

Graphical Evidence

Anindya Bhadra

www.stat.purdue.edu/~bhadra

Purdue University

Overview

- Marginal likelihood or *evidence* is fundamental to Bayesian statistics.
- Used for empirical Bayes tuning of hyperparameters, model selection using Bayes factors.
- There is no dearth of generic approaches, yet calculation of evidence is mostly unresolved in Gaussian graphical models (GGMs), except for very specific priors such as the Wishart or G-Wishart.
- **Goal: To provide a tractable approach for evidence calculation in GGMs under mild requirements.**
- *Joint work with Ksheera Sagar (Purdue), Sayantan Banerjee (IIM Indore) and Jyotishka Datta (Virginia Tech).*

Evidence in GGMs

- Suppose $\mathbf{y}_{n \times p} \sim \mathcal{N}(\mathbf{0}, \mathbf{I}_n \otimes \boldsymbol{\Omega}_{p \times p}^{-1})$. Evidence calculation is simple in principle:

$$f(\mathbf{y}) = \int_{\boldsymbol{\Omega} \in \mathcal{M}_p^+} f(\mathbf{y} \mid \boldsymbol{\Omega}) f(\boldsymbol{\Omega}) d\boldsymbol{\Omega}.$$

- The restriction of the integral to the space of positive definite matrices causes a lot of difficulties, except for Wishart and specific instances of G-Wishart (Uhler et al., 2018, AoS).
- For the same reason, a “default” covering density is very hard to design: difficulties for importance, bridge or path sampling.

Generic approaches for estimating evidence

- Harmonic mean estimates and variants (Newton and Raftery, 1994, JRSSB; Gelfand and Dey, 1994, JRSSB)
- Importance sampling approaches:
 - Bridge sampling and variants (path, warped bridge) (Gelman and Meng, 1998, Stats. Sci.; Meng and Wong, 1996, Sinica; Meng and Schilling, 2002, JCGS),
 - Annealed importance sampling (Neal, 2001, Stats. Comput.)
- Nested sampling (Skilling, 2006, BA).
- Chib (1995, JASA) and Chib and Jeliazkov (2001, JASA) based on MCMC posterior draws.
- Excellent review article by Llorente et al. (2022, SIAM Review).

Do generic approaches work in GGMs?

- HM estimates can have unbounded variance: limit distribution is α stable (Wolpert and Schmeider, 2012).
- We are not aware of any principled way of choosing an importance or bridge density **under a positive definite restriction**.
- Nested sampling requires sampling from a progressively higher likelihood region: very hard to implement in high dimensions.
- **A case in point**: the specialized Monte Carlo method of Atay-Kayis and Massam (2005, Biometrika) for G-Wishart marginals appeared a good 10 years after these generic approaches.

Chib (1995)

- Recall the fundamental Bayesian identity:

$$f(\mathbf{y}) = \frac{f(\mathbf{y} | \theta)f(\theta)}{f(\theta | \mathbf{y})},$$

- The likelihood and the prior can typically be evaluated at some $\theta = \theta^*$, the trouble is *evaluating* $f(\theta | \mathbf{y})$.
- Chib's strategy:
 - Decompose $\Omega = (z, \theta) = (\text{nuisance parameter}, \text{parameter of interest})$.
 - Run a Gibbs sampler iterating between $f(z | \theta, \mathbf{y})$ and $f(\theta | z, \mathbf{y})$. Converges to $f(z, \theta | \mathbf{y})$. Correct marginals for $(z | \mathbf{y})$ and $(\theta | \mathbf{y})$.
 - Estimate using the Gibbs draws:

$$\hat{f}(\theta^* | \mathbf{y}) = M^{-1} \sum_{i=1}^M f(\theta^* | z^{(i)}, \mathbf{y}), \quad z^{(i)} \sim f(z | \mathbf{y}).$$

- Need the constants only for $f(\theta | z, \mathbf{y})$; not for $f(z | \theta, \mathbf{y})$.

Pros and cons of Chib (1995)

- Chib's approach is automatic in the same way a Gibbs sampler is automatic: a covering (importance, bridge) density is not required.
- But applying Chib's method requires designing a suitable $f(\theta | z, \mathbf{y})$ that can be **evaluated** (merely **sampling** from it is not enough).
- Application is a matter of art and not generic in a way the harmonic mean estimate is generic.
- Some known difficulties in finite mixture models (Neal, 1999).

Chib's approach for GGMs: the telescoping block decomposition

- Apply the decomposition:

$$\Omega_{p \times p} = \begin{bmatrix} \Omega_{(p-1) \times (p-1)} & \omega_{\bullet p} \\ \omega_{\bullet p}^T & \omega_{pp} \end{bmatrix}.$$

- Let $\theta_p = (\omega_{\bullet p}, \omega_{pp})$ and $z =$ collection of all other latent variables.
- Wang (2012, BA) showed **in the context of sampling** that $f(\theta_p | \mathbf{y}, z) = f(\omega_{\bullet p}, \omega_{pp} | \mathbf{y}, z) = f(\omega_{\bullet p} | \mathbf{y}, z) f(\omega_{pp} | \omega_{\bullet p}, \mathbf{y}, z)$ decomposes as **(normal \times gamma)** under suitable priors on $\Omega_{p \times p}$.
- We will use this for **density evaluation**, since the normalizing constants for both normal and gamma densities are available!

Chib's approach for GGMs: the telescoping block decomposition

- We have

$$\log f(\mathbf{y}_{1:p}) = \log f(\mathbf{y}_{1:p} | \boldsymbol{\theta}_p) + \log f(\boldsymbol{\theta}_p) - \log f(\boldsymbol{\theta}_p | \mathbf{y}_{1:p}).$$

- Slightly rewrite:

$$\begin{aligned} \log f(\mathbf{y}_{1:p}) &= \log f(\mathbf{y}_p | \mathbf{y}_{1:p-1}, \boldsymbol{\theta}_p) + \log f(\mathbf{y}_{1:p-1} | \boldsymbol{\theta}_p) + \log f(\boldsymbol{\theta}_p) - \log f(\boldsymbol{\theta}_p | \mathbf{y}_{1:p}) \\ &:= \text{I}_p + \text{II}_p + \text{III}_p - \text{IV}_p. \end{aligned}$$

- We can evaluate the partial likelihood I_p using

$$\mathbf{y}_p | \mathbf{y}_{1:p-1}, \boldsymbol{\theta}_p \sim \mathcal{N}(-\mathbf{y}_{1:p-1} \boldsymbol{\omega}_{\bullet p} / \omega_{pp}, 1/\omega_{pp}),$$

- Assume III_p can be evaluated and Wang's result from the previous slide will be used for evaluating IV_p . There remains II_p to deal with.

Chib's approach for GGMs: the telescoping block decomposition

- BUT! The term II is telescoping. We have:

$$\begin{aligned}\text{II}_p &= \log f(\mathbf{y}_{1:p-1} \mid \boldsymbol{\theta}_p) \\ &= \log f(\mathbf{y}_{p-1} \mid \mathbf{y}_{1:p-2}, \boldsymbol{\theta}_p, \boldsymbol{\theta}_{p-1}) + \log f(\mathbf{y}_{1:p-2} \mid \boldsymbol{\theta}_p, \boldsymbol{\theta}_{p-1}) \\ &\quad + \log f(\boldsymbol{\theta}_{p-1} \mid \boldsymbol{\theta}_p) - \log f(\boldsymbol{\theta}_{p-1} \mid \mathbf{y}_{1:p-1}, \boldsymbol{\theta}_p) \\ &:= \text{I}_{p-1} + \text{II}_{p-1} + \text{III}_{p-1} - \text{IV}_{p-1}.\end{aligned}$$

- We use a form of iterative proportional scaling (IPS). Define $\tilde{\boldsymbol{\Omega}}_{(p-1) \times (p-1)}$ as:

$$\tilde{\boldsymbol{\Omega}}_{(p-1) \times (p-1)} = \boldsymbol{\Omega}_{(p-1) \times (p-1)} - \frac{\boldsymbol{\omega}_{\bullet p} \boldsymbol{\omega}_{\bullet p}^T}{\omega_{pp}} := \begin{bmatrix} \tilde{\boldsymbol{\Omega}}_{(p-2) \times (p-2)} & \tilde{\boldsymbol{\omega}}_{\bullet(p-1)} \\ \tilde{\boldsymbol{\omega}}_{\bullet(p-1)}^T & \tilde{\omega}_{(p-1)(p-1)} \end{bmatrix}.$$

Then $\tilde{\boldsymbol{\Omega}}_{(p-1) \times (p-1)}$ is p.d. and $(\mathbf{y}_{1:p-1} \mid \boldsymbol{\theta}_p, \boldsymbol{\Omega}_{(p-1) \times (p-1)}) \sim \mathcal{N}(\mathbf{0}, \tilde{\boldsymbol{\Omega}}_{(p-1) \times (p-1)}^{-1})$.

- Thus, I_{p-1} can be evaluated using:

$$\mathbf{y}_{p-1} \mid \mathbf{y}_{1:p-2}, \boldsymbol{\theta}_p, \boldsymbol{\theta}_{p-1} \sim \mathcal{N}(-\mathbf{y}_{1:p-2} \tilde{\boldsymbol{\omega}}_{\bullet(p-1)} / \tilde{\omega}_{(p-1)(p-1)}, 1 / \tilde{\omega}_{(p-1)(p-1)}).$$

Overall strategy

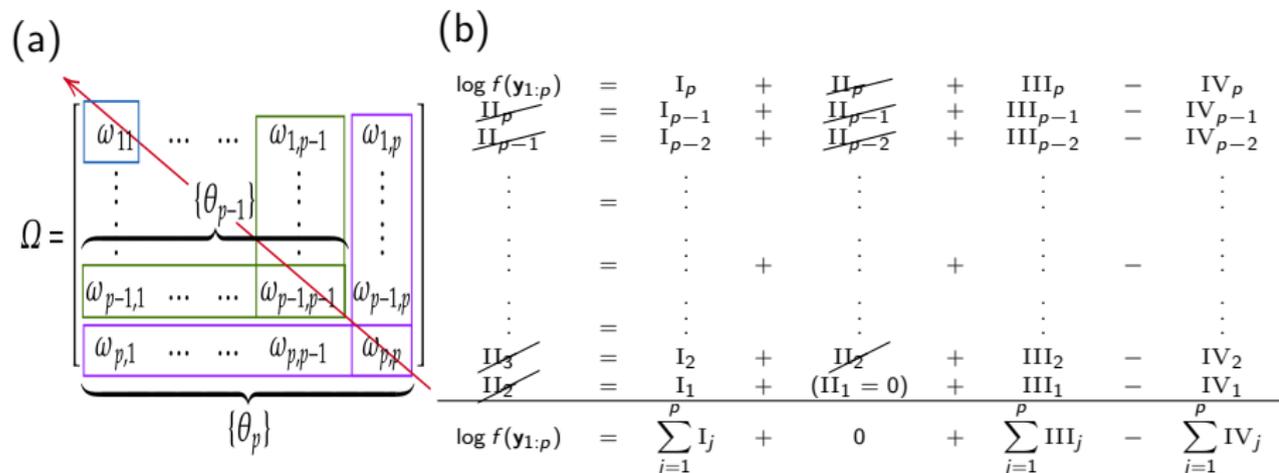


Figure: (a) Decomposition of $\Omega_{p \times p}$. Purple, green and blue blocks denote θ_p , θ_{p-1} and finally $\theta_1 = \omega_{11}$. Red arrow denotes how the algorithm proceeds, fixing one row/column at a time, and (b) the telescoping sum giving the log-marginal $\log f(\mathbf{y}_{1:p})$.

- Run Chib p times, adjusting Ω each time. In each equation, evaluate only I, III and IV. Eliminate II via telescoping sum.

A demonstration on Wishart (known marginal)

- Suppose $\boldsymbol{\Omega} \sim \mathcal{W}_p(\mathbf{V}, \alpha)$. Then,

$$\begin{aligned} \log f(\mathbf{y}_{1:p}) &= -\frac{np}{2} \log(\pi) + \log \Gamma_p \left(\frac{\alpha + n}{2} \right) - \log \Gamma_p \left(\frac{\alpha}{2} \right) \\ &\quad + \frac{(\alpha + n)}{2} \log \left| \mathbf{I}_p + \mathbf{V}^{1/2} \mathbf{S} \mathbf{V}^{1/2} \right|. \end{aligned}$$

- The closed form expression for the marginal provides an *oracle*.

Computing $\text{III}_p (= \log f(\boldsymbol{\theta}_p))$

- Recall, $\boldsymbol{\theta}_p = (\boldsymbol{\omega}_{\bullet,p}, \omega_{pp})$.
- If $\boldsymbol{\Omega} \sim \mathcal{W}_p(\mathbf{I}_p, \alpha)$ then $f(\boldsymbol{\omega}_{\bullet,p}, \omega_{pp}) = f(\boldsymbol{\omega}_{\bullet,p} \mid \omega_{pp})f(\omega_{pp})$, where,
 $\boldsymbol{\omega}_{\bullet,p} \mid \omega_{pp} \sim \mathcal{N}(0, \omega_{pp}\mathbf{I}_{p-1})$, $\omega_{pp} \sim \text{Gamma}(\text{shape} = \alpha/2, \text{rate} = 1/2)$.
- Computing III_p is easy: normal \times gamma.

Computing $IV_p (= \log f(\boldsymbol{\theta}_p | \mathbf{y}_{1:p}))$

- Decompose $\mathbf{S} = \mathbf{y}^T \mathbf{y}$ analogous to $\boldsymbol{\Omega}$ and reparameterize $(\boldsymbol{\omega}_{\bullet,p}, \omega_{pp}) \mapsto (\boldsymbol{\beta}_{\bullet,p}, \gamma_{pp})$:

$$\mathbf{S} = \begin{bmatrix} \mathbf{S}_{(p-1) \times (p-1)} & \mathbf{s}_{\bullet,p} \\ \mathbf{s}_{\bullet,p}^T & s_{pp} \end{bmatrix}, \boldsymbol{\beta}_{\bullet,p} = \boldsymbol{\omega}_{\bullet,p}, \gamma_{pp} = \omega_{pp} - \boldsymbol{\omega}_{\bullet,p}^T \boldsymbol{\Omega}_{(p-1) \times (p-1)}^{-1} \boldsymbol{\omega}_{\bullet,p}.$$

- Key result of Wang (2012, BA):

$$f(\boldsymbol{\beta}_{\bullet,p}, \gamma_{pp} | \text{rest}) = \mathcal{N}(\boldsymbol{\beta}_{\bullet,p} | -\mathbf{C}\mathbf{s}_{\bullet,p}, \mathbf{C}) \times \text{G}\left(\gamma_{pp} \mid \frac{n + \alpha - p - 1}{2} + 1, \frac{s_{pp} + 1}{2}\right).$$

where $\mathbf{C} = \{(s_{pp} + 1)\boldsymbol{\Omega}_{(p-1) \times (p-1)}^{-1}\}^{-1}$

- Allows computation of IV_p using Chib's two block strategy.

Computing $\text{III}_{p-1}, \dots, \text{III}_1$ and $\text{IV}_{p-1}, \dots, \text{IV}_1$

- Same strategy as in going from I_p to I_{p-1} .
- Proceed backwards starting from the p th row. At each step, adjust the upper left sub-matrix $\Omega_{j \times j}$ via IPS:

for ($j=p-1, \dots, 1$) **do**

$$\text{Update } \Omega_{j \times j} \leftarrow \Omega_{j \times j} - \frac{\omega_{\bullet(j+1)} \omega_{\bullet(j+1)}^T}{\omega_{j+1,j+1}}.$$

end for

- Calculate III_j and IV_j with the updated $\Omega_{j \times j}$.

Results for Wishart

Dimension and Parameters	Truth	Proposed	AIS	Nested
$(p = 5, n = 10, \alpha = 7)$	-84.13	-84.13 (-0.02)	-84.3 (0.68)	-84.26 (0.57)
$(p = 10, n = 20, \alpha = 13)$	-365.11	-365.12 (0.06)	-397.64 (6.1)	-392.2 (6.04)
$(p = 15, n = 30, \alpha = 20)$	-837.7	-837.67 (0.13)	-1000.45 (13.5)	-994.87 (13.7)
$(p = 25, n = 50, \alpha = 33)$	-2417.65	-2417.14 (1.11)	$-\infty$	$-\infty$
$(p = 30, n = 60, \alpha = 39)$	-3553.62	-3548.02 (3.04)	$-\infty$	$-\infty$

Table: Mean (sd) of estimated log marginal for Wishart for the proposed approach, AIS, nested sampling; under 25 random permutations of the nodes $\{1, \dots, p\}$ using 5000 samples.

Evidence under element-wise priors

- Clearly, we did not get into all this trouble just for Wishart!
- Consider the element-wise prior:

$$f(\mathbf{\Omega} \mid \lambda) = C^{-1} \prod_{i < j} f(\omega_{ij} \mid \lambda) \prod_{j=1}^p f(\omega_{jj} \mid \lambda) \mathbb{1}(\mathbf{\Omega} \in \mathcal{M}_p^+).$$

- Two examples (global-local shrinkage priors):
 - Bayesian graphical lasso (BGL):

$$f(\omega_{ij} \mid \lambda) = (\lambda/2) \exp(-\lambda|\omega_{ij}|)$$

\Updownarrow (Andrews and Mallows, 1974)

$$\omega_{ij} \mid \tau_{ij}, \lambda \sim \mathcal{N}(0, \tau_{ij}), \quad \tau_{ij} \mid \lambda \sim \text{Exp}(\lambda^2/2).$$

- Graphical horseshoe (GHS):

$$\omega_{ij} \mid \tau_{ij}, \lambda \sim \mathcal{N}(0, \tau_{ij}), \quad \tau_{ij} \mid \lambda \sim \mathcal{C}^+(0, \lambda).$$

Evidence under element-wise priors

- The off-diagonal ω_{ij} terms are normal **conditional on** τ_{ij} .
- Similarly, the diagonal ω_{jj} are exponential.
- The presence of these mixing τ_{ij} variables is the **ONLY** difference with the Wishart case for our purposes.
- The τ_{ij} terms can be sampled easily.
- **MAIN IDEA:** Absorb the τ_{ij} terms into Chib's latent z (they are sampled, but their densities are not evaluated). Conditional on these, evaluate the normal and gamma densities exactly as in Wishart.

Computing $IV_p (= \log f(\boldsymbol{\theta}_p \mid \mathbf{y}_{1:p}))$

- We have

$$f(\boldsymbol{\beta}_{\bullet,p}, \gamma_{pp} \mid \boldsymbol{\tau}_{\bullet,p}, \boldsymbol{\Omega}_{(p-1) \times (p-1)}, \mathbf{y}_{1:p}) = \mathcal{N}(\boldsymbol{\beta}_{\bullet,p} \mid -\mathbf{C}\boldsymbol{s}_{\bullet,p}, \mathbf{C}) \\ \times \text{Gamma}\left(\gamma_{pp} \mid \frac{n}{2} + 1, \frac{s_{pp} + \lambda}{2}\right),$$

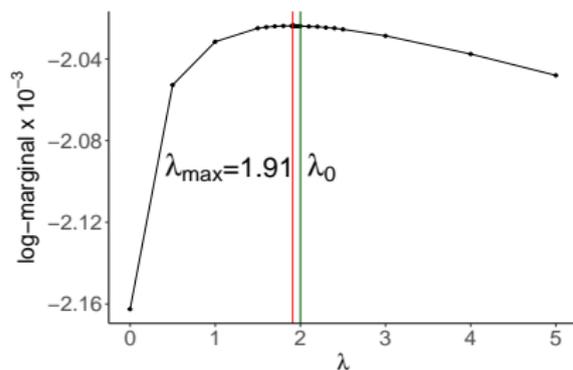
$$\text{where } \mathbf{C} = \{\text{diag}^{-1}(\boldsymbol{\tau}_{\bullet,p}) + (s_{pp} + \lambda)\boldsymbol{\Omega}_{(p-1) \times (p-1)}^{-1}\}^{-1}$$

- Recall, for Wishart we had

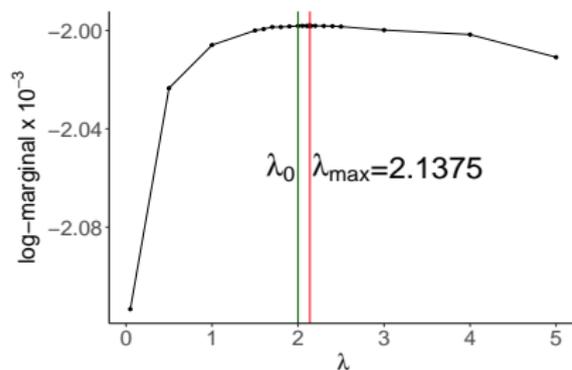
$$f(\boldsymbol{\beta}_{\bullet,p}, \gamma_{pp} \mid \boldsymbol{\Omega}_{(p-1) \times (p-1)}, \mathbf{y}_{1:p}) = \mathcal{N}(\boldsymbol{\beta}_{\bullet,p} \mid -\mathbf{C}\boldsymbol{s}_{\bullet,p}, \mathbf{C}) \\ \times \text{G}\left(\gamma_{pp} \mid \frac{n + \alpha - p - 1}{2} + 1, \frac{s_{pp} + 1}{2}\right).$$

$$\text{where } \mathbf{C} = \{(s_{pp} + 1)\boldsymbol{\Omega}_{(p-1) \times (p-1)}^{-1}\}^{-1}.$$

Results



(a) BGL



(b) GHS

Figure: Log marginal vs. λ under (a) BGL and (b) GHS ($p = 10, n = 150$).

	λ	0.05	1	2 ($= \lambda_0$)	3	4	5
log BF	BGL	138.84	7.86	0.18	4.98	13.9	24.34
	GHS	115.31	7.89	0.12	1.79	3.63	12.83

Table: Logarithm of Bayes factors.

Additional results and applications

- The strategy also works for calculating evidence under G-Wishart priors.
- Results are quite competitive with current state of the art (Atay-Kayis and Massam, 2005)
- As a by product, we are also able to develop a **new row-wise sampler** for G-Wishart that does not require a maximal clique decomposition.
- Details in the paper.

Concluding remarks

- The strategy developed will work whenever: (a) the priors on the off-diagonals of Ω are scale mixtures of normal and (b) the diagonals of Ω are scale mixtures of exponential.
- These are very mild requirements and can handle a broad class of priors.
- Although we did not do so in this paper, one may also shift focus from the prior to likelihood that are mixtures of normal! Consider $\mathbf{y} \sim t_\nu(\boldsymbol{\mu}, \Omega^{-1})$
- This is equivalent to $\mathbf{y} \mid \tau \sim \mathcal{N}(\boldsymbol{\mu}, \tau^{-1}\Omega^{-1})$, $\tau \sim \text{Gamma}(\nu/2, \nu/2)$.
- Should be possible to absorb the τ in the likelihood into Chib's z and proceed.

Main references

- **Bhadra, A.**, Sagar, K., Banerjee, S. and Datta, J. (2022+). Graphical Evidence. (submitted). [arXiv:2205.01016]
 - Code: https://github.com/sagarknk/Graphical_Evidence
- Chib, S. (1995). Marginal likelihood from the Gibbs output. *Journal of the American Statistical Association* **90**, 1313–1321.
- Wang, H. (2012). Bayesian graphical lasso models and efficient posterior computation. *Bayesian Analysis* **7**, 867–886.