

Supplementary material for Network structure learning under uncertain interventions

Federico Castelletti

Department of Statistical Sciences, Università Cattolica del Sacro Cuore, Milan

and

Stefano Peluso

Department of Statistics and Quantitative Methods, Università degli Studi di
Milano-Bicocca, Milan

Sections 1 and 2 contain proofs of Propositions 3.1 and 3.2. In Section 3 we show and discuss Proposition 3.3. Sections 4 and 5 provide details about the proposed MCMC scheme and additional simulated results. Finally, Sections 6 and 7 contain sensitivity analyses to hyperparameter choices and further results for the two real data applications in the main text.

1 Proof of Proposition 3.1 and related discussion

Proposition 3.1. *For $\mathcal{K}_j(\mathcal{D}) := a_j^{\mathcal{D}}/U_{jj|\prec j \succ}$ we have*

$$\gamma_j^{(k)}(\mathbf{X} | \mathcal{A}_j^k)/n^{(k)} \xrightarrow{d} \mathcal{N}\left(\frac{1}{2}\left(\mathcal{K}_j(\mathcal{D})\phi_{0j}^{(k)} - 1 + \text{tr}\Sigma_{0\prec j \succ}/g - \ln\left(\mathcal{K}_j(\mathcal{D})\phi_{0j}^{(k)}\right) + C_1\right),\right. \\ \left.\left(\phi_{0j}^{(k)}\mathcal{K}_j(\mathcal{D}) - 1\right)^2/(2n^{(k)}) + \text{tr}\Sigma_{0\prec j \succ}^2/(2n^{(k)})\right),$$

where $C_1 = \ln(a_j^{\mathcal{D}}/2) - \psi\left(a_j^{\mathcal{D}}/2\right)$. Furthermore $\gamma_j^{(k)}(\mathbf{X} | \mathcal{A}_j^k) \xrightarrow{a.s.} +\infty$.

Proof. First we can write

$$\begin{aligned}
\tilde{\gamma}_j^{(k)}(\mathbf{X}, \boldsymbol{\theta}) &= \ln \frac{\pi_k(j)}{1 - \pi_k(j)} + \sum_{i=1}^{n^{(k)}} \ln \frac{\varphi(x_{ij}^{(k)} | 0, \phi_j^{(k)})}{\varphi(x_{ij}^{(k)} | -\mathbf{L}_{\prec j}^\top \mathbf{x}_{i, \text{pa}_{\mathcal{D}}(j)}, \sigma_j^2)} \\
&= \ln \frac{\pi_k(j)}{1 - \pi_k(j)} - \frac{1}{2} (\phi_j^{(k)})^{-1} \sum_{i=1}^{n^{(k)}} (x_{ij}^{(k)})^2 + \frac{n^{(k)}}{2} \ln \frac{\sigma_j^2}{\phi_j^{(k)}} \\
&\quad + \frac{1}{2} \sigma_j^{-2} \sum_{i=1}^{n^{(k)}} (x_{ij}^{(k)} + \mathbf{L}_{\prec j}^\top \mathbf{x}_{i, \text{pa}_{\mathcal{D}}(j)})^2.
\end{aligned}$$

To derive the quantity of interest $\gamma_j^{(k)}(\mathbf{X} | \mathcal{A}_j^k)$ note that $(\mathbf{D}, \mathbf{L}) | \mathbf{X}^{(k)}, \mathcal{A}_j^k \stackrel{d}{=} (\mathbf{D}, \mathbf{L})$ as in (6), and

$$\phi_j^{(k)} | \mathbf{X}^{(k)}, \mathcal{A}_j^k \sim \text{I-Ga} \left(a_j^{(k)} + n^{(k)}/2, b_j^{(k)} + \sum_i (x_{ij}^{(k)})^2/2 \right).$$

Therefore we can derive $\mathbb{E}_{\boldsymbol{\theta} | \mathbf{X}, \mathcal{A}_j^k} [1/\sigma_j^2] = \mathcal{K}_j(\mathcal{D})$, $\mathbb{E}_{\boldsymbol{\theta} | \mathbf{X}, \mathcal{A}_j^k} [\ln \sigma_j^2] = \ln(\mathbf{U}_{jj | \prec j} / 2) - \psi(a_j^{\mathcal{D}}/2)$, where ψ is the digamma function, and, recalling $\psi(x) \asymp \ln x$ for large x ,

$$\begin{aligned}
\mathbb{E}_{\boldsymbol{\theta} | \mathbf{X}, \mathcal{A}_j^k} [1/\phi_j^{(k)}] &= [2a_j^{(k)} + n^{(k)}] / \left[2b_j^{(k)} + \sum_i (x_{ij}^{(k)})^2 \right] \xrightarrow{\text{a.s.}} 1/\phi_{0j}^{(k)}, \\
\mathbb{E}_{\boldsymbol{\theta} | \mathbf{X}, \mathcal{A}_j^k} [\ln \phi_j^{(k)}] &= \ln \left(b_j^{(k)} + \sum_i (x_{ij}^{(k)})^2/2 \right) - \psi(a_j^{(k)} + n^{(k)}/2) \xrightarrow{\text{a.s.}} \ln \phi_{0j}^{(k)}, \\
\mathbb{E}_{\boldsymbol{\theta} | \mathbf{X}, \mathcal{A}_j^k} [\mathbf{L}_{\prec j} / \sigma_j^2] &= -\mathcal{K}_j(\mathcal{D}) \mathbf{U}_{\prec j}^{-1} \mathbf{U}_{\prec j} = \mathbf{0}, \\
\mathbb{E}_{\boldsymbol{\theta} | \mathbf{X}, \mathcal{A}_j^k} [\mathbf{L}_{\prec j} \mathbf{L}_{\prec j}' / \sigma_j^2] &= \mathbf{U}_{\prec j}^{-1} + \mathcal{K}_j(\mathcal{D}) \mathbf{U}_{\prec j}^{-1} \mathbf{U}_{\prec j} (\mathbf{U}_{\prec j}^{-1} \mathbf{U}_{\prec j})^\top = \mathbf{I}_{|\text{pa}_{\mathcal{D}}(j)|} / g,
\end{aligned}$$

from which, with $\bar{X}_{2j}^{(k)} := \sum_i (x_{ij}^{(k)})^2 / n^{(k)}$ and $\bar{\mathbf{X}}_{2\prec j}^{(k)} := \sum_i \mathbf{x}_{i, \text{pa}_{\mathcal{D}}(j)}^{(k)} (\mathbf{x}_{i, \text{pa}_{\mathcal{D}}(j)}^{(k)})^\top / n^{(k)}$,

$$\begin{aligned}
\gamma_j^{(k)}(\mathbf{X} | \mathcal{A}_j^k) / n^{(k)} &\asymp \frac{1}{2} [\mathcal{K}_j(\mathcal{D}) - 1/\phi_{0j}^{(k)}] \bar{X}_{2j}^{(k)} + \frac{1}{2g} \text{tr} \bar{\mathbf{X}}_{2\prec j}^{(k)} - \frac{1}{2} \ln 2 \mathbf{U}_{jj | \prec j}^{-1} \phi_{0j}^{(k)} - \frac{1}{2} \psi(a_j^{\mathcal{D}}/2) \\
&\xrightarrow{\text{a.s.}} \frac{1}{2} \left(\mathcal{K}_j(\mathcal{D}) \phi_{0j}^{(k)} - 1 + \text{tr} \boldsymbol{\Sigma}_{0\prec j} / g - \ln \mathcal{K}_j(\mathcal{D}) \phi_{0j}^{(k)} + C_1 \right). \tag{1}
\end{aligned}$$

Since $\psi(x) < \ln x$ for all $x > 1/2$, it is true that $C_1 > 0$, since $a_j^{\mathcal{D}} \geq 1$ from the discussion following Equation (6). Also, $\mathcal{K}_j(\mathcal{D}) \phi_{0j}^{(k)} - 1 - \ln \mathcal{K}_j(\mathcal{D}) \phi_{0j}^{(k)} \geq 0$ for all possible values of the terms involved, with equality holding for $\mathcal{K}_j(\mathcal{D})^{-1} = \phi_{0j}^{(k)}$. Then the quantity in (1) is strictly positive, showing $\gamma_j^{(k)}(\mathbf{X} | \mathcal{A}_j^k) \xrightarrow{\text{a.s.}} +\infty$.

To prove asymptotic normality, denote vec the matrix vectorization (and vec^{-1} its inverse function), and note that from the above discussion $\gamma_j^{(k)}(\mathbf{X} | \mathcal{A}_j^k) / n^{(k)} \asymp h(\bar{X}_{2j}^{(k)}, \text{vec} \bar{\mathbf{X}}_{2\prec j}^{(k)})$, where $h(x, y) = (\mathcal{K}_j(\mathcal{D})x - 1 + \text{tr} \text{vec}^{-1} y - \ln \mathcal{K}_j(\mathcal{D})x + C_1) / 2$, with derivatives computed as $h'_x(x, y) = \partial h(x, y) / \partial x = (\mathcal{K}_j(\mathcal{D}) - 1/x) / 2$ and $h'_y(x, y) = \partial h(x, y) / \partial y = \text{vec} \mathbf{I}_{|\text{pa}_{\mathcal{D}}(j)|} / 2$. Since $\bar{X}_{2j}^{(k)} \xrightarrow{d} \mathcal{N}(\phi_{0j}^{(k)}, 2(\phi_{0j}^{(k)})^2 / n^{(k)})$ and $\text{vec} \bar{\mathbf{X}}_{2\prec j}^{(k)} \xrightarrow{d} \mathcal{N}(\boldsymbol{\Sigma}_{0\prec j}, 2(\boldsymbol{\Sigma}_{0\prec j} \otimes \boldsymbol{\Sigma}_{0\prec j}) / n^{(k)})$ (see for

instance Ghazal and Neudecker 2000), with null covariance between $\bar{X}_{2j}^{(k)}$ and $\text{vec } \bar{\mathbf{X}}_{2\prec j \succ}^{(k)}$, from the Delta method we have

$$\begin{aligned} \gamma_j^{(k)}(\mathbf{X} | \mathcal{A}_j^k)/n^{(k)} &\xrightarrow{d} \mathcal{N} \left(h \left(\phi_{0j}^{(k)}, \text{vec } \Sigma_{0\prec j \succ} \right), 2(\phi_{0j}^{(k)})^2/n^{(k)} \cdot h'_x \left(\phi_{0j}^{(k)}, \text{vec } \Sigma_{0\prec j \succ} \right)^2 + \right. \\ &\quad \left. 2/n^{(k)} h'_y \left(\phi_{0j}^{(k)}, \text{vec } \Sigma_{0\prec j \succ} \right)' (\Sigma_{0\prec j \succ} \otimes \Sigma_{0\prec j \succ}) h'_y \left(\phi_{0j}^{(k)}, \text{vec } \Sigma_{0\prec j \succ} \right) \right) \end{aligned}$$

corresponding to the distribution in the statement. \square

The proposition states that, if $j \in I_k$, i.e. node j is a target under intervention k , this will be detected with sample size large enough and for any given graph \mathcal{D} , therefore with a false negative rate eventually zero. The scaled log-odds of the (correct) target classification has asymptotic Gaussian distribution. From the normality and by replacing $\phi_{0j}^{(k)}$ and $\Sigma_{0\prec j \succ}$ by some estimators, we can further provide a reasonable range within which the posterior log-odds is expected to be for a given fixed sample size. Also note that $\phi_{0j}^{(k)} \rightarrow 0$ corresponds to the limit case of deterministic (non stochastic) interventions. In our result this scenario clearly coincides with a degenerate normal distribution with infinite mean, since there is no uncertainty in the log-odds of an intervention.

Note that the mean of the asymptotic distribution increases when $a_j^{\mathcal{D}}$ is large. This is typical of a node with many parents in a large graph, that, coherently to intuition, makes easier the identification of an intervention, since the latter will suppress many dependence relations. With $C_1 \approx 0$ ($a_j^{\mathcal{D}}$ large), we would have that for $\mathcal{K}_j(\mathcal{D})^{-1} = \phi_{0j}^{(k)}$, $\gamma_j^{(k)}(\mathbf{X} | \mathcal{A}_j^k) \xrightarrow{a.s.} \text{tr} \Sigma_{0\prec j \succ} / (2g)$: with many parents and equal variances with or without intervention, target discrimination is still feasible thanks to the variability of the node parents. Then the extreme case when it is not possible to identify an intervention is when $a_j^{\mathcal{D}}$ large, pre- and post-intervention variances coincide, and either (a) the node has no parents or (b) the parents have a degenerate distribution with null variances.

2 Proof of Proposition 3.2 and related discussion

Proposition 3.2. For $\delta_{jk} := \frac{a_j^{(k)}}{b_j^{(k)}} \sigma_{0j}^2 - 1$ we have

$$\begin{aligned} \gamma_j^{(k)}(\mathbf{X} | \bar{\mathcal{A}}_j^k)/n^{(k)} &\xrightarrow{d} \mathcal{N} \left(\frac{1}{2} \left(\frac{a_j^{(k)}}{b_j^{(k)}} \Sigma_{0jj} - 1 - \ln \frac{a_j^{(k)}}{b_j^{(k)}} \sigma_{0j}^2 + C_2 \right), \right. \\ &\quad \left. \left(\frac{\Sigma_{0jj}}{\sigma_{0j}^2} \delta_{jk} \right)^2 / (2n^{(k)}) + \frac{\Sigma_{0jj}^2 - \sigma_{0j}^4}{\sigma_{0j}^4} \left[\delta_{jk}/n^{(k)} + \left(\frac{5\Sigma_{0jj} - 3\sigma_{0j}^2}{\Sigma_{0jj} + \sigma_{0j}^2} \right) / (2n^{(k)}) \right] \right) \end{aligned}$$

where $C_2 = \ln(a_j^{(k)}) - \psi(a_j^{(k)})$. Furthermore $\gamma_j^{(k)}(\mathbf{X} | \bar{\mathcal{A}}_j^k) \xrightarrow{a.s.} +\infty$.

Proof. From Ben-David et al. (2015) we know that

$$\begin{aligned}\sigma_j^2 | \mathbf{X}^{(k)} &\sim \text{I-Ga} \left(\frac{1}{2} \left(a_j^{\mathcal{D}} + n^{(k)} \right), \frac{1}{2} \mathbf{U}_{jj|\prec j \succ}^{(k)} \right), \\ \mathbf{L}_{\prec j} | \sigma_j^2, \mathbf{X}^{(k)} &\sim \mathcal{N}_{|\text{pa}_{\mathcal{D}}(j)|} \left(- \left(\mathbf{U}_{\prec j \succ}^{(k)} \right)^{-1} \mathbf{U}_{\prec j}^{(k)}, \sigma_j^2 \left(\mathbf{U}_{\prec j \succ}^{(k)} \right)^{-1} \right),\end{aligned}$$

where $\mathbf{U}^{(k)} := \mathbf{U} + (\mathbf{X}^{(k)})' \mathbf{X}^{(k)}$, and we have that $\Phi^{(k)} | \mathbf{X}^{(k)}, \bar{\mathcal{A}}_j^k \stackrel{d}{=} \Phi^{(k)}$ as in (7) in the main text. Therefore we can first compute $\mathbb{E}_{\theta | \mathbf{X}, \bar{\mathcal{A}}_j^k} [1/\phi_j^{(k)}] = a_j^{(k)}/b_j^{(k)}$ and $\mathbb{E}_{\theta | \mathbf{X}, \bar{\mathcal{A}}_j^k} [\ln \phi_j^{(k)}] = \ln b_j^{(k)} - \psi(a_j^{(k)})$. Furthermore, with $\mathcal{K}_j^{(k)}(\mathcal{D}) := (a_j^{\mathcal{D}} + n^{(k)})/\mathbf{U}_{jj|\prec j \succ}^{(k)}$, we can derive $\mathbb{E}_{\theta | \mathbf{X}, \bar{\mathcal{A}}_j^k} [1/\sigma_j^2] = \mathcal{K}_j^{(k)}(\mathcal{D}) \xrightarrow{a.s.} \sigma_{0j}^{-1}$, and

$$\begin{aligned}\mathbb{E}_{\theta | \mathbf{X}, \bar{\mathcal{A}}_j^k} [\ln \sigma_j^2] &= \ln \left(\mathbf{U}_{jj|\prec j \succ}^{(k)} / 2 \right) - \psi \left(\frac{1}{2} \left(a_j^{\mathcal{D}} + n^{(k)} \right) \right) \xrightarrow{a.s.} \ln \sigma_{0j}^2, \\ \mathbb{E}_{\theta | \mathbf{X}, \bar{\mathcal{A}}_j^k} [\mathbf{L}_{\prec j} / \sigma_j^2] &= -\mathcal{K}_j^{(k)}(\mathcal{D}) \left(\mathbf{U}_{\prec j \succ}^{(k)} \right)^{-1} \mathbf{U}_{\prec j}^{(k)} \xrightarrow{a.s.} -\sigma_{0j}^{-2} \Sigma_{0\prec j \succ}^{-1} \Sigma_{0\prec j}, \\ \mathbb{E}_{\theta | \mathbf{X}, \bar{\mathcal{A}}_j^k} [\mathbf{L}_{\prec j} \mathbf{L}'_{\prec j} / \sigma_j^2] &= \left(\mathbf{U}_{\prec j \succ}^{(k)} \right)^{-1} + \mathcal{K}_j^{(k)}(\mathcal{D}) \left(\mathbf{U}_{\prec j \succ}^{(k)} \right)^{-1} \mathbf{U}_{\prec j}^{(k)} \mathbf{U}_{\prec j}^{(k)'} \left(\mathbf{U}_{\prec j \succ}^{(k)} \right)^{-1} \\ &\xrightarrow{a.s.} \sigma_{0j}^{-2} \Sigma_{0\prec j \succ}^{-1} \Sigma_{0\prec j} \Sigma'_{0\prec j} \Sigma_{0\prec j}^{-1}\end{aligned}$$

Then, with $\bar{X}_{2j}^{(k)} := \sum_i (x_{ij}^{(k)})^2 / n^{(k)} \xrightarrow{a.s.} \Sigma_{0jj}$, $\text{vec } \bar{\mathbf{X}}_{2\prec j \succ}^{(k)} := \text{vec } \sum_i \mathbf{x}_{i, \text{pa}_{\mathcal{D}}(j)}^{(k)} \left(\mathbf{x}_{i, \text{pa}_{\mathcal{D}}(j)}^{(k)} \right)^\top / n^{(k)} \xrightarrow{a.s.} \text{vec } \Sigma_{0\prec j \succ}$ and $\bar{\mathbf{X}}_{2[j \succ]}^{(k)} := \sum_{i=1}^{n^{(k)}} x_{ij}^{(k)} \left(\mathbf{x}_{i, \text{pa}_{\mathcal{D}}(j)}^{(k)} \right)^\top / n^{(k)} \xrightarrow{a.s.} \Sigma_{0[j \succ]}$,

$$\begin{aligned}\bar{\gamma}_j^{(k)}(\mathbf{X} | \bar{\mathcal{A}}_j^k) / n^{(k)} &\asymp \frac{1}{2} \bar{X}_{2j}^{(k)} \mathbb{E}_{\theta | \mathbf{X}, \bar{\mathcal{A}}_j^k} [1/\phi_j^{(k)}] - \frac{1}{2} \left(\mathbb{E}_{\theta | \mathbf{X}, \bar{\mathcal{A}}_j^k} [\ln \sigma_j^2] - \mathbb{E}_{\theta | \mathbf{X}, \bar{\mathcal{A}}_j^k} [\ln \phi_j^{(k)}] \right) \\ &\quad - \frac{1}{2} \bar{X}_{2j}^{(k)} \mathbb{E}_{\theta | \mathbf{X}, \bar{\mathcal{A}}_j^k} [1/\sigma_j^2] - \frac{1}{2} \text{tr} \left\{ \bar{\mathbf{X}}_{2\prec j \succ}^{(k)} \mathbb{E}_{\theta | \mathbf{X}, \bar{\mathcal{A}}_j^k} [\mathbf{L}_{\prec j} \mathbf{L}'_{\prec j} / \sigma_j^2] \right\} \\ &\quad - \bar{\mathbf{X}}_{2[j \succ]}^{(k)} \mathbb{E}_{\theta | \mathbf{X}, \bar{\mathcal{A}}_j^k} [\mathbf{L}_{\prec j} / \sigma_j^2] \\ &\xrightarrow{a.s.} \frac{1}{2} \frac{a_j^{(k)}}{b_j^{(k)}} \Sigma_{0jj} - \frac{1}{2} \left(\ln \sigma_{0j}^2 - \ln b_j^{(k)} + \psi(a_j^{(k)}) \right) - \frac{1}{2} \frac{\Sigma_{0jj}}{\sigma_{0j}^2} \\ &\quad - \frac{1}{2} \sigma_{0j}^{-2} \text{tr} \left\{ \Sigma_{0\prec j} \Sigma'_{0\prec j} \Sigma_{0\prec j}^{-1} \right\} + \sigma_{0j}^{-2} \Sigma'_{0\prec j} \Sigma_{0\prec j}^{-1} \Sigma_{0\prec j} \\ &= \frac{1}{2} \left(\frac{a_j^{(k)}}{b_j^{(k)}} \Sigma_{0jj} - 1 - \ln \frac{a_j^{(k)}}{b_j^{(k)}} \sigma_{0j}^2 + C_2 \right),\end{aligned}\tag{2}$$

where we have used $\Sigma'_{0\prec j} \Sigma_{0\prec j}^{-1} \Sigma_{0\prec j} = \Sigma_{0jj} - \sigma_{0j}^2$. Since $\psi(x) < \ln x$ for all $x > 1/2$, it is true that $C_2 > 0$, since $a_j^{(k)} \geq 1$ from the discussion following Equation (7) in the main text. Also, $a_j^{(k)}/b_j^{(k)} \Sigma_{0jj} - 1 - \ln a_j^{(k)}/b_j^{(k)} \sigma_{0j}^2 \geq 0$ for all possible values of the terms involved, with equality holding for $a_j^{(k)}/b_j^{(k)} = \Sigma_{0jj} = \sigma_{0j}^2$. Then (2) is strictly positive, showing $\bar{\gamma}_j^{(k)}(\mathbf{X} | \bar{\mathcal{A}}_j^k) \xrightarrow{a.s.} +\infty$.

To show the asymptotic distribution, first note that we can write the quantity of interest as $\bar{\gamma}_j^{(k)}(\mathbf{X} | \bar{\mathcal{A}}_j^k) / n^{(k)} = \bar{h} \left(\bar{X}_{2j}^{(k)}, \text{vec } \bar{\mathbf{X}}_{2\prec j \succ}^{(k)}, \bar{\mathbf{X}}_{2\prec j}^{(k)} \right)$, where vec is the matrix vectorization (and

vec^{-1} its inverse function) and $\bar{\mathbf{X}}_{2 \prec j}^{(k)} = (\bar{\mathbf{X}}_{2[j \succ]}^{(k)})^\top$, and where the function \bar{h} is defined as

$$\begin{aligned}\bar{h}(x, y, z) &= \frac{1}{2} \left(\frac{a_j^{(k)}}{b_j^{(k)}} - u(x, y, z)^{-1} \right) x - \frac{1}{2} \text{tr} \left(\frac{zz' \omega(y)^{-1}}{u(x, y, z)} \right) \\ &\quad + \frac{1}{u(x, y, z)} z' \omega(y)^{-1} z - \frac{1}{2} \left(\ln u(x, y, z) - \ln b_j^{(k)} + \psi \left(\frac{a_j^{(k)}}{b_j^{(k)}} \right) \right) \\ &= \frac{1}{2} \left(\frac{a_j^{(k)}}{b_j^{(k)}} x - 1 - \ln \frac{a_j^{(k)}}{b_j^{(k)}} u(x, y, z) + C_2 \right),\end{aligned}\quad (3)$$

with $u(x, y, z) = x - z^\top \omega(y)^{-1} z$ and $\omega(y) = \text{vec}^{-1}(y)$. The partial derivatives of \bar{h} can be recovered as $\bar{h}'_x(x, y, z) = \partial \bar{h}(x, y, z) / \partial x = (a_j^{(k)} / b_j^{(k)} - 1 / u(x, y, z)) / 2$, $\bar{h}'_y(x, y, z) = \partial \bar{h}(x, y, z) / \partial y = -\text{vec}(\omega(y)^{-1} z z^\top \omega(y)^{-1} / u(x, y, z)) / 2$ and $\bar{h}'_z(x, y, z) = \partial \bar{h}(x, y, z) / \partial z = \omega(y)^{-1} z / u(x, y, z)$. It is then easy to see, using $\mathbf{L}_{0 \prec j} = -\Sigma_{0 \prec j \succ}^{-1} \Sigma_{0 \prec j}$, that

$$\nabla \bar{g} \left(\bar{\mathbf{X}}_{2j}^{(k)}, \text{vec } \bar{\mathbf{X}}_{2 \prec j \succ}^{(k)}, \bar{\mathbf{X}}_{2 \prec j}^{(k)} \right)^\top \xrightarrow{a.s.} \frac{1}{2\sigma_{0j}^2} \left(\frac{a_j^{(k)}}{b_j^{(k)}} \sigma_{0j}^2 - 1, -\text{vec} \{ \mathbf{L}_{0 \prec j} \mathbf{L}_{0[j \succ]} \}^\top, -2\mathbf{L}_{0[j \succ]} \right)^\top \quad (4)$$

Furthermore, from the properties of higher-order moments of the multivariate Gaussian distribution, see for instance Ghazal and Neudecker (2000), we know that

$$\begin{aligned}\bar{\mathbf{X}}_{2j}^{(k)} &\xrightarrow{d} \mathcal{N}_1 \left(\Sigma_{0jj}, \frac{2}{n^{(k)}} \Sigma_{0jj}^2 \right), \\ \text{vec } \bar{\mathbf{X}}_{2 \prec j \succ}^{(k)} &\xrightarrow{d} \mathcal{N}_{|\text{pa}_{\mathcal{D}}(j)|^2} \left(\text{vec } \Sigma_{0 \prec j \succ}, \frac{2}{n^{(k)}} \Sigma_{0 \prec j \succ} \otimes \Sigma_{0 \prec j \succ} \right), \\ \bar{\mathbf{X}}_{2 \prec j}^{(k)} &\xrightarrow{d} \mathcal{N}_{|\text{pa}_{\mathcal{D}}(j)|} \left(\Sigma_{0 \prec j}, \frac{1}{n^{(k)}} [\Sigma_{0jj} \Sigma_{0 \prec j \succ} + \Sigma_{0 \prec j} \Sigma_{0[j \succ]}] \right),\end{aligned}\quad (5)$$

where \otimes is the Kronecker product. The cross-covariance matrices can be derived as follows:

$$\begin{aligned}\mathbb{C} \left(\bar{\mathbf{X}}_{2j}^{(k)}, \text{vec } \bar{\mathbf{X}}_{2 \prec j \succ}^{(k)} \right) &= \frac{1}{n^{(k)}} \mathbb{C} \left(\left(x_{ij}^{(k)} \right)^2, \text{vec} \left\{ \mathbf{x}_{i, \text{pa}_{\mathcal{D}}(j)}^{(k)} (\mathbf{x}_{i, \text{pa}_{\mathcal{D}}(j)}^{(k)})^\top \right\} \right) \\ &= \frac{1}{n^{(k)}} \mathbb{C} \left(\mathbb{E} \left[\left(x_{ij}^{(k)} \right)^2 \middle| \mathbf{x}_{i, \text{pa}_{\mathcal{D}}(j)}^{(k)} \right], \text{vec} \left\{ \mathbf{x}_{i, \text{pa}_{\mathcal{D}}(j)}^{(k)} (\mathbf{x}_{i, \text{pa}_{\mathcal{D}}(j)}^{(k)})^\top \right\} \right) \\ &= \frac{1}{n^{(k)}} \mathbb{C} \left(\mathbf{L}_{0[j \succ]} \mathbf{x}_{i, \text{pa}_{\mathcal{D}}(j)}^{(k)} (\mathbf{x}_{i, \text{pa}_{\mathcal{D}}(j)}^{(k)})^\top \mathbf{L}_{0 \prec j}, \text{vec} \left\{ \mathbf{x}_{i, \text{pa}_{\mathcal{D}}(j)}^{(k)} (\mathbf{x}_{i, \text{pa}_{\mathcal{D}}(j)}^{(k)})^\top \right\} \right) \\ &= \frac{1}{n^{(k)}} \mathbb{C} \left((\mathbf{L}_{0[j \succ]} \otimes \mathbf{L}_{0[j \succ]}) \text{vec} \left\{ \mathbf{x}_{i, \text{pa}_{\mathcal{D}}(j)}^{(k)} (\mathbf{x}_{i, \text{pa}_{\mathcal{D}}(j)}^{(k)})^\top \right\}, \text{vec} \left\{ \mathbf{x}_{i, \text{pa}_{\mathcal{D}}(j)}^{(k)} (\mathbf{x}_{i, \text{pa}_{\mathcal{D}}(j)}^{(k)})^\top \right\} \right) \\ &= \frac{1}{n^{(k)}} (\mathbf{L}_{0[j \succ]} \otimes \mathbf{L}_{0[j \succ]}) \mathbb{V} \left(\text{vec} \left\{ \mathbf{x}_{i, \text{pa}_{\mathcal{D}}(j)}^{(k)} (\mathbf{x}_{i, \text{pa}_{\mathcal{D}}(j)}^{(k)})^\top \right\} \right) \\ &= \frac{2}{n^{(k)}} (\mathbf{L}_{0[j \succ]} \Sigma_{0 \prec j \succ}) \otimes (\mathbf{L}_{0[j \succ]} \Sigma_{0 \prec j \succ}) = \frac{2}{n^{(k)}} \Sigma_{0[j \succ]} \otimes \Sigma_{0[j \succ]},\end{aligned}\quad (6)$$

where we have used $\text{vec}(\mathbf{ABC}) = (\mathbf{C}^\top \otimes \mathbf{A}) \text{vec } \mathbf{B}$ and $(\mathbf{A} \otimes \mathbf{B})(\mathbf{C} \otimes \mathbf{D}) = (\mathbf{AC}) \otimes (\mathbf{BD})$ for

conformable matrices \mathbf{A} , \mathbf{B} , \mathbf{C} and \mathbf{D} . Similarly, we see that

$$\begin{aligned}
\mathbb{C} \left(\bar{\mathbf{X}}_{2 \prec j}^{(k)}, \text{vec } \bar{\mathbf{X}}_{2 \prec j}^{(k)} \right) &= \frac{1}{n^{(k)}} \mathbb{C} \left(x_{ij}^{(k)} \mathbf{x}_{i, \text{pa}_{\mathcal{D}}(j)}^{(k)}, \text{vec} \left\{ \mathbf{x}_{i, \text{pa}_{\mathcal{D}}(j)}^{(k)} (\mathbf{x}_{i, \text{pa}_{\mathcal{D}}(j)}^{(k)})^\top \right\} \right) \\
&= -\frac{1}{n^{(k)}} \mathbb{C} \left(\mathbf{x}_{i, \text{pa}_{\mathcal{D}}(j)}^{(k)} (\mathbf{x}_{i, \text{pa}_{\mathcal{D}}(j)}^{(k)})^\top \mathbf{L}_{0 \prec j}, \text{vec} \left\{ \mathbf{x}_{i, \text{pa}_{\mathcal{D}}(j)}^{(k)} (\mathbf{x}_{i, \text{pa}_{\mathcal{D}}(j)}^{(k)})^\top \right\} \right) \\
&= -\frac{1}{n^{(k)}} \mathbb{C} \left((\mathbf{L}_{0[j]} \otimes \mathbf{I}_{|\text{pa}_{\mathcal{D}}(j)|}) \text{vec} \left\{ \mathbf{x}_{i, \text{pa}_{\mathcal{D}}(j)}^{(k)} (\mathbf{x}_{i, \text{pa}_{\mathcal{D}}(j)}^{(k)})^\top \right\}, \right. \\
&\quad \left. \text{vec} \left\{ \mathbf{x}_{i, \text{pa}_{\mathcal{D}}(j)}^{(k)} (\mathbf{x}_{i, \text{pa}_{\mathcal{D}}(j)}^{(k)})^\top \right\} \right) \\
&= -\frac{2}{n^{(k)}} (\mathbf{L}_{0[j]} \Sigma_{0 \prec j}) \otimes \Sigma_{0 \prec j} = \frac{2}{n^{(k)}} \Sigma_{0[j]} \otimes \Sigma_{0 \prec j}, \tag{7}
\end{aligned}$$

where we have used the relation $\text{vec}(\mathbf{AB}) = (\mathbf{B}^\top \otimes \mathbf{I}) \text{vec } \mathbf{A}$, and

$$\begin{aligned}
\mathbb{C} \left(\bar{X}_{2j}^{(k)}, \bar{\mathbf{X}}_{2 \prec j}^{(k)} \right) &= \frac{1}{n^{(k)}} \mathbb{C} \left(\left(x_{ij}^{(k)} \right)^2, x_{ij}^{(k)} \mathbf{x}_{i, \text{pa}_{\mathcal{D}}(j)}^{(k)} \right) \\
&= \frac{1}{n^{(k)}} \mathbb{E} \left[\mathbb{C} \left(\left(x_{ij}^{(k)} \right)^2, x_{ij}^{(k)} \mathbf{x}_{i, \text{pa}_{\mathcal{D}}(j)}^{(k)} \middle| \mathbf{x}_{i, \text{pa}_{\mathcal{D}}(j)}^{(k)} \right) \right] \\
&\quad + \frac{1}{n^{(k)}} \mathbb{C} \left(\mathbb{E} \left[\left(x_{ij}^{(k)} \right)^2 \middle| \mathbf{x}_{i, \text{pa}_{\mathcal{D}}(j)}^{(k)} \right], \mathbf{x}_{i, \text{pa}_{\mathcal{D}}(j)}^{(k)} \mathbb{E} \left[x_{ij}^{(k)} \middle| \mathbf{x}_{i, \text{pa}_{\mathcal{D}}(j)}^{(k)} \right] \right) \\
&= \frac{1}{n^{(k)}} \left\{ -2\sigma_{0j}^2 \mathbf{L}_{0[j]} \mathbb{E} \left[\mathbf{x}_{i, \text{pa}_{\mathcal{D}}(j)}^{(k)} (\mathbf{x}_{i, \text{pa}_{\mathcal{D}}(j)}^{(k)})^\top \right] + \right. \\
&\quad \left. + \mathbb{C} \left(\mathbf{L}_{0[j]} \mathbf{x}_{i, \text{pa}_{\mathcal{D}}(j)}^{(k)} (\mathbf{x}_{i, \text{pa}_{\mathcal{D}}(j)}^{(k)})^\top \mathbf{L}_{0 \prec j}, -\mathbf{x}_{i, \text{pa}_{\mathcal{D}}(j)}^{(k)} (\mathbf{x}_{i, \text{pa}_{\mathcal{D}}(j)}^{(k)})^\top \mathbf{L}_{0 \prec j} \right) \right\} \\
&= \frac{1}{n^{(k)}} \left\{ 2\sigma_{0j}^2 \Sigma_{0[j]} - \mathbb{C} \left((\mathbf{L}_{0[j]} \otimes \mathbf{L}_{0[j]}) \text{vec} \left\{ \mathbf{x}_{i, \text{pa}_{\mathcal{D}}(j)}^{(k)} (\mathbf{x}_{i, \text{pa}_{\mathcal{D}}(j)}^{(k)})^\top \right\}, \right. \right. \\
&\quad \left. \left. (\mathbf{L}_{0[j]} \otimes \mathbf{I}_{|\text{pa}_{\mathcal{D}}(j)|}) \text{vec} \left\{ \mathbf{x}_{i, \text{pa}_{\mathcal{D}}(j)}^{(k)} (\mathbf{x}_{i, \text{pa}_{\mathcal{D}}(j)}^{(k)})^\top \right\} \right) \right\} \\
&= \frac{2}{n^{(k)}} \left\{ \sigma_{0j}^2 \Sigma_{0[j]} - (\mathbf{L}_{0[j]} \otimes \mathbf{L}_{0[j]}) (\Sigma_{0 \prec j} \otimes \Sigma_{0 \prec j}) (\mathbf{L}_{0 \prec j} \otimes \mathbf{I}_{|\text{pa}_{\mathcal{D}}(j)|}) \right\} \\
&= \frac{2}{n^{(k)}} \left\{ \sigma_{0j}^2 \Sigma_{0[j]} - (\mathbf{L}_{0[j]} \Sigma_{0 \prec j} \mathbf{L}_{0 \prec j}) (\mathbf{L}_{0[j]} \Sigma_{0 \prec j} \mathbf{I}_{|\text{pa}_{\mathcal{D}}(j)|}) \right\} \\
&= \frac{2}{n^{(k)}} \Sigma_{0jj} \Sigma_{0[j]}, \tag{8}
\end{aligned}$$

using $(\mathbf{A} \otimes \mathbf{B})^\top = \mathbf{A}^\top \otimes \mathbf{B}^\top$ and $\Sigma_{0jj} - \sigma_{0j}^2 = \mathbf{L}_{0[j]} \Sigma_{0 \prec j} \mathbf{L}_{0 \prec j}$. Putting together the results in (3)-(8), by Delta method that we have that

$$\begin{aligned}
\bar{h} \left(\bar{X}_{2j}^{(k)}, \text{vec } \bar{\mathbf{X}}_{2 \prec j}^{(k)}, \bar{\mathbf{X}}_{2 \prec j}^{(k)} \right) &\xrightarrow{d} \mathcal{N} \left(\bar{h} \left(\Sigma_{0jj}, \text{vec } \Sigma_{0 \prec j}, \Sigma_{0 \prec j} \right), \right. \\
&\quad \left. \frac{2}{n^{(k)}} \nabla \bar{h} \left(\Sigma_{0jj}, \text{vec } \Sigma_{0 \prec j}, \Sigma_{0 \prec j} \right)^\top \Lambda \nabla \bar{h} \left(\Sigma_{0jj}, \text{vec } \Sigma_{0 \prec j}, \Sigma_{0 \prec j} \right) \right),
\end{aligned}$$

where

$$\Lambda = \begin{pmatrix} \Sigma_{0jj}^2 & \Sigma_{0[j]} \otimes \Sigma_{0[j]} & \Sigma_{0[j]} \otimes \Sigma_{0jj} \\ \Sigma_{0 \prec j} \otimes \Sigma_{0 \prec j} & \Sigma_{0 \prec j} \otimes \Sigma_{0 \prec j} & \Sigma_{0 \prec j} \otimes \Sigma_{0 \prec j} \\ \Sigma_{0 \prec j} \otimes \Sigma_{0jj} & \Sigma_{0[j]} \otimes \Sigma_{0 \prec j} & \frac{1}{2} (\Sigma_{0 \prec j} \Sigma_{0jj} + \Sigma_{0 \prec j} \Sigma_{0[j]}) \end{pmatrix}$$

and this distribution corresponds to the stated result. \square

Proposition 3.2 tells that a true negative case of no intervention, i.e. $j \notin I_k$, will be eventually detected with sample size large enough. Note that the mean of the asymptotic distribution is closer to zero when node j is independent from any other node, that is $|\text{pa}_{\mathcal{D}}(j)| = \emptyset$, and therefore $\Sigma_{0jj} = \sigma_{0j}^2$. Intuitively, it is more difficult to understand the absence of an intervention since there are no parent-child relations that are removed by the intervention on node j , which makes the intervention effect less apparent. In this special case, the second addend in the asymptotic variance disappears, and the whole asymptotic distribution becomes very similar, *mutatis mutandis*, to the one recovered in Proposition 3.1.

3 Graph Consistency

In the current section we prove model selection consistency of the true DAG observational equivalence class; the latter, combined with consistent estimation of targets, allows to identify the group-specific intervention graphs. In the following lemma and proposition, we first extend the conjugacy result on the DAG-Wishart prior of Ben-David et al. (2015) to interventional Gaussian multivariate data from multiple groups; then we prove, following Cao et al. (2019) and Peluso and Consonni (2020), its Bayes factor and posterior ratio consistency outside $[\mathcal{D}_0]$, the equivalence class of the true DAG, and its asymptotic compatibility within $[\mathcal{D}_0]$. Let \mathbf{P}^k be a $q \times q$ diagonal matrix such that $\mathbf{P}_{jj}^k = 1$ if $j \notin I_k$, that is node j is not intervened in group k .

Lemma 3.1. *Let \mathcal{D}_0 be the true DAG. Assume $(\mathbf{D}, \mathbf{L}) | \mathcal{D}$ follows a DAG-Wishart distribution with hyperparameters \mathbf{U} and $\mathbf{a}^{\mathcal{D}}$ and consider the likelihood function $f(\mathbf{X} | \boldsymbol{\theta}, I_1, \dots, I_K, \mathcal{D})$ defined as in (5) of the main text. A posteriori we have that $(\mathbf{D}, \mathbf{L}) | \mathcal{D}, \mathbf{X}, I_1, \dots, I_K$ is also DAG-Wishart with hyperparameters $n\tilde{\mathbf{S}}$ and $\mathbf{a} + \mathbf{n}^*$, where $\tilde{\mathbf{S}} := \mathbf{U}/n + \sum_k \mathbf{P}^k \mathbf{X}^{(k)\top} \mathbf{X}^{(k)} \mathbf{P}^k / n$ and $\mathbf{n}^* := \sum_k \text{diag}\{\mathbf{P}^k\} n^{(k)}$.*

Proof. The DAG-Wishart (prior) distribution has density

$$\frac{1}{\mathcal{Z}_{\mathcal{D}}(\mathbf{U}, \mathbf{a}^{\mathcal{D}})} \exp \left\{ -\frac{1}{2} \text{tr} \left(\left(\mathbf{L} \mathbf{D}^{-1} \mathbf{L}^{\top} \right) \mathbf{U} \right) \right\} \prod_{j \in V} \mathbf{D}_{jj}^{-\frac{\mathbf{a}_j^{\mathcal{D}}}{2}}$$

for all (\mathbf{D}, \mathbf{L}) in the Cholesky space, with $\mathcal{Z}_{\mathcal{D}}(\mathbf{U}, \mathbf{a}^{\mathcal{D}})$ being the normalizing constant. In the posterior distribution, we can write the likelihood $f(\mathbf{X} | \boldsymbol{\theta}, I_1, \dots, I_K, \mathcal{D})$ as proportional to

$$\begin{aligned} \prod_{k=1}^K p(\mathbf{X}^{(k)} | \mathbf{D}, \mathbf{L}, I_1, \dots, I_K) &\propto \prod_{k=1}^K \exp \left\{ -\frac{1}{2} \text{tr} \left(\left(\mathbf{L} \mathbf{D}^{-1} \mathbf{L}^{\top} \right) \mathbf{P}^k \mathbf{S}^{(k)} \mathbf{P}^k \right) \right\} \prod_{j \in V} \mathbf{D}_{jj}^{-\frac{\mathbf{P}_{jj}^k n^{(k)}}{2}} \\ &= \exp \left\{ -\frac{1}{2} \text{tr} \left(\left(\mathbf{L} \mathbf{D}^{-1} \mathbf{L}^{\top} \right) \sum_k \mathbf{P}^k \mathbf{S}^{(k)} \mathbf{P}^k \right) \right\} \prod_{j \in V} \mathbf{D}_{jj}^{-\sum_k \frac{\mathbf{P}_{jj}^k n^{(k)}}{2}}, \end{aligned}$$

with $\mathbf{S}^{(k)} = (\mathbf{X}^{(k)})^{\top} \mathbf{X}^{(k)}$, where we omitted the interventional terms in the likelihood not dependent on \mathbf{D} and \mathbf{L} . The DAG-Wishart (posterior) distribution of $(\mathbf{D}, \mathbf{L}) | \mathcal{D}, \mathbf{X}, I_1, \dots, I_K$ follows immediately from Bayes theorem by combining prior and likelihood. \square

We have *Bayes factor consistency* if, for all $\mathcal{D} \neq \mathcal{D}_0$, the Bayes factor

$$BF_{\mathcal{D}, \mathcal{D}_0} = \frac{m(\mathbf{X} | \mathcal{D}, I_1, \dots, I_K)}{m(\mathbf{X} | \mathcal{D}_0, I_1, \dots, I_K)} \xrightarrow{\bar{P}} 0,$$

whenever \mathcal{D}_0 is the true DAG generating \mathbf{X} , where $\xrightarrow{\bar{P}}$ denotes convergence in probability, \bar{P} is the probability measure under the true DAG \mathcal{D}_0 , and $m(\mathbf{X} | \mathcal{D}, I_1, \dots, I_K)$ is the marginal (or integrated) likelihood. Define also $\tilde{\mathbf{n}} = \mathbf{n} - \mathbf{n}^*$. On the other hand, we have *posterior ratio consistency* if, with \mathcal{D}_0 being the true DAG, it holds that, as $n \rightarrow \infty$,

$$\max_{\mathcal{D} \neq \mathcal{D}_0} \frac{p(\mathcal{D} | \mathbf{X}, I_1, \dots, I_K)}{p(\mathcal{D}_0 | \mathbf{X}, I_1, \dots, I_K)} = \max_{\mathcal{D} \neq \mathcal{D}_0} BF_{\mathcal{D}, \mathcal{D}_0}(\mathbf{X} | I_1, \dots, I_K) \frac{p(\mathcal{D})}{p(\mathcal{D}_0)} \xrightarrow{\bar{P}} 0. \quad (9)$$

Proposition 3.3. *Let \mathcal{D}_0 be the true DAG. Assume $(\mathbf{D}, \mathbf{L}) | \mathcal{D}$ follows a DAG-Wishart distribution with hyperparameters \mathbf{U} and $\mathbf{a}^{\mathcal{D}}$ as in Equation (8) of the main text, and consider the likelihood function $f(\mathbf{X} | \boldsymbol{\theta}, I_1, \dots, I_K, \mathcal{D})$ defined as in (5) of the main text. If (a) $a_j^{\mathcal{D}} = a + |\text{pa}_{\mathcal{D}}(j)| - q + 1$, (b) $\tilde{n}_j = o(n_j^*)$ for all $j \in V$, and (c) for all $j \neq l \in V$ there exists a k such that $j \notin I_k$ and $l \notin I_k$ hold, then as $n \rightarrow \infty$,*

$$\begin{aligned} i) \quad & \max_{\mathcal{D} \notin [\mathcal{D}_0]} \frac{p(\mathcal{D} | \mathbf{X}, I_1, \dots, I_K)}{p(\mathcal{D}_0 | \mathbf{X}, I_1, \dots, I_K)} \xrightarrow{\bar{P}} 0, \\ ii) \quad & \frac{p(\mathcal{D} | \mathbf{X}, I_1, \dots, I_K)}{p(\mathcal{D}_0 | \mathbf{X}, I_1, \dots, I_K)} \xrightarrow{\bar{P}} \frac{p(\mathcal{D})}{p(\mathcal{D}_0)} \text{ for all } \mathcal{D} \in [\mathcal{D}_0]. \end{aligned}$$

Proof. From the conjugacy of the DAG-Wishart distribution (Lemma 3.1), we can write the marginal likelihood of DAG \mathcal{D} as

$$m(\mathbf{X} | \mathcal{D}, I_1, \dots, I_K) = (2\pi)^{-n/2} \mathcal{Z}_{\mathcal{D}}(n\tilde{\mathbf{S}}, \mathbf{a}^{\mathcal{D}} + \mathbf{n}^*) / \mathcal{Z}_{\mathcal{D}}(\mathbf{U}, \mathbf{a}^{\mathcal{D}}).$$

Accordingly, the Bayes factor between any two DAGs \mathcal{D} and \mathcal{D}_0 is

$$BF_{\mathcal{D}, \mathcal{D}_0}(\mathbf{X} | I_1, \dots, I_K) = \frac{\mathcal{Z}_{\mathcal{D}}(n\tilde{\mathbf{S}}, \mathbf{a}^{\mathcal{D}} + \mathbf{n}^*)}{\mathcal{Z}_{\mathcal{D}}(\mathbf{U}, \mathbf{a}^{\mathcal{D}})} \frac{\mathcal{Z}_{\mathcal{D}_0}(\mathbf{U}, \mathbf{a}^{\mathcal{D}_0})}{\mathcal{Z}_{\mathcal{D}_0}(n\tilde{\mathbf{S}}, \mathbf{a}^{\mathcal{D}} + \mathbf{n}^*)}.$$

Since $\tilde{n}_j = o(n_j^*)$ for all $j \in V$, then $(a + n_j)/(a + n_j^*) = (a + n_j^* + \tilde{n}_j)/(a + n_j^*) \rightarrow 1$. Furthermore,

$$\begin{aligned} \sum_{k=1}^K \mathbf{P}^k \mathbf{X}^{(k)\top} \mathbf{X}^{(k)} \mathbf{P}^k &= \left(\sum_{k=1}^K \sum_{i=1}^{n^{(k)}} \mathbf{x}_{ij}^{(k)} \mathbf{x}_{il}^{(k)} \mathbb{1}_{\{j \notin I_k\}} \mathbb{1}_{\{l \notin I_k\}} \right)_{jl} \\ &\asymp \left(\sum_{k=1}^K \sum_{i=1}^{n^{(k)}} \mathbf{x}_{ij}^{(k)} \mathbf{x}_{il}^{(k)} \right)_{jl} = \sum_{k=1}^K \mathbf{X}^{(k)\top} \mathbf{X}^{(k)} = \mathbf{X}^\top \mathbf{X}, \end{aligned}$$

where the asymptotic equivalence in probability holds from the assumption (c) that for all nodes there exists a setting (group) k where both nodes are not intervened. Then, from Continuity Mapping theorem

$$BF_{\mathcal{D}, \mathcal{D}_0}(\mathbf{X} | I_1, \dots, I_K) \asymp \frac{\mathcal{Z}_{\mathcal{D}}(\mathbf{U} + \mathbf{X}'\mathbf{X}, \mathbf{a}^{\mathcal{D}} + \mathbf{n})}{\mathcal{Z}_{\mathcal{D}}(\mathbf{U}, \mathbf{a}^{\mathcal{D}})} \frac{\mathcal{Z}_{\mathcal{D}_0}(\mathbf{U}, \mathbf{a}^{\mathcal{D}_0})}{\mathcal{Z}_{\mathcal{D}_0}(\mathbf{U} + \mathbf{X}'\mathbf{X}, \mathbf{a}^{\mathcal{D}} + \mathbf{n})},$$

in the same form provided by Cao et al. (2019). Result *i*) comes therefore from Cao et al. (2019, Theorem 4.1), whilst result *ii*) comes from assumption b) together with Peluso and Consonni (2020, Proposition 3.5). \square

Proposition 3.3 shows that posterior ratio and Bayes factor consistency under DAG-Wishart prior holds *outside* the Markov equivalence class of the true generating DAG \mathcal{D}_0 . On the other hand, the posterior ratio tends to the prior ratio (Bayes factor equal to one) within the true equivalence class. This result is coherent with Peluso and Consonni (2020), in the context of a single-group observational data. Assumption (a) guarantees *compatibility* within equivalence classes (Geiger and Heckerman, 2002; Peluso and Consonni, 2020), to have Bayes factors between equivalent graphs always equal to one. The removal of this assumption does not break posterior consistency outside $[\mathcal{D}_0]$, but the posterior ratio within $[\mathcal{D}_0]$ would converge to some value dependent on the true parameters, losing compatibility. With assumption (b) we require that, for each node, the number of non-interventional observations for node j will eventually dominate the number of intervened records for the node. Finally, the last assumption (c) intuitively says that, for each pair of nodes, we can estimate their dependence if there exists at least a group in which these two nodes are not intervened.

We further note that assumption (b) on the asymptotic dominance of purely observational data over interventional ones is particularly convenient since it permits to reconduct the problem of graph consistency to the one in Cao et al. (2019) and in Peluso and Consonni (2020) for the comparison, respectively, of (observationally) non-equivalent and equivalent graphs. In this way we guarantee the identification of the correct observational equivalence class and, together with the correct identification of target nodes, we are able to consistently estimate the interventional DAG, up to relations of interventional equivalences. We conjecture that the removal of this assumption, substituted by $n_j^* = \alpha n$ for some $\alpha \in (0, 1)$, that is observational and interventional sample sizes diverging at the same speed, can lead to the more precise result of graph posterior consistency by breaking, within the observational equivalence class of the true DAG, the compatibility among graphs that are not anymore equivalent after the interventions. We refer the reader to Hauser and Bühlmann (2012) for the formal definition of interventional Markov equivalence. Moreover, by Chickering (1995, Theorem 2), there exists a sequence of observationally Markov equivalent DAGs differing only for the reversal of a *covered edge*, that is for an edge (u, v) of \mathcal{D} for which $\text{pa}_v(\mathcal{D}) = \text{pa}_u(\mathcal{D}) \cup \{u\}$. A comparison between \mathcal{D} and \mathcal{D}_0 , differing only by one covered edge reversal, could show, when an intervention occurs on u or on v and breaks the equivalence, that compatibility is replaced by posterior ratio consistency.

4 Sampler for posterior inference

4.1 Marginal likelihood derivation

In this section we provide further details on the MCMC scheme adopted for posterior inference on DAGs and intervention targets. Our algorithm is based on a collapsed Metropolis-Hastings sampler (Metropolis et al., 1953), where the DAG parameters $\boldsymbol{\theta} = \{\mathbf{D}, \mathbf{L}, \boldsymbol{\Phi}^{(1)}, \dots, \boldsymbol{\Phi}^{(K)}\}$ are integrated out. Accordingly, we first focus on the integrated likelihood

$$m(\mathbf{X} | I_1, \dots, I_K, \mathcal{D}) = \int_{\boldsymbol{\theta} \in \Theta} f(\mathbf{X} | \boldsymbol{\theta}, I_1, \dots, I_K, \mathcal{D}) p(\boldsymbol{\theta}) d\boldsymbol{\theta}. \quad (10)$$

To integrate out all parameters in $\boldsymbol{\theta}$, we exploit the structure of the prior $p(\boldsymbol{\theta})$ and we first re-write Equation (5) in the main text as

$$f(\mathbf{X} | \boldsymbol{\theta}, I_1, \dots, I_K, \mathcal{D}) = \prod_{j=1}^q \left\{ \varphi_{n_j^*}(\mathbf{X}_j^* | -\mathbf{X}_{\text{pa}_{\mathcal{D}}(j)}^* \mathbf{L}_{\prec j}, \sigma_j^2 \mathbf{I}_{n_j^*}) \prod_{k:j \in I_k} \varphi_{n^{(k)}}(\mathbf{X}_j^{(k)} | \mathbf{0}, \phi_j^{(k)} \mathbf{I}_{n_j^{(k)}}) \right\},$$

where \mathbf{X}_j^* is the $(n_j^*, 1)$ vector collecting all the observations $\mathbf{X}_j^{(k)}$ such that $j \notin I_k$ while $\mathbf{X}_{\text{pa}_{\mathcal{D}}(j)}^*$ is the corresponding $n_j^* \times |\text{pa}_{\mathcal{D}}(j)|$ matrix of parent-nodes observations. Accordingly, $n_j^* = \sum_{k:j \notin I_k} n^{(k)}$. Now notice that, because of prior parameter independence, the likelihood admits the same factorization of (8) (see our main text). Therefore, the integrated likelihood can be obtained as $m(\mathbf{X} | I_1, \dots, I_K, \mathcal{D})$

$$\begin{aligned} &= \prod_{j=1}^q \left\{ \int_0^\infty \int_{\mathbb{R}^{|\text{pa}_{\mathcal{D}}(j)|}} \left\{ \varphi_{n_j^*}(\mathbf{X}_j^* | -\mathbf{X}_{\text{pa}_{\mathcal{D}}(j)}^* \mathbf{L}_{\prec j}, \sigma_j^2 \mathbf{I}_{n_j^*}) \cdot p(\mathbf{L}_{\prec j} | \sigma_j^2) p(\sigma_j^2) d\mathbf{L}_{\prec j} d\sigma_j^2 \right\} \right. \\ &\quad \left. \cdot \prod_{k:j \in I_k} \int_0^\infty \left\{ \varphi_{n^{(k)}}(\mathbf{X}_j^{(k)} | \mathbf{0}, \phi_j^{(k)} \mathbf{I}_{n_j^{(k)}}) \cdot p(\phi_j^{(k)}) d\phi_j^{(k)} \right\} \right\}. \quad (11) \end{aligned}$$

Also, because of conjugacy each integral in the product is available in closed-form as the ratio of prior and posterior normalizing constants. Therefore we obtain

$$m(\mathbf{X} | I_1, \dots, I_K, \mathcal{D}) = \prod_{j=1}^q \left\{ m(\mathbf{X}_j^* | \mathbf{X}_{\text{pa}_{\mathcal{D}}(j)}^*, \mathcal{D}) \cdot \prod_{k:j \in I_k} m(\mathbf{X}_j^{(k)}) \right\}, \quad (12)$$

where

$$\begin{aligned} m(\mathbf{X}_j^* | \mathbf{X}_{\text{pa}_{\mathcal{D}}(j)}^*, \mathcal{D}) &= (2\pi)^{-\frac{n_j^*}{2}} \cdot \frac{|g\mathbf{I}_{|\text{pa}_{\mathcal{D}}(j)|}|^{1/2}}{|g\mathbf{I}_{|\text{pa}_{\mathcal{D}}(j)|} + \mathbf{S}_{\text{pa}_{\mathcal{D}}(j)}^*|^{1/2}} \cdot \frac{\Gamma\left(\frac{a_j^{\mathcal{D}}}{2} + \frac{n_j^*}{2}\right)}{\Gamma\left(\frac{a_j^{\mathcal{D}}}{2}\right)} \cdot \frac{\left(\frac{g}{2}\right)^{a_j^{\mathcal{D}}/2}}{\left(\frac{g+r_j^*}{2}\right)^{(a_j^{\mathcal{D}}+n_j^*)/2}}, \\ m(\mathbf{X}_j^{(k)}) &= (2\pi)^{-\frac{n_j^{(k)}}{2}} \cdot \frac{\Gamma\left(\frac{a_j^{(k)}}{2} + \frac{n_j^{(k)}}{2}\right)}{\Gamma\left(\frac{a_j^{(k)}}{2}\right)} \cdot \frac{\left(\frac{g}{2}\right)^{a_j^{(k)}/2}}{\left(\frac{g+s_j^{(k)}}{2}\right)^{(a_j^{(k)}+n_j^{(k)})/2}}, \end{aligned}$$

and

$$\begin{aligned} \mathbf{S}_{\text{pa}_{\mathcal{D}}(j)}^* &= \mathbf{X}_{\text{pa}_{\mathcal{D}}(j)}^{*\top} \mathbf{X}_{\text{pa}_{\mathcal{D}}(j)}^*, \\ r_j^* &= \mathbf{X}_j^{*\top} \mathbf{X}_j^* - \mathbf{X}_j^{*\top} \mathbf{X}_{\text{pa}_{\mathcal{D}}(j)}^* \left(g \mathbf{I}_{|\text{pa}_{\mathcal{D}}(j)|} + \mathbf{S}_{\text{pa}_{\mathcal{D}}(j)}^* \right)^{-1} \mathbf{X}_{\text{pa}_{\mathcal{D}}(j)}^{*\top} \mathbf{X}_j^*, \\ s_j^{(k)} &= \mathbf{X}_j^{(k)\top} \mathbf{X}_j^{(k)}. \end{aligned}$$

We then construct a collapsed Metropolis-Hastings sampler (Metropolis et al., 1953) to approximate

$$p(I_1, \dots, I_K, \mathcal{D} \mid \mathbf{X}) \propto m(\mathbf{X} \mid I_1, \dots, I_K, \mathcal{D}) \cdot p(I_1, \dots, I_K) p(\mathcal{D}), \quad (13)$$

the posterior distribution of DAGs and intervention targets. The output is a collection of DAGs $\{\mathcal{D}^{(s)}\}_{s=1}^S$ and targets $\{I_1^{(s)}, \dots, I_K^{(s)}\}_{s=1}^S$ approximately sampled from (13), where S is the number of finally kept MCMC iterations. We summarize our MCMC scheme in Algorithm 1 and provide details about the update of DAGs and targets in the following sections.

4.2 Update of \mathcal{D}

The full conditional distribution of \mathcal{D} is $p(\mathcal{D} \mid I_1, \dots, I_K, \mathbf{X}) \propto m(\mathbf{X} \mid I_1, \dots, I_K, \mathcal{D}) p(\mathcal{D})$. Update of DAG \mathcal{D} can be performed through a Metropolis Hastings step, where, given the current DAG, a new DAG $\tilde{\mathcal{D}}$ is proposed from $q(\tilde{\mathcal{D}} \mid \mathcal{D})$ and accepted with probability

$$\alpha = \min \left\{ 1; \frac{m(\mathbf{X} \mid I_1, \dots, I_K, \tilde{\mathcal{D}})}{m(\mathbf{X} \mid I_1, \dots, I_K, \mathcal{D})} \cdot \frac{p(\tilde{\mathcal{D}})}{p(\mathcal{D})} \cdot \frac{q(\mathcal{D} \mid \tilde{\mathcal{D}})}{q(\tilde{\mathcal{D}} \mid \mathcal{D})} \right\}. \quad (14)$$

The proposal distribution $q(\tilde{\mathcal{D}} \mid \mathcal{D})$ follows the lines of Chickering (2002). We consider three types of operators that locally modify the input DAG \mathcal{D} : insert a directed edge (InsertD $u \rightarrow v$ for short), delete a directed edge (DeletedD $u \rightarrow v$) and reverse a directed edge (ReverseD $u \rightarrow v$). For a given \mathcal{D} we then construct the set of valid operators $\mathcal{O}_{\mathcal{D}}$, that is operators whose resulting graph is a DAG. A DAG $\tilde{\mathcal{D}}$ is then called a *direct successor* of \mathcal{D} if it can be reached by applying an operator in $\mathcal{O}_{\mathcal{D}}$ to \mathcal{D} . Given the current \mathcal{D} we propose $\tilde{\mathcal{D}}$ by uniformly sampling an element in $\mathcal{O}_{\mathcal{D}}$ and applying it to \mathcal{D} . Since there is a one-to-one correspondence between each operator and resulting DAG, the probability of transition is $q(\tilde{\mathcal{D}} \mid \mathcal{D}) = 1/|\mathcal{O}_{\mathcal{D}}|$, for each $\tilde{\mathcal{D}}$ direct successor of \mathcal{D} .

4.3 Update of I_1, \dots, I_K

Conditionally on DAG \mathcal{D} we update the K targets I_1, \dots, I_K (equivalently, the indicator vectors $\mathbf{h}_1, \dots, \mathbf{h}_K$) sequentially; see also Section 2.4 in the main text. For a given k , the full conditional of I_k is $p(I_k \mid \{I_s\}_{s \neq k}, \mathcal{D}, \mathbf{X}) \propto m(\mathbf{X} \mid I_1, \dots, I_K, \mathcal{D}) \cdot p(I_1, \dots, I_K)$. Again, update of I_k conditionally on $\{I_s\}_{s \neq k}$ and DAG \mathcal{D} can be performed through a Metropolis Hastings step, where a

new target \tilde{I}_k proposed from $q(\tilde{I}_k | I_k)$ is accepted with probability

$$\beta_k = \min \left\{ 1; \frac{m(\mathbf{X} | \tilde{I}_k, \{I_s\}_{s \neq k}, \mathcal{D})}{m(\mathbf{X} | I_k, \{I_s\}_{s \neq k}, \mathcal{D})} \cdot \frac{p(\tilde{I}_k)}{p(I_k)} \cdot \frac{q(I_k | \tilde{I}_k)}{q(\tilde{I}_k | I_k)} \right\}. \quad (15)$$

Starting from a given target I_k , a candidate target \tilde{I}_k is obtained by randomly drawing a node $j \in \{1, \dots, q\}$ which is included in I_k if $j \notin I_k$, otherwise removed. The structure of this proposal is such that $q(I_k | \tilde{I}_k)/q(\tilde{I}_k | I_k) = 1$ for any I_k and \tilde{I}_k differing in one j , and equal to 0 otherwise.

Algorithm 1: MCMC scheme to sample from the posterior (13)

Input: K datasets $\mathbf{X}^{(1)}, \dots, \mathbf{X}^{(K)}$
Output: S samples from the posterior (13)

- 1 Initialize $\mathcal{D}^{(0)}$, e.g. the empty DAG, and the targets I_1, \dots, I_K , e.g. $I_k = \emptyset$ for $k = 1, \dots, K$;
- 2 **for** $s = 1, \dots, S$ **do**
- 3 Sample $\tilde{\mathcal{D}}$ from $q(\tilde{\mathcal{D}} | \mathcal{D}^{(s-1)})$ and set $\mathcal{D}^{(s)} = \tilde{\mathcal{D}}$ with probability α in (14),
- 4 otherwise $\mathcal{D}^{(s)} = \mathcal{D}^{(s-1)}$;
- 5 **for** $k = 1, \dots, K$ **do**
- 6 Propose \tilde{I}_k from $q(\tilde{I}_k | I_k)$ and set $I_k^{(s)} = \tilde{I}_k$ with probability β_k in (15),
- 7 otherwise $I_k^{(s)} = I_k^{(s-1)}$;
- 8 **end**
- 9 **end**

4.4 Computational and convergence issues

We first investigate the computational time of our method as a function of the number of variables q and sample size $n^{(k)}$. Figure 1 summarizes the computational time (averaged over 12 repetitions) *per iteration* for $q \in \{5, 10, 20, 50, 100, 200\}$ and for sample sizes $n^{(k)} \in \{10, 20, 50, 100, 200, 500, 1000\}$. Each dotted line in the left (right) panel of the figure corresponds to a fixed value of q ($n^{(k)}$); moreover, higher curves are associated to higher values of q and $n^{(k)}$ respectively. The behavior of all curves suggests a polynomial dependence of the computational time from both q and $n^{(k)}$. Results were obtained on a PC Intel(R) Core(TM) i7-8550U 1,80 GHz.

Computational challenges are related the huge dimension of the graph space. More precisely, our proposal distribution for DAG moves in the MCMC is *uniform* over the set of direct successors of the current DAG \mathcal{D} , a set consisting of all those DAGs obtained from \mathcal{D} by applying single-edge removals, insertions or reversals. This kind of proposal results in a DAG acceptance rate which is relatively small, especially when the posterior over DAGs is highly concentrated (typically when n is large), meaning that many of the proposed DAGs are discarded by the MCMC. As a consequence, convergence to the target posterior over DAGs can require a higher number of iterations. One possible solution not investigated in the current manuscript

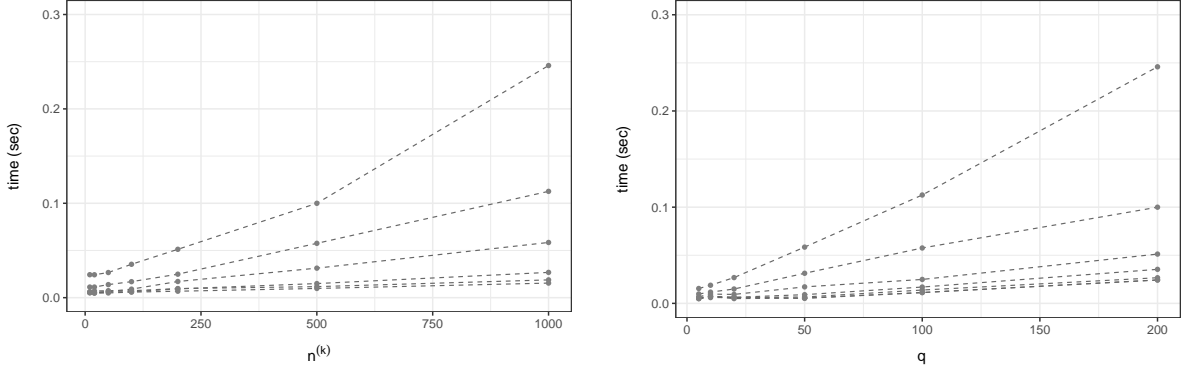


Figure 1: Computational time (in seconds) *per iteration*, as a function of the sample size $n^{(k)}$ (left panel, for increasing values of q) and as a function of q (right panel, for increasing values of $n^{(k)}$), averaged over 12 simulated multiple datasets.

is represented by Feature Inclusion Stochastic Search (FINCS) methods, based on the notion of (marginal) probability of inclusion for a *potential* edge; see for instance Scott and Carvalho (2008) and Altomare et al. (2013) who adopt FINCS for undirected graphs and DAGs with fixed ordering of the nodes respectively. In brief, the intuition is that insertion of edges with a higher posterior probability of inclusion (computed up to time t) is more likely to be accepted. Accordingly, at each step t of the MCMC, moves are proposed with different probabilities, each reflecting the probability of (non) inclusion of those edges which are involved in the move.

In a similar way, alternative choices of the proposal distribution for target nodes can speed up MCMC convergence. More precisely, in the current version of our MCMC scheme, starting from a given target I_k , a candidate target \tilde{I}_k is obtained by randomly drawing a node $j \in \{1, \dots, q\}$ which is included in I_k if $j \notin I_k$, otherwise removed from I_k if j already belongs to I_k . The structure of this proposal is such that $q(I_k | \tilde{I}_k) / q(\tilde{I}_k | I_k) = 1$ for any I_k and \tilde{I}_k differing in one j , and equal to 0 otherwise. Accordingly, our proposal modifies one target at a time. Alternatively, one can randomly choose a node $j \in \{1, \dots, q\}$ and propose a new allocation of the node in the K datasets. More precisely, for $j = 1, \dots, q$, let $\boldsymbol{\xi}_j = (\xi_j(1), \dots, \xi_j(K))^T$ be a K -th vector of indicators such that for each $k = 1, \dots, K$, $\xi_j(k) = 1$ if and only if $j \in I_k$. Similarly as before, we can then sample $j \in \{1, \dots, q\}$ and propose a new configuration of vector $\boldsymbol{\xi}_j$ given the current one. This kind of move applied for the *simultaneous* update of a set of nodes can result in an improved MCMC mixing and convergence to the posterior of DAGs and targets.

Finally, we assess convergence of the MCMC algorithm by investigating how selected features of the model space behave across iterations, and by evaluating the agreement between independent MCMC chains with different starting points. To this end, we considered some pilot simulations for each setting defined by $q = 20$ and $q = 40$ variables for a number of iterations $S = 25000$ and $S = 50000$ respectively. Specifically, under each simulation, we ran two inde-

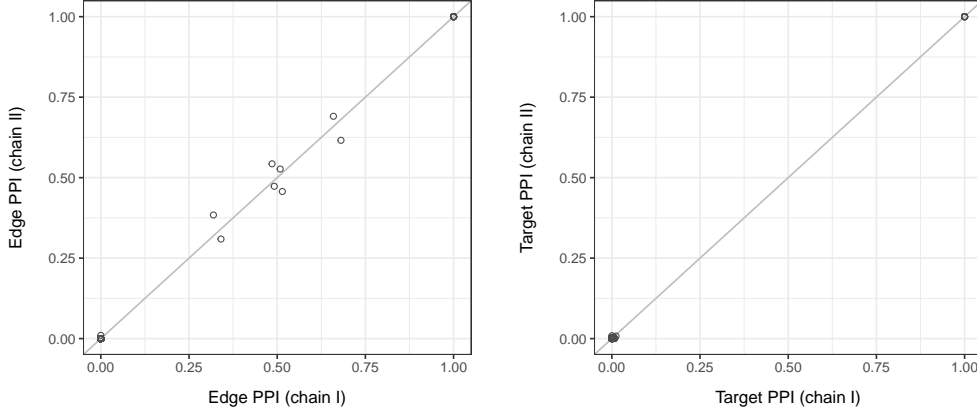


Figure 2: Simulations. Scatter plots of estimated marginal posterior probabilities of edge inclusion (left panel) and marginal posterior probabilities of intervention (target inclusion, right panel) obtained from two independent MCMC chains, chain I and chain II.

pendent chains of same length, with random starting points. We then compare the estimated marginal posterior probabilities of edge inclusion obtained from the two chains as well as the marginal posterior probabilities of intervention (target inclusion); see also Equations (13) and (14) in the main text. As an instance, Figure 2 shows the scatter plots of the estimated posterior probabilities of edge inclusion and intervention obtained from the two chains for one pilot simulation with $q = 20$. By visual inspection, we see that the agreement between the chains is highly satisfactory, since points are clustered around the main diagonal of the plot.

5 More simulated scenarios

We present extensive simulation studies to evaluate the performance of our methodology and comparisons with alternative approaches. We first include the complete results for the simulation study introduced in the main text, with all interventional datasets. Next, we consider additional simulation scenarios including both observational *and* interventional data. Under both cases, we fix the number of groups (interventions) $K = 4$, while we vary the number of variables $q \in \{20, 40\}$ and the number of observations for each group $k = 1, \dots, K$, $n^{(k)} \in \{10, 20, 50, 100, 200, 500\}$. Under each scenario defined by $(q, n^{(k)})$, we perform 40 simulations, each corresponding to a true DAG \mathcal{D} , family of targets $\{I_1, \dots, I_K\}$ and resulting in a (multiple with $K = 4$ groups) dataset. Specifically, we first randomly generate a topologically ordered DAG \mathcal{D} with probability of edge inclusion $p_{edge} = 2/q$. Generation of targets I_1, \dots, I_K is specific to each of the two following settings and therefore detailed in the next. Given DAG \mathcal{D} and family of targets $\{I_1, \dots, I_K\}$, parameters \mathbf{D} , \mathbf{L} and $\Phi^{(k)} = \{\phi_j^{(k)}, j \in I_k\}$ are generated by fixing $\mathbf{D} = \mathbf{I}_q$, and $\phi_j^{(k)} = 0.1$ for each $j \in I_k$ and $k = 1, \dots, K$; non zero elements of \mathbf{L} are

instead uniformly chosen in the interval $[-1, -0.1] \cup [0.1, 1]$; see also Equation (4) of the main text.

In the first case with all interventional datasets, a family of intervention targets I_1, \dots, I_K is generated under two scenarios resembling different degrees of “sparsity” in the targets. Scenario *Sparse* is characterized by a moderate number of interventions, with each target I_k obtained by drawing without replacement $s \in \{2, 4\}$ nodes, respectively for $q = 20$ or $q = 40$. On the other hand, in Scenario *Diffuse* we assign each node to one of the K targets. As a consequence, each variable is involved in one of the K interventions, with an overall larger number of simultaneous interventions (sizes of the targets I_k). Next, for each $k = 1 \dots, K$, $n^{(k)}$ i.i.d. interventional data collected in the $n^{(k)} \times q$ data matrix $\mathbf{X}^{(k)}$ are generated as in Equation (4) of our paper. Each dataset is therefore a collection of K interventional data matrices $\mathbf{X}^{(1)}, \dots, \mathbf{X}^{(K)}$.

In the following we present results relative to our Bayesian method (*Bayes*) and some benchmark methods. In particular we include the *Unknown Target Interventional Greedy Sparsest Permutation* algorithm of Squires et al. (2020), a non-parametric strategy for target estimation and structure learning based on permutation tests that we implement at significance level $\alpha \in \{0.1\%, 0.001\%\}$ (IGSP 0.1% and IGSP 0.001% respectively), as recommended in the original paper. Other significance levels are not included since result in worse performances. We also include Algorithm 1 of He and Geng (2016), a method for DAG learning from interventional data with unknown targets, that builds on the PC algorithm of Spirtes et al. (2000). We implement it at significance level 0.05; results obtained at different significance levels and not included led to worse performances. As a further benchmark, we also construct a baseline node-wise regression approach, by adapting to our interventional setting the two-stage adaptive lasso method of Han et al. (2016): we first recover a baseline DAG structure by applying the adaptive lasso to the union of the K data matrices; next, we apply adaptive lasso to each dataset $\mathbf{X}^{(k)}$ separately, and estimate intervention targets by comparing overall and group-specific DAGs; we call this benchmark *Node-wise*. Finally, we include the Greedy Interventional Equivalence Search (GIES) method of Hauser and Bühlmann (2012), a search-and-score method based on maximum likelihood estimation, which provides an estimate of the graph representing the interventional Markov equivalence class of the true DAG. GIES was developed for known intervention targets, that we input, as in an *oracle* setting, using the *true* intervened nodes. We implement GIES using the Extended Bayesian Information Criterion (EBIC) with tuning coefficient $\gamma \in \{0.5, 1\}$ (GIES 0.5 and GIES 1 respectively) as also recommended in Foygel and Drton (2010). All methods can be adopted for DAG structure learning, whilst only IGSP and *Node-wise* can perform target estimation. Finally, notice that IGSP requires $n^{(k)} > q$, so that results of IGPS for $n^{(k)} \in \{10, 20\}$ are missing.

All methods can be adopted for DAG structure learning, whilst only IGSP 0.1, IGSP 0.001% and *Node-wise* can perform target estimation; see the main text for more details. Moreover,

while all the benchmarks directly output single estimates of DAGs (and targets), we adapt our Bayesian method to provide point estimates as described in our paper (Section 5.2). Finally, notice that IGSP requires $n^{(k)} > q$, so that results of IGSP for $n^{(k)} \in \{10, 20\}$ are missing.

A summary for $q = 20$ and $q = 40$ is reported in the box-plots of Figures 3 and 4. This reports the distribution of FPR and FNR constructed across the $N = 40$ simulations for three methods under evaluation and increasing sample sizes $n^{(k)}$ under *Sparse* and *Diffuse* Scenarios. We notice that, coherently with the theoretical results of Section 4 in the main text, for our method both sources of error vanish as sample size increases. This tendency is more evident for FNR that rapidly goes to zero already at moderate sample sizes, e.g. $n^{(k)} = 20$. It is clear the outperformance of our proposal, relative to the benchmarks, with *Node-wise* performing equally well only in terms of FNR.

We then evaluate the performance of each methodology in recovering the DAG structure. We compare each DAG estimate $\hat{\mathcal{D}}$ with the true DAG \mathcal{D} , by measuring the Structural Hamming Distance (SHD, Tsamardinos et al. 2006) between the two graphs; see also the main paper for details. Results are summarized in the box-plots of Figures 9 and 10, where each plot reports the distribution of SHD across the $N = 40$ simulated datasets for the various methods and increasing sample sizes $n^{(k)} \in \{10, \dots, 500\}$ under *Sparse* and *Diffuse* Scenarios. It is clear the tendency of a better and better recovery of the true graphical structure as we increase the amount of available data, and an overall better performance relative to all the benchmarks. The only exception is GIES 0.5 which was however implemented with input the true targets and outperforms our Bayesian method in few settings characterized by small sample sizes, where indeed target identification was more difficult for our method. However, it performs worse than *Bayes* as $n^{(k)}$ increases, especially under Scenario *Sparse*.

In the second case, we repeat all simulation settings above, but we impose $I_1 = \emptyset$, implying no interventions for the first data matrix $\mathbf{X}^{(1)}$, which therefore consists of *observational* data. Targets I_2, I_3, I_4 are again generated under two different scenarios, *Sparse* and *Diffuse*. We present results relative to both target estimation (Figures 7 and 8) and DAG structure learning (Figures 9 and 10). We note no substantial differences relative to the case with only interventional data.

6 Sensitivity analyses

For the two real data problems presented in the main text (Section 6), we perform sensitivity analyses with respect to the hyperparameter η , which represents the prior probability of edge inclusion in the DAG. Therefore, tuning of η can regulate the inclusion of specific links in the graph space. Since interventions modify the DAG structure through edge removals, choice of η may affect intervention targets' identification. Accordingly, we evaluate the impact of η

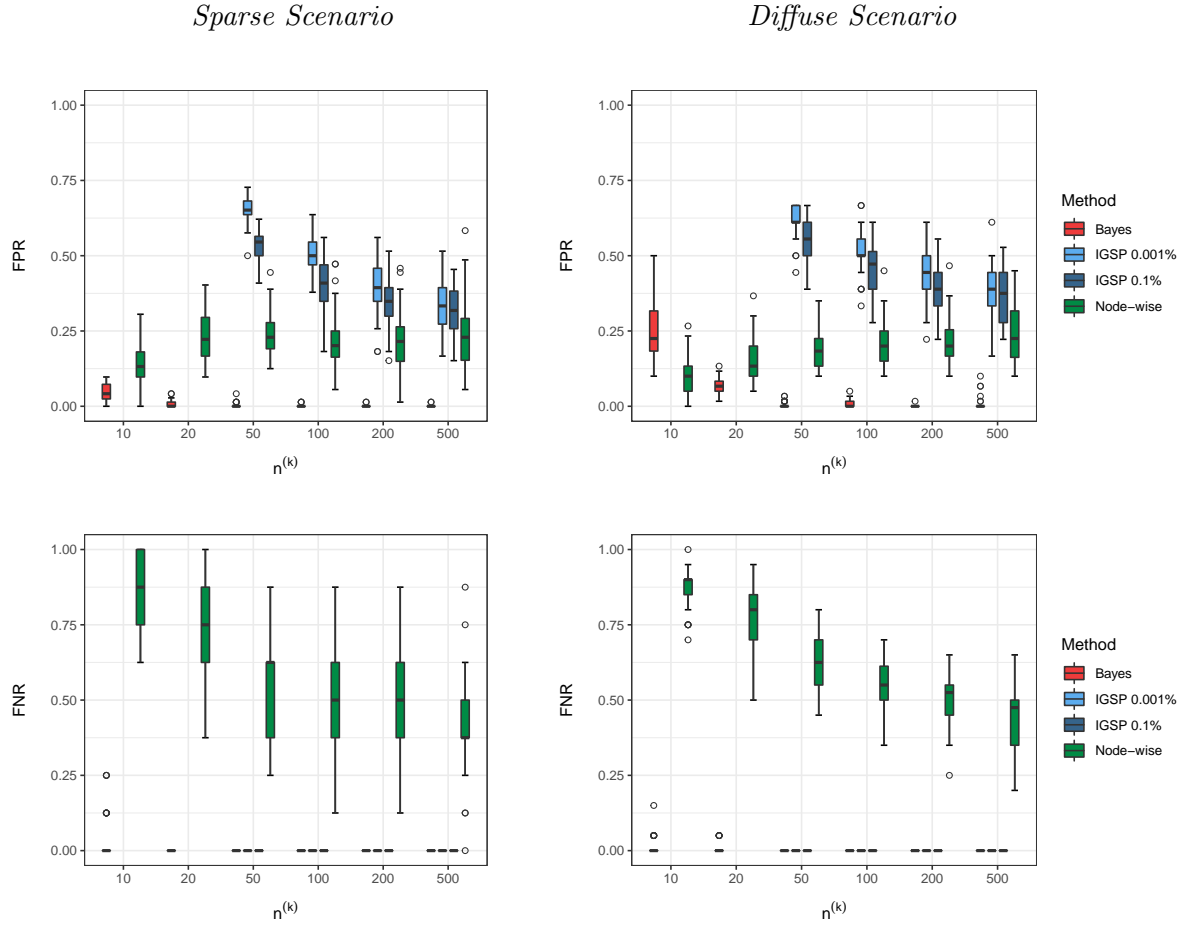


Figure 3: Simulation 1, interventional data. Distribution of the False Positive Rate (FPR, first row) and False Negative Rate (FNR, second row) across $N = 40$ simulated datasets under *Sparse* (first column) and *Diffuse* (second column) Scenarios, for number of nodes $q = 20$ and increasing sample sizes $n^{(k)}$. Methods under comparison are: our Bayesian methodology (Bayes), the *Unknown Target Interventional Greedy Sparsest Permutation* algorithm implemented at significance level $\alpha \in \{0.1\%, 0.001\%\}$ (IGSP 0.1% and IGSP 0.001%) and node-wise regression (Node-wise).

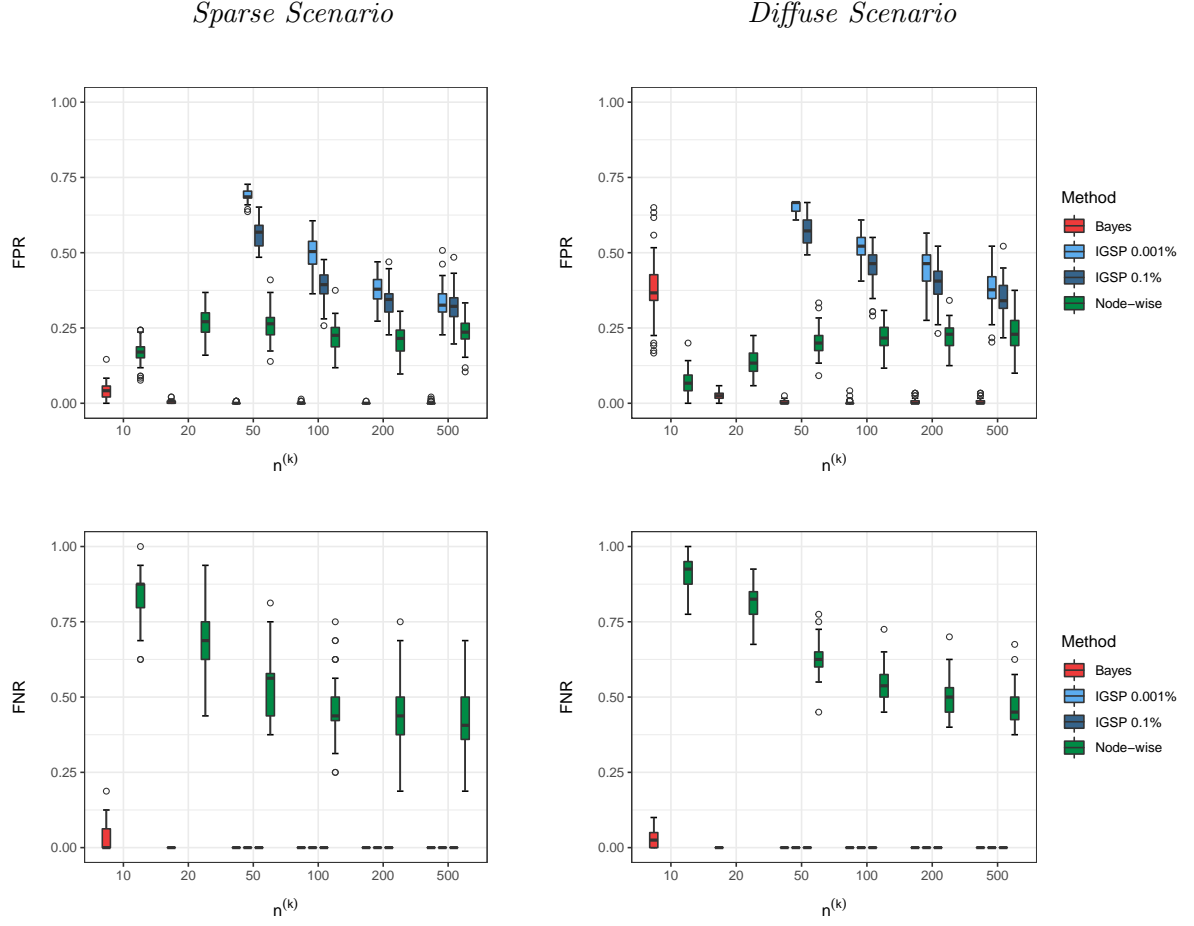


Figure 4: Simulation 1, interventional data. Distribution of the False Positive Rate (FPR, first row) and False Negative Rate (FNR, second row) across $N = 40$ simulated datasets under *Sparse* (first column) and *Diffuse* (second column) Scenarios, for number of nodes $q = 40$ and increasing sample sizes $n^{(k)}$. Methods under comparison are: our Bayesian methodology (Bayes), the *Unknown Target Interventional Greedy Sparsest Permutation* algorithm implemented at significance level $\alpha \in \{0.1\%, 0.001\%\}$ (IGSP 0.1%, IGSP 0.001%) and node-wise regression (Node-wise).

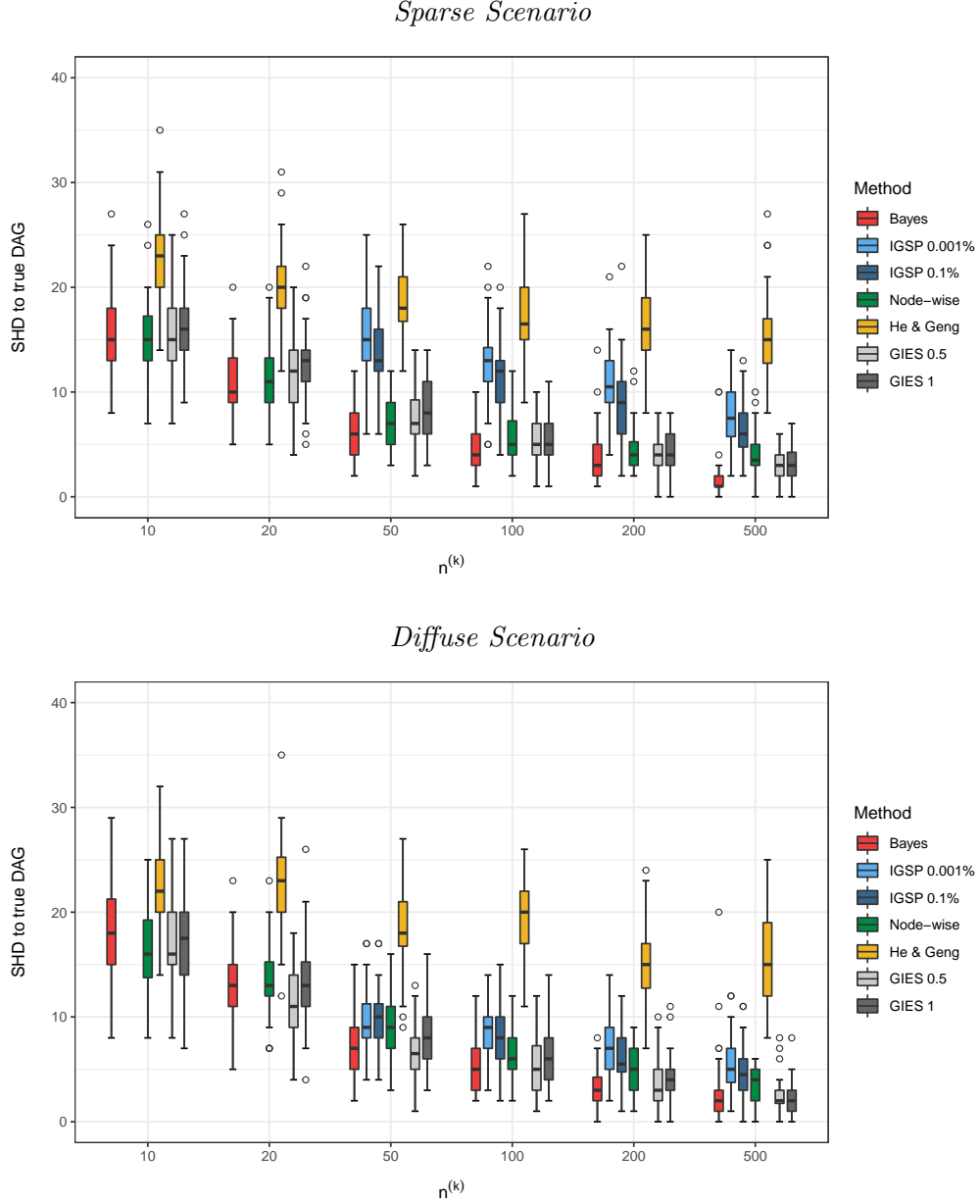


Figure 5: Simulation 1, interventional data. Distribution across $N = 40$ simulated datasets of the Structural Hamming Distance (SHD) between estimated and true DAG under *Sparse* (top) and *Diffuse* (bottom) Scenarios, for number of nodes $q = 20$ and increasing sample sizes $n^{(k)}$. Methods under comparison are: our Bayesian methodology (Bayes), the *Unknown Target Interventional Greedy Sparsest Permutation* algorithm implemented at significance level $\alpha \in \{0.1\%, 0.001\%\}$ (IGSP 0.1%, IGSP 0.001%), node-wise regression (Node-wise), Algorithm 1 of He and Geng (2016) (He & Geng) and the *Greedy Interventional Equivalence Search* method with tuning coefficient $\gamma \in \{0.5, 1\}$ (GIES 0.5, GIES 1).

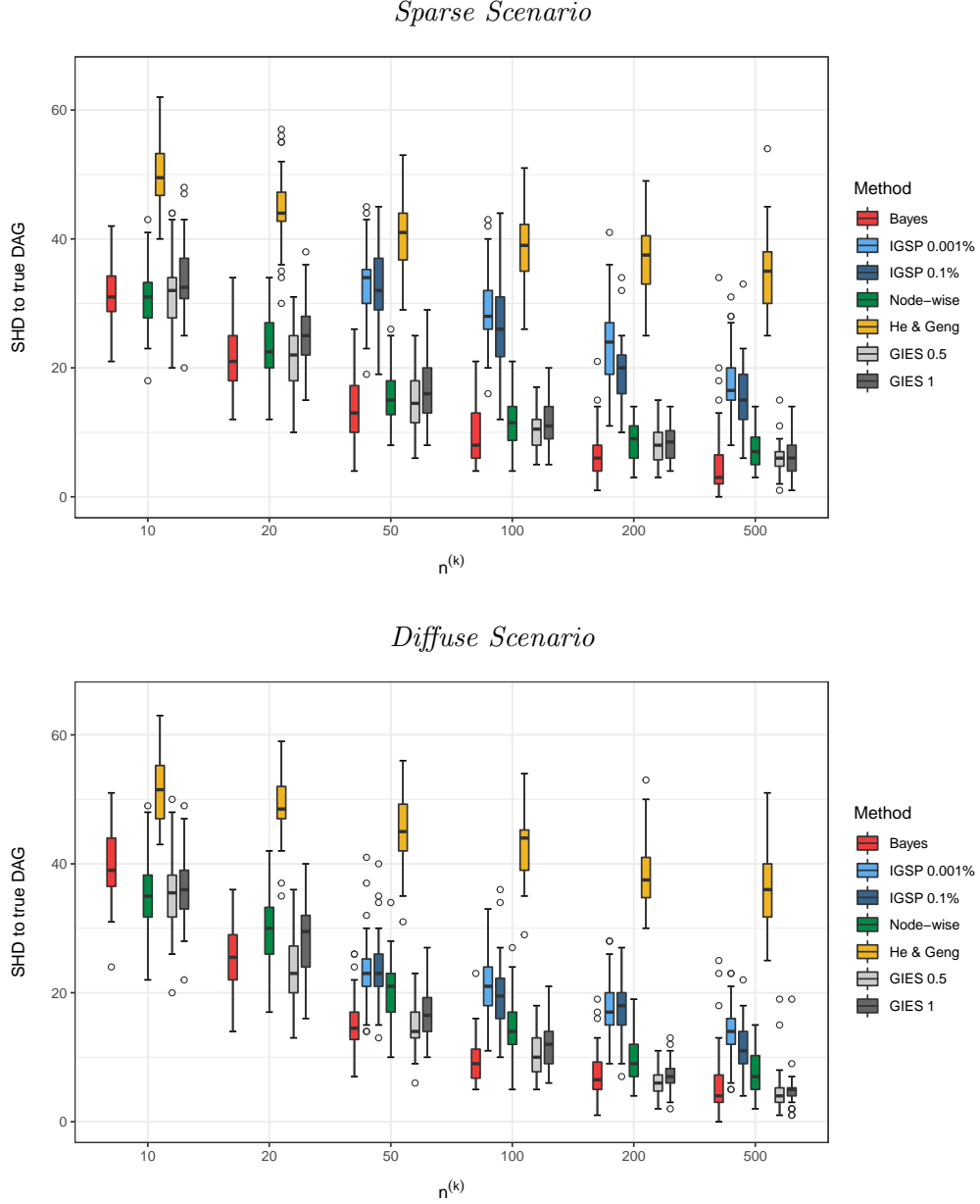


Figure 6: Simulation 1, interventional data. Distribution across $N = 40$ simulated datasets of the Structural Hamming Distance (SHD) between estimated and true DAG under *Sparse* (top) and *Diffuse* (bottom) Scenarios, for number of nodes $q = 40$ and increasing sample sizes $n^{(k)}$. Methods under comparison are: our Bayesian methodology (Bayes), the *Unknown Target Interventional Greedy Sparsest Permutation* algorithm implemented at significance level $\alpha \in \{0.1\%, 0.001\%\}$ (IGSP 0.1% IGSP 0.001%), node-wise regression (Node-wise), Algorithm 1 of He and Geng (2016) (He & Geng) and the *Greedy Interventional Equivalence Search* method with tuning coefficient $\gamma \in \{0.5, 1\}$ (GIES 0.5, GIES 1).

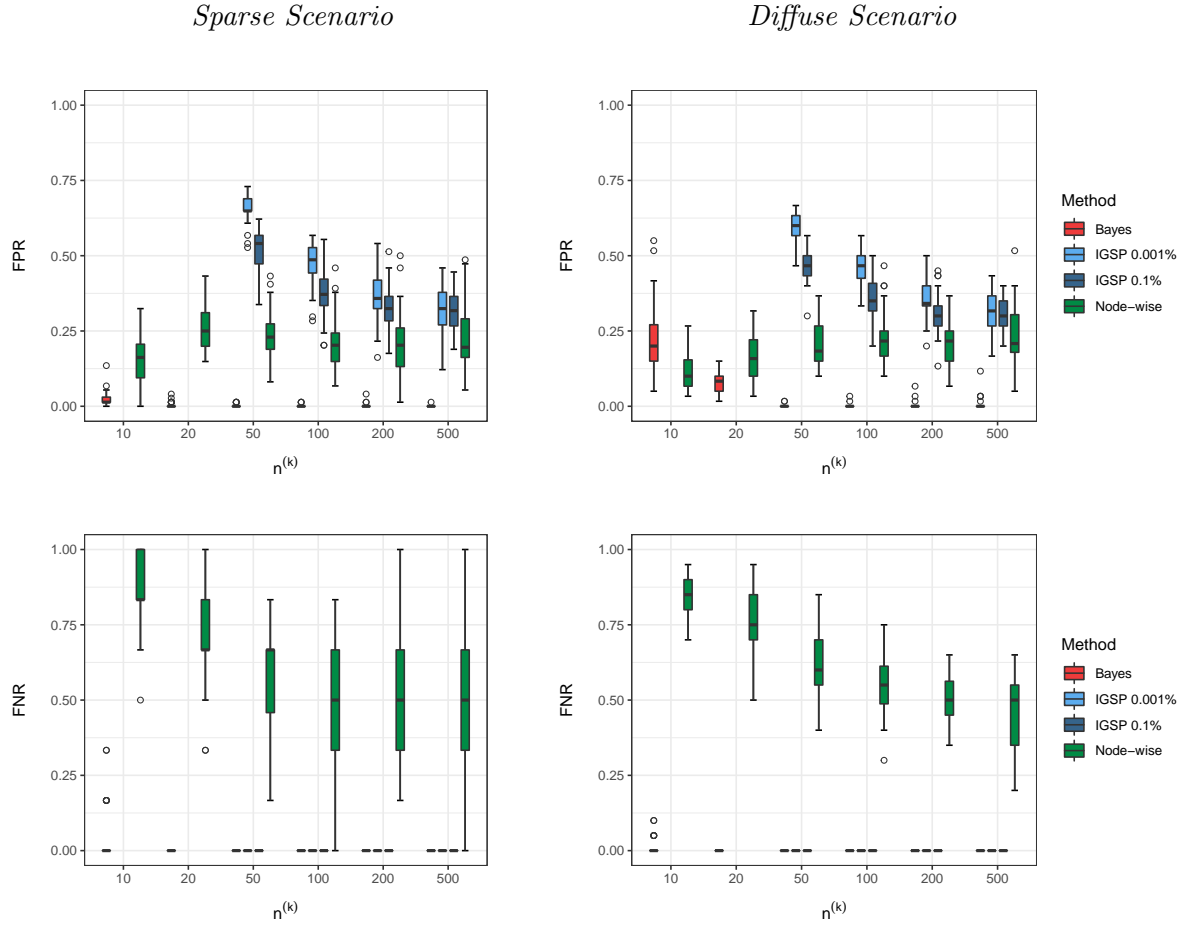


Figure 7: Simulation 2, observational and interventional data. Distribution of the False Positive Rate (FPR, first row) and False Negative Rate (FNR, second row) across $N = 40$ simulated datasets under *Sparse* (first column) and *Diffuse* (second column) Scenarios, for number of nodes $q = 20$ and increasing sample sizes $n^{(k)}$. Methods under comparison are: our Bayesian methodology (Bayes), the *Unknown Target Interventional Greedy Sparsest Permutation* algorithm implemented at significance level $\alpha \in \{0.1\%, 0.001\%\}$ (IGSP 0.1%, IGSP 0.001%) and node-wise regression (Node-wise).

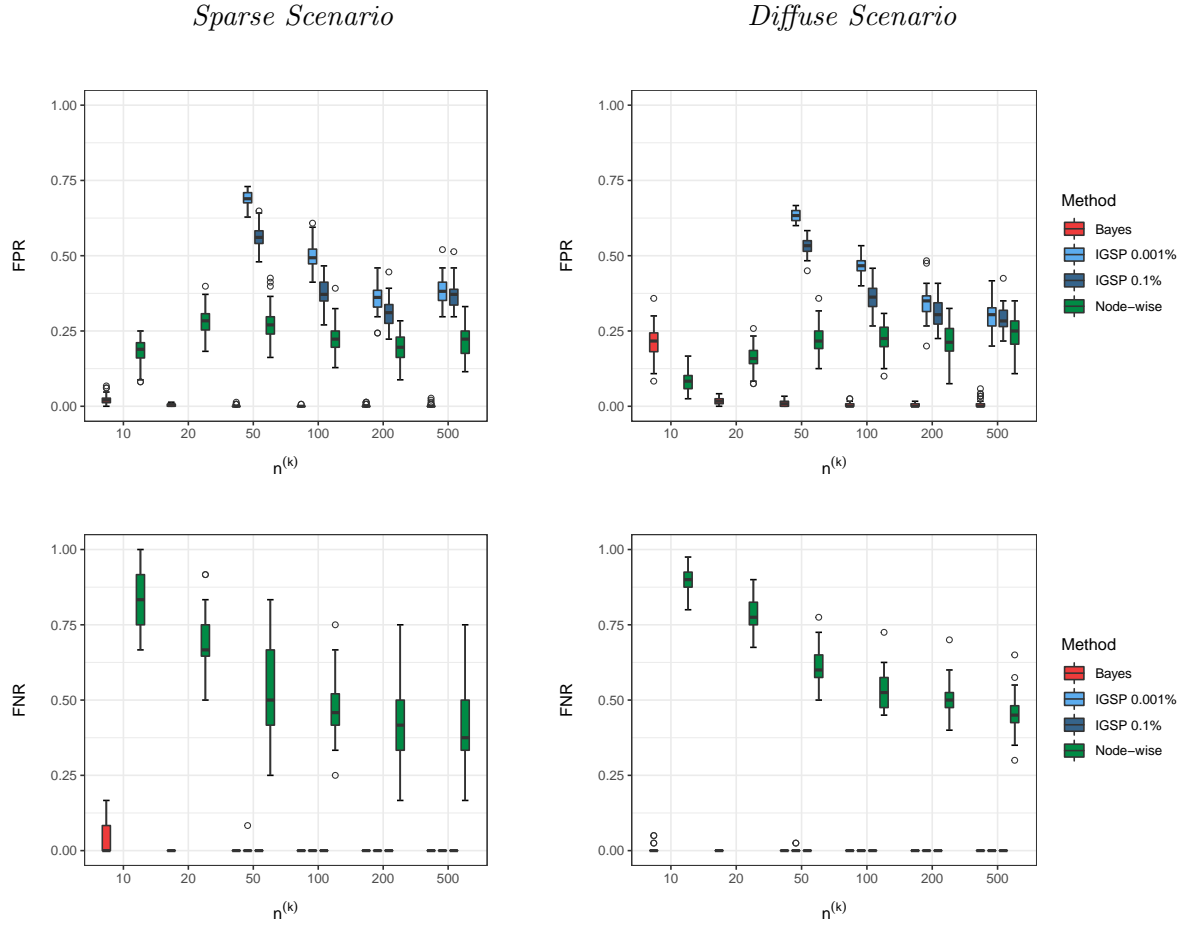
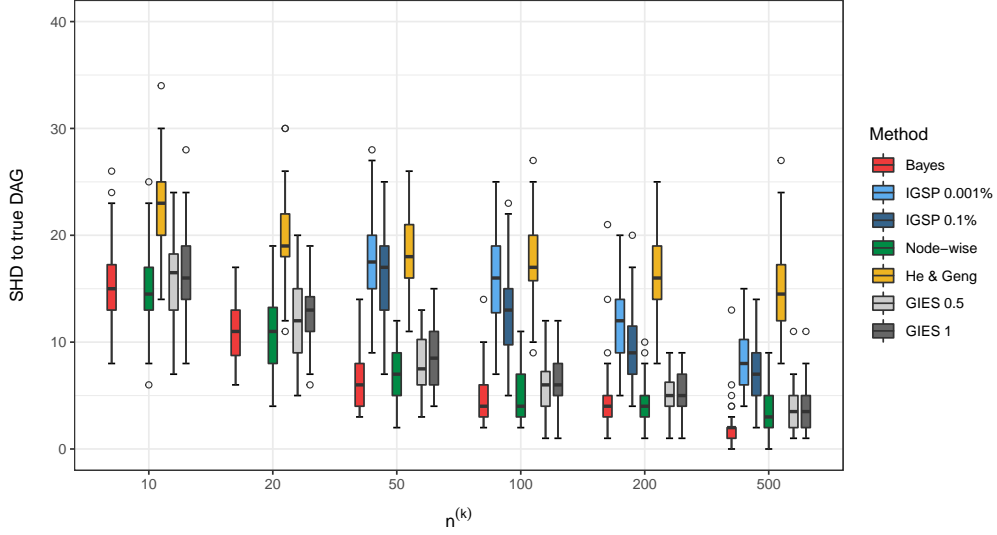


Figure 8: Simulation 2, observational and interventional data. Distribution of the False Positive Rate (FPR, first row) and False Negative Rate (FNR, second row) across $N = 40$ simulated datasets under *Sparse* (first column) and *Diffuse* (second column) Scenarios, for number of nodes $q = 40$ and increasing sample sizes $n^{(k)}$. Methods under comparison are: our Bayesian methodology (Bayes), the *Unknown Target Interventional Greedy Sparsest Permutation* algorithm implemented at significance level $\alpha \in \{0.1\%, 0.001\%\}$ (IGSP 0.1%, IGSP 0.001%) and node-wise regression (Node-wise).

Sparse Scenario



Diffuse Scenario

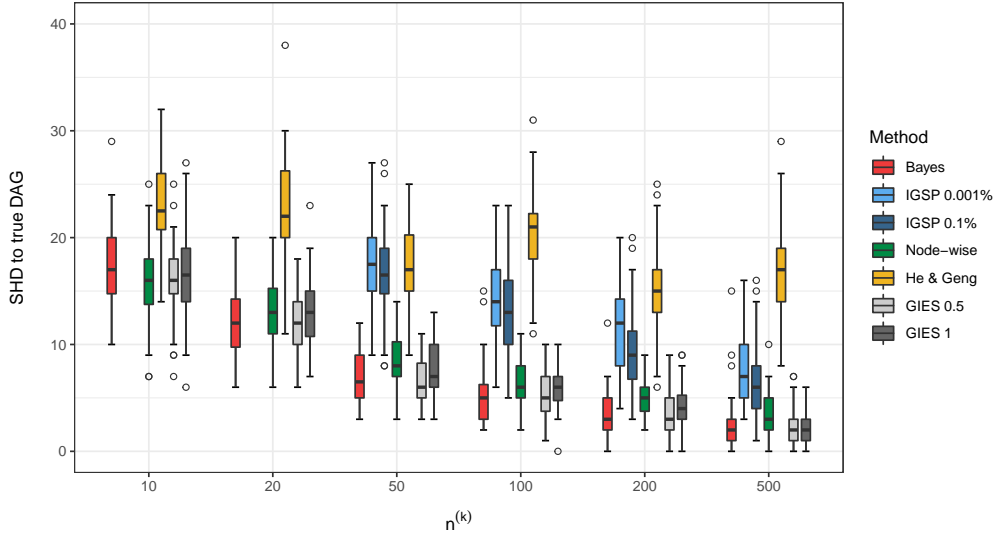
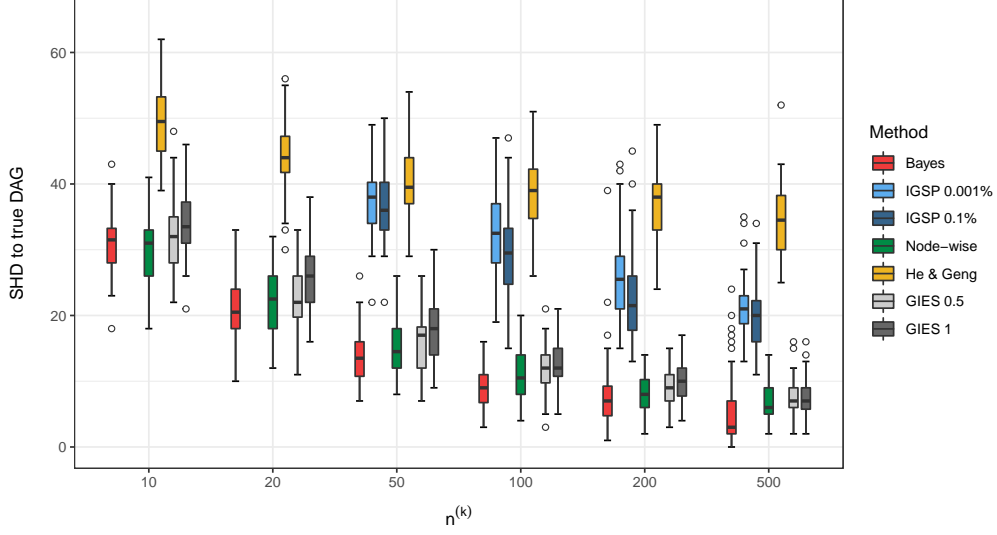


Figure 9: Simulation 2, observational and interventional data. Distribution across $N = 40$ simulated datasets of the Structural Hamming Distance (SHD) between estimated and true DAG under *Sparse* (top) and *Diffuse* (bottom) Scenarios, for number of nodes $q = 20$ and increasing sample sizes $n^{(k)}$. Methods under comparison are: our Bayesian methodology (Bayes), the *Unknown Target Interventional Greedy Sparsest Permutation* algorithm implemented at significance level $\alpha \in \{0.1\%, 0.001\%\}$ (IGSP 0.1%, IGSP 0.001%), node-wise regression (Node-wise), Algorithm 1 of He and Geng (2016) (He & Geng) and the *Greedy Interventional Equivalence Search* method with tuning coefficient $\gamma \in \{0.5, 1\}$ (GIES 0.5, GIES 1).

Sparse Scenario



Diffuse Scenario

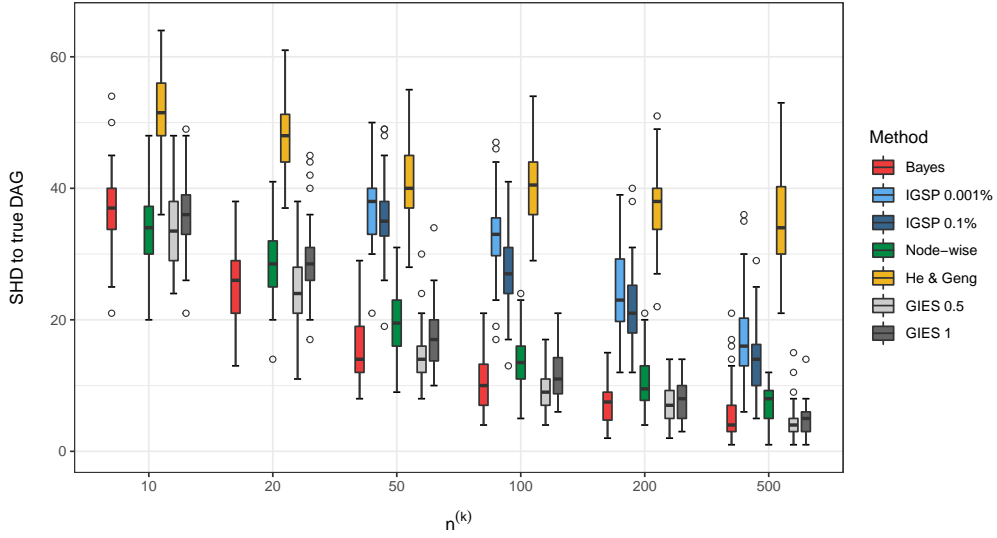


Figure 10: Simulation 2, observational and interventional data. Distribution across $N = 40$ simulated datasets of the Structural Hamming Distance (SHD) between estimated and true DAG under *Sparse* (top) and *Diffuse* (bottom) Scenarios, for number of nodes $q = 40$ and increasing sample sizes $n^{(k)}$. Methods under comparison are: our Bayesian methodology (Bayes), the *Unknown Target Interventional Greedy Sparsest Permutation* algorithm implemented at significance level $\alpha \in \{0.1\%, 0.001\%\}$ (IGSP 0.1%, IGSP 0.001%) node-wise regression (Node-wise), Algorithm 1 of He and Geng (2016) (He & Geng) and the *Greedy Interventional Equivalence Search* method with tuning coefficient $\gamma \in \{0.5, 1\}$ (GIES 0.5, GIES 1).

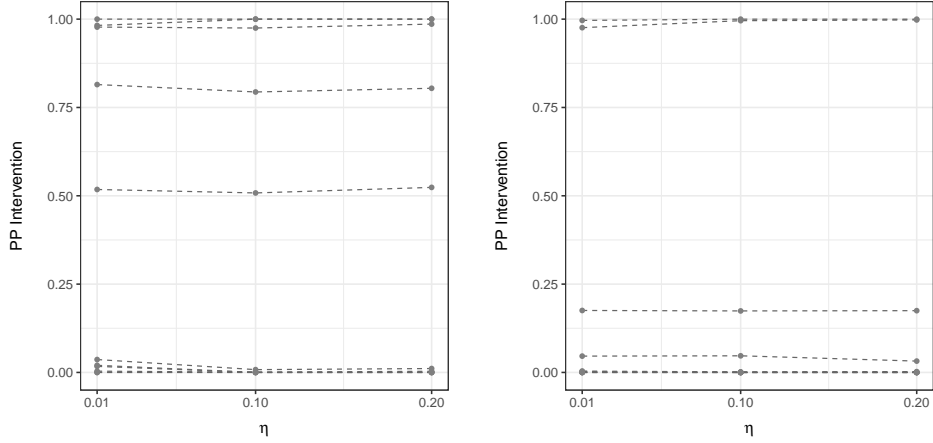


Figure 11: Sachs data. Posterior probabilities of intervention $\hat{p}_{j \in I_k}$ as a function of $\eta \in \{0.01, 0.1, 0.2\}$ for each node $j = 1, \dots, q$, for the two datasets with unknown targets, I_8 (left) and I_9 (right).

on the posterior probabilities of intervention $\hat{p}_{j \in I_k}$ as defined in Equation (13) of our paper, for each node $j = 1, \dots, q$ and target I_k , $k = 1, \dots, K$. For both the applications, we vary $\eta \in \{0.01, 0.10, 0.20\}$.

Results for the Sachs data are reported in Figure 11 where each plot refers to one of the two datasets with unknown targets. It is clear that in this setting, characterized by large group sample sizes $n^{(k)}$, results are quite insensitive to the choice of η . Similar results were obtained for the epilepsy data: from Figure 12, which reports the behaviour of $\hat{p}_{j \in I_k}$ as a function of η , for $j = 1, \dots, q$ and each of the three unknown targets (drugs) I_k , $k = 2, 3, 4$, posterior probabilities of inclusion are stable around zero for most of the nodes and targets. Differently, results are more affected by the choice of η for nodes with non-zero probabilities of inclusion. However, largest variations do not exceed a range of 0.25; in addition, estimated intervention targets, obtained with a threshold for inclusion of 0.5, never change, with the only exception being one node in the last group/intervention, included in the estimated target only for $\eta = 0.10$, with a probability of inclusion just above 0.5.

7 Real data analyses: additional results and figures

In this section we include additional results and plots for the application to real data. Specifically, for the epilepsy dataset we report in Figure 13 the heat map collecting the posterior probabilities of intervention computed for each node $v \in \{1, \dots, 100\}$ under each of the three drug therapies (interventions) as in Equation (15) of our paper. Sparsity in the map reveals that there are few genes exhibiting a high posterior probability of intervention under some of the treatments. In particular, only six genes are associated with probabilities of intervention exceeding 0.5. The same result is also clear from the posterior distribution of the number of intervened nodes (size

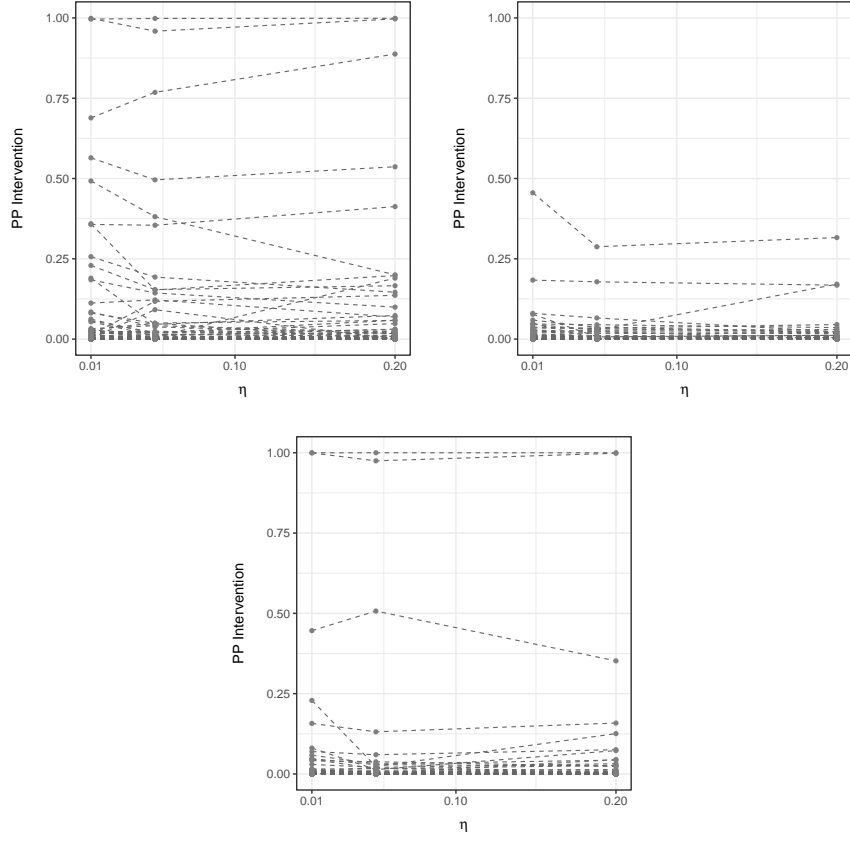


Figure 12: Epilepsy data. Posterior probabilities of intervention $\hat{p}_{j \in I_k}$ as a function of $\eta \in \{0.01, 0.1, 0.2\}$ for each node $j = 1, \dots, q$ and each dataset (targets I_2, I_3, I_4) corresponding to one of the three administered drugs.

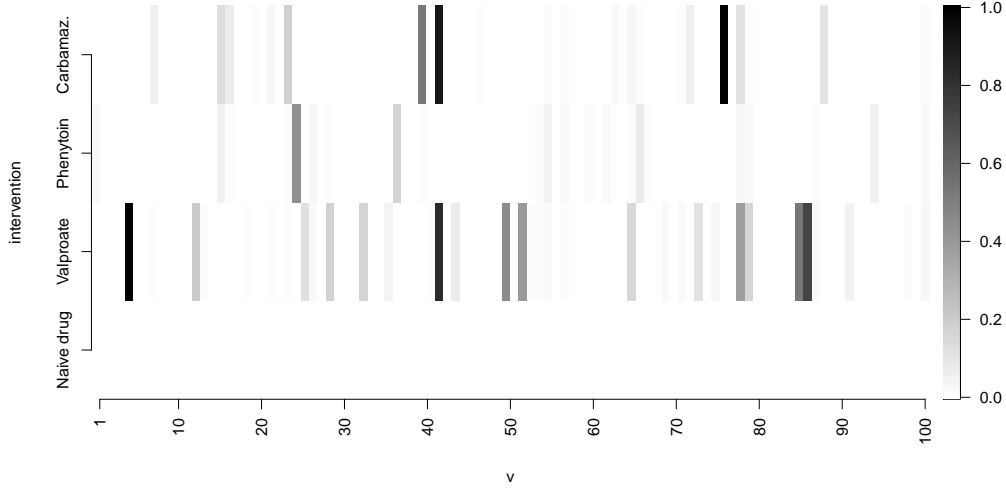


Figure 13: Epilepsy data. Heat map with estimated posterior probabilities of intervention computed for each drug-therapy (intervention) and each node/gene v ($v = 1, \dots, 100$).

of the intervention target) obtained for each of the three drug therapies. The approximated posterior distribution resulting from our MCMC output is reported in Figure 14.

For the Sachs data, Figure 15 reports our estimated (median probability) DAG $\hat{\mathcal{D}}$, with blue and pink circles representing nodes which are included in the estimated intervention targets \hat{I}_8 and \hat{I}_9 respectively. The same figure contains the intervention DAGs (Definition 2.1 in the main text) of $\hat{\mathcal{D}}$ given the intervention targets \hat{I}_8 and \hat{I}_9 , $\hat{\mathcal{D}}^{\hat{I}_8}$ and $\hat{\mathcal{D}}^{\hat{I}_9}$ respectively. The two graphs help understanding how the original DAG structure modifies after one of the reagents is applied.

We also report in Figure 16 three DAG estimates obtained from benchmark methods *Node-wise*, *He & Geng* and *IGSP 0.1%* applied to Sachs data; see also Section 5 within our paper for a brief presentation of these methods. Again, blue (pink) circles represent nodes included in the (estimated) intervention target \hat{I}_8 (\hat{I}_9); grey circles represent instead nodes which are included in both the two targets. Notice that, because *He & Geng* does not provide target estimates, this information is missing from the corresponding estimated graph. Also notice that its output is not in general a DAG but a *partially* directed graph; accordingly we name it $\hat{\mathcal{G}}_{H\&G}$. Two of these graphs, $\hat{\mathcal{D}}_{NW}$ and $\hat{\mathcal{G}}_{H\&G}$, are equal to our estimate $\hat{\mathcal{D}}$ (Figure 15) in terms of skeleton. In addition, $\hat{\mathcal{D}}$ and $\hat{\mathcal{D}}_{NW}$ also share the same edge orientations, with the only exception of link $p38 \leftarrow JNK$ in $\hat{\mathcal{D}}$ which is reversed in $\hat{\mathcal{D}}_{NW}$. By converse, orientation of most edges is not recovered in $\hat{\mathcal{G}}_{H\&G}$ and only two oriented edges, $PLC \rightarrow PIP2$ and $PIP3 \rightarrow PIP2$, are also present in $\hat{\mathcal{D}}$ and $\hat{\mathcal{D}}_{NW}$. Differently, $\hat{\mathcal{D}}_{IGSP}$ presents a different skeleton, because of the absence of two links, and also differs from $\hat{\mathcal{D}}$ by the orientation of edges $PKA \rightarrow Erk$ and $p38 \rightarrow JNK$.

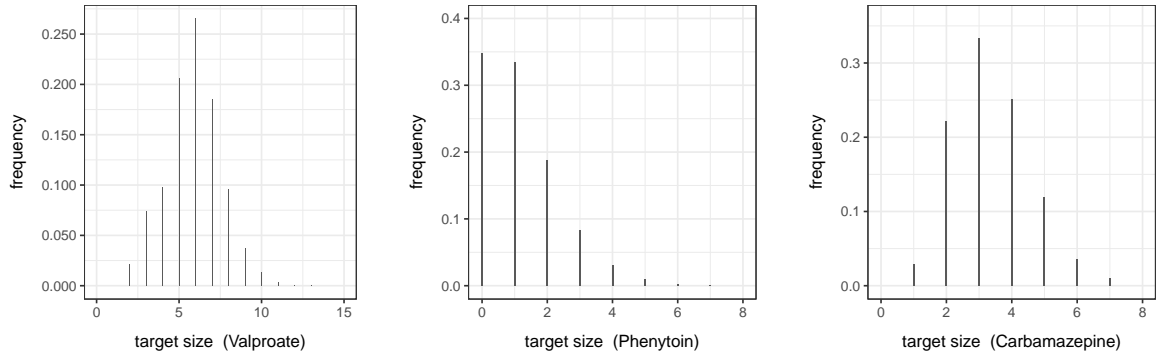


Figure 14: Epilepsy data. Posterior distribution of the number of intervened nodes (size of the intervention target) for each drug-treatment included in the study.

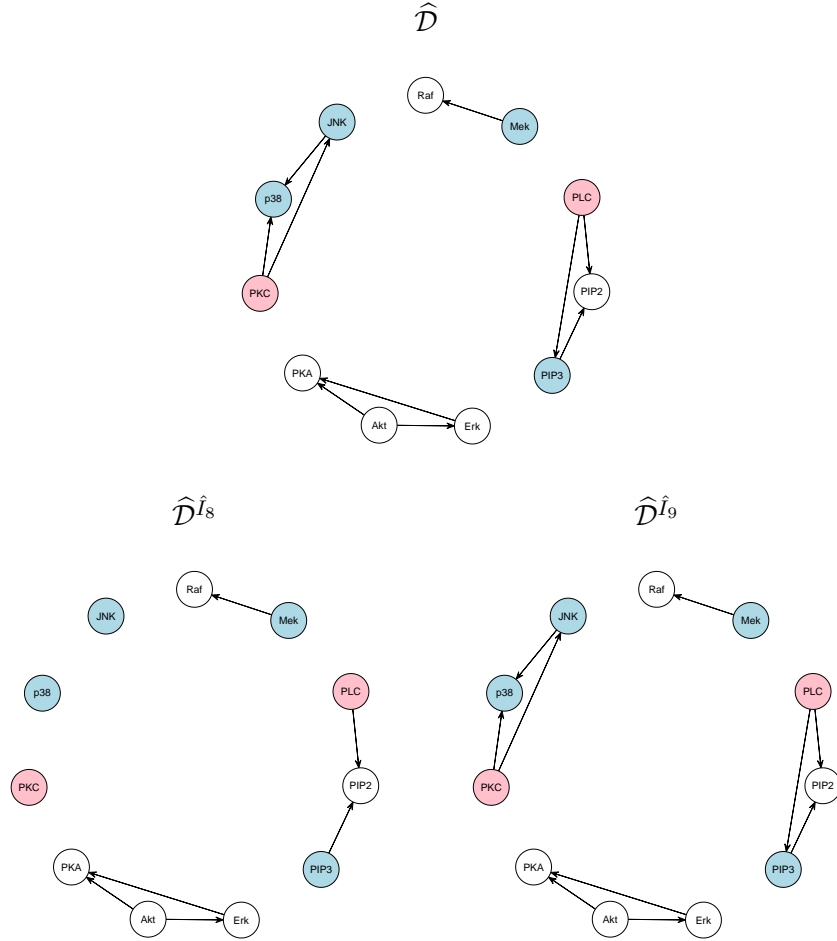


Figure 15: Sachs data. Estimated (median probability) DAG $\hat{\mathcal{D}}$ obtained from our method with blue (pink) circles representing nodes included in targets \hat{I}_8 and \hat{I}_9 ; two corresponding intervention DAGs, $\hat{\mathcal{D}}^{\hat{I}_8}$ and $\hat{\mathcal{D}}^{\hat{I}_9}$.

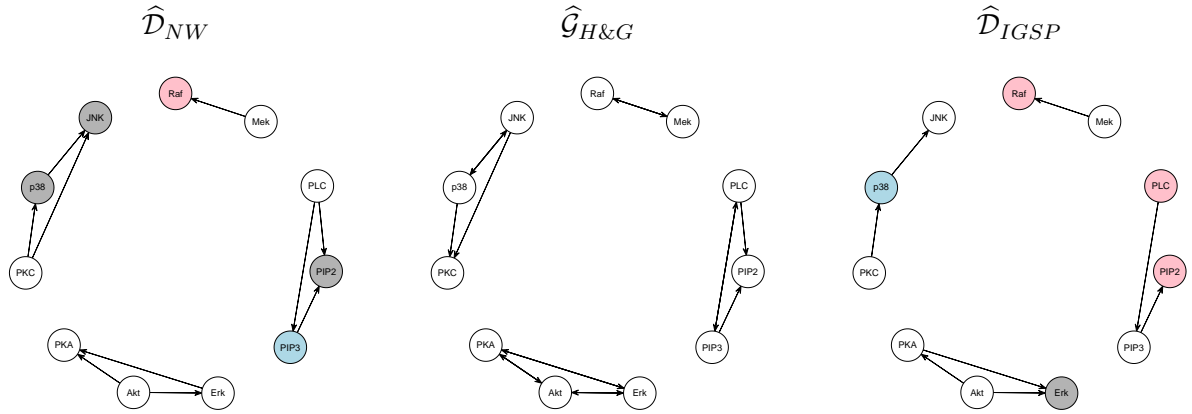


Figure 16: Sachs data. Estimated graphs obtained from *Node-wise* regression, Algorithm 1 of He and Geng (2016) and IGSP method, with blue (pink) circles representing nodes included in targets \hat{I}_8 and \hat{I}_9 (grey nodes included in both).

References

- Altomare, D., G. Consonni, and L. La Rocca (2013). Objective Bayesian search of Gaussian directed acyclic graphical models for ordered variables with non-local priors. *Biometrics* 69(2), 478–487.
- Ben-David, E., T. Li, H. Massam, and B. Rajaratnam (2015). High dimensional Bayesian inference for Gaussian directed acyclic graph models. *arXiv pre-print*.
- Cao, X., K. Khare, and M. Ghosh (2019). Posterior graph selection and estimation consistency for high-dimensional Bayesian DAG models. *The Annals of Statistics* 47(1), 319–348.
- Chickering, D. M. (1995). A transformational characterization of equivalent Bayesian network structures. In *Proceedings of the Eleventh conference on Uncertainty in artificial intelligence*, pp. 87–98. Morgan Kaufmann Publishers Inc.
- Chickering, D. M. (2002). Learning equivalence classes of Bayesian-network structures. *Journal of Machine Learning Research* 2(3), 445–498.
- Foygel, R. and M. Drton (2010). Extended Bayesian information criteria for Gaussian graphical models. In *Advances in Neural Information Processing Systems 23*, pp. 2020–2028.
- Geiger, D. and D. Heckerman (2002). Parameter priors for directed acyclic graphical models and the characterization of several probability distributions. *The Annals of Statistics* 30(5), 1412–1440.

- Ghazal, G. A. and H. Neudecker (2000). On second-order and fourth-order moments of jointly distributed random matrices: a survey. *Linear Algebra and its Applications* 321(1-3), 61–93.
- Han, S. W., G. Chen, M.-S. Cheon, and H. Zhong (2016). Estimation of directed acyclic graphs through two-stage adaptive lasso for gene network inference. *Journal of the American Statistical Association* 111(515), 1004–1019.
- Hauser, A. and P. Bühlmann (2012). Characterization and greedy learning of interventional Markov equivalence classes of directed acyclic graphs. *Journal of Machine Learning Research* 13(79), 2409–2464.
- He, Y. and Z. Geng (2016). Causal network learning from multiple interventions of unknown manipulated targets. *arXiv preprint arXiv:1610.08611*.
- Metropolis, N., A. W. Rosenbluth, M. N. Rosenbluth, A. H. Teller, and E. Teller (1953). Equation of state calculations by fast computing machines. *The journal of Chemical Physics* 21(6), 1087–1092.
- Peluso, S. and G. Consonni (2020). Compatible priors for model selection of high-dimensional Gaussian DAGs. *Electronic Journal of Statistics* 14(2), 4110–4132.
- Scott, J. G. and C. M. Carvalho (2008). Feature-inclusion stochastic search for Gaussian graphical models. *Journal of Computational and Graphical Statistics* 17(4), 790–808.
- Spirtes, P., C. N. Glymour, R. Scheines, and D. Heckerman (2000). *Causation, prediction, and search*. MIT press.
- Squires, C., Y. Wang, and C. Uhler (2020). Permutation-based causal structure learning with unknown intervention targets. In *Conference on Uncertainty in Artificial Intelligence*, pp. 1039–1048. PMLR.
- Tsamardinos, I., L. E. Brown, and C. F. Aliferis (2006). The max-min hill-climbing Bayesian network structure learning algorithm. *Machine learning* 65(1), 31–78.