

Causal Inference and Machine Learning

Session 7: Bayesian causal discovery and causal effect estimation

Alessandro Mascaro

February 18th, 2025

Barcelona School of Economics,
Master's degree in Data Science Methodology

Bayesian inference and model selection in a few slides

Bayesian inference in one slide

- Suppose we observe a (n, p) data matrix of observations \mathbf{X} , containing i.i.d. samples of a random vector $X := (X_1, \dots, X_p)$, and that X follows parametric distribution with unknown parameter Θ ;
- In the Bayesian setting, Θ is a random variable and inference on it is done using Bayes theorem to derive the **posterior distribution**

$$p(\Theta \mid \mathbf{X}) \propto p(\mathbf{X} \mid \Theta) p(\Theta)$$

where $p(\Theta)$ is the **prior distribution**, specified by the statistician;

- Sometimes the posterior has a known form, sometimes it has to be approximated, mainly through **sampling methods**;

Bayesian Model Selection in one slide

- Bayesian inference is very flexible and it can be seamlessly used to do **model selection**: one just specifies a prior on the model M_k and its parameters Θ_k and derives the joint posterior distribution:

$$p(M_k, \Theta_k | \mathbf{X}) \propto p(\mathbf{X} | \Theta_k, M_k) p(\Theta_k | M_k) p(M_k)$$

- Usually, Bayesian Model Selection (BMS) procedures directly target the marginal posterior distribution of the model:

$$p(M_k | \mathbf{X}) \propto p(\mathbf{X} | M_k) p(M_k);$$

where

$$p(\mathbf{X} | \mathcal{M}_k) = \int p(\mathbf{X} | \Theta_k, \mathcal{M}_k) p(\Theta_k | \mathcal{M}_k) d\Theta_k$$

is the **marginal** (i.e. integrated w.r.t. Θ_k) **likelihood** of \mathcal{M}_k

Marginal likelihood

- The **marginal likelihood** is a fundamental quantity in BMS;
- You can informally think of it as a score assigned to model M_k ;
- The marginal likelihood heavily depends on the specification of the prior $p(\Theta_k | \mathcal{M}_k)$, which is chosen by the analyst. As it influences the "score", specific care must be paid in the choice of the parameter prior distributions!
- Typically, $p(\Theta_k | \mathcal{M}_k)$ should satisfy some *compatibility* requirements. For instance, we would like two unidentifiable model to be assigned the same marginal likelihood (**score equivalence**)

Bayesian Causal Discovery

- In the Bayesian setting, causal discovery can be tackled as a Bayesian model selection problem, where the models considered are DAG models
- Suppose we have a collection of possible DAGs $(\mathcal{D}_1, \dots, \mathcal{D}_q)$. If \mathcal{S}_q is the space of "all" DAGs on q nodes/variables, our target is

$$p(\mathcal{D} | \mathbf{X}) = \frac{m(\mathbf{X} | \mathcal{D}) p(\mathcal{D})}{\sum_{\mathcal{D} \in \mathcal{S}_q} m(\mathbf{X} | \mathcal{D}) p(\mathcal{D})} \propto m(\mathbf{X} | \mathcal{D}) p(\mathcal{D})$$

i.e. the posterior distribution over DAG models

Bayesian Causal Discovery: How to

- The first step is to specify a Bayesian model that will define the posterior distribution. This consists of:
 - $p(X \mid \Theta, \mathcal{D})$: the statistical model;
 - $p(\Theta \mid \mathcal{D})$: the parameter prior;
 - $p(\mathcal{D})$: the model prior;
- Once the model is specified, we can derive the posterior distribution. Usually, this is not calculated exactly, but approximated via **sampling methods**.
 - For example, in our setting as p grows the number of DAGs grows exponentially, and evaluating the posterior probability of each of them becomes infeasible
- We will now tackle both steps, focusing on the Gaussian setting;

Bayesian causal discovery: Model specification in the Gaussian setting

- Suppose X is generated by a linear Gaussian Structural Equation Model (SEM) with independent error components and causal structure represented by the DAG \mathcal{D} .
- As we saw in **L1**, this implies that X is distributed as a Gaussian DAG model, i.e.

$$X_1, \dots, X_q \mid \Sigma_{\mathcal{D}} \sim \mathcal{N}_p(\mathbf{0}, \Sigma_{\mathcal{D}})$$
$$\Sigma_{\mathcal{D}} = \mathbf{L}^{-T} \mathbf{D} \mathbf{L}^{-1}$$

where $\mathbf{L} = (\mathbf{I}_p - \mathbf{B})$ and \mathbf{B} is the matrix of coefficients of the SEM;
coefficients

- The joint pdf of X factorizes as

$$p(x \mid (L, D), \mathcal{D}) = \prod_{j=1}^p d\mathcal{N}(x_j; -L_{\text{pa}_{\mathcal{D}}(j),j}^{\top} x_{\text{pa}_j(\mathcal{D})}, D_{jj})$$

- Consequently, we can write the likelihood of the (n, p) data matrix \mathbf{X} , containing i.i.d. samples from X as

$$p(\mathbf{X} \mid (L, D), \mathcal{D}) = \prod_{j=1}^p d\mathcal{N}_n(\mathbf{X}_{.j}; -\mathbf{X}_{.\text{pa}_j(\mathcal{D})} L_{\text{pa}_j(\mathcal{D}),j}, D_{jj} \mathbf{I}_n)$$

- Ben-David et al. (2011) developed the **DAG-Wishart** distribution, which is defined on the space of matrices (\mathbf{L}, \mathbf{D}) of a DAG.
 \implies It is a natural candidate for our parameter prior specification problem!
- The DAG-Wishart distribution is parameterized by a shape parameter $\alpha(\mathcal{D}) := (\alpha_1(\mathcal{D}), \dots, \alpha_p(\mathcal{D}))$ and a rate hyperparameter \mathbf{U} , a (p, p) s.p.d. matrix;
- In general terms, we thus specify the parameter prior as:

$$(\mathbf{L}, \mathbf{D}) \mid \mathcal{D} \sim \text{DAG-Wishart}(\alpha(\mathcal{D}), \mathbf{U})$$

The DAG-Wishart prior

- The pdf of a DAG-Wishart distribution has the following form:

$$p(\mathbf{L}, \mathbf{D} \mid \mathcal{D}) = \frac{1}{\mathcal{Z}_{\mathcal{D}}(\mathbf{U}, \boldsymbol{\alpha}(\mathcal{D}))} \exp \left\{ -\frac{1}{2} \text{tr} \left(\left(\mathbf{L} \mathbf{D}^{-1} \mathbf{L}^T \right) \mathbf{U} \right) \right\} \\ \cdot \prod_{j=1}^p \mathbf{D}_{jj}^{-\frac{\alpha_j(\mathcal{D})}{2}}$$

- The normalizing constant $\mathcal{Z}_{\mathcal{D}}(\mathbf{U}, \boldsymbol{\alpha}(\mathcal{D}))$ is

$$\mathcal{Z}_{\mathcal{D}}(\mathbf{U}, \boldsymbol{\alpha}) = \prod_{j=1}^p 2^{\left(\frac{\alpha_j(\mathcal{D})-2}{2}\right)} \pi^{\frac{|\text{pa}_j(\mathcal{D})|}{2}} \Gamma \left(\frac{\alpha_j(\mathcal{D}) - |\text{pa}_j(\mathcal{D})| - 2}{2} \right) \\ \cdot \frac{|\mathbf{U}_{\text{pa}_j(\mathcal{D})}|^{\frac{\alpha_j(\mathcal{D}) - \text{pa}_j(\mathcal{D}) - 3}{2}}}{|\mathbf{U}_{\text{fa}_j(\mathcal{D})}|^{\frac{\alpha_j(\mathcal{D}) - \text{pa}_j(\mathcal{D}) - 2}{2}}}$$

The DAG-Wishart prior

- The DAG-Wishart distribution presents many useful features
- First, it is **conjugate** to the likelihood of a Gaussian DAG model, which means that, if \mathbf{X} contains i.i.d. samples from a Gaussian DAG Model \mathcal{D} , we have

$$(\mathbf{L}, \mathbf{D}) | \mathbf{X}, \mathcal{D} \sim \text{DAG-Wishart} \left(\boldsymbol{\alpha}(\mathcal{D}) + n, \mathbf{U} + \mathbf{X}^T \mathbf{X} \right)$$

i.e., the posterior distribution has the same form as the prior distribution;

- Moreover, its normalizing constant, although involved, is available in closed form, which allows the calculation of the marginal likelihood;

The DAG-Wishart prior

- Finally, the DAG-Wishart distribution implies a set of local distributions on the non-null elements of (\mathbf{L}, \mathbf{D}) that are node-wise independent. In other words, it holds that

$$p(\mathbf{L}, \mathbf{D} \mid \mathcal{D}) = \prod_{j=1}^p p(\mathbf{L}_{\text{pa}_j(\mathcal{D}),j} \mid \mathbf{D}_{jj}, \mathcal{D}) p(\mathbf{D}_{jj} \mid \mathcal{D})$$

- In particular, we have that:

$$\begin{aligned} \mathbf{D}_{jj} \mid \mathcal{D} &\sim \text{Inv-Ga} \left(\frac{a_j(\mathcal{D}) - |\text{pa}_j(\mathcal{D})|}{2} - 1, \frac{1}{2} \mathbf{U}_{j|\text{pa}_j(\mathcal{D})} \right), \\ \mathbf{L}_{\text{pa}_j(\mathcal{D}),j} \mid \mathbf{D}_{jj}, \mathcal{D} &\sim \mathcal{N}_{|\text{pa}_j(\mathcal{D})|} \left(-\mathbf{U}_{\text{pa}_j(\mathcal{D})}^{-1} \mathbf{U}_{\text{pa}_j(\mathcal{D}),j}, \mathbf{D}_{jj} \mathbf{U}_{\text{pa}_j(\mathcal{D})}^{-1} \right), \end{aligned}$$

where $\mathbf{U}_{j|\text{pa}_j(\mathcal{D})} := \mathbf{U}_{jj} - \mathbf{U}_{j,\text{pa}_j(\mathcal{D})} (\mathbf{U}_{\text{pa}_j(\mathcal{D}),\text{pa}_j(\mathcal{D})})^{-1} \mathbf{U}_{\text{pa}_j(\mathcal{D}),j}$

- Nice consequence: sampling from the posterior of (\mathbf{L}, \mathbf{D}) is easy!

The DAG-Wishart prior

- The most important effect of this last property of the DAG-Wishart distribution emerges when computing the marginal likelihood of a DAG \mathcal{D} :

$$\begin{aligned} m(\mathbf{X} \mid \mathcal{D}) &= \int p(\mathbf{X} \mid (\mathbf{L}, \mathbf{D}), \mathcal{D}) p(\mathbf{L}, \mathbf{D} \mid \mathcal{D}) d(\mathbf{L}, \mathbf{D}) \\ &= \int \prod_{j=1}^p p(\mathbf{X}_{\cdot j} \mid \mathbf{X}_{\cdot \text{pa}_j(\mathcal{D})}, \mathbf{L}_{\text{pa}_j(\mathcal{D}), j}, \mathbf{D}_{jj}, \mathcal{D}) \\ &\quad p(\mathbf{L}_{\text{pa}_j(\mathcal{D}), j} \mid \mathbf{D}_{jj}, \mathcal{D}) p(\mathbf{D}_{jj} \mid \mathcal{D}) d(\mathbf{L}, \mathbf{D}) \\ &= \prod_{j=1}^p \int p(\mathbf{X}_{\cdot j} \mid \mathbf{X}_{\cdot \text{pa}_j(\mathcal{D})}, \mathbf{L}_{\text{pa}_j(\mathcal{D}), j}, \mathbf{D}_{jj}, \mathcal{D}) \\ &\quad p(\mathbf{L}_{\text{pa}_j(\mathcal{D}), j} \mid \mathbf{D}_{jj}, \mathcal{D}) p(\mathbf{D}_{jj} \mid \mathcal{D}) d(\mathbf{L}_{\text{pa}_j(\mathcal{D})}, \mathbf{D}_{jj}) \\ &= \prod_{j=1}^p m(\mathbf{X}_{\cdot j} \mid \mathbf{X}_{\cdot \text{pa}_j(\mathcal{D})}, \mathcal{D}) \end{aligned}$$

The DAG-Wishart prior

- Using a DAG-Wishart prior, the marginal likelihood follows the same factorization that the Gaussian DAG model implies on the likelihood;
- In this case, we say that the marginal likelihood is decomposable, a property that will turn out to be very useful when designing sampling algorithms to sample from the posterior over the DAG space!
- Each element of that factorisation can be easily computed from the normalising constant of the DAG-Wishart distribution. Denoting with $\tilde{U} = U + X^T X$, and with $\tilde{a}_j = a_j + n$ we have

$$m(\mathbf{X}_j | \mathbf{X}_{\text{pa}(j)}, \mathcal{D}) = (2\pi)^{-\frac{n}{2}} \cdot \frac{|\mathbf{U}_{\text{pa}_j, \text{pa}_j}|^{\frac{1}{2}}}{|\tilde{\mathbf{U}}_{\text{pa}_j, \text{pa}_j}|^{\frac{1}{2}}} \cdot \frac{\Gamma(\frac{1}{2}\tilde{a}_j)}{\Gamma(\frac{1}{2}a_j)} \cdot \frac{\left(\frac{1}{2}\mathbf{U}_j | \text{pa}_j\right)^{\frac{1}{2}a_j}}{\left(\frac{1}{2}\tilde{\mathbf{U}}_j | \text{pa}_j\right)^{\frac{1}{2}\tilde{a}_j}}$$

Score equivalence

- As we discussed at the beginning, the parameter prior we specify influences the marginal likelihood, i.e. the "score" we assign to each model;
- In the Gaussian setting, two Markov equivalent DAGs are unidentifiable given data alone;
- However, their marginal likelihood can be different because of the prior information we included via our prior distribution on the parameters, which is undesirable;
- Peluso & Consonni (2020) showed that we can ensure score equivalence with DAG-Wishart prior by setting, for each $j \in [q]$:

$$a_j(\mathcal{D}) = a - p + 2|\text{pa}_j(\mathcal{D})| + 3$$

where $a > p - 1$ ensures that the prior is proper;

- $p(\mathcal{D})$ can be assigned through a collection of Bernoulli distributions on 0-1 elements indicating absence/presence of edges in DAG \mathcal{D}
- Let $S^{\mathcal{D}}$ be the 0-1 *adjacency matrix* of the skeleton of \mathcal{D} :

$$S_{u,v}^{\mathcal{D}} = \begin{cases} 1 & \text{if } u \rightarrow v \in \mathcal{D} \text{ or } u \leftarrow v \in \mathcal{D} \\ 0 & \text{otherwise} \end{cases}$$

- We assign a prior on the DAGs based on their skeleton. In particular:

$$S_{u,v}^{\mathcal{D}} \mid \pi \stackrel{\text{iid}}{\sim} \text{Ber}(\pi) \quad u < v, \pi \in (0, 1)$$

- The prior probability assigned to each DAG is thus:

$$p(\mathcal{D} \mid \pi) = \pi^{|S^{\mathcal{D}}|} (1 - \pi)^{\frac{q(q-1)}{2} - |S^{\mathcal{D}}|}$$

Model specified

- We have now specified the whole Bayesian model for our Bayesian causal discovery problem!
- We just need to compute the posterior distribution;

**Bayesian causal discovery:
Sampling from the posterior
distribution**

Sampling from the posterior

- The denominator of $p(\mathcal{D} | \mathbf{X})$ involves a sum over a finite, but *very large*, number of DAGs \implies we cannot compute the posterior in exact form:
- The only thing we need to do is to approximate it using, for example, Markov Chain Monte Carlo (MCMC) sampling schemes;
- In particular, we will use a **Metropolis-Hastings** algorithm, which is based on the following steps
 - Start from an (arbitrary) initial DAG $\mathcal{D}^{(0)}$
 - Given a current DAG \mathcal{D} propose a new candidate DAG $\tilde{\mathcal{D}}$
 - Accept/reject $\tilde{\mathcal{D}}$ with "some" probability
 - Iterate the previous steps for a number of times S

Proposing a new DAG

- Suppose \mathcal{D} is the current DAG. We **propose** a new candidate DAG by **inserting**, **deleting** or **reversing** at random an edge in \mathcal{D} and **checking that the resulting graph is a DAG!**
- In practice, we build the set $\mathcal{O}_{\mathcal{D}}$ of all possible DAGs that can be reached from \mathcal{D} in one of the three moves above and propose uniformly one DAG among them
- The probability of transition from \mathcal{D} to $\tilde{\mathcal{D}}$ is then

$$q(\tilde{\mathcal{D}} | \mathcal{D}) = \frac{1}{|\mathcal{O}_{\mathcal{D}}|}$$

with $|\mathcal{O}_{\mathcal{D}}|$ number of DAGs obtained from \mathcal{D} by one of the local moves above, i.e. the number of *direct successors* DAGs of \mathcal{D} ;

Acceptance/Rejection step

- Given a current DAG \mathcal{D} , a new DAG $\tilde{\mathcal{D}}$ drawn from the proposal $q(\tilde{\mathcal{D}} | \mathcal{D})$ is accepted with probability

$$\alpha_{\tilde{\mathcal{D}}, \mathcal{D}} = \min \left\{ 1; \frac{m(\mathbf{X} | \tilde{\mathcal{D}})}{m(\mathbf{X} | \mathcal{D})} \cdot \frac{p(\tilde{\mathcal{D}})}{p(\mathcal{D})} \cdot \frac{q(\mathcal{D} | \tilde{\mathcal{D}})}{q(\tilde{\mathcal{D}} | \mathcal{D})} \right\}$$

which depends on:

- The marginal likelihood ratio;
- The prior ratio;
- The proposal ratio;

- The **proposal ratio**

$$\frac{q(\mathcal{D} | \tilde{\mathcal{D}})}{q(\tilde{\mathcal{D}} | \mathcal{D})} = \frac{|\mathcal{O}_{\mathcal{D}}|}{|\mathcal{O}_{\tilde{\mathcal{D}}}|}$$

requires the enumeration of all operators that can be applied to \mathcal{D} and lead to a valid graph (i.e. a DAG).

- It is usually computationally expensive, but for p large it can be approximated to 1;

Acceptance/Rejection step

- As we are using local moves and thanks to the decomposability of the marginal likelihood, the **marginal likelihood** ratio simplifies to the components which are affected by the local move.
- If, for instance, $\tilde{\mathcal{D}}$ differs from \mathcal{D} for the addition of an edge pointing towards node t , we have:

$$\frac{m(\mathbf{X} \mid \tilde{\mathcal{D}})}{m(\mathbf{X} \mid \mathcal{D})} = \frac{m(\mathbf{X}_{.t} \mid \mathbf{X}_{.\text{pa}_t(\tilde{\mathcal{D}})} \tilde{\mathcal{D}})}{m(\mathbf{X}_{.t} \mid \mathbf{X}_{.\text{pa}_t(\mathcal{D})} \mathcal{D})},$$

which significantly speeds up computations!

Algorithm 1: Collapsed MCMC to sample from $p(\mathcal{D} \mid \mathbf{X})$

Input: \mathbf{X} (n, q) dataset; S number of MCMC iterations; prior hyperparameters

Output: S samples from the posterior $p(\mathcal{D} \mid \mathbf{X})$

```
1 Initialize  $\mathcal{D}^{(0)}$ , e.g. the empty DAG;  
2 for  $s = 1, \dots, S$  do  
3   Sample  $\tilde{\mathcal{D}}$  from  $q(\tilde{\mathcal{D}} \mid \mathcal{D}^{(s-1)})$ ;  
4   Compute the acceptance probability  $\alpha_{\tilde{\mathcal{D}}, \mathcal{D}}$ ;  
5   Set  $\mathcal{D}^{(s)} = \tilde{\mathcal{D}}$  with probability  $\alpha_{\tilde{\mathcal{D}}, \mathcal{D}}$ , otherwise  $\mathcal{D}^{(s)} = \mathcal{D}^{(s-1)}$ ;  
6 end  
7 return  $\{\mathcal{D}^{(1)}, \dots, \mathcal{D}^{(S)}\}$ 
```

- The resulting algorithm is a *collapsed* sampler over the space of DAGs, since we integrated out the parameter;
- Its output is a collection of DAGs approximately drawn from the posterior $p(\mathcal{D} \mid \mathbf{X})$

Posterior inference

- Output of the MH algorithm is a collection of DAGs $\{\mathcal{D}^{(1)}, \dots, \mathcal{D}^{(S)}\}$
- We can provide an estimate of the posterior probability of $\mathcal{D} \in \mathcal{S}_q$ as

$$\hat{p}(\mathcal{D} | \mathbf{X}) = \frac{1}{S} \sum_{s=1}^S \mathbf{1} \left\{ \mathcal{D}^{(s)} = \mathcal{D} \right\}$$

i.e. the proportion of DAGs, visited by the MCMC, equal to \mathcal{D}

- Other summaries:
 - Estimate of the (marginal) posterior probability of edge inclusion for each $u \rightarrow v$

$$\hat{p}(u \rightarrow v | \mathbf{X}) = \frac{1}{S} \sum_{s=1}^S \mathbf{1} \left\{ u \rightarrow v \in \mathcal{D}^{(s)} \right\}$$

computed as the proportion of DAGs, visited by the MCMC, containing $u \rightarrow v$;

- DAG point estimates from the posterior over DAGs can be obtained by:
 - including those edges whose posterior probability is higher than some threshold, e.g. 0.5 (Median Probability DAG Model, MPM)

$$\hat{S}_{u,v} = \begin{cases} 1 & \text{if } \hat{p}(u \rightarrow v \mid \mathbf{X}) \geq 0.5 \\ 0 & \text{otherwise} \end{cases}$$

- selecting the DAG having the highest posterior probability (Maximum A Posteriori DAG, MAP)

$$\hat{\mathcal{D}}_{MAP} = \underset{\mathcal{D}}{\operatorname{argmax}} \hat{p}(\mathcal{D} \mid \mathbf{X})$$

- Algorithm 1, specialized to Gaussian DAGs, is implemented in the R package BCDAG (Bayesian structure and Causal learning of Gaussian DAGs) within the function `learnDAG` and under the setting `collapse = TRUE`
- Inputs of the function are:
 - `data` : the (n, q) data matrix \mathbf{X}
 - `S` : the number of MCMC iterations
 - `burn` : a burn-in period
 - `a, U` : hyperparameters of the DAG-Wishart prior
 - `w` : prior probability of edge inclusion for $p(\mathcal{D})$
- Function for posterior summaries and MCMC diagnostics are also provided within the package. See Castelletti & Mascaro (2022) for full details

Bayesian estimation of causal effects

Estimating causal effects for fixed DAG

- As we saw yesterday discussing Maximum Likelihood Estimation, a causal effect can be identified and estimated directly from $\Sigma^{\mathcal{D}}$ as

$$\gamma_{ty} = [(\Sigma_{\tilde{Z}, \tilde{Z}}^{\mathcal{D}})^{-1}(\Sigma_{\tilde{Z}, y}^{\mathcal{D}})]_1$$

where $\tilde{Z} := (X_t, \mathbf{Z})$ and \mathbf{Z} is an adjustment set;

- In the Bayesian setting, the "estimate" of $\Sigma^{\mathcal{D}}$ is its posterior distribution $p(\Sigma^{\mathcal{D}} \mid \mathbf{X}, \mathcal{D})$;
- In the same way, the "estimate" of the causal effect will be a posterior distribution $p(\gamma_{ty} \mid \mathbf{X})$
- We cannot compute these posteriors exactly, but again we can approximate them via sampling!

Estimating causal effects for unknown DAG

- As both γ_{ty} and $\Sigma_{\mathcal{D}}$ are functions of the DAG parameters (L, D) , we can first try to approximate the joint posterior distribution

$$p(L, D, \mathcal{D} | X) = p(L, D | \mathcal{D}, X) \cdot p(\mathcal{D} | X);$$

- As we saw before, $p(L, D | \mathcal{D}, X)$ is just a DAG-Wishart and it can be easily sampled from once the DAG is known, and we can sample from $p(\mathcal{D} | X)$ using the same MH as before! Hence:

Algorithm 2: MCMC to sample from $p(\theta, \mathcal{D} | X)$

Input: X (n, q) dataset; S number of MCMC iterations; prior hyperparameters

Output: S samples from the posterior $p(\theta, \mathcal{D} | X)$

```
1 Initialize  $\mathcal{D}^{(0)}$ , e.g. the empty DAG;
2 for  $s = 1, \dots, S$  do
3   Sample  $\tilde{\mathcal{D}}$  from  $q(\tilde{\mathcal{D}} | \mathcal{D}^{(s-1)})$ ;
4   Compute the acceptance probability  $\alpha_{\tilde{\mathcal{D}}, \mathcal{D}}$ ;
5   Set  $\mathcal{D}^{(s)} = \tilde{\mathcal{D}}$  with probability  $\alpha_{\tilde{\mathcal{D}}, \mathcal{D}}$ , otherwise  $\mathcal{D}^{(s)} = \mathcal{D}^{(s-1)}$ ;
6   Sample  $\theta^{(s)}$  from its full conditional;
7 end
8 return  $\{(\theta^{(1)}, \mathcal{D}^{(1)}), \dots, (\theta^{(S)}, \mathcal{D}^{(S)})\}$ 
```

Posterior sampler for DAGs and parameters

- Output of Algorithm 2 is a collection of draws from the posterior $p(\mathbf{L}, \mathbf{D}, \mathcal{D} \mid \mathbf{X})$ of the form

$$\left\{ \left(\mathbf{L}^{(1)}, \mathbf{D}^{(1)}, \mathcal{D}^{(1)} \right), \dots, \left(\mathbf{L}^{(S)}, \mathbf{D}^{(S)}, \mathcal{D}^{(S)} \right) \right\}$$

- We are interested in the causal effect γ_{ty}
- We can obtain a sample from it by, for each $s \in [S]$:
 - Calculating $\Sigma_{\mathcal{D}}^{(s)} = (\mathbf{L}^{(s)})^{-T} \mathbf{D}^{(s)} (\mathbf{L}^{(s)})^{-1}$;
 - From $\Sigma_{\mathcal{D}}^{(s)}$ deriving $\gamma_{ty}^{(s)}$ using the adjustment formula above;
- The collection $\left\{ \gamma_{ty}^{(1)}, \dots, \gamma_{ty}^{(S)} \right\}$ provides an approximation of $p(\gamma_{ty} \mid \mathbf{X})$ which naturally accounts for DAG-model uncertainty, since each draw potentially depends on a different underlying DAG;

Posterior inference on causal effects

- A point estimate of γ_{ty} is then

$$\hat{\gamma}_{ty}^{BMA} = \frac{1}{S} \sum_{s=1}^S \gamma_{ty}^{(s)},$$

which implicitly performs **Bayesian Model Averaging** (BMA) through the MCMC frequencies of the visited DAGs;

- Other summaries/queries of interest are:
 - $\hat{p}(\gamma_{ty} < 0 \mid \mathbf{X}) = \frac{1}{S} \sum_{s=1}^S \mathbf{1}\{\gamma_{ty}^{(s)} < 0\};$
 - $\hat{p}(\gamma_{ty} = 0 \mid \mathbf{X}) = \frac{1}{S} \sum_{s=1}^S \mathbf{1}\{\gamma_{ty}^{(s)} = 0\};$
 - $\hat{p}(\gamma_{ty} > 0 \mid \mathbf{X}) = \frac{1}{S} \sum_{s=1}^S \mathbf{1}\{\gamma_{ty}^{(s)} > 0\}$

- Algorithm 2 is also implemented in BCDAG within the function `learn_DAG` and under the setting `collapse = FALSE`;
- Inputs of the function are the same as of Algorithm 1:
 - `data` : the (n, q) data matrix \mathbf{X} ;
 - `S` : the number of MCMC iterations;
 - `burn` : a burn-in period;
 - `a, U` : hyperparameters of the DAG-Wishart prior;
 - `w` : prior probability of edge inclusion for $p(\mathcal{D})$;

- Function `get_causaleffect` recovers the posterior of causal effect parameters of interest from the output of `learnDAG`;
- Inputs of `get_causaleffect` are:
 - `learnDAG_output` : output of `learnDAG`, an object of class `bcdag`;
 - `targets` : numerical label of variable X_v (target);
 - `response` : numerical label of variable Y (response)
- Functions for MCMC diagnostics and summaries are also available (see `get_diagnostics` function);

Thank you!

References:

- Ben-David, E., Li, T., Massam, H., Rajaratnam, B. (2011). *High dimensional Bayesian inference for Gaussian directed acyclic graph models*. arXiv preprint arXiv:1109.4371.
- Castelletti, F., & Consonni, G. (2021). *Bayesian inference of causal effects from observational data in Gaussian graphical models*. Biometrics, 77(1), 136-149.
- Castelletti, F., & Mascaro, A. (2022). *BCDAG: An R package for Bayesian structure and causal learning of Gaussian DAGs*. arXiv preprint arXiv:2201.12003.
- Peluso, S., & Consonni, G. (2020). *Compatible priors for model selection of high-dimensional Gaussian DAGs*. Electronic Journal of Statistics, 14(2), 4110-4132.