

VIII. Model selection

A. Marginal likelihood

Suppose we're trying to choose among a series of models:

Model 1: $p(\mathbf{y}|\theta_1)$

⋮

Model M : $p(\mathbf{y}|\theta_M)$

where θ_m are possibly of different dimension

The Bayesian might think in terms of an unobserved random variable:

$s = 1$ if Model 1 is true

⋮

$s = M$ if Model M is true

and assign prior probabilities

$$\pi_1 = \Pr(s = 1)$$

⋮

$$\pi_M = \Pr(s = M)$$

with associated priors on the parameters

$$p(\boldsymbol{\theta}_1|s = 1)$$

⋮

$$p(\boldsymbol{\theta}_M|s = M)$$

From such a perspective, the probability that Model m is true given the data is

$$\begin{aligned} p(s = m|\mathbf{y}) &= \frac{\pi_m \int p(\mathbf{y}|\boldsymbol{\theta}_m)p(\boldsymbol{\theta}_m|s = m)d\boldsymbol{\theta}_m}{\sum_{j=1}^M \pi_j \int p(\mathbf{y}|\boldsymbol{\theta}_j)p(\boldsymbol{\theta}_j|s = j)d\boldsymbol{\theta}_j} \\ &\equiv \frac{\pi_m p_m(\mathbf{y})}{\sum_{j=1}^M \pi_j p_j(\mathbf{y})} \end{aligned}$$

The expression

$$p_m(\mathbf{y}) = \int p(\mathbf{y}|\boldsymbol{\theta}_m)p(\boldsymbol{\theta}_m|s = m)d\boldsymbol{\theta}_m$$

is sometimes called the "marginal likelihood" of Model m

The Bayesian would say that the data favor the model for which $p(s = m|\mathbf{y})$ is biggest. With diffuse priors ($\pi_m = 1/M$) this is equivalent to choosing the model with the highest marginal likelihood.

VIII. Model selection

- A. Marginal likelihood
- B. Schwarz criterion

First let's examine the behavior of $p_m(\mathbf{y}) = \int p(\mathbf{y}|\boldsymbol{\theta}_m)p(\boldsymbol{\theta}_m|s = m)d\boldsymbol{\theta}_m$ as the sample size T gets large

Suppose

$$\log p(\mathbf{y}|\boldsymbol{\theta}) = \sum_{t=1}^T \log p(\mathbf{y}_t|\boldsymbol{\theta})$$

and let $\hat{\boldsymbol{\theta}}_T$ denote the MLE

$$\hat{\boldsymbol{\theta}}_T = \arg \max \log p(\mathbf{y}|\boldsymbol{\theta})$$

Recall Taylor's Theorem:

$$\log p(\mathbf{y}|\boldsymbol{\theta}) = \log p(\mathbf{y}|\hat{\boldsymbol{\theta}}_T) - \frac{1}{2} \sqrt{T} (\boldsymbol{\theta} - \hat{\boldsymbol{\theta}}_T)' \times \\ T^{-1} \sum_{t=1}^T \mathbf{H}_t(\tilde{\boldsymbol{\theta}}_T) \sqrt{T} (\boldsymbol{\theta} - \hat{\boldsymbol{\theta}}_T)$$

$$\tilde{\boldsymbol{\theta}}_T = \lambda_T \boldsymbol{\theta} + (1 - \lambda_T) \hat{\boldsymbol{\theta}}_T$$

$$\mathbf{H}_t(\boldsymbol{\theta}) \equiv -\frac{\partial^2 \log p(\mathbf{y}_t|\boldsymbol{\theta})}{\partial \boldsymbol{\theta} \partial \boldsymbol{\theta}'}$$

Let λ_2 denote smallest eigenvalue of $T^{-1} \sum_{t=1}^T \mathbf{H}_t(\tilde{\boldsymbol{\theta}}_T)$ in neighborhood around true value $\boldsymbol{\theta}^*$, so that in this neighborhood,

$$\log p(\mathbf{y}|\boldsymbol{\theta}) = \log p(\mathbf{y}|\hat{\boldsymbol{\theta}}_T) - \frac{1}{2} \sqrt{T} (\boldsymbol{\theta} - \hat{\boldsymbol{\theta}}_T)' \times \\ T^{-1} \sum_{t=1}^T \mathbf{H}_t(\tilde{\boldsymbol{\theta}}_T) \sqrt{T} (\boldsymbol{\theta} - \hat{\boldsymbol{\theta}}_T) \\ \leq \log p(\mathbf{y}|\hat{\boldsymbol{\theta}}_T) - \frac{T\lambda_2}{2} (\boldsymbol{\theta} - \hat{\boldsymbol{\theta}}_T)' (\boldsymbol{\theta} - \hat{\boldsymbol{\theta}}_T)$$

If also there exists B such that

$B \geq p(\theta)$, we should have

$$p(\mathbf{y}) = \int p(\mathbf{y}|\theta)p(\theta)d\theta$$

$$\leq \int \exp\left\{\log p(\mathbf{y}|\hat{\theta}_T) - \frac{T\lambda_2}{2}(\theta - \hat{\theta}_T)'(\theta - \hat{\theta}_T)\right\} B d\theta$$

$$= Bp(\mathbf{y}|\hat{\theta}_T) \int \exp\left\{-\frac{T\lambda_2}{2}(\theta - \hat{\theta}_T)'(\theta - \hat{\theta}_T)\right\} d\theta$$

But

$$\int \exp\left\{-\frac{T\lambda_2}{2}(\theta - \hat{\theta}_T)'(\theta - \hat{\theta}_T)\right\} d\theta$$

$$= \left[\frac{2\pi}{T\lambda_2}\right]^{-k/2} \times$$

$$\int \left[\frac{T\lambda_2}{2\pi}\right]^{k/2} \exp\left\{-\frac{T\lambda_2}{2}(\theta - \hat{\theta}_T)'(\theta - \hat{\theta}_T)\right\} d\theta$$

$$= \left[\frac{2\pi}{T\lambda_2}\right]^{-k/2}$$

for k the dimension of θ

Conclusion:

$$p(\mathbf{y}) = \int p(\mathbf{y}|\theta)p(\theta)d\theta$$

$$\leq Bp(\mathbf{y}|\hat{\theta}_T) \left[\frac{2\pi}{T\lambda_2}\right]^{-k/2}$$

$$\log p(\mathbf{y}) \leq \log p(\mathbf{y}|\hat{\theta}_T) - (k/2) \log T + R_{2k}$$

$$R_{2k} = \log \left\{ B \left[\frac{2\pi}{\lambda_2}\right]^{-k/2} \right\}$$

R_{2k} does not depend on T

Similar argument reasons that

$$\begin{aligned} \log p(\mathbf{y}|\boldsymbol{\theta}) &= \log p(\mathbf{y}|\hat{\boldsymbol{\theta}}_T) - \frac{1}{2} \sqrt{T} (\boldsymbol{\theta} - \hat{\boldsymbol{\theta}}_T)' \times \\ &\quad T^{-1} \sum_{t=1}^T \mathbf{H}_t(\tilde{\boldsymbol{\theta}}_T) \sqrt{T} (\boldsymbol{\theta} - \hat{\boldsymbol{\theta}}_T) \\ &\geq \log p(\mathbf{y}|\hat{\boldsymbol{\theta}}_T) - \frac{T\lambda_1}{2} (\boldsymbol{\theta} - \hat{\boldsymbol{\theta}}_T)' (\boldsymbol{\theta} - \hat{\boldsymbol{\theta}}_T) \end{aligned}$$

for λ_1 biggest eigenvalue of

$T^{-1} \sum_{t=1}^T \mathbf{H}_t(\tilde{\boldsymbol{\theta}}_T)$ in neighborhood
around true value $\boldsymbol{\theta}^*$

$$\log p(\mathbf{y}) \geq \log p(\mathbf{y}|\hat{\boldsymbol{\theta}}_T) - (k/2) \log T + R_{1k}$$

$$\log p(\mathbf{y}) \leq \log p(\mathbf{y}|\hat{\boldsymbol{\theta}}_T) - (k/2) \log T + R_{2k}$$

Implication: for large T we have the
approximation

$$\log p(\mathbf{y}) \simeq \log p(\mathbf{y}|\hat{\boldsymbol{\theta}}_T) - (k/2) \log T$$

for k the dimension of $\boldsymbol{\theta}$

Choosing the model m for which

$$\log p(\mathbf{y}|\hat{\boldsymbol{\theta}}_{mT}) - (k_m/2) \log T$$

is biggest is known as the using the

Schwarz Information Criterion (SIC)

or Bayesian Information Criterion (BIC)

Since it is asymptotically a Bayesian decision rule, SIC inherits the properties of being asymptotically admissible and consistent

However, note that this result required the same regularity conditions needed to get asymptotic Normality of MLE

VIII. Model selection

- A. Marginal likelihood
- B. Schwarz criterion
- C. Calculating the marginal likelihood with the Gibbs sampler

Goal: calculate

$$p(\mathbf{y}) = \int p(\mathbf{y}|\boldsymbol{\theta})p(\boldsymbol{\theta})d\boldsymbol{\theta}$$

Couldn't we get this by drawing $\boldsymbol{\theta}^{(j)}$ $j = 1, \dots, J$ from $p(\boldsymbol{\theta})$ and then

$$\hat{p}(\mathbf{y}) = J^{-1} \sum_{j=1}^J p(\mathbf{y}|\boldsymbol{\theta}^{(j)})?$$

Answer: no, this algorithm is badly behaved numerically.

Chib's idea: think of evaluating at a point with a lot of mass (say the posterior mean $\boldsymbol{\theta}^*$).

Note that for any $\boldsymbol{\theta}^*$ we have the identity

$$p(\boldsymbol{\theta}^*|\mathbf{y})p(\mathbf{y}) = p(\mathbf{y}|\boldsymbol{\theta}^*)p(\boldsymbol{\theta}^*)$$

$$p(\mathbf{y}) = p(\mathbf{y}|\boldsymbol{\theta}^*)p(\boldsymbol{\theta}^*)/p(\boldsymbol{\theta}^*|\mathbf{y})$$

In many applications, we know $p(\mathbf{y}|\boldsymbol{\theta}^*)$ and $p(\boldsymbol{\theta}^*)$ analytically (evaluating the likelihood and prior at posterior mean, respectively), but couldn't calculate $p(\boldsymbol{\theta}^*|\mathbf{y})$ explicitly

Suppose we've generated draws from a two-block Gibbs sampler:

$$p(\theta_1|\theta_2, \mathbf{y}) \text{ and } p(\theta_2|\theta_1, \mathbf{y})$$

The object of interest is given by

$$p(\theta_1^*, \theta_2^*|\mathbf{y}) = p(\theta_1^*|\theta_2^*, \mathbf{y})p(\theta_2^*|\mathbf{y})$$

where we may know $p(\theta_1^*|\theta_2^*, \mathbf{y})$ analytically.

We know that

$$p(\theta_2^*|\mathbf{y}) = \int p(\theta_2^*|\theta_1, \mathbf{y})p(\theta_1|\mathbf{y})d\theta_1$$

and can therefore estimate

$$\hat{p}(\theta_2^*|\mathbf{y}) = G^{-1} \sum_{g=1}^G p(\theta_2^*|\theta_1^{(g)}, \mathbf{y})$$

that is, the average value of

$p(\theta_2^*|\theta_1, \mathbf{y})$ across Gibbs simulated draws for θ_1

Conclusion: for a two-block Gibbs sampler, we can estimate the marginal likelihood from

$$\hat{p}(\mathbf{y}) = \frac{p(\mathbf{y}|\theta_1^*, \theta_2^*)p(\theta_1^*)p(\theta_2^*)}{p(\theta_1^*|\theta_2^*, \mathbf{y})G^{-1} \sum_{g=1}^G p(\theta_2^*|\theta_1^{(g)}, \mathbf{y})}$$

How about 3 blocks? Now we want to estimate the denominator of

$$p(\mathbf{y}) = p(\mathbf{y}|\boldsymbol{\theta}^*)p(\boldsymbol{\theta}^*)/p(\boldsymbol{\theta}_1^*, \boldsymbol{\theta}_2^*, \boldsymbol{\theta}_3^*|\mathbf{y})$$

$$p(\boldsymbol{\theta}_1^*, \boldsymbol{\theta}_2^*, \boldsymbol{\theta}_3^*|\mathbf{y})$$

$$= p(\boldsymbol{\theta}_1^*|\mathbf{y})p(\boldsymbol{\theta}_2^*|\boldsymbol{\theta}_1^*, \mathbf{y})p(\boldsymbol{\theta}_3^*|\boldsymbol{\theta}_1^*, \boldsymbol{\theta}_2^*, \mathbf{y})$$

$$p(\boldsymbol{\theta}_1^*, \boldsymbol{\theta}_2^*, \boldsymbol{\theta}_3^*|\mathbf{y})$$

$$= p(\boldsymbol{\theta}_1^*|\mathbf{y})p(\boldsymbol{\theta}_2^*|\boldsymbol{\theta}_1^*, \mathbf{y})p(\boldsymbol{\theta}_3^*|\boldsymbol{\theta}_1^*, \boldsymbol{\theta}_2^*, \mathbf{y})$$

First term can be estimated as before:

$$\hat{p}(\boldsymbol{\theta}_1^*|\mathbf{y}) = G^{-1} \sum_{g=1}^G p(\boldsymbol{\theta}_1^*|\boldsymbol{\theta}_2^{(g)}, \boldsymbol{\theta}_3^{(g)}, \mathbf{y})$$

Third term $p(\boldsymbol{\theta}_3^*|\boldsymbol{\theta}_1^*, \boldsymbol{\theta}_2^*, \mathbf{y})$ is known analytically

$$p(\boldsymbol{\theta}_1^*, \boldsymbol{\theta}_2^*, \boldsymbol{\theta}_3^*|\mathbf{y})$$

$$= p(\boldsymbol{\theta}_1^*|\mathbf{y})p(\boldsymbol{\theta}_2^*|\boldsymbol{\theta}_1^*, \mathbf{y})p(\boldsymbol{\theta}_3^*|\boldsymbol{\theta}_1^*, \boldsymbol{\theta}_2^*, \mathbf{y})$$

Second term:

$$p(\boldsymbol{\theta}_2^*|\boldsymbol{\theta}_1^*, \mathbf{y}) = \int p(\boldsymbol{\theta}_2^*|\boldsymbol{\theta}_1^*, \boldsymbol{\theta}_3, \mathbf{y})p(\boldsymbol{\theta}_3|\boldsymbol{\theta}_1^*, \mathbf{y})d\boldsymbol{\theta}_3$$

But how do we generate a sample from $p(\boldsymbol{\theta}_3|\boldsymbol{\theta}_1^*, \mathbf{y})$?

Suppose we do a 2-block Gibbs sampler between θ_2 and θ_3 with θ_1^* fixed throughout:

$$p(\theta_2|\theta_1^*, \theta_3, \mathbf{y})$$

$$p(\theta_3|\theta_1^*, \theta_2, \mathbf{y})$$

The ergodic distribution of θ_3 determined by this Markov chain ($q = 1, \dots, Q$) is $p(\theta_3|\theta_1^*, \mathbf{y})$

$$\hat{p}(\theta_2^*|\theta_1^*, \mathbf{y}) = Q^{-1} \sum_{q=1}^Q p(\theta_2^*|\theta_1^*, \theta_3^{(q)}, \mathbf{y})$$

So we estimate $p(\mathbf{y})$ from

$$\frac{p(\mathbf{y}|\theta_1^*, \theta_2^*, \theta_3^*)p(\theta_1^*)p(\theta_2^*)p(\theta_3^*)}{\hat{p}(\theta_1^*, \theta_2^*, \theta_3^*|\mathbf{y})}$$

$$\hat{p}(\theta_1^*, \theta_2^*, \theta_3^*|\mathbf{y}) =$$

$$G^{-1} \sum_{g=1}^G p(\theta_1^*|\theta_2^{(g)}, \theta_3^{(g)}, \mathbf{y}) \times$$

$$Q^{-1} \sum_{q=1}^Q p(\theta_2^*|\theta_1^*, \theta_3^{(q)}, \mathbf{y}) \times$$

$$p(\theta_3^*|\theta_1^*, \theta_2^*, \mathbf{y})$$
