

# Tutorial and Lab Problems # 11

## MATH3871/MATH5970

1. **Bayes Factors.** Consider model 1 with likelihood  $\pi(x|\theta, m=1)$  and model 2 with  $\pi(x|\theta, m=2) = \pi(x|\psi^{-1}(\theta), m=1)$  for an invertible function  $\psi$ . If we adopt Jeffrey's priors  $\pi_J(\theta)$  and  $\pi_J(\psi(\theta))$  for model 1 and 2 respectively, show that the Bayes factor  $B_{1|2} = B_{2|1} = 1$ .

**Answer:** From the definition of Jeffrey's prior we know that

$$\begin{aligned}\pi_J(\theta) &\propto \sqrt{\mathbb{E} \left( \frac{\partial L}{\partial \theta} \right)^2} \\ &= \sqrt{\mathbb{E} \left( \frac{\partial L}{\partial \psi} \frac{d\psi}{d\theta} \right)^2} \\ &= \left| \frac{d\psi}{d\theta} \right| \sqrt{\mathbb{E} \left( \frac{\partial L}{\partial \psi} \right)^2} \\ &\propto \pi_J(\psi) |d\psi/d\theta| .\end{aligned}$$

Therefore, by a change of variable, we have:

$$\begin{aligned}\pi(x|m=1) &= \int \pi(x|\theta, m=1) \pi_J(\theta) d\theta \\ &= \int \pi(x|\psi(\theta), m=2) \pi_J(\psi(\theta)) \frac{d\psi}{d\theta} d\theta \\ &= \int \pi(x|\psi(\theta), m=2) \pi_J(\psi(\theta)) d\psi \\ &= \pi_J(x|m=2) .\end{aligned}$$

2. **Tute/Lab # 10 revisited (this question not examinable).** Assume that  $\mathbf{D} = \alpha^2 \mathbf{I}_p$ , where  $\mathbf{I}_p \in \mathbb{R}^{p \times p}$ , and that  $\sigma^2$  is unknown. Let  $\mathbf{C} := \mathbf{I}_n + \alpha^2 \mathbf{X}\mathbf{X}^\top$ . Show that (use Woodberry matrix identity, see Wikipedia) an alternative formula for  $\hat{\sigma}^2$  is

$$\hat{\sigma}^2 = \mathbf{y}^\top \mathbf{C}^{-1} \mathbf{y} / n .$$

Also, an alternative formula for the log-determinant of  $\Sigma = \alpha^2(\mathbf{I}_p + \alpha^2 \mathbf{X}^\top \mathbf{X})^{-1}$  is

$$\ln |\Sigma| = p \ln(\alpha^2) - \ln |\mathbf{C}| .$$

Then, for a model with  $p$  predictors, we have

$$\begin{aligned} -2 \ln g(\mathbf{y} | p) &= p \ln(\alpha^2) - 2 \ln \Gamma(n/2) + n \ln(\pi n \hat{\sigma}^2) - \ln |\Sigma| \\ &= n \ln(\pi \mathbf{y}^\top \mathbf{C}^{-1} \mathbf{y}) + \ln |\mathbf{C}| - 2 \ln \Gamma(n/2) . \end{aligned}$$

Now suppose that  $m$  is the maximum number of predictors that we can include in the linear model and denote them by  $\mathbf{v}_1, \dots, \mathbf{v}_m (\in \mathbb{R}^n)$ . Let  $\mathbf{z}$  be a binary vector of length  $m$  that encodes a particular model under consideration.

For example, with  $m = 5$  and  $\mathbf{z} = (0, 1, 0, 1, 1)^\top$ , this means that  $\mathbf{X} := [\mathbf{v}_2, \mathbf{v}_4, \mathbf{v}_5]$  has three columns ( $p = 3$ ) and includes only predictors 2, 4, 5. Thus, for a fixed  $n$ , the model (as determined by the binary vector  $\mathbf{z}$ ) with the highest evidence will be the one minimizing the loss over  $\mathbf{z} \in \{0, 1\}^m$ :

$$\ell(\mathbf{z}) := n \ln(\mathbf{y}^\top [\mathbf{C}(\mathbf{z})]^{-1} \mathbf{y}) + \ln |\mathbf{C}(\mathbf{z})| ,$$

where

$$\mathbf{C}(\mathbf{z}) := \mathbf{I}_n + \alpha^2 \sum_{k=1}^m \mathbb{I}\{z_k = 1\} \mathbf{v}_k \mathbf{v}_k^\top .$$

In fact, assuming uniform prior, the posterior density on the space of all  $2^m$  models with  $m$  predictors can be written as:

$$g(\mathbf{z}) \propto e^{-\ell(\mathbf{z})}, \quad \mathbf{z} \in \{0, 1\}^m .$$

(We may exclude the all zero vector  $\mathbf{z} = (0, \dots, 0)^\top$ , if we want at least one predictor.)

To simulate from this posterior we can use Gibbs sampling, whereby sampling from the conditional  $g(z_k | \mathbf{z}_{-k})$  consists of flipping a coin to decide if the  $k$ -th component is  $z_k = 1$  or  $z_k = 0$  (equivalent to keeping or removing the  $k$ -th predictor). The probability of success of the coin flip is:

$$g(z_k = 1 | \mathbf{z}_{-k}) = \frac{e^{-\ell(\mathbf{z}; z_k=1)}}{e^{-\ell(\mathbf{z}; z_k=1)} + e^{-\ell(\mathbf{z}; z_k=0)}} .$$

To be able to implement the Gibbs sampler efficiently we must be able to update the value of  $\ell(\mathbf{z})$  after a single change to one of the  $z$ 's.

In other words, given the inverse  $\mathbf{C}^{-1}$  and log-determinant  $\ln |\mathbf{C}|$ , we want to compute very fast the inverse and determinant of the new updated matrix  $\mathbf{C} \pm \mathbf{v} \mathbf{v}^\top$ . This is a rank-1 update/change of  $\mathbf{C}$ .

Instead of computing  $(\mathbf{C} \pm \mathbf{v}\mathbf{v}^\top)^{-1}$  and  $\ln |\mathbf{C} \pm \mathbf{v}\mathbf{v}^\top|$  from scratch (this will take  $\mathcal{O}(n^3)$  cost), we compute these using the formulas (these take  $\mathcal{O}(n^2)$  cost):

$$\begin{aligned}(\mathbf{C} \pm \mathbf{v}\mathbf{v}^\top)^{-1} &= \mathbf{C}^{-1} \mp \frac{\mathbf{C}^{-1}\mathbf{v}\mathbf{v}^\top\mathbf{C}^{-1}}{1 \pm \mathbf{v}^\top\mathbf{C}^{-1}\mathbf{v}} \\ \ln |\mathbf{C} \pm \mathbf{v}\mathbf{v}^\top| &= \ln |\mathbf{C}| + \ln(1 \pm \mathbf{v}^\top\mathbf{C}^{-1}\mathbf{v}) .\end{aligned}$$

This then suggests the following  $\mathcal{O}(n^2)$ -cost algorithm for updating  $\ell$  when we add a predictor (similarly for removing one).

---

**Algorithm 1 :** Updating  $\ell(\mathbf{z})$  when we add/remove the predictor  $\mathbf{v}$

---

**Require:** Current values in memory:  $\mathbf{y}^\top\mathbf{C}^{-1}\mathbf{y}, \mathbf{C}^{-1}, \mathbf{v}, \ln |\mathbf{C}|$

---

$\mathbf{v} \leftarrow \alpha\mathbf{v}$   
 $\mathbf{a} \leftarrow \mathbf{C}^{-1}\mathbf{v}$   
 $b \leftarrow 1 \pm \mathbf{a}^\top\mathbf{v}$  (minus when removing a predictor)  
 $\ell \leftarrow \mp(\mathbf{a}^\top\mathbf{y})^2/b$   
 $\ell \leftarrow \ell + \mathbf{y}^\top\mathbf{C}^{-1}\mathbf{y}$   
 $\ell \leftarrow n \ln(\ell) + \ln |\mathbf{C}| + \ln(b)$   
**return** updated value  $\ell$

---

I can give an numerical experiment/illustration of the full Bayesian model selection algorithm in the tute/lab.