf{1}) \quad (14)$$

$$= [\mathbf{I} - \Phi_{\mathbf{X}}(i)] \mathbf{r}_j(\hat{\mathcal{P}}_j \setminus i). \quad (15)$$

We define the *normalized difference in squared residual norms* as

$$\Delta_{ij} \doteq \left( \|\mathbf{r}_j(\hat{\mathcal{P}}_j \setminus i)\|_2^2 - \|\mathbf{r}_j(\hat{\mathcal{P}}_j)\|_2^2 \right) / t_j(a_j), \quad (16)$$

where  $t_j(a_j)$  denotes the number of time slots when node  $j$  is under intervention  $a_j$ .

**Lemma 1.** If  $\mathcal{P}_j \subseteq \hat{\mathcal{P}}_j$  and  $i \in \hat{\mathcal{P}}_j \setminus \mathcal{P}_j$ , the following upper bound on  $\Delta_{ij}$  holds almost surely for sufficiently large  $t_j$ .

$$\Delta_{ij} \leq \overline{\Delta}_{ij}(\text{F parent}) \doteq \frac{\left( \epsilon_j^2 \sum_{\mathbf{a}} \pi(\mathbf{a}) [\mathbf{C}_a]_{ji} \right)^2}{\sum_{ii} \pi_{ii} + (\mu_i^\pi)^2}.$$

*Proof:* In both cases, the difference in squared norms can be reformulated as

$$\begin{aligned} & \|\mathbf{r}_j(\hat{\mathcal{P}}_j \setminus i)\|_2^2 - \|\mathbf{r}_j(\hat{\mathcal{P}}_j)\|_2^2 \\ &= [\mathbf{r}_j(\hat{\mathcal{P}}_j \setminus i) + \mathbf{r}_j(\hat{\mathcal{P}}_j)]^\top [\mathbf{r}_j(\hat{\mathcal{P}}_j \setminus i) - \mathbf{r}_j(\hat{\mathcal{P}}_j)] \end{aligned} \quad (17)$$

$$\stackrel{(a)}{=} \mathbf{r}_j(\hat{\mathcal{P}}_j \setminus i)^\top [2\mathbf{I} - \Phi_{\mathbf{X}}(i)] \Phi_{\mathbf{X}}(i) \mathbf{r}_j(\hat{\mathcal{P}}_j \setminus i) \quad (18)$$

$$\stackrel{(b)}{=} \mathbf{r}_j(\hat{\mathcal{P}}_j \setminus i)^\top \Phi_{\mathbf{X}}(i) \mathbf{r}_j(\hat{\mathcal{P}}_j \setminus i), \quad (19)$$

where (a) comes from the substitution of (15) and (b) from the fact that  $\Phi_{\mathbf{X}}(i)$  is a projection matrix.

Given the fact that  $i \in \hat{\mathcal{P}}_j \setminus \mathcal{P}_j$ , we can rewrite the residual vector as

$$\mathbf{r}_j(\hat{\mathcal{P}}_j \setminus i) = [\mathbf{I} - \Phi_{\mathbf{X}}(\hat{\mathcal{P}}_j \setminus i)](\mathbf{X}_j - \nu_j \mathbf{1}) \quad (20)$$

$$\stackrel{(a)}{=} [\mathbf{I} - \Phi_{\mathbf{X}}(\hat{\mathcal{P}}_j \setminus i)] [\mathbf{X}_{\mathcal{P}_j} \mathbf{B}_{j, \mathcal{P}_j} + \mathbf{e}_j - \nu_j \mathbf{1}] \quad (21)$$

$$\stackrel{(b)}{=} [\mathbf{I} - \Phi_{\mathbf{X}}(\hat{\mathcal{P}}_j \setminus i)] \tilde{\mathbf{e}}_j, \quad (22)$$

where  $\tilde{\mathbf{e}}_j \doteq \mathbf{e}_j - \nu_j \mathbf{1}$  denotes the vector of centered exogenous variable. Equation (a) is a direct subsequence of (3) while (b) holds because value vectors of parent nodes  $\mathbf{X}_{\mathcal{P}_j}$  exist exactly in the projection space. Based on the new formulation, the normalized difference can be bounded from above as

$$\Delta_{ij} = \frac{\tilde{\mathbf{e}}_j^\top [\mathbf{I} - \Phi_{\mathbf{X}}(\hat{\mathcal{P}}_j \setminus i)] \mathbf{X}_i \mathbf{X}_i^\top [\mathbf{I} - \Phi_{\mathbf{X}}(\hat{\mathcal{P}}_j \setminus i)] \tilde{\mathbf{e}}_j}{t_j \mathbf{X}_i^\top \mathbf{X}_i} \quad (23)$$

$$= \frac{(\tilde{\mathbf{e}}_j^\top [\mathbf{I} - \Phi_{\mathbf{X}}(\hat{\mathcal{P}}_j \setminus i)] \mathbf{X}_i / t_j)^2}{\mathbf{X}_i^\top \mathbf{X}_i / t_j} \quad (24)$$

$$\stackrel{(a)}{\leq} \frac{\left[ \sum_{\tau=1}^{t_j} (e_j^\tau - \nu_j) (\sum_k [C_{\mathbf{a}^\tau}]_{ki} e_k^\tau) / t_j \right]^2}{\sum_{\tau=1}^{t_j} (x_i^\tau)^2 / t_j} \quad (25)$$

$$\doteq \tilde{\Delta}_{ij}. \quad (26)$$

Inequality (a) holds because magnitude of the inner product is smaller in the projection space. As the number of samples  $t_j$  approaches infinity, the strong law of large numbers guarantees the following convergence result:

$$\frac{\sum_{\tau=1}^{t_j(0)} (e_j^\tau - \nu_j) \sum_k [C_{\mathbf{a}^\tau}]_{ki} e_k^\tau}{t_j} \xrightarrow{\text{a.s.}} \mathbb{E} \left[ \tilde{\mathbf{e}}_j \sum_k [C_{\mathbf{a}}]_{ki} e_k \right] \quad (27)$$

$$\stackrel{(a)}{=} \sum_{\mathbf{a}} \pi(\mathbf{a}) \mathbb{E} \left[ \tilde{\mathbf{e}}_j [C_{\mathbf{a}}]_{ji} e_j \mid \mathbf{a} \right] \quad (28)$$

$$= \epsilon_j^2 \sum_{\mathbf{a}} \pi(\mathbf{a}) [C_{\mathbf{a}}]_{ji}, \quad (29)$$

where  $\xrightarrow{\text{a.s.}}$  denotes almost sure convergence and equality (a) comes from the independence among exogenous variables. Similarly, the denominator of (25) converges as

$$\frac{\sum_{\tau=1}^{t_j} (x_i^\tau)^2}{t_j} \xrightarrow{\text{a.s.}} \mathbb{E}[x_i^2] = \Sigma_{ii}^\pi + (\mu_i^\pi)^2. \quad (30)$$

Lastly, according to the continuous mapping theorem [1] and the fact that multiplication preserves almost sure convergence [2], we have

$$\tilde{\Delta}_{ij} \xrightarrow{\text{a.s.}} \frac{(\epsilon_j^2 \sum_{\mathbf{a}} \pi(\mathbf{a}) [C_{\mathbf{a}}]_{ji})^2}{\Sigma_{ii}^\pi + (\mu_i^\pi)^2}, \text{ as } t_j \rightarrow \infty, \quad (31)$$

which completes the proof.  $\square$

**Lemma 2.** If  $\mathcal{P}_j \subseteq \hat{\mathcal{P}}_j$  and  $i \in \mathcal{P}_j$ , the following lower bound on  $\Delta_{ij}$  holds almost surely for sufficiently large  $t_j$ .

$$\Delta_{ij} \geq \underline{\Delta}_{ij}(\text{T parent}) \doteq \frac{B_{ij}^2}{\Sigma_{ii}^\pi + (\mu_i^\pi)^2} \times \max \left[ 0, \Sigma_{ii}^\pi + (\mu_i^\pi)^2 - \frac{\sum_{k \in \hat{\mathcal{P}}_j \setminus i} (\Sigma_{ik}^\pi + \mu_i^\pi \mu_k^\pi)^2}{\lambda_{\min}(\Sigma_{\hat{\mathcal{P}}_j \setminus i}^\pi + \boldsymbol{\mu}_{\hat{\mathcal{P}}_j \setminus i}^\pi (\boldsymbol{\mu}_{\hat{\mathcal{P}}_j \setminus i}^\pi)^\top)} \right]^2, \quad (32)$$

where  $\lambda_{\min}$  denotes the minimum eigenvalue of a matrix.

*Proof:* Since  $i \in \mathcal{P}_j$ , we can rewrite the normalized difference based on (19) as

$$\Delta_{ij} = \frac{((\mathbf{X}_j - \nu_j \mathbf{1})^\top [\mathbf{I} - \Phi_{\mathbf{X}}(\hat{\mathcal{P}}_j \setminus i)] \mathbf{X}_i)^2}{t_j \mathbf{X}_i^\top \mathbf{X}_i} \quad (33)$$

$$\stackrel{(a)}{=} \frac{[(B_{ij} \mathbf{X}_i + \tilde{\mathbf{e}}_j)^\top [\mathbf{I} - \Phi_{\mathbf{X}}(\hat{\mathcal{P}}_j \setminus i)] \mathbf{X}_i / t_j]^2}{\mathbf{X}_i^\top \mathbf{X}_i / t_j} \quad (34)$$

$$= \frac{1}{\mathbf{X}_i^\top \mathbf{X}_i / t_j} \left( \frac{\tilde{\mathbf{e}}_j^\top [\mathbf{I} - \Phi_{\mathbf{X}}(\hat{\mathcal{P}}_j \setminus i)] \mathbf{X}_i}{t_j} + B_{ij} \frac{\mathbf{X}_i^\top \mathbf{X}_i - \mathbf{X}_i^\top \Phi_{\mathbf{X}}(\hat{\mathcal{P}}_j \setminus i) \mathbf{X}_i}{t_j} \right)^2. \quad (35)$$

Equation (a) holds because the value vectors of all parents except for node  $i$  lie within the projection space. To understand the asymptotic behavior of  $\Delta_{ij}$ , we need to study its building components.

First, the convergence result of  $\mathbf{X}_i^\top \mathbf{X}_i / t_j$  is provided in (30). The third component can be bounded as

$$\frac{\mathbf{X}_i^\top \Phi_{\mathbf{X}}(\hat{\mathcal{P}}_j \setminus i) \mathbf{X}_i}{t_j} = \frac{\mathbf{X}_i^\top \mathbf{X}_{\hat{\mathcal{P}}_j \setminus i} (\mathbf{X}_{\hat{\mathcal{P}}_j \setminus i}^\top \mathbf{X}_{\hat{\mathcal{P}}_j \setminus i})^{-1} \mathbf{X}_{\hat{\mathcal{P}}_j \setminus i}^\top \mathbf{X}_i}{t_j} \stackrel{(a)}{\leq} \frac{\lambda_{\max}((\mathbf{X}_{\hat{\mathcal{P}}_j \setminus i}^\top \mathbf{X}_{\hat{\mathcal{P}}_j \setminus i})^{-1}) \|\mathbf{X}_i^\top \mathbf{X}_{\hat{\mathcal{P}}_j \setminus i}\|_2^2}{t_j} \quad (36)$$

$$= \frac{\sum_{k \in \hat{\mathcal{P}}_j \setminus i} (\mathbf{X}_i^\top \mathbf{X}_k / t_j)^2}{\lambda_{\min}(\mathbf{X}_{\hat{\mathcal{P}}_j \setminus i}^\top \mathbf{X}_{\hat{\mathcal{P}}_j \setminus i} / t_j)} \quad (37)$$

$$\stackrel{\text{a.s.}}{\rightarrow} \frac{\sum_{k \in \hat{\mathcal{P}}_j \setminus i} (\Sigma_{ik}^\pi + \mu_i^\pi \mu_k^\pi)^2}{\lambda_{\min}(\Sigma_{\hat{\mathcal{P}}_j \setminus i}^\pi + \boldsymbol{\mu}_{\hat{\mathcal{P}}_j \setminus i}^\pi (\boldsymbol{\mu}_{\hat{\mathcal{P}}_j \setminus i}^\pi)^\top)}. \quad (38)$$

Inequality (a) follows directly from the definition of maximum eigenvalue. The convergence relies on the strong law of large numbers, the fact that eigenvalues are continuous functions of the matrix [3] and the continuous mapping theorem [1]. Since magnitude of the last component is bounded from above by zero asymptotically,

$$\left| \frac{\tilde{\mathbf{e}}_j^\top [\mathbf{I} - \Phi_{\mathbf{X}}(\hat{\mathcal{P}}_j \setminus i)] \mathbf{X}_i}{t_j} \right| \leq \left| \frac{\tilde{\mathbf{e}}_j^\top \mathbf{X}_i}{t_j} \right| \xrightarrow{\text{a.s.}} |\mathbb{E}[\tilde{\mathbf{e}}_j x_i]| = 0, \quad (39)$$

the first component converges to zero almost surely.

Eventually, since addition and multiplication preserve almost sure convergence [2], combination of the aforementioned results completes the proof.  $\square$

## REFERENCES

- [1] A. W. Van der Vaart, *Asymptotic statistics*. Cambridge university press, 2000, vol. 3.
- [2] G. Casella and R. Berger, *Statistical inference*. CRC press, 2024.
- [3] R. A. Horn and C. R. Johnson, *Matrix analysis*. Cambridge university press, 2012.