

The full model:

$$p(y | \beta, \omega, \Lambda) \sim N(X\beta, (\omega\Lambda)^{-1})$$

$$\begin{aligned} \Lambda &= \text{Diag}(\lambda_1, \dots, \lambda_n) \\ \lambda_i &\sim \text{Gamma}\left(\frac{h}{2}, \frac{h}{2}\right) \quad \text{where } h \text{ is a fixed hyperparameter.} \\ (\beta | \omega) &\sim N(m, (\omega K)^{-1}) \\ \omega &\sim \text{Gamma}\left(\frac{d}{2}, \frac{n}{2}\right) \end{aligned}$$

A) Under this model, what is the implied conditional distribution  $p(y_i | x_i, \beta, \omega)$ ?  
Notice that  $\lambda_i$  has been marginalized out. This should look familiar.

The conditional distribution of  $y_i$

$$P(y_i | x_i, \beta, \omega) = P(y_i | x_i, \beta, \omega, \lambda_i) P(\lambda_i)$$

$$\propto \int_0^\infty \sqrt{\omega \lambda_i} \exp\left(-\frac{\omega \lambda_i}{2} (y_i - x_i^T \beta)^2\right) \lambda_i^{\frac{h}{2}-1} \exp\left(-\frac{\lambda_i h}{2}\right) d\lambda_i$$

$$\propto \int_0^\infty \lambda_i^{\frac{h}{2}-1} \exp\left(-\frac{\lambda_i}{2} (\omega (y_i - x_i^T \beta)^2 + h)\right) d\lambda_i$$

$$\propto \left(\frac{1}{2} (\omega (y_i - x_i^T \beta)^2 + h)\right)^{-\frac{(h+1)}{2}}$$

Kernel of  $\text{gamma}\left(\frac{h+1}{2}, \frac{1}{2} (\omega (y_i - x_i^T \beta)^2 + h)\right)$

$$\propto \left(1 + \frac{1}{h} \frac{(y_i - x_i^T \beta)^2}{(\sqrt{\omega^{-1}})^2}\right)^{-\frac{(h+1)}{2}}$$

$$\sim T_h(x_i^T \beta, \sqrt{\omega^{-1}}), \text{ a } t \text{ dist'n with } h \text{ degrees of freedom, centered at } x_i^T \beta, \text{ with } \sqrt{\omega^{-1}} \text{ scale}$$

B) What is the conditional posterior distribution  $p(\lambda_i | y_i, \beta, \omega)$ ?

To find the conditional posterior, I will first find the full posterior. We have

$$p(\lambda_i, \beta, \omega | y_i) \propto p(y_i | \beta, \omega, \lambda_i) p(\beta | \omega) p(\omega) p(\lambda_i)$$

$$\propto \sqrt{\frac{\omega \lambda_i}{2\pi}} \exp\left(-\frac{\omega \lambda_i}{2} (y_i - x_i^T \beta)^2\right) \sqrt{\frac{\omega K}{2\pi}} \exp\left(-\frac{\omega K}{2} (\beta - m)^2\right) \omega^{\frac{d}{2}-1} \exp\left(-\frac{\omega n}{2}\right) \lambda_i^{\frac{h}{2}-1} \exp\left(-\frac{\lambda_i h}{2}\right) \quad (1)$$

To find  $p(\lambda_i | y_i, \beta, \omega)$ , we can set  $y_i, \beta$ , and  $\omega$  to constants and only work with the pieces of the posterior that involve  $\lambda_i$ . Thus, (1) will simplify to

$$p(\lambda_i | y_i, \beta, \omega) \propto \sqrt{\frac{\omega \lambda_i}{2\pi}} \exp\left(-\frac{\omega \lambda_i}{2} (y_i - x_i^T \beta)^2\right) \lambda_i^{\frac{h}{2}-1} \exp\left(-\frac{\lambda_i h}{2}\right)$$

$$\propto \lambda_i^{\frac{h}{2}-1} \exp\left(-\lambda_i \left(\frac{1}{2} (\omega (y_i - x_i^T \beta)^2 + h)\right)\right)$$

$$\sim \text{gamma}\left(\frac{h+1}{2}, \frac{1}{2} (\omega (y_i - x_i^T \beta)^2 + h)\right)$$

C) Code up a gibbs sampler that repeatedly cycles through

c) Code up a gibbs sampler that repeatedly cycles through

$$p(\beta | y, \omega, \lambda)$$

$$p(\omega | y, \lambda)$$

$$p(\lambda_i | y, \beta, \omega)$$

The first two should follow identically from your previous results, except that we are explicitly conditioning on  $\lambda$ , which is a random variable rather than a fixed hyperparameter.

$$p(\lambda_i | \beta, \omega | y_i) \propto \sqrt{\frac{\omega \lambda_i}{2\pi}} \exp\left(-\frac{\omega \lambda_i}{2} (y_i - x_i^T \beta)^2\right) \sqrt{\frac{\omega K}{2\pi}} \exp\left(-\frac{\omega K}{2} (\beta - m)^2\right) \omega^{\frac{d}{2}-1} \exp\left(-\frac{\omega \eta_i}{2}\right) \lambda_i^{\frac{h}{2}-1} \exp\left(-\frac{\lambda_i h}{2}\right) \quad (1)$$

From here, we can see that if we were to find the conditional posterior of  $\beta$  and  $\omega$ , we can simply disregard the prior for  $\lambda_i$  since both the  $p(\beta | \omega)$  and  $p(\omega)$  are independent of  $\lambda_i$ . Thus, the conditional posterior  $p(\beta, \omega | y, \lambda)$  is:

$$p(\beta, \omega | y, \lambda) \propto \frac{\omega^{\frac{d}{2}-1}}{1(\omega \lambda)^{\frac{1}{2}} 1(\omega K)^{\frac{1}{2}}} \exp\left(-\frac{\omega}{2} (\beta - A^{-1}b)^T A (\beta - A^{-1}b)\right) \exp\left(-\frac{\omega}{2} (\eta - b^T A^{-1}b + y^T \lambda y + m^T K m)\right)$$

$$\text{where } A = (X^T X + K) \\ b = (y \lambda X + m^T K)$$

$$\text{Thus, } p(\beta | y, \omega, \lambda) \propto \exp\left(-\frac{\omega}{2} (\beta - A^{-1}b)^T A (\beta - A^{-1}b)\right)$$

$$\sim \text{MVN}(A^{-1}b, (\omega A)^{-1})$$

$$p(\omega | y, \beta, \lambda) \propto \omega^{\frac{d+n+p}{2}-1} \exp\left[-\frac{\omega}{2} ((\beta - A^{-1}b)^T A (\beta - A^{-1}b) + (\eta + b^T A^{-1}b + y^T \lambda y + m^T K m))\right]$$

$$\sim \text{Gam}\left(\frac{d+n+p}{2}, \frac{1}{2}((\beta - A^{-1}b)^T A (\beta - A^{-1}b) + (\eta + b^T A^{-1}b + y^T \lambda y + m^T K m))\right)$$

$$p(\lambda_i | y, \beta, \omega) \propto \lambda_i^{\frac{h+1}{2}-1} \exp\left(-\lambda_i \left(\frac{1}{2} (\omega (y_i - x_i^T \beta)^2 + h)\right)\right)$$

$$\sim \text{gamma}\left(\frac{h+1}{2}, \frac{1}{2} (\omega (y_i - x_i^T \beta)^2 + h)\right)$$

The Gibbs sampler algorithm is:

1-) Start with some  $(\beta, \omega, \lambda)^{(0)}$

2-) At each iteration  $t$ , for each  $j=1, \dots, p$ , sample  $(\beta^{(t)}, \omega^{(t)}, \lambda^{(t)})$  from

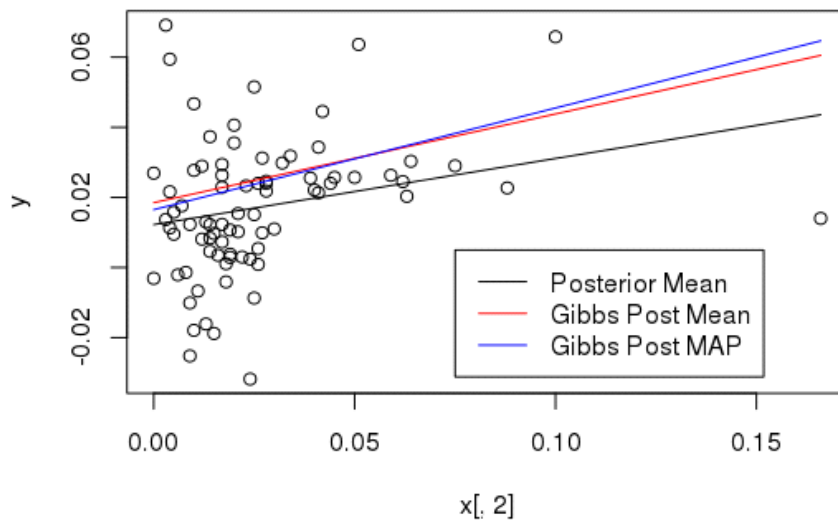
$$p(\beta^{(t)} | \omega^{(t-1)}, \lambda^{(t-1)})$$

$$p(\omega^{(t)} | \beta^{(t)}, \lambda^{(t-1)})$$

$$p(\lambda_i^{(t)} | \beta^{(t)}, \omega^{(t)}, \lambda_1^{(t)}, \dots, \lambda_{i-1}^{(t)}, \lambda_{i+1}^{(t-1)}, \dots, \lambda_p^{(t-1)})$$

We can discard the first thousand draws as the burn in and draw an additional 3000 samples.

## Bayesian Linear Regression plot



The hyperparameters I used for this simulation were:

$$\begin{aligned}
 d &= 2 \\
 \eta &= 2 \\
 h &= 2 \\
 m &= (0.02, 0.1) \\
 K &= \begin{bmatrix} 0.01 & 0 \\ 0 & 0.01 \end{bmatrix}
 \end{aligned}$$

After burning 1000 samples, I thinned the remaining 3000 samples by a factor of 3. To compute the red fitted line, I took the average of my sampled betas. As you can see from the figure above, the red line has a slightly higher slope and intercept than the original Bayesian linear model. Considering the true fit of the model should have a very steep slope, the averaged gibbs result is a little bit better. However, it still is influenced by the outliers a lot. The MAP estimates for beta again produce a slightly better result but it still has the same problem. In general, the MAP estimate was slightly better in all cases. Tuning the hyperparameters showed that the model is very sensitive to changes in  $m$ . This makes sense as  $m$  holds our prior beliefs about what  $\beta_0$  and  $\beta_1$  should be.