

# Bayesian constrained-based structure learning

## Background

### 1. Multivariate Gaussian data

$$\begin{aligned} \underline{x}_1, \dots, \underline{x}_n \mid \Omega &\stackrel{\text{iid}}{\sim} \mathcal{N}_q(\underline{0}, \Omega^{-1}) & \Omega \in \rho \quad \text{s.p.d. matrices} \\ \Omega &\sim W_q(a, U) \end{aligned}$$

$$p(\underline{x}_1, \dots, \underline{x}_n \mid \Omega) = (2\pi)^{-\frac{np}{2}} |\Omega|^{\frac{n}{2}} \exp\left\{-\frac{1}{2} \text{tr}(S\Omega)\right\}$$

$$S = \sum_{i=1}^n \underline{x}_i \underline{x}_i^T$$

$$p(\Omega) = c(a, U) |\Omega|^{\frac{a-p-1}{2}} \exp\left\{-\frac{1}{2} \text{tr}(U\Omega)\right\}$$

$$c(a, U) = \frac{|U|^{\frac{a}{2}}}{2^{\frac{aq}{2}} \Gamma_q\left(\frac{a}{2}\right)}$$

prior normalizing constant  
↳ multivariate

The marginal data distribution (i.e. the marginal likelihood) is :

$$p(x_1, \dots, x_n) = \int p(x_1, \dots, x_n | \Omega) p(\Omega) d\Omega$$

$$= \frac{c(a, U)}{c(a+n, U+S)} \cdot (2\pi)^{-\frac{nq}{2}} := p(X)$$

$c(a+n, U+S)$  posterior normalizing constant

Now consider  $A \subseteq \{1, \dots, q\}$ . We have:

$$x_1^A, \dots, x_n^A \mid \Omega_{A|\bar{A}} \stackrel{iid}{\sim} \mathcal{N}_{|A|}(\underline{0}, (\Omega_{A|\bar{A}})^{-1})$$

$$\Omega_{A|\bar{A}} \sim W_{|A|}(a - |\bar{A}|, U_{AA})$$

with  $x_i^A = (x_{ij})_{j \in A}$   $U_{AA} = [U_{jj}]_{j \in A}$

$$\Omega_{A|\bar{A}} = \Omega_{AA} - \Omega_{A\bar{A}} (\Omega_{\bar{A}\bar{A}})^{-1} \Omega_{\bar{A}A}$$

$$= (\Sigma_{AA})^{-1} \quad \text{i.e. } (\Omega_{A|\bar{A}})^{-1} = \Sigma_{AA}$$

the marginal  
covariance matrix

Therefore we obtain :

$$\begin{aligned} p(\underline{x}_1^A, \dots, \underline{x}_n^A) &= \int p(\underline{x}_1^A, \dots, \underline{x}_n^A | \Omega_{A|\bar{A}}) p(\Omega_{A|\bar{A}}) d\Omega_{A|\bar{A}} \\ &= \frac{c(a - |\bar{A}|, U_{AA})}{c(a - |\bar{A}| + n, U_{AA} + S_{AA})} \cdot (2\pi)^{-\frac{n|A|}{2}} \\ &:= p(X_A) \quad (*) \end{aligned}$$

Refs : Press (1982) "Applied Multivariate Analysis"  
Cousonni & La Rocca (2012, SJS)

## 2. Bayesian DAG-model selection

$$\mathcal{D} = (V, E) \quad \text{a DAG with} \quad V = \{1, \dots, p\}$$
$$E \subseteq V \times V$$

$$\text{pa}_{\mathcal{D}}(j) = \{u : (u, j) \in E\} \quad \text{parents of } j \text{ in } \mathcal{D}$$

$$\text{fa}_{\mathcal{D}}(j) = j \cup \text{pa}_{\mathcal{D}}(j) \quad \text{family of } j \text{ in } \mathcal{D}$$

Under  $\mathcal{D}$  we have,

$$p(x_1, \dots, x_p | \mathcal{D}) = \prod_{j=1}^p p(x_j | \underline{x}_{\text{pa}_{\mathcal{D}}(j)})$$

In a Gaussian DAG model:

$$\underline{x}_1, \dots, \underline{x}_n | \Omega_{\mathcal{D}} \sim \mathcal{N}_p(\underline{0}, \Omega_{\mathcal{D}}^{-1})$$

$$\Omega_{\mathcal{D}} \sim p(\Omega_{\mathcal{D}})$$

$\Omega_{\mathcal{D}} \in \mathcal{P}_{\mathcal{D}}$  precision matrices Markov w.r.t.  $\mathcal{D}$

Which prior for  $\Omega_{\mathcal{D}}$ ?

If the goal is to compute the DAG marginal likelihood

$$p(X|D) = \int p(x_1, \dots, x_n | \Omega_D) p(\Omega_D) d\Omega_D$$

then there is no need to specify "directly"  $p(\Omega_D)$  but rather a prior on  $\Omega$  **just s.p.d.** and some assumptions under which we recover  $p(X|D)$  as:

$$p(X|D) = \prod_{j=1}^n \left\{ \frac{p(X_{fa_D(j)})}{p(X_{pa_D(j)})} \right\}$$

and  $p(X_{fa_D(j)})$ ,  $p(X_{pa_D(j)})$  are as in

with  $A = fa_D(j)$  ,  $A = pa_D(j)$  

Refs : Geiger & Heckerman (2002, AOS)

Coussonni & La Rocca (2012, STS)

### 3. The PC Algorithm for DAG estimation

Refs : Kalish & Buhlmann (2007, JMLR)  
Slides on Graphical Models

Idea in PC Algorithm is to recover a CPDAG  
(Completed Partially DAG)  
through a sequence of Conditional  
Independence (CI) tests

Specifically, in Step 1 (skeleton estimation)  
the PC algorithm remove an edge  $u-v$  if

$$X_u \perp\!\!\!\perp X_v \mid X_S \quad \text{for some } S \subseteq V \setminus \{u, v\}$$

i.e. for at least one set  $S$  not including  
 $u$  and  $v$

Idea in practice is to start from a set  $S$   
of size 0 (i.e.  $\emptyset$ ), then increase it by one  
and so on; we stop when we find a set  $S$   
for which  $X_u \perp\!\!\!\perp X_v \mid X_S$ .

Which CI test ?

For Gaussian data, a test based on partial correlation coefficients  $\rho_{uv|s}$ .

What is  $\rho_{uv|s}$  ?

It is the  $\text{Corr}(X_u, X_v | X_s)$  that is the correlation coefficient between  $X_u$  and  $X_v$  in the joint, conditional distribution

$$(X_u, X_v) | X_s$$

Start from  $X_A = (X_j)_{j \in A}$  with  $A = \{u, v, s\}$

for which we know that

$$X_A | \Omega_{A|\bar{A}} \sim \mathcal{N}_{|A|}(\mathbf{0}, (\Omega_{A|\bar{A}})^{-1})$$

$$\text{and } \Omega_{A|\bar{A}}^{-1} = \Sigma_{AA}$$

Then, we partition  $A$  into  $\{u, v\}$  and  $s$  and  $\Sigma_{AA}$  accordingly

We consider the conditional distribution of  $X_{\{u,v\}}$  given  $X_s$  :

$$X_{\{u,v\}} | X_s \sim N_2 (\mu_{\{u,v\}|s}, \Sigma_{\{u,v\}|s})$$

with  $\mu_{\{u,v\}|s} = \Sigma_{s,\{u,v\}} (\Sigma_{\{u,v\}\{u,v\}})^{-1} \Sigma_{\{u,v\},s}$

$$\Sigma_{\{u,v\}|s} = \Sigma_{\{u,v\}\{u,v\}} - \Sigma_{\{u,v\},s} (\Sigma_{s,s})^{-1} \Sigma_{s,\{u,v\}}$$

Remark: the out-diagonal element of  $\Sigma_{\{u,v\}|s}$  is the covariance between  $X_u$  and  $X_v$  in the conditional distribution of  $(X_u, X_v)$  given  $X_s$  :  $\text{Cov}(X_u, X_v | X_s)$

The partial correlation coefficient is then :

$$\text{Corr}(X_u, X_v | X_s) \stackrel{\text{def.}}{=} \frac{[\Sigma_{\{u,v\}|s}]_{u,v}^{1,2}}{\sqrt{[\Sigma_{\{u,v\}|s}]_{u,u}^{1,1} [\Sigma_{\{u,v\}|s}]_{v,v}^{2,2}}}$$

$\rho_{\{u,v\}|s}$



$$\rho_{\{u,v\}|s} = 0 \quad \text{iff} \quad X_u \perp\!\!\!\perp X_v \mid X_s$$

$$\text{and } \rho_{\{u,v\}|s} = 0 \quad \text{iff} \quad [\Sigma_{\{u,v\}|s}]_{u,v} = 0$$

Now, the joint distribution of  $X_{\{u,v\}} \mid X_s$  can be equivalently parameterise as:

$$\Sigma_{\{u,v\}|s} \mapsto \{L_{u|v,s} ; D_{u|v,s}, D_{v|s}\}$$

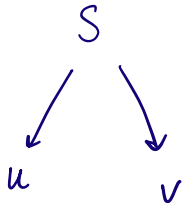
where  $L_{u|v,s}$  is the regression coefficient in the regression of  $X_u$  on  $X_v \cup X_s$ , while  $D_{u|v,s}$  and  $D_{v|s}$  are the conditional variances.

$$L_{u|v,s} = [\Sigma_{\{u,v\}|s}]_{u, \{v,s\}} [\Sigma_{\{u,v\}|s}]_{\{v,s\} \{v,s\}}^{-1}$$

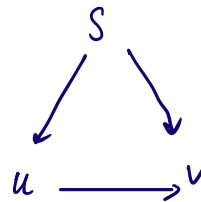
This should be equivalent as doing the following.

Consider

$\mathcal{D}_0$



$\mathcal{D}_1$



$$\begin{aligned}
 BF_{01} &= \frac{m(X | \mathcal{D}_0)}{m(X | \mathcal{D}_1)} = \frac{\cancel{p(X_u | X_s)} p(X_v | X_s) \cancel{p(X_s)}}{\cancel{p(X_u | X_s)} p(X_v | X_u, X_s) \cancel{p(X_s)}} = \\
 &= \frac{p(X_v, X_s) / p(X_s)}{p(X_v, X_u, X_s) / p(X_u, X_s)} = \\
 &= \frac{p(X_{\text{fa}_{\mathcal{D}_0}}(v)) / p(X_{\text{pa}_{\mathcal{D}_0}}(v))}{p(X_{\text{fa}_{\mathcal{D}_1}}(v)) / p(X_{\text{pa}_{\mathcal{D}_1}}(v))}
 \end{aligned}$$

Similarly to what is given in Geiger & Heckermann (2002) it should coincide with the ratio of two marginal likelihoods of two complete DAGs.

# Uncertainty quantification

PCalg does not allow for uncertainty quantification underlying the estimated graphical structure.

Here instead we have  $p(u \rightarrow v | X)$

Let  $p(u \rightarrow v)$  and  $p(u \nrightarrow v)$  be prior probabilities of having or not an edge between  $u$  and  $v$ .

$$\begin{aligned} p(u \rightarrow v | X) &= \frac{\cancel{m(X | u \rightarrow v)} p(u \rightarrow v)}{m(X | u \rightarrow v) p(u \rightarrow v) + m(X | u \nrightarrow v) p(u \nrightarrow v)} \times \frac{m(X | u \rightarrow v) p(u \rightarrow v)}{\cancel{m(X | u \rightarrow v)} p(u \rightarrow v)} = \\ &= \frac{1}{\frac{m(X | u \rightarrow v) p(u \rightarrow v) + m(X | u \nrightarrow v) p(u \nrightarrow v)}{m(X | u \rightarrow v) p(u \rightarrow v)}} = \\ &= \frac{1}{1 + BF_{01} \frac{p(u \nrightarrow v)}{p(u \rightarrow v)}} \end{aligned}$$

ie given prior probabilities of edges and the BF we can compute  $p(u \rightarrow v | X)$

If :  $\bullet BF_{01} = 0 \Rightarrow m(X | \mathcal{D}_0) = 0 \Rightarrow p(u \rightarrow v | X) = 1$

$\bullet BF_{01} \rightarrow +\infty \Rightarrow m(X | \mathcal{D}_1) = 0 \Rightarrow p(u \rightarrow v | X) = 0$