



Marginal Likelihood from the Gibbs output  
- Chib (1995)

Journal of American Statistical Association

**Research Question** - The problem of calculating marginal distribution (Normalizing constant) of the posterior density, which is an input in calculating Bayes factor, has proved extremely challenging as it involves integrating the likelihood function with respect to the prior density.

**Formulae :**  $m(y|M_k) = \int \underbrace{f(y|\theta_k, M_k)}_{\text{likelihood or density function of } y \text{ for } \theta_k \text{ and model } k} \underbrace{\pi(\theta_k|M_k)}_{\text{prior distribution for } \theta_k \text{ for model } k} d\theta_k$

Integrating over  $\theta_k$  to get marginal

$$f(y|\theta_k, M_k) \pi(\theta_k|M_k) = \underbrace{\pi(y, \theta_k|M_k)}_{\text{joint distribution}}$$

**Previous studies solutions**

1) Newton and Raftery (1994)

$$\hat{m}_{NR} = \left\{ \frac{1}{G} \sum_{i=1}^G \left( \frac{1}{f(y|\theta_k^{(i)}, M_k)} \right) \right\}^{-1}$$

Harmonic mean of the likelihood functions  
evaluated at  $\theta_k^{(1)}, \theta_k^{(2)} \dots \theta_k^{(g)}$

limitations - not stable as inverse likelihood does not have finite variance.

## 2) Gelfand and Dey (1993)

density with thin tails than the denominator. (tuning function)

$$\hat{m}_{GD} = \left\{ \frac{1}{G} \sum_{g=1}^G \left( \frac{\overbrace{p(\theta_k^{(g)})}}{\underbrace{f(y|\theta_k^{(g)}, M_k)} \pi(\theta_k^{(g)}|M_k)} \right) \right\}^{-1}$$

limitation: Tuning function can be hard to find in high dimensional problems.

## Derivation of the Approach

Let's set up the problem

$\theta$  = parameters of interest

$y$  = data

$z$  = latent variable (To allow for data augmentation).

Data augmentation — ??

It refers to the scheme augmenting the observed data. We do it when there are missing values or truncated data. Sometimes it is difficult or impossible to sample  $y$  directly, but there exist

a latent variable " $y^*$ " s.t it is possible to conditionally sample  $y|y^*$  and  $y^*|y$ .

Example: Probit regression explained later.

### Basic marginal likelihood identity (BML)

$$\text{Bayes rule : } \pi(\theta|y) = \frac{f(y|\theta) \pi(\theta)}{m(y)}$$

$$m(y) = \frac{f(y|\theta) \pi(\theta)}{\pi(\theta|y)} \quad - \text{BML}$$

Now, in order to calculate LHS, we need to evaluate RHS. We can evaluate RHS at some fixed value of  $\theta$  say  $\theta^*$  as we already know likelihood, prior and posterior. For that we can exploit the already existing information that we have from Gibbs sampler -  $\{\pi(\theta_r|y, \theta_s(s \neq r), z)\}_{r=1}^B$

On the logarithmic scale

$$\ln \hat{m}(y) = \ln f(y|\theta^*) + \ln \pi(\theta^*) - \ln \pi(\theta^*|y)$$

where  $\hat{m}(y)$  is the estimator for  $m(y)$ .

The identity is valid for any  $\theta^*$  but it is recommended to take  $\theta^*$  from high density for efficiency.

## Estimating $\pi(\theta^*|y)$ using Gibbs Sampler

### Two vector Block

Full conditionals :  $\pi(\theta|y, z)$  ,  $p(z|y, \theta)$

Gibbs sampler output :  $\{\theta^{(g)}, z^{(g)}\}_{g=1}^G$

Posterior density which we are interested in (does not include  $z$ )

$$\begin{aligned}\pi(\theta|y) &= \int \overbrace{\pi(\theta, z|y)}^{\text{target density in Gibbs sampler}} dz \\ &= \int \overbrace{\pi(\theta|y, z)}^{\text{we know the full conditional}} \underbrace{\pi(z|y)} dz\end{aligned}$$

Monte Carlo Estimate of  $\pi(\theta|y)$  at  $\theta^*$

$$\hat{\pi}(\theta^*|y) = G^{-1} \sum_{i=1}^G \pi(\theta^*|y, z^{(g)})$$

$z^{(g)}$  is distributed from  $\pi(z|y)$  as we are drawing  $\theta$  and  $z$  from target distribution which is  $\pi(\theta, z|y)$  which can be written as  $\pi(\theta|y, z) \pi(z|y)$ . Therefore,  $\theta$  is coming from left value &  $z$  is coming from right value. We could have written it as  $\pi(z|y, \theta)$  as well.

Under regularity conditions. [ Similar to gauss Markov assumption for OLS)

$$\hat{\pi}(\theta^*|y) \xrightarrow{a.s.} \pi(\theta^*|y)$$

Recap of law of large numbers

Convergence in probability :  $z_n \xrightarrow{P} z$

$$\lim_{n \rightarrow \infty} P(|z_n - z| \leq \delta) = 1 \quad \forall \delta$$

Convergence almost surely :  $z_n \xrightarrow{a.s.} z$

$$P\left(\lim_{n \rightarrow \infty} |z_n - z| \leq \delta\right) = 1$$

Regularity conditions

- 1)  $y_i$  is i.i.d with density  $f(y; \theta)$
- 2) parameter space is compact (closed and bounded)
- 3)  $\theta_0 = \max_{\theta \in \Theta} E_{\theta_0} \log f(Y_i; \theta)$  where  $\theta_0$  = true parameter
- 4) Likelihood function is continuous
- 5)  $E_{\theta_0} \log f(Y_i; \theta)$  exist
- 6)  $\sup_{\theta \in \Theta} \left| \frac{1}{n} L(y; \theta) - E_{\theta_0} \log f(Y_i; \theta) \right| < \delta$  almost surely  $\forall \delta$ .

∴ estimation of marginal likelihood will look like.

$$\ln \hat{m}(y) = \ln f(y|\theta^*) + \ln \pi(\theta^*) - \ln \left\{ G^{-1} \sum_{j=1}^G \pi(\theta^* | y, z^{(j)}) \right\}$$

### Three vector Block

Full conditionals:  $\pi(\theta_1 | y, \theta_2, z)$ ,  $\pi(\theta_2 | y, \theta_1, z)$ ,  $p(z | y, \theta)$   
 $\theta = \{\theta_1, \theta_2\}$

Again, the objective is to find  $\pi(\theta^* | y)$  which will go in the derivation of marginal likelihood.

Posterior density which we are interested in (does not include  $z$ )

$$\pi(\theta_1^*, \theta_2^* | y) = \pi(\overset{\textcircled{1}}{\theta_1^*} | y) \pi(\overset{\textcircled{2}}{\theta_2^*} | \theta_1^*, y)$$

Now,  $\pi(\overset{\textcircled{1}}{\theta_1^*} | y) = \iint \underbrace{\pi(\theta_1^*, \theta_2, z | y)}_{\text{target density in Gibbs sampler}} dz d\theta_2$

$$= \iint \underbrace{\pi(\theta_1^* | \theta_2, z, y)}_{\text{we know the full conditionals}} \underbrace{\pi(z, \theta_2 | y)} dz d\theta_2$$

$$\therefore \hat{\pi}(\theta_1^* | y) = G^{-1} \sum_{j=1}^G \pi(\theta_1^* | \theta_2^{(j)}, z^{(j)}, y)$$

evaluated at  $\theta_1^*$  and  $\theta_2$  &  $z$  are sampled from  $\pi(z, \theta_2 | y)$  ←

Now Consider (2)

$$\pi(\theta_2^* | \theta_1^*, y) = \int \pi(\theta_2^* | \theta_1^*, y, z) \pi(z | \theta_1^*, y) dz$$

Draws of  $z$  from Gibbs sampler is from  $z|y$  and not  $z|\theta_1, y$

$$\pi(\theta_1^*, \theta_2^*, z | y) = \pi(\theta_1^* | \theta_2^*, z, y) \pi(\theta_2^* | y, z) \pi(z | y)$$

In Gibbs sampler we are sampling  $\theta_1, \theta_2$  and  $z$  from  $\pi(\theta_1^*, \theta_2^*, z | y)$  or in other words sampling  $\theta_1^*, \theta_2^*$  and  $z$  from RHS components respectively. RHS can be written in other sequences of  $\theta_1, \theta_2$  &  $z$  as well but our problem set up demands us to make the RHS in this form only.

How to draw  $z$  from  $f(z | \theta_1^*, y)$ ?

Continue the original gibbs sampler for  $G$  runs with full conditionals.

$$\pi(\theta_2 | y, \theta_1^*, z) \quad \& \quad p(z | y, \theta_1^*, \theta_2) \quad - (*)$$

Then, Sampling will be from  $\pi(\theta_2, z | y, \theta_1^*)$   
 $= \pi(\theta_2 | z, y, \theta_1^*) \pi(z | y, \theta_1^*)$



$$\therefore \hat{\pi}(\theta_2^* | y, \theta_1^*) = G^{-1} \sum_{i=1}^G \pi(\theta_2^* | y, \theta_1^*, z^{(i)})$$

where  $z^{(j)}$  are sampled from  $(*)$

$$\ln \hat{m}(y) = \ln f(y | \theta^*) + \ln \pi(\theta^*) - \ln \hat{\pi}(\theta_1^* | y) - \ln(\theta_2^* | y, \theta_1^*)$$

### Multiple Blocks

$$\begin{aligned} \text{posterior density} &= \pi(\theta^* | y) = \pi(\theta_1^*, \theta_2^* \dots \theta_B^* | y) \\ &= \underbrace{\pi(\theta_1^* | y)} \pi(\theta_2^* | y, \theta_1^*) \dots \pi(\theta_B^* | y, \theta_1^* \dots \theta_{B-1}^*) \end{aligned}$$

This can be estimated  
using usual Gibbs sampler  
draws from full conditionals

$$\pi(\theta_1^* | y) = \iint \dots \int \underbrace{\pi(\theta_1^*, \theta_2 \dots \theta_B, z | y)}_{\text{usual draw from gibbs sampler}} d\theta_2 d\theta_3 \dots d\theta_B dz$$

$$\pi(\theta_1^* | y) = \iint \dots \int \pi(\theta_1^* | \theta_2 \dots \theta_B, z, y) \pi(\theta_2 | \theta_3 \dots \theta_B, z, y) \dots \pi(z | y) d\theta_2 d\theta_3 \dots dz$$

$$\hat{\pi}(\theta_1^* | y) = G^{-1} \sum_{i=1}^G \pi(\theta_1^* | \theta_2^{(i)}, \theta_3^{(i)} \dots, \theta_B^{(i)}, z^{(i)}, y)$$

↓  
Sampled from  
 $\pi(\theta_2 | y, \theta_3 \dots \theta_B, z, y)$

↓  
Sampled from  
 $\pi(\theta_B | z, y)$   
Sampled from  $\pi(z | y)$

→ Sampled from  $\pi(\theta_3 | y, \theta_4 \dots \theta_B, z)$

Any other general term is written as

$$\pi(\theta_r^* | y, \theta_1^*, \theta_2^*, \dots, \theta_{r-1}^*) = \iiint \pi(\theta_r^*, \theta_x(l > r), z | y, \theta_1^* \dots \theta_{r-1}^*)$$

conditional density

$$= \iiint \pi(\theta_r^* | \theta_x(l > r), z, y, \theta_1^* \dots \theta_{r-1}^*)$$

$$\pi(\theta_x(l > r) | z, y, \dots \theta_{r-1}^*) \dots \pi(z | y, \theta_1^* \dots \theta_{r-1}^*) d\theta_x(l > r) dz$$

$\theta_x(l > r)$  and  $z$  are not sample from the above distributions.

$z$  is sampled from  $z | y$  rather than  $z | y, \theta_1^* \dots \theta_{r-1}^*$ .  
 $\theta_{r+1}$  is sampled from  $\theta_{r+1} | r+2, \dots, r_B, z, y$  rather than  $\theta_{r+1} | r+2, \dots, r_B, z, y, \theta_1^* \dots \theta_{r-1}^*$ .

So, to estimate the term, continue sampling with conditional densities of  $\{\theta_r, \theta_{r+1}, \dots, \theta_B, z\}$  and substitute  $\theta_1^*, \theta_2^* \dots \theta_{r-1}^*$  instead of  $\theta_1, \theta_2 \dots \theta_{r-1}$ . (\*\*)

$$\hat{\pi}(\theta_r | y, \theta_s^*(s < r)) = G^{-1} \sum_{i=1}^G \pi(\theta_r^* | y, \theta_1^*, \theta_2^* \dots \theta_{r-1}^*, \theta_x^{(ij)}(l > r), z^{(ij)})$$

where  $\theta_x^{(ij)}(l > r)$  and  $z^{(ij)}$  are the results from extra sampling done in (\*\*)

$$\ln \hat{m}(y) = \ln f(y | \theta^*) + \ln \pi(\theta^*) - \sum_{r=1}^B \ln \hat{\pi}(\theta_r^* | y, \theta_s^*(s < r))$$

## Bayes factor

Bayes factor to compare model  $k$  and  $j$

$$\hat{B}_{kj} = \exp \{ \ln \hat{\pi}(y | M_k) - \ln \hat{\pi}(y | M_j) \}$$

## Application: Binary Probit model

$$y = \begin{cases} 1 & \text{if cancer has spread} \\ 0 & \text{otherwise} \end{cases}$$

$x_1$  = age

$x_2$  = level of serum acid phosphatase

$x_3$  = X ray result (Dummy)

$x_4$  = Size of tumor (Dummy)

$x_5$  = Grade of tumor (Dummy)

53 observations

Examined 9 different models

$$\text{Probit model: } \Pr(y=1 | M_k) = \underbrace{\phi(x_{ik}' \beta_k)}_{\text{CDF of standard normal distribution}}$$

as  $y$  is binary, the likelihood of  $y$  is multiplying iid bernoulli distribution for each  $i$ .

$$f(y|M_k, \beta_k) = \prod_{i=1}^{53} [\Phi(x_i' \beta_k)]^{y_i} [1 - \Phi(x_i' \beta_k)]^{1-y_i}$$

Setting up the gibbs sampler (for any model k)

Define a latent variable  $z$  s.t

$$z_i \sim N(x_i' \beta, 1) ; \quad y_i = \underbrace{I(z_i > 0)}_{\text{indicator function}}$$

$$\Pr(z > 0) = \Phi(x_i' \beta)$$

$\pi(\beta)$  - prior distribution

$\pi(\beta | y, z)$  - posterior distribution (calculated)

$$p(z_i | y, \beta) \propto \phi(z_i | x_i' \beta, 1) I[0, \infty] \quad \text{if } y_i = 1$$

$$\propto \phi(z_i | x_i' \beta, 1) I(-\infty, 0] \quad \text{if } y_i = 0$$

where  $\phi(\cdot | \mu, 1) I[a, b]$  is the normal density truncated to the interval  $[a, b]$ .

then log of marginal likelihood of model  $k$  is

$$\ln f(y|M_k) = \ln f(y|M_k, \beta_k^*) + \ln \pi(\beta_k)$$

$$= \ln \left\{ \frac{1}{5000} \sum_{g=1}^{5000} \phi(\beta_k^* | y, z^{(g)}) \right\}$$

Bayes factor can be calculated to compare models