
Structure Learning with Adaptive Random Neighborhood Informed MCMC (Supplementary Material)

Anonymous Author(s)

Affiliation

Address

email

1 A Derivation of the posterior model probability of DAG model

2 In this section, we will introduce the companion form for the prior on W_γ and derive the marginal
3 likelihood of γ starting from a spike-and-slab prior, as described in Section 2. In addition, under the
4 *Markov property*, we show that the posterior distribution $\pi(\cdot)$ can be factorised into a product form
5 [Koller and Friedman, 2009].

6 Let γ_j be the j -th column of γ , γ_{ij} is the i -th component of γ_j (equivalent to the i, j -th entry
7 of indicator matrix γ). Let $\gamma_{-i,j}$ be γ_j but excluding the i -th component γ_{ij} . Let X_{γ_j} denote the
8 “active” variables in γ_j (those variables such that $\gamma_{ij} = 1$). The variables X_{γ_j} are the parents of X_j
9 implied by γ . Let $W_{\gamma,j}$ be the coefficients of j -th column in W_γ selected by γ_j . In addition, let
10 $\mathbf{X}_i = \{X_{j,i}\}_{j=1}^n$ be a vector of nodes for a fixed i -th observation, with X_j being the j -th node, while
11 let $x_j = \{X_{j,i}\}_{i=1}^N$ be the vector of the N observations corresponding to the j -th node. The matrix
12 x_{γ_j} is the collection of observations on the set of variables defined by γ_j .

13 Under γ , the parent set of X_j ($\text{Pa}(X_j)$) is equivalent to set of nodes collected in γ_j (i.e., X_{γ_j}). The
14 joint prior distribution on W_γ then is

$$p(W_\gamma | \gamma, \{\sigma_j^2\}) = \prod_{j=1}^n p(W_{\gamma,j} | \gamma_j, \sigma_j^2) \quad (\text{S.1})$$

15 where

$$W_{\gamma,j} | \gamma_j, \sigma_j^2 \sim N_{q_j}(0, g\sigma_j^2 I_{q_j}) \quad (\text{S.2})$$

16 and q_j is the the parent size of X_j .

17 **Remark A.1.** The identity matrix I_{q_j} in the prior (S.2) can be replaced by the inverse of the Gram
18 matrix defined as $(x_{\gamma_j}^T x_{\gamma_j})^{-1}$. The resulting prior is then an analog, in the context of DAG sampling,
19 to *g*-prior in the Bayesian variable selection context.

20 Note that in a DAG model, each node X_j is independent of its non-descendants given its parents X_{γ_j} .
21 This property implies a factorisation of the joint likelihood over data \mathcal{D} given by

$$p(\mathcal{D} | W_\gamma, \{\sigma_j^2\}, \gamma) = \prod_{j=1}^n p(X_j | X_{\gamma_j}, W_{\gamma,j}, \sigma_j^2, \gamma_j),$$

22 which translates into

$$p(\mathcal{D} | W_\gamma, \{\sigma_j^2\}, \gamma) = \prod_{j=1}^n N(x_j | x_{\gamma_j}^T W_{\gamma,j}, \sigma_j^2 I_N)$$

in the context of Gaussian graphical models. Due to the factorisation above, it is possible to exchange the integral of W_γ and $\{\sigma_j^2\}$ with the product in j , and integrate out $W_{\gamma,j}$ and σ_j^2 from the joint probabilities

$$p(X_j, W_{\gamma,j}, \sigma_j^2 | X_{\gamma_j}, \gamma_j) = p(X_j | X_{\gamma_j}, W_{\gamma,j}, \sigma_j^2, \gamma_j) p(W_{\gamma,j} | \sigma_j^2, \gamma_j) p(\sigma_j^2) \quad (\text{S.3})$$

for each individual node X_j . In what follows, we will show how to derive $p(X_j | X_{\gamma_j}, \gamma_j)$ from the above, by solving these integrals.

We start from the integration of the coefficients $W_{\gamma,j}$. We can marginalise out $W_{\gamma,j}$ simply due to conjugacy as below

$$\begin{aligned} p(X_j | X_{\gamma_j}, \sigma_j^2, \gamma_j) &= \int p(X_j | X_{\gamma_j}, W_{\gamma,j}, \sigma_j^2, \gamma_j) p(W_{\gamma,j} | \sigma_j^2, \gamma_j) dW_{\gamma,j} \\ &= \int N(x_j | x_{\gamma_j}^T W_{\gamma,j}, \sigma_j^2 I_n) N(W_{\gamma,j} | 0, g \sigma_j^2 I_{q_j}) dW_{\gamma,j} \\ &= (2\pi \sigma_j^2)^{-\frac{N}{2}} g^{-\frac{q_j}{2}} \det(\Sigma_j)^{-\frac{1}{2}} \exp \left\{ -\frac{1}{2\sigma_j^2} (x_j^T x_j - x_j^T x_{\gamma_j} \Sigma_j^{-1} x_{\gamma_j}^T x_j) \right\} \end{aligned}$$

where $\Sigma_j = x_{\gamma_j}^T x_{\gamma_j} + \frac{1}{g} I_{q_j}$. Recall that the prior for σ_j^2 is specified by $p(\sigma_j^2) \propto \sigma_j^{-2}$, and thus we can marginalise out σ_j^2 in a similar way

$$\begin{aligned} p(X_j | X_{\gamma_j}, \gamma_j) &= \int p(X_j | X_{\gamma_j}, \sigma_j^2, \gamma_j) p(\sigma_j^2) d\sigma_j^2 \\ &= \int_{\mathbb{R}^+} (2\pi \sigma_j^2)^{-\frac{N}{2}} g^{-\frac{q_j}{2}} \det(\Sigma_j)^{-\frac{1}{2}} \exp \left\{ -\frac{1}{2\sigma_j^2} (x_j^T x_j - x_j^T x_{\gamma_j} \Sigma_j^{-1} x_{\gamma_j}^T x_j) \right\} \cdot \sigma_j^{-2} d\sigma_j^2 \\ &= (2\pi)^{-\frac{N}{2}} g^{-\frac{q_j}{2}} \det(\Sigma_j)^{-\frac{1}{2}} \Gamma\left(\frac{N}{2}\right) \left[\frac{1}{2} (x_j^T x_j - x_j^T x_{\gamma_j} \Sigma_j^{-1} x_{\gamma_j}^T x_j) \right]^{-\frac{N}{2}} \quad (\text{S.4}) \end{aligned}$$

where $\Sigma_j = x_{\gamma_j}^T x_{\gamma_j} + \frac{1}{g} I_{p_j}$ and $\Gamma(\cdot)$ is the Gamma function. Eventually, putting all the steps above together, we can write down the model's marginal likelihood, conditional on γ only, as follows

$$\begin{aligned} p(\mathcal{D} | \gamma) &= \int \int p(\mathcal{D} | W_\gamma, \{\sigma_j^2\}, \gamma) dW_\gamma d\{\sigma_j^2\} \\ &= \prod_{j=1}^n \int \int p(X_j | X_{\gamma_j}, \beta_j, \gamma_j, \sigma_j^2) dW_{\gamma,j} d\sigma_j^2 \\ &= \prod_{j=1}^n p(X_j | X_{\gamma_j}, \gamma_j) \\ &\propto \prod_{j=1}^n g^{-\frac{q_j}{2}} \det(\Sigma_j)^{-\frac{1}{2}} \left[\frac{1}{2} (x_j^T x_j - x_j^T x_{\gamma_j} \Sigma_j^{-1} x_{\gamma_j}^T x_j) \right]^{-\frac{N}{2}} \quad (\text{S.5}) \end{aligned}$$

where $\Sigma_j = x_{\gamma_j}^T x_{\gamma_j} + \frac{1}{g} I_{p_j}$. From the above derivation, we can also conclude the following property for the posterior distribution $\pi(\cdot)$:

Property 1 (Markov property.). *The posterior distribution of indicator variable can be factorised as the following*

$$\pi(\gamma) \propto \prod_{j=1}^n p(\gamma_j | \mathcal{D}) \times \mathbb{I}\{\mathcal{G}_\gamma \text{ is a DAG}\} \quad (\text{S.6})$$

where $p(\gamma_j | \mathcal{D}) = p(X_j | X_{\gamma_j}, \gamma_j) p^u(\gamma_j)$ and $p^u(\cdot)$ is the unconstrained prior

$$p^u(\gamma_j) = \left(\frac{h}{1-h} \right)^{d_{\gamma_j}}. \quad (\text{S.7})$$

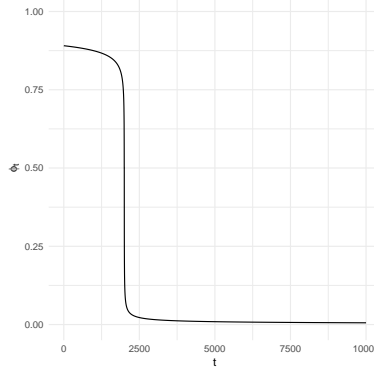


Figure S.1: The plot of diminishing sequence ϕ_t for 10,000 iterations with burn-in period of 2,000 iterations.

39 B The adaptive ϕ_t parameter

40 The decaying parameter ϕ_t is used for the adaptation of the algorithmic tuning parameter η as given
 41 in (9), and controls the trade-off between the pre-processing approximation $\tilde{\pi}$ and the ergodic average
 42 $\hat{\pi}^{(t)}$. Intuitively, one should rely more on the approximation $\tilde{\pi}$ during the burn-in period, since the
 43 ergodic average $\hat{\pi}^{(t)}$ has not fully converged. After the burn-in period though, the proposal should
 44 rely more heavily on the ergodic average. Based on the arguments, we use the following sequence of
 45 ϕ_t to adapt η :

$$\phi_t = \begin{cases} 1 - \frac{1}{2} \left(\frac{1}{N_b - t + 1} \right)^{0.2} & \text{if } t \leq N_b \\ \frac{1}{2} \left(\frac{1}{t - N_b} \right)^{0.5} & \text{if } t > N_b. \end{cases} \quad (\text{S.8})$$

46 In particular, we have $\phi_t = 1/2$ when $t = N_b$. An example plot of decaying ϕ_t is given in Figure
 47 S.1, where we considered 10,000 iterations, with the first 2,000 iterations discarded as burn-in.

48 In addition, in the following lemma we derive the decaying rate of ϕ_t , and thus of $\eta^{(t)}$ as well:

49 **Lemma S1.** *After the burn-in period, the diminishing rate of ϕ_t defined by (S.8) is $\mathcal{O}(t^{-1})$. And thus*
 50 *the diminishing rate of tuning parameter $\hat{\eta}^{(t)}$ in (9) is $\mathcal{O}(t^{-1})$.*

51 *Proof.* Suppose $t > N_b$, we start from the definition of η_t and consider the diminishing rate

$$\begin{aligned} |\phi_{t+1} - \phi_t| &= \left| \frac{1}{2} \left(\frac{1}{t+1-N_b} \right)^{0.5} - \frac{1}{2} \left(\frac{1}{t-N_b} \right)^{0.5} \right| \\ &\leq \frac{1}{2} \left(\frac{(t+1-N_b)^{0.5} - (t-N_b)^{0.5}}{(t+1-N_b)^{0.5}(t-N_b)^{0.5}} \right) \\ &\leq \frac{1}{2} \frac{1}{t-N_b} \\ &= \mathcal{O}(t^{-1}) \end{aligned}$$

52 as required.

53 We can also show that the diminishing rate of the ergodic average $\hat{\pi}_{ij}^t$ has the same decaying rate

$$\begin{aligned} |\hat{\pi}_{ij}^{(t+1)} - \hat{\pi}_{ij}^{(t)}| &= \left| \frac{t}{t+1} \hat{\pi}_{ij}^{(t)} + \frac{1}{t+1} \gamma_{ij}^{(t+1)} - \hat{\pi}_{ij}^{(t)} \right| \\ &\leq \frac{1}{t+1} \hat{\eta}_{ij}^{(t)} + \frac{1}{t+1} \gamma_{ij}^{(t+1)} \\ &\leq \frac{2}{t+1} = \mathcal{O}(t^{-1}) \end{aligned}$$

54 as both $\gamma_{ij}^{(t)}$ and $\gamma_{ij}^{(t+1)}$ are bounded between 0 and 1.

55 We complete the proof by showing that also the tuning parameter $\hat{\eta}_{ij}^{(t)}$ has the same decaying rate

$$\begin{aligned} |\hat{\eta}_{ij}^{(t+1)} - \hat{\eta}_{ij}^{(t)}| &= |(\phi_t - \phi_{t+1})\tilde{\pi}_{ij} + (\phi_{t+1} - \phi_t)\hat{\pi}_{ij}^{(t+1)}| \\ &\leq |\phi_t - \phi_{t+1}|\tilde{\pi}_{ij} + |\phi_{t+1} - \phi_t|\hat{\pi}_{ij}^{(t+1)} \\ &= \mathcal{O}(t^{-1}) + \mathcal{O}(t^{-1}) = \mathcal{O}(t^{-1}) \end{aligned}$$

56 as $\tilde{\pi}_{ij}$ is between 0 and 1.

57

□

58 Given the above lemma, and combining Lemmas 1 and 2 in Section 4.2.3 of Liang et al. [2022], we
59 can establish the similar ergodicity and strong law of large numbers as in Theorem 2 of Liang et al.
60 [2022].

61 C The warm-start approximation $\tilde{\pi}_{ij}$

62 The main idea of how to compute the warm-start approximation $\tilde{\pi}_{ij}$ is explained in Section 3.2 of the
63 main body. In this section, we will fill in with more details on calculation of $\tilde{\pi}_{ij}$.

64 Consider the extended collection of permissible parent set \mathbf{h}_+^j , formally defined as

$$\mathbf{h}_+^j = \mathbf{h}^j \cup \{m \cup \{j\} \mid m \in \mathbf{h}^j, j \in (h^j)^c \setminus \{j\}\}. \quad (\text{S.9})$$

65 The extended set \mathbf{h}_+^j includes the parent set implied by \mathbf{h}^j , but also the additional parent sets obtained
66 by including at least one more parent not included in h^j . Trivially thus the extended set \mathbf{h}_+^j contains
67 more permissible parents and minimises the risk of leaving out important parents when reducing the
68 search space to a skeleton graph for scalability reasons [Kuipers et al., 2022].

69 We will show now how we calculate the marginal probability π_{ij}^u under the unconstrained posterior
70 π^u and the extended sets $\{\mathbf{h}_+^j\}$. We define $\Gamma_j \subset \{0, 1\}^p$ to be the sub-space of Γ such that $\gamma_{jj} = 0$ if
71 $\gamma_j \in \Gamma_j$. Given Markov property, thanks to which the columns of γ can be factorised, we can define
72 the unconstrained posterior distribution as

$$\pi^u(\gamma) \propto \prod_{j=1}^n p(\gamma_j | \mathcal{D}), \quad (\text{S.10})$$

73 the marginal distribution of γ_j under the unconstrained posterior π^u is

$$\pi^u(\gamma_j) = \frac{p(\gamma_j | \mathcal{D})}{\sum_{\gamma_j \in \Gamma_j} p(\gamma_j | \mathcal{D})}. \quad (\text{S.11})$$

74 However, the sub-space Γ_j contains 2^{n-1} candidate parent sets and the full enumeration is computa-
75 tionally infeasible when n is big. So we replace the subspace by the extended collection of parent
76 sets \mathbf{h}_+^j which contains much less models than Γ_j . So this approximate the above as

$$\pi^u(\gamma_j) \approx \frac{p(\gamma_j | \mathcal{D})}{\sum_{m \in \mathbf{h}_+^j} p(\rho(m) | \mathcal{D})}. \quad (\text{S.12})$$

77 where $\rho : [p] \rightarrow \{0, 1\}^p$ is a one-to-one mapping that converts the parent set m to γ_j which lies in
78 the binary space Γ_j . Then we approximate the PEP with the following formula:

$$\pi_{ij}^u \stackrel{\text{def}}{=} \pi^u(\gamma_{ij} = 1) \approx \sum_{m \in \mathbf{h}_+^j : i \in m} \pi^u(\rho(m)). \quad (\text{S.13})$$

79 The set $\{m \in \mathbf{h}_+^j \mid i \in m\}$ is equivalent to the set $\{\gamma_j \in \Gamma_j \mid \gamma_{ij} = 1\}$ of which X_i must be a parent
80 of X_j from all these permissible parent sets. We also guarantee $\pi_{jj}^u = 0$ since j is not one of the
81 permissible parents of j and not included in \mathbf{h}_+^j .

After the approximations, the π_{ij}^u are calculated and stored, then we can compute the warm-start approximations $\tilde{\pi}_{ij}$, followed by the equations (11) and (12) described in the main body. The calculations of $\tilde{\pi}_{ij}$ are equivalent to

$$\tilde{\pi}_{ij} = \frac{\pi_{ij}^u(1 - \pi_{ji}^u)}{\pi_{ij}^u(1 - \pi_{ji}^u) + (1 - \pi_{ij}^u)\pi_{ji}^u + (1 - \pi_{ij}^u)(1 - \pi_{ji}^u)}. \quad (\text{S.14})$$

The last step combines $\tilde{\pi}_{ij}$ with the ergodic average $\hat{\pi}_{ij}$ to form the algorithmic tuning parameters $\hat{\eta}_{ij}$ used in the PARNI-DAG proposal.

D The full PARNI-DAG proposal

We first recall the random neighbourhood construction specified in Section 3.1 of the main body. The neighbourhood indicator $k \in \mathcal{K} = \{0, 1\}^{n \times n}$ indicates the positions that can be potentially flipped in the current DAG γ , and k follows the conditional distribution as given in (3) of the main body. The neighbourhood $\mathcal{N}(\gamma, k)$ given in (4) of the main body consists of 2^{d_k} DAG models, and the full enumeration over it is computationally efficient (and sometimes infeasible). For this reason we consider a pointwise implementation consisting in smaller sub-proposals. The first step in the pointwise implementation is to convert k into a set of variables $K = \{K_r\}_{r=1}^R$, which can take two categories:

- if $k_{ij} = 1$ but $k_{ji} = 0$, K_r only consists of one position $K_r = \{(i, j)\}$;
- if $k_{ij} = 1$ and $k_{ji} = 1$, K_r then contains two positions $K_r = \{(i, j), (j, i)\}$.

The order of K_r is random. To fully specify the sub-neighbourhoods $\mathcal{N}(\gamma(r-1), K_r)$ at each time r , we first define a new mapping that “flips” the components of γ according to the set K ,

$$\gamma' = \text{flip}(\gamma, K), \quad (\text{S.15})$$

with $\gamma'_{ij} = 1 - \gamma_{ij}$ if $(i, j) \in K$ and $\gamma'_{ij} = \gamma_{ij}$ otherwise. The resulting neighbourhoods $\mathcal{N}(\gamma(r-1), K_r)$ are defined as follows:

- If K only contains one position (i, j) , then

$$\mathcal{N}(\gamma(r-1), K_r) = \{\gamma(r-1), \text{flip}(\gamma(r-1), K_r)\}. \quad (\text{S.16})$$

- If K contains two positions (i, j) and (j, i) , we construct a neighbourhood with reversal move given by

$$\mathcal{N}(\gamma(r-1), K_r) = \{\gamma(r-1), \text{flip}(\gamma(r-1), (i, j)), \text{flip}(\gamma(r-1), (j, i)), \text{flip}(\gamma(r-1), K_r)\}. \quad (\text{S.17})$$

The next components in PARNI-DAG are the normalising constants of the sub-proposals q_{g, K_r} and of the reversal sub-proposals q_{g, K'_r} , used to facilitate the calculations of Metropolis-Hastings acceptance probability. As mentioned in the main body, we construct the auxiliary variables K' in the reversal move where K' consists of the same elements of K but with different order. One property of such a design is that the intermediate DAGs in the reversal move are identical to the those ones in the proposal move but only with opposite order. We then can re-use the posterior model probabilities calculated during the proposal move to compute the reversal probabilities directly. In addition, the Proposition 3 of Liang et al. [2022] implies that the acceptance probability in (7) can be expressed in terms of the product of the normalising constants. Given the normalising constants

$$Z(r) = \sum_{\gamma^* \in \mathcal{N}(\gamma(r-1), K_r)} g \left(\frac{\pi(\gamma^*)p(K_r|\gamma^*)}{\pi(\gamma(r-1))p(K_r|\gamma(r-1))} \right) \quad (\text{S.18})$$

$$Z(R-r+1)' = \sum_{\gamma^* \in \mathcal{N}(\gamma(r), K_r)} g \left(\frac{\pi(\gamma^*)p(K_r|\gamma^*)}{\pi(\gamma(r))p(K_r|\gamma(r))} \right), \quad (\text{S.19})$$

the MH acceptance probability is then simplified as follows

$$\alpha(\gamma, \gamma') = \min \left\{ 1, \prod_{r=1}^R \frac{Z(r)}{Z'(R-r+1)} \right\}. \quad (\text{S.20})$$

To summarise the PARNI-DAG proposal, we now give its full algorithmic pseudo-code is given in **Algorithm 1**.

Algorithm 1 The PARNI-DAG proposal

Compute $\{\tilde{\pi}_{ij}\}$ as described in **Section 3.2** and **Appendix C**
Initialise $\gamma^{(1)}$, $\omega^{(1)}$ and $\hat{\eta}^{(1)} = \{\hat{\eta}_{ij}^{(1)}\}$
for $t = 1$ to $t = T$ **do**
 Sample $k \sim p_{\hat{\eta}^{(t)}}(\cdot | \gamma^{(t)})$ as in (3)
 Construct $K = \{K_r\}$ as described in **Appendix D**
 Set $\gamma(0) = \gamma^{(i)}$ and $\tilde{N} = 0$
 for $r = 1$ to $r = R$ **do**
 if K_r has one element **then**
 Construct $\mathcal{N}(\gamma(r-1), K_r)$ as in (S.16)
 else
 Construct $\mathcal{N}(\gamma(r-1), K_r)$ as in (S.17)
 end if
 Sample $U_1 \sim \text{Unif}(0, 1)$
 if $U_1 < \omega^{(t)}$ **then**
 Sample $\gamma(r) \sim q_{g, K_r}(\gamma(r-1), \cdot)$ as in (5)
 Calculate $Z(r)$ as in (S.18)
 Calculate $Z'(R-r+1)$ as in (S.19)
 else
 Set $\gamma(r) = \gamma(r-1)$, $Z(r) = Z(R-r+1)' = 1$.
 end if
 end for
 Set $\gamma' = \gamma(R)$, sample $U_2 \sim \text{Unif}(0, 1)$ and compute $\alpha(\gamma^{(i)}, \gamma')$ as in (S.20);
 if $U_2 < \alpha(\gamma^{(i)}, \gamma')$ **then**
 Set $\gamma^{(t+1)} = \gamma'$
 else
 Set $\gamma^{(t+1)} = \gamma^{(t)}$
 end if
 Update $\hat{\eta}^{(t+1)} = \{\hat{\eta}_{ij}^{(t+1)}\}$ according to (9).
 Update $\omega^{(t+1)}$ according to (13).
end for

117 E Implementation of other MCMC algorithms

118 In this section, we will briefly describe the Add-Delete-Reverse (ADR) proposal, the order MCMC
119 and the partition MCMC implemented in the experimental Section 4.

120 E.1 The Add-Delete-Reverse proposal

121 The ADR proposal is another example of random neighbourhood proposal. After the random
122 neighbourhood is generated, unlike the PARNI-DAG proposal, ADR uniformly proposes a model
123 within that neighbourhood.

124 In ADR, there are three possible moves, “Addition”, “Deletion” and “Reversal”, which defines the
 125 space \mathcal{K} of indicator k . The indicator k is then uniformly drawn from one of these three possible
 126 moves. We construct the neighbourhood $\mathcal{N}(\gamma, k)$ according to the current DAG γ and indicator k as
 127 the following

$$\begin{aligned}\mathcal{N}(\gamma, \text{“Addition”}) &= \{\gamma^* \in \Gamma \mid d_H(\gamma, \gamma^*) = 1, d_{\gamma^*} = d_{\gamma} + 1, \gamma_{ji}^* = 0 \text{ if } \gamma_{ij} = 1\} \\ \mathcal{N}(\gamma, \text{“Deletion”}) &= \{\gamma^* \in \Gamma \mid d_H(\gamma, \gamma^*) = 1, d_{\gamma^*} = d_{\gamma} - 1\} \\ \mathcal{N}(\gamma, \text{“Reversal”}) &= \{\gamma^* \in \Gamma \mid d_H(\gamma, \gamma^*) = 2, d_{\gamma^*} = d_{\gamma}, \exists i, j: \gamma_{ij} = \gamma_{ji}^* = 0 \text{ if } \gamma_{ji} = \gamma_{ij}^* = 1\}.\end{aligned}$$

128 where $d_H(\cdot, \cdot)$ denotes the Hamming distance and d_{γ} is the number of edges induced by γ . A new
 129 DAG model $\gamma' \in \Gamma$ is proposed according to the uniform density:

$$q_k^{\text{ADR}}(\gamma, \gamma') = \begin{cases} \frac{1}{|\mathcal{N}(\gamma, k)|}, & \text{if } \gamma \in \mathcal{N}(\gamma, k), \\ 0, & \text{otherwise.} \end{cases} \quad (\text{S.21})$$

130 As ADR is not guaranteed (unlike PARNI-DAG) to return a DAG, we carry out acyclicity diagnostics
 131 check when γ' is sampled. If γ does not lead to a valid DAG, we will accept the current state γ as the
 132 next state of the Markov chain. If γ' leads to a valid DAG, to preserve π -reversibility, we consider k'
 133 in the reverse move where k' is the opposite move to k (e.g. the addition is opposite to the deletion
 134 move, and vice versa). To complete the ADR proposal, the new DAG γ' is accepted as the next state
 135 of the Markov chain with probability

$$\begin{aligned}\alpha_k(\gamma, \gamma') &= \min \left\{ 1, \frac{\pi(\gamma')q_{k'}(\gamma', \gamma)}{\pi(\gamma)q_k(\gamma, \gamma')} \right\} \\ &= \min \left\{ 1, \frac{\pi(\gamma')|\mathcal{N}(\gamma, k)|}{\pi(\gamma)|\mathcal{N}(\gamma', k')|} \right\}.\end{aligned}$$

136 E.2 Order MCMC and Partition MCMC

137 Order and Partition MCMC samplers rely on DAG scoring measures where $p(\mathcal{G}|\mathcal{D}) =$
 138 $\prod_{j=1}^n S(X_j, \text{Pa}(X_j)|\mathcal{D})$, and on restricting the search space of DAGs to some lower capacity set.
 139 Friedman and Koller [2003] first proposed order MCMC, which operates in the space of orders
 140 (\mathbf{X}, \prec) instead of the DAGs one. This search space is smoother, so that order MCMC generally
 141 achieves faster convergence with respect to classic neighborhood MCMC. The method considers
 142 permutations (j_1, \dots, j_n) of nodes that imply a linear order $j_1 \prec \dots \prec j_n$. Each linear order is
 143 then given a score $R(\prec|\mathcal{D})$ equal to the sum of scores of all the compatible DAGs with the order
 144 $\prec \in (\mathbf{X}, \prec)$. To restrict the space capacity even more, sum and product operations can be swapped
 145 so that the order’s score is:

$$\begin{aligned}R(\prec|\mathcal{D}) &= \sum_{\mathcal{G} \in \prec} p(\mathcal{G}|\mathcal{D}) \propto \sum_{\mathcal{G} \in \prec} \prod_{j=1}^n S(X_j, \text{Pa}(X_j)|\mathcal{D}) \\ &\propto \prod_{j=1}^n \sum_{\text{Pa}(X_j) \in \prec} S(X_j, \text{Pa}(X_j)|\mathcal{D}),\end{aligned}$$

146 where $S(\cdot)$ is a DAG scoring function, such as the Bayesian Gaussian equivalent (BGe) score
 147 for continuous variables [Geiger and Heckerman, 2002]. Constructing a MCMC in the space of
 148 orders means moving from order \prec to order \prec' with MH acceptance probability $\alpha_{\prec}(\prec, \prec') =$
 149 $\min \left\{ 1, \frac{R(\prec'|\mathcal{D})}{R(\prec|\mathcal{D})} \right\}$, according to following moves: a) “Local Swap” swaps adjacent nodes in \prec ; b)
 150 “Global Swap” swaps random nodes in \prec ; c) “Relocate” picks a node and replaces it in all possible
 151 positions in \prec , while keeping the rest of the order unchanged. Order MCMC is efficient, but has two
 152 main drawbacks compared to neighborhood-like structure MCMC: i) it outputs a sample of orders,
 153 not a DAG directly; sample of DAGs can be obtained by sampling, for each node independently, a
 154 parent set compatible with the relative order \prec , thanks to score decomposability (then get the DAG
 155 with maximum score in the sample as MAP estimate); ii) Order MCMC places a non-uniform prior
 156 on some DAGs by over-representing them in more than one order \prec , which causes the sample to be
 157 biased [Ellis and Wong, 2008].

158 In the attempt to correct sampling bias, Kuipers and Moffa [2017] proposes Partition MCMC, which
 159 operates in the space of ordered partitions instead. A labelled partition Λ is a pair of order \prec and

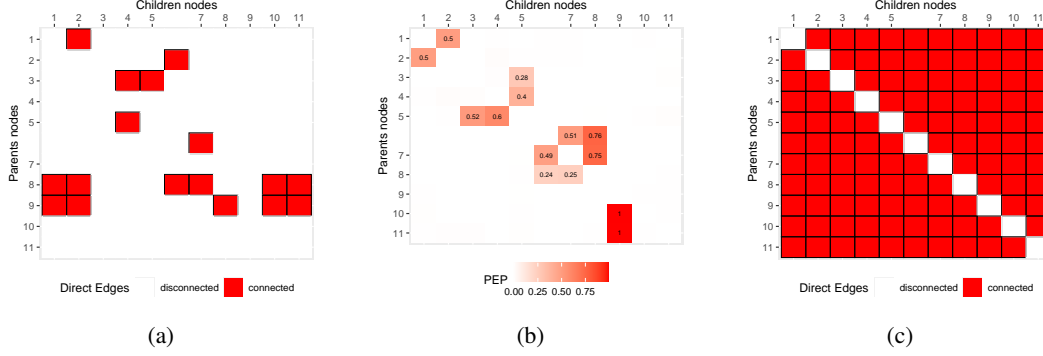


Figure S.2: Protein dataset: (a) The ground-truth DAG. (b) The true PEPs estimated by a long chain of partition MCMC algorithm. (c) The full skeleton used in Section 5.1.

partitioning vector $\delta = (\delta_1, \dots, \delta_p)$ where $p < n$, and $\sum_{j=1}^p \delta_j = n$. Vector δ divides the order into p parts v_1, \dots, v_p such that v_1 contains the first δ_1 nodes of the order \prec and so on. A DAG is compatible with partition $\Lambda = (\prec, \delta)$ if for every node $X_j \in v_j$: i) X_j has at least one parent in part v_{j+1} ; ii) all $\text{Pa}(X_j)$ belong to parts with indices higher than j ; iii) $\text{Pa}(X_j) = \emptyset$ only if $j = p$. The possible moves of partition MCMC are: a) swap 2 nodes from different parts; b) swap 2 nodes from adjacent parts; c) split or join two parts; d) move a node in another existing part or create a new part with that node. Partition MCMC is unbiased in terms of DAG sampling; however it is extremely slow as it has very high computational complexity, and therefore can only be used to sample DAGs with very few nodes [Kuipers et al., 2022].

The orderMCMC and partitionMCMC schemes in the experiments are implemented via the BiDAG¹ package [Suter et al., 2023] (available from **CRAN**) written in **rcpp** [Eddelbuettel and François, 2011].

F Additional experimental results

F.1 Addition to Section 5.1

Protein data. We provide additional information on the real-world protein-signalling dataset ($n = 11$, $N = 853$) studied in **Section 4.1**. The ground-truth DAG constructed through expert knowledge is provided in Figure S.2 (a). The ground-truth estimate of PEP estimated by running a long chain of the partition MCMC algorithm is provided in Figure S.2 (b). The full skeleton used to implement different MCMC schemes is provided in S.2 (c).

gsim100 data. We also provide additional information on the simulated **gsim100** dataset ($n = 100$, $N = 100$) studied in **Section 4.1**, and featured in Suter et al. [2023]. The ground-truth DAG, containing 161 true edges, used to simulate data is presented in Figure S.3 (a). The ground-truth PEP estimates obtained from a long chain run of the PARNI-DAG proposal are provided in S.3 (b). The skeleton learned via the PC-algorithm (with a significance rate of 5%) \mathcal{H}_{PC} and the skeleton learned from the iterative MCMC $\mathcal{H}_{\text{iter}}$ are provided in Figure S.3 (c) and (d) respectively.

F.2 Additional studies on the effects of tuning parameter of the PC-algorithm

In addition to **Section 4.1** results, we perform another batch of numerical studies to study the PARNI-DAG proposal run over different starting skeletons, by comparing their convergence behavior. The skeletons are mainly obtained by running the PC-algorithm with different levels of significance. We consider 4 different values of significance level, 0, 0.01, 0.05 and 0.1. The resulting skeletons are provided in Figure S.4 (The skeleton with significance level of 0.05 is provided in Figure S.3 (c)). As the significance level of the PC-algorithm increases, the resulting skeleton includes more edges

¹More information on: <https://cran.r-project.org/web/packages/BiDAG/BiDAG.pdf>.

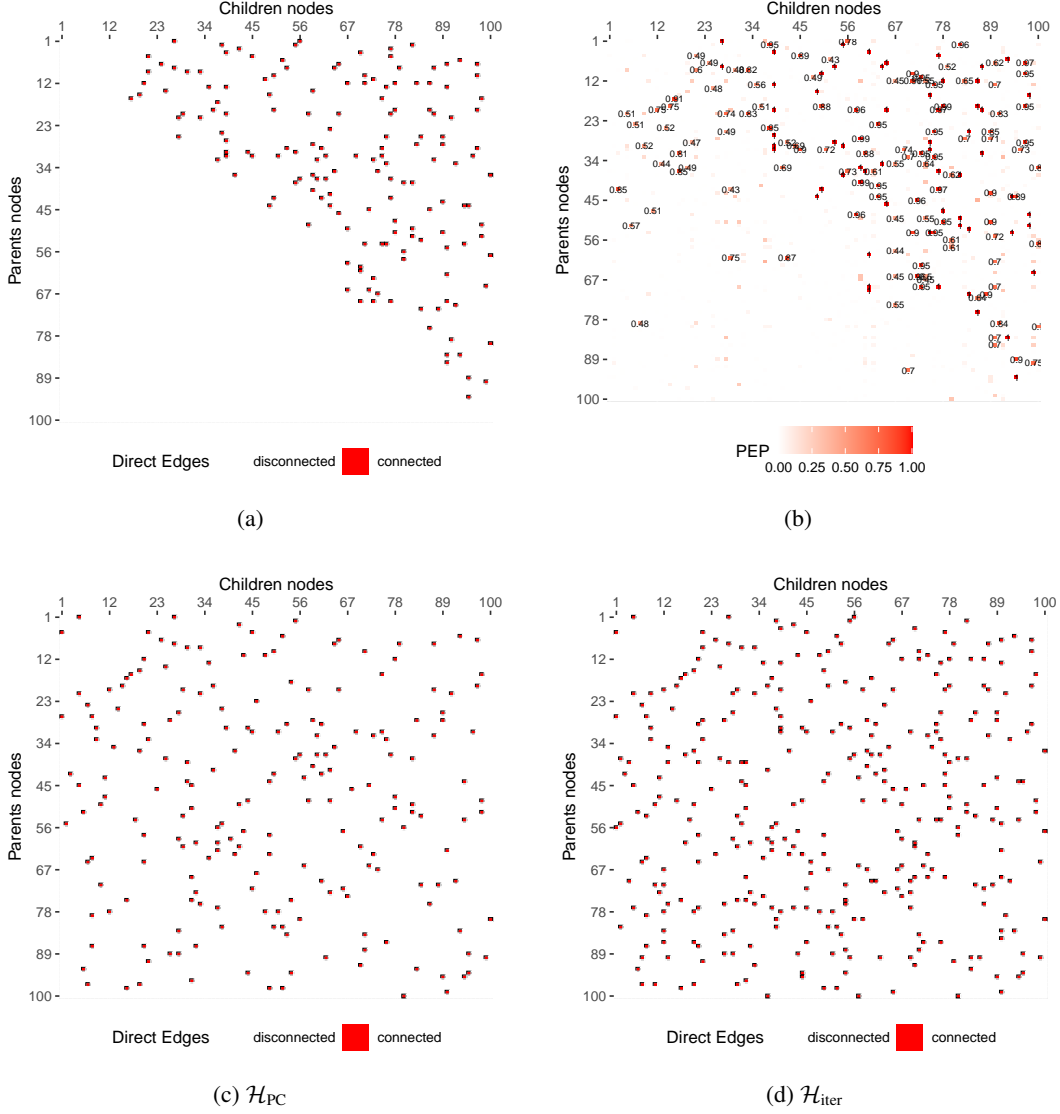


Figure S.3: gsim100 dataset: (a) The ground-truth DAG. (b) The true PEPs estimated by a long chain of the PARNI-DAG proposal. (c) \mathcal{H}_{PC} with significance level of 0.05. (d) \mathcal{H}_{iter} .

whilst the skeleton is empty when significance level is 0. In addition, we also compare these results to the skeleton from iteration MCMC as in **Section 4.1**.

Figure S.5 presents the trace plots of the log posterior model probabilities from the PARNI-DAG proposal under different skeletons. Using the skeleton \mathcal{H}_{PC} with significance level of 0 trivially leads to the slowest convergence rate, while the skeleton \mathcal{H}_{iter} (which includes more possible parent sets) results in the fastest convergence. The behaviour of the chains is similar. We can draw similar conclusions from the median MSEs presented in Table 1. The skeleton \mathcal{H}_{iter} leads to the minimum MSE, while \mathcal{H}_{PC} with significance level of 0 again trivially results in the largest MSE. The MSEs from other three options of significance levels are not significant different between each other, so using any one of them will give similar level of accuracy.

F.3 Additional studies on DAG learning performance

In addition to the results provided in Section 4.2 of the main body, in a similar way we compare PARNI-DAG with ADR and the same other methods in terms of accuracy in recovering the underlying DAG structure on another example. We consider a randomly generated DAG structure featuring

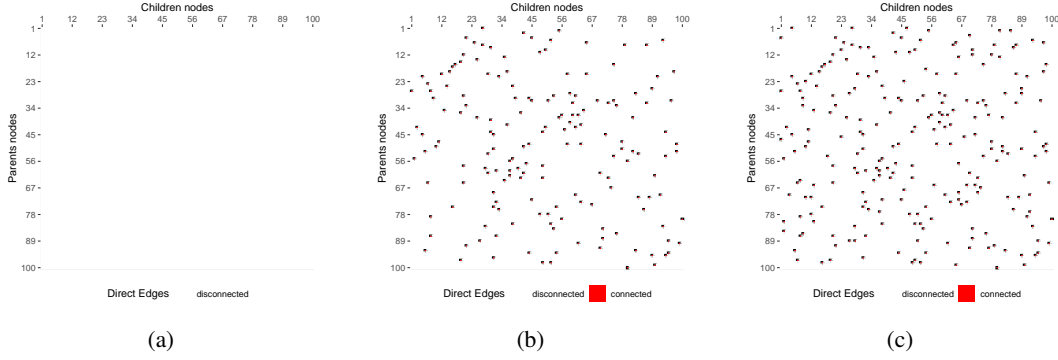


Figure S.4: gsim100 dataset: (a) \mathcal{H}_{PC} with significance level of 0. (b) \mathcal{H}_{PC} with significance level of 0.01. (c) \mathcal{H}_{PC} with significance level of 0.1

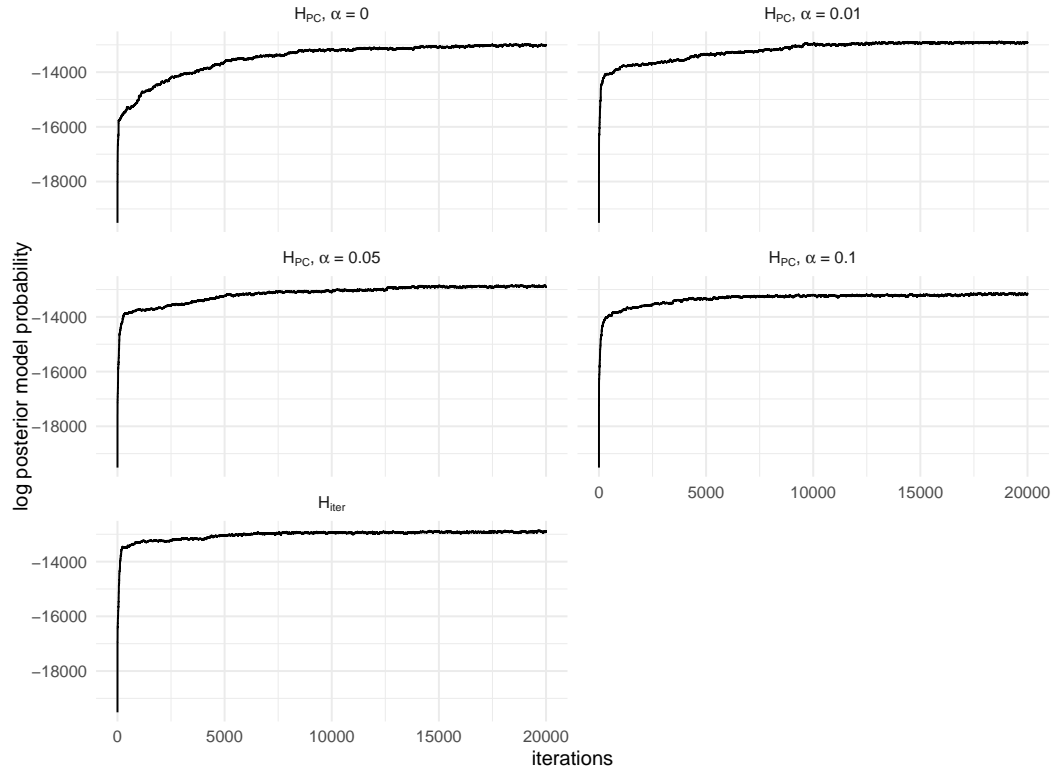


Figure S.5: gsim100 dataset: trace plots of log posterior model probabilities for the PARNI-DAG proposal under different skeletons.

206 50 nodes and 119 edges. Given this DAG structure, we simulate Gaussian data for three different
 207 sample sizes $N \in \{50, 100, 150\}$.

208 To recall, the models compared are: i) PC algorithm [Spirtes et al., 2000]; ii) GES algorithm
 209 [Chickering, 1996, 2002]; iii) Order MCMC [Friedman and Koller, 2003]; iv) Iterative MCMC
 210 [Kuipers et al., 2022]; v) Partition MCMC [Kuipers and Moffa, 2017]; vi) ADR; vii) PARNI-DAG.

211 Models' performance is compared according to same metric, the *Structural Hamming Distance*
 212 $d_H(\hat{\gamma}, \gamma)$ between the estimated and the true DAG. The experiments are replicated for 10 times,
 213 for each of the sample sizes above. Results on the distributions of SHD are depicted in Figure
 214 S.6. We can see how Iterative MCMC and PARNI-DAG are the best competing models out of
 215 this example. Note that this is a medium size graph, and while we have argued in the main body

| skeleton | significance rate (α) | MSE |
|-----------------------------|--------------------------------|------|
| \mathcal{H}_{PC} | 0 | 2.40 |
| | 0.01 | 1.92 |
| | 0.05 | 1.95 |
| | 0.1 | 1.95 |
| $\mathcal{H}_{\text{iter}}$ | - | 1.85 |

Table 1: gsim100 dataset: time-normalised median MSE ($\times 10^{-7}$) on estimating posterior edge probabilities. (*Lower is better.*)

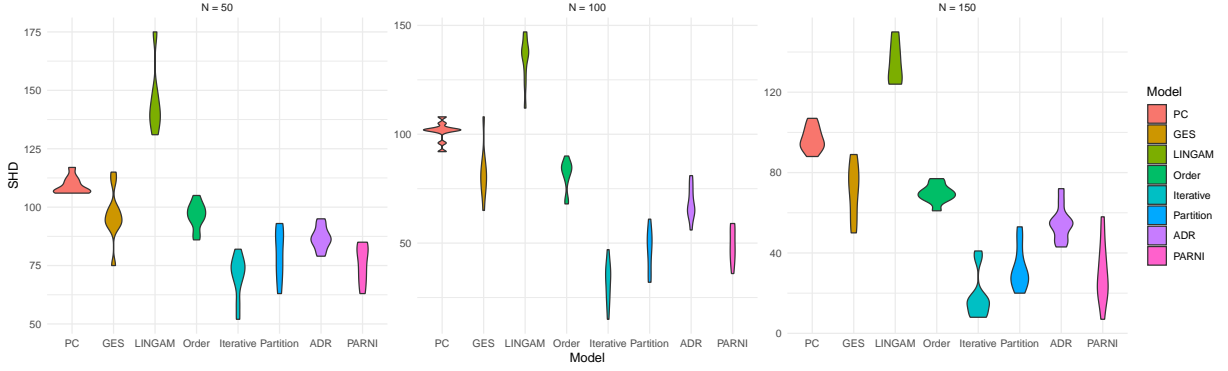


Figure S.6: Distribution of SHD of each compared model, for different sample sizes $N \in \{50, 100, 150\}$, over 10 replications.

216 that PARNI-DAG particularly excels in settings characterized by larger size graphs ($|\mathcal{V}| \geq N$), it
217 is nonetheless competitive with the best performing method, Iterative MCMC here, also on lower
218 dimensional DAGs.

References

- David Maxwell Chickering. Learning Bayesian networks is NP-complete. *Learning from data: Artificial intelligence and statistics V*, pages 121–130, 1996.
- David Maxwell Chickering. Optimal structure identification with greedy search. *Journal of machine learning research*, 3(Nov):507–554, 2002.
- Dirk Eddelbuettel and Romain François. Rcpp: Seamless R and C++ Integration. *Journal of Statistical Software*, 40(8):1–18, 2011. doi: 10.18637/jss.v040.i08.
- Byron Ellis and Wing Hung Wong. Learning causal Bayesian network structures from experimental data. *Journal of the American Statistical Association*, 103(482):778–789, 2008.
- Nir Friedman and Daphne Koller. Being Bayesian about network structure. A Bayesian approach to structure discovery in Bayesian networks. *Machine learning*, 50:95–125, 2003.
- Dan Geiger and David Heckerman. Parameter priors for directed acyclic graphical models and the characterization of several probability distributions. *The Annals of Statistics*, 30(5):1412–1440, 2002.
- Daphne Koller and Nir Friedman. *Probabilistic graphical models: principles and techniques*. MIT press, 2009.
- Jack Kuipers and Giusi Moffa. Partition MCMC for Inference on Acyclic Digraphs. *Journal of the American Statistical Association*, 112(517):282–299, 2017.
- Jack Kuipers, Polina Suter, and Giusi Moffa. Efficient Sampling and Structure Learning of Bayesian Networks. *Journal of Computational and Graphical Statistics*, 31(3):639–650, 2022.
- Xitong Liang, Samuel Livingstone, and Jim Griffin. Adaptive random neighbourhood informed Markov chain Monte Carlo for high-dimensional Bayesian variable selection. *Statistics and Computing*, 32(5):84, 2022.
- Peter Spirtes, Clark N Glymour, Richard Scheines, and David Heckerman. *Causation, prediction, and search*. MIT press, 2000.
- Polina Suter, Jack Kuipers, Giusi Moffa, and Niko Beerenwinkel. Bayesian Structure Learning and Sampling of Bayesian Networks with the R Package BiDAG. *Journal of Statistical Software*, 105(9):1–31, 2023.