

## Normal Linear Regression Model with Conjugate Priors

*R* scripts:

mod6s2a, mod6s2b, mod6s2c, mod6s2d

“Conjugate” refers to the property of a prior to generate, when combined with the likelihood function, a posterior that has the same density as the prior itself. Conjugate priors are mainly chosen for computational convenience, as they lead to exact analytical expressions for the posterior and the marginal likelihood. In practice, they are very restrictive, and in few real applications will we be able to formalize our prior understanding of the world in conjugate terms. However, this particular model is very useful for pedagogical purposes, and is thus a good starting point.

Importantly, it illustrates that Bayesian modeling is not inherently synonymous with "simulation techniques". The latter is simply a computational tool that broadens the choice and thus applicability of Bayesian models.

### Model Specification

Step 1:

Start by writing down the regression model:

$$\mathbf{y} = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\varepsilon} \quad \boldsymbol{\varepsilon} \sim n(\mathbf{0}, \sigma^2 \mathbf{I}) \quad (1)$$

For this model with conjugate priors it is more convenient to work with the “precision”  $h = \sigma^{-2}$  instead of the variance for the error term. Thus, our model parameters about which we wish to learn are

$$\boldsymbol{\theta} = [\boldsymbol{\beta}' \quad h]'$$

Step 2:

Next, write down the likelihood function (LHF) for this model. It’s the same you would use for Maximum Likelihood Estimation (MLE), except perhaps for the switch from  $\sigma^2$  to  $h$ . Nothing Bayesian yet.

$$p(\mathbf{y} | \boldsymbol{\theta}, \mathbf{X}) = (2\pi)^{-n/2} h^{n/2} \exp\left(-\frac{h}{2}(\mathbf{y} - \mathbf{X}\boldsymbol{\beta})'(\mathbf{y} - \mathbf{X}\boldsymbol{\beta})\right) \quad (2)$$

For this model we need to use a lot of “tricks” to get convenient and intuitive expressions for interim and final results. First, we will express the LHF in terms of the OLS estimator  $\mathbf{b}$ , and the summed squared errors (*SSE*). Note the following equivalence:

$$\begin{aligned}
(\mathbf{y} - \mathbf{X}\boldsymbol{\beta})'(\mathbf{y} - \mathbf{X}\boldsymbol{\beta}) &= (\mathbf{y} - \mathbf{X}\boldsymbol{\beta} + \mathbf{X}\mathbf{b} - \mathbf{X}\mathbf{b})'(\mathbf{y} - \mathbf{X}\boldsymbol{\beta} + \mathbf{X}\mathbf{b} - \mathbf{X}\mathbf{b}) = \\
&= (\mathbf{y} - \mathbf{X}\mathbf{b} + (\mathbf{X}(\mathbf{b} - \boldsymbol{\beta})))'(\mathbf{y} - \mathbf{X}\mathbf{b} + (\mathbf{X}(\mathbf{b} - \boldsymbol{\beta}))) = \\
&= (\mathbf{y} - \mathbf{X}\mathbf{b})'(\mathbf{y} - \mathbf{X}\mathbf{b}) + 2(\mathbf{X}(\mathbf{b} - \boldsymbol{\beta}))'(\mathbf{y} - \mathbf{X}\mathbf{b}) + (\mathbf{b} - \boldsymbol{\beta})' \mathbf{X}'\mathbf{X}(\mathbf{b} - \boldsymbol{\beta}) = \\
&= (\mathbf{y} - \mathbf{X}\mathbf{b})'(\mathbf{y} - \mathbf{X}\mathbf{b}) + (\mathbf{b} - \boldsymbol{\beta})' \mathbf{X}'\mathbf{X}(\mathbf{b} - \boldsymbol{\beta}) = SSE + (\mathbf{b} - \boldsymbol{\beta})' \mathbf{X}'\mathbf{X}(\mathbf{b} - \boldsymbol{\beta})
\end{aligned} \tag{3}$$

since

$$(\mathbf{X}(\mathbf{b} - \boldsymbol{\beta}))'(\mathbf{y} - \mathbf{X}\mathbf{b}) = (\mathbf{b} - \boldsymbol{\beta})' \mathbf{X}'(\mathbf{y} - \mathbf{X}\mathbf{b}) = (\mathbf{b} - \boldsymbol{\beta})' (\mathbf{X}'\mathbf{y} - \mathbf{X}'\mathbf{X}(\mathbf{X}'\mathbf{X})^{-1} \mathbf{X}'\mathbf{y}) = 0$$

Note that only in the last line do we use the explicit form of the OLS estimator. Thus, we can express the LHF as

$$p(\mathbf{y} | \boldsymbol{\theta}, \mathbf{X}) = (2\pi)^{-n/2} h^{n/2} \exp\left(-\frac{h}{2} \left( SSE + (\mathbf{b} - \boldsymbol{\beta})' \mathbf{X}'\mathbf{X}(\mathbf{b} - \boldsymbol{\beta}) \right)\right) \tag{4}$$

### Step 3:

Next, we need to define the prior distributions for our parameters. Here, we choose “conjugate priors” that, when combined with the likelihood, yield a posterior from the exact same conjugate family. A trademark of conjugate priors is that the prior for one set of parameters depends on the prior for another set. In other words, the parameters of interest do NOT have independent priors. Here, the prior of  $\boldsymbol{\beta}$ , conditional on  $h$ , is normal, and the prior of  $h$  is gamma.

$$p(\boldsymbol{\beta}, h) = p(\boldsymbol{\beta} | h) p(h) \quad \text{where}$$

$$\boldsymbol{\beta} | h \sim n(\boldsymbol{\mu}_0, h^{-1}\mathbf{V}_0), \quad h \sim g(\tau_0, \nu_0)$$

$$p(\boldsymbol{\beta} | h) = (2\pi)^{-k/2} |h^{-1}\mathbf{V}_0|^{-1/2} \exp\left(-\frac{1}{2}(\boldsymbol{\beta} - \boldsymbol{\mu}_0)'(h^{-1}\mathbf{V}_0)^{-1}(\boldsymbol{\beta} - \boldsymbol{\mu}_0)\right) \tag{5}$$

$$p(h) = \left(\frac{2\tau_0}{\nu_0}\right)^{-\nu_0/2} \Gamma\left(\frac{\nu_0}{2}\right)^{-1} h^{\left(\frac{\nu_0-2}{2}\right)} \exp\left(-\frac{h\nu_0}{2\tau_0}\right) \quad \text{with} \quad \Gamma(x) = \int_0^{\infty} t^{x-1} \exp(-t) dt$$

where  $\boldsymbol{\mu}_0$  is the prior mean for  $\boldsymbol{\beta}$ ,  $\mathbf{V}_0$  is the unscaled prior variance-covariance matrix for  $\boldsymbol{\beta}$ , and  $\tau_0$  and  $\nu_0$  are the mean and degrees of freedom (DoF) for the gamma distribution. (Note: In later examples and models we might express the gamma density in form of its shape and scale instead of mean and DoF. Here we're sticking with the version chosen in KPT, ch. 10). Note:

$$\begin{aligned}
|h^{-1}\mathbf{V}_0| &= h^{-k} |\mathbf{V}_0| && \text{(e.g. Greene, A-48)} \\
(h^{-1}\mathbf{V}_0)^{-1} &= h(\mathbf{V}_0)^{-1}
\end{aligned} \tag{6}$$

**R** script `mod6s2a` illustrates the synergies between these conjugate priors. The upshot is that with conjugate priors, the prior settings for one parameter will affect the prior density of another. A bit of trial and error (and density plotting) will be needed to arrive at the desired prior densities for all parameters.

Step 4:

Write down the complete posterior distribution.

$$\begin{aligned}
 p(\boldsymbol{\beta}, h | \mathbf{y}, \mathbf{X}) &= \frac{p(\boldsymbol{\beta}, h) p(\mathbf{y} | \boldsymbol{\beta}, h, \mathbf{X})}{p(\mathbf{y} | \mathbf{X})} \propto p(\boldsymbol{\beta}, h) p(\mathbf{y} | \boldsymbol{\beta}, h, \mathbf{X}) = p(\boldsymbol{\beta} | h) p(h) p(\mathbf{y} | \boldsymbol{\beta}, h, \mathbf{X}) = \\
 & (2\pi)^{-k/2} h^{k/2} |\mathbf{V}_0|^{-1/2} \exp\left(-\frac{h}{2}(\boldsymbol{\beta} - \boldsymbol{\mu}_0)' \mathbf{V}_0^{-1} (\boldsymbol{\beta} - \boldsymbol{\mu}_0)\right) * \\
 & \left(\frac{2\tau_0}{v_0}\right)^{-v_0/2} \Gamma\left(\frac{v_0}{2}\right)^{-1} h^{\left(\frac{v_0-2}{2}\right)} \exp\left(-\frac{h v_0}{2\tau_0}\right) * \\
 & (2\pi)^{-n/2} h^{n/2} \exp\left(-\frac{h}{2}\left(SSE + (\mathbf{b} - \boldsymbol{\beta})' \mathbf{X}'\mathbf{X}(\mathbf{b} - \boldsymbol{\beta})\right)\right)
 \end{aligned} \tag{7}$$

Note the switch from equality (“=”) to proportionality ( $\propto$ ) in the first line. This simplifies our posterior computation since we don’t need the marginal likelihood  $p(\mathbf{y} | \mathbf{X})$  to find the posterior densities for our parameters.

At this point it is customary (for ANY Bayesian model) to further ignore (= drop) all terms that are multiplicatively unrelated to our parameters of interest. Thus:

$$\begin{aligned}
 & p(\boldsymbol{\beta} | h) p(h) p(\mathbf{y} | \boldsymbol{\beta}, h, \mathbf{X}) \propto \\
 & h^{k/2} \exp\left(-\frac{h}{2}(\boldsymbol{\beta} - \boldsymbol{\mu}_0)' \mathbf{V}_0^{-1} (\boldsymbol{\beta} - \boldsymbol{\mu}_0)\right) * h^{\left(\frac{v_0-2}{2}\right)} \exp\left(-\frac{h v_0}{2\tau_0}\right) * \\
 & h^{n/2} \exp\left(-\frac{h}{2}\left(SSE + (\mathbf{b} - \boldsymbol{\beta})' \mathbf{X}'\mathbf{X}(\mathbf{b} - \boldsymbol{\beta})\right)\right).
 \end{aligned} \tag{8}$$

Collecting terms the posterior kernel further simplifies to

$$\begin{aligned}
 & p(\boldsymbol{\beta} | h) p(h) p(\mathbf{y} | \boldsymbol{\beta}, h, \mathbf{X}) \propto \\
 & h^{\frac{k+v_0-2+n}{2}} \exp\left(-\frac{h}{2}\left(v_0\tau_0^{-1} + SSE + (\boldsymbol{\beta} - \boldsymbol{\mu}_0)' \mathbf{V}_0^{-1} (\boldsymbol{\beta} - \boldsymbol{\mu}_0) + (\mathbf{b} - \boldsymbol{\beta})' \mathbf{X}'\mathbf{X}(\mathbf{b} - \boldsymbol{\beta})\right)\right).
 \end{aligned} \tag{9}$$

Here comes another trick:

Let  $\mathbf{V}_1 = (\mathbf{V}_0^{-1} + \mathbf{X}'\mathbf{X})^{-1}$  and  $\boldsymbol{\mu}_1 = \mathbf{V}_1(\mathbf{V}_0^{-1}\boldsymbol{\mu}_0 + \mathbf{X}'\mathbf{Xb}) = \mathbf{V}_1(\mathbf{V}_0^{-1}\boldsymbol{\mu}_0 + \mathbf{X}'\mathbf{y})$

Then:

$$\begin{aligned}
& (\boldsymbol{\beta} - \boldsymbol{\mu}_0)' \mathbf{V}_0^{-1} (\boldsymbol{\beta} - \boldsymbol{\mu}_0) + (\mathbf{b} - \boldsymbol{\beta})' \mathbf{X}'\mathbf{X} (\mathbf{b} - \boldsymbol{\beta}) = \\
& \boldsymbol{\beta}' \mathbf{V}_0^{-1} \boldsymbol{\beta} - \boldsymbol{\mu}_0' \mathbf{V}_0^{-1} \boldsymbol{\beta} - \boldsymbol{\beta}' \mathbf{V}_0^{-1} \boldsymbol{\mu}_0 + \boldsymbol{\mu}_0' \mathbf{V}_0^{-1} \boldsymbol{\mu}_0 + \boldsymbol{\beta}' \mathbf{X}'\mathbf{X} \boldsymbol{\beta} - \boldsymbol{\beta}' \mathbf{X}'\mathbf{X} \mathbf{b} - \mathbf{b}' \mathbf{X}'\mathbf{X} \boldsymbol{\beta} + \mathbf{b}' \mathbf{X}'\mathbf{X} \mathbf{b} = \\
& \boldsymbol{\beta}' (\mathbf{V}_0^{-1} + \mathbf{X}'\mathbf{X}) \boldsymbol{\beta} - \boldsymbol{\beta}' (\mathbf{V}_0^{-1} \boldsymbol{\mu}_0 + \mathbf{X}'\mathbf{Xb}) - (\boldsymbol{\mu}_0' \mathbf{V}_0^{-1} + \mathbf{b}' \mathbf{X}'\mathbf{X}) \boldsymbol{\beta} + \boldsymbol{\mu}_0' \mathbf{V}_0^{-1} \boldsymbol{\mu}_0 + \mathbf{b}' \mathbf{X}'\mathbf{X} \mathbf{b} = \\
& \boldsymbol{\beta}' \mathbf{V}_1^{-1} \boldsymbol{\beta} - \boldsymbol{\beta}' \mathbf{V}_1^{-1} \boldsymbol{\mu}_1 - \boldsymbol{\mu}_1' \mathbf{V}_1^{-1} \boldsymbol{\beta} + \boldsymbol{\mu}_0' \mathbf{V}_0^{-1} \boldsymbol{\mu}_0 + \mathbf{b}' \mathbf{X}'\mathbf{X} \mathbf{b}, \quad \text{since}
\end{aligned} \tag{10}$$

$$\mathbf{V}_1^{-1} \boldsymbol{\mu}_1 = \mathbf{V}_1^{-1} \mathbf{V}_1 (\mathbf{V}_0^{-1} \boldsymbol{\mu}_0 + \mathbf{X}'\mathbf{Xb}) = (\mathbf{V}_0^{-1} \boldsymbol{\mu}_0 + \mathbf{X}'\mathbf{Xb}) \quad \text{and}$$

$$(\boldsymbol{\mu}_0' \mathbf{V}_0^{-1} + \mathbf{b}' \mathbf{X}'\mathbf{X}) = (\mathbf{V}_0^{-1} \boldsymbol{\mu}_0 + \mathbf{X}'\mathbf{Xb})' = (\mathbf{V}_1^{-1} \boldsymbol{\mu}_1)' = \boldsymbol{\mu}_1' \mathbf{V}_1^{-1}.$$

Simplifying further:

$$\begin{aligned}
& \boldsymbol{\beta}' \mathbf{V}_1^{-1} \boldsymbol{\beta} - \boldsymbol{\beta}' \mathbf{V}_1^{-1} \boldsymbol{\mu}_1 - \boldsymbol{\mu}_1' \mathbf{V}_1^{-1} \boldsymbol{\beta} + \boldsymbol{\mu}_0' \mathbf{V}_0^{-1} \boldsymbol{\mu}_0 + \mathbf{b}' \mathbf{X}'\mathbf{X} \mathbf{b} = \\
& \boldsymbol{\beta}' \mathbf{V}_1^{-1} \boldsymbol{\beta} - \boldsymbol{\beta}' \mathbf{V}_1^{-1} \boldsymbol{\mu}_1 - \boldsymbol{\mu}_1' \mathbf{V}_1^{-1} \boldsymbol{\beta} + \boldsymbol{\mu}_1' \mathbf{V}_1^{-1} \boldsymbol{\mu}_1 - \boldsymbol{\mu}_1' \mathbf{V}_1^{-1} \boldsymbol{\mu}_1 + \boldsymbol{\mu}_0' \mathbf{V}_0^{-1} \boldsymbol{\mu}_0 + \mathbf{b}' \mathbf{X}'\mathbf{X} \mathbf{b} = \\
& (\boldsymbol{\beta} - \boldsymbol{\mu}_1)' \mathbf{V}_1^{-1} (\boldsymbol{\beta} - \boldsymbol{\mu}_1) - \boldsymbol{\mu}_1' \mathbf{V}_1^{-1} \boldsymbol{\mu}_1 + \boldsymbol{\mu}_0' \mathbf{V}_0^{-1} \boldsymbol{\mu}_0 + \mathbf{b}' \mathbf{X}'\mathbf{X} \mathbf{b}.
\end{aligned} \tag{11}$$

The last three terms can be further “massaged” to yield:

$$\begin{aligned}
& \boldsymbol{\mu}_0' \mathbf{V}_0^{-1} \boldsymbol{\mu}_0 + \mathbf{b}' \mathbf{X}'\mathbf{X} \mathbf{b} - \boldsymbol{\mu}_1' \mathbf{V}_1^{-1} \boldsymbol{\mu}_1 = \\
& \boldsymbol{\mu}_0' \mathbf{V}_0^{-1} \boldsymbol{\mu}_0 + \mathbf{b}' \mathbf{X}'\mathbf{X} \mathbf{b} - \left( (\mathbf{V}_0^{-1} \boldsymbol{\mu}_0 + \mathbf{X}'\mathbf{Xb})' \mathbf{V}_1 (\mathbf{V}_0^{-1} \boldsymbol{\mu}_0 + \mathbf{X}'\mathbf{Xb}) \right) = \\
& \boldsymbol{\mu}_0' \mathbf{V}_0^{-1} \boldsymbol{\mu}_0 + \mathbf{b}' \mathbf{X}'\mathbf{X} \mathbf{b} - \left( \boldsymbol{\mu}_0' \mathbf{V}_0^{-1} \mathbf{V}_1 \mathbf{V}_0^{-1} \boldsymbol{\mu}_0 + \mathbf{b}' \mathbf{X}'\mathbf{X} \mathbf{V}_1 \mathbf{V}_0^{-1} \boldsymbol{\mu}_0 + \boldsymbol{\mu}_0' \mathbf{V}_0^{-1} \mathbf{V}_1 \mathbf{X}'\mathbf{X} \mathbf{b} + \mathbf{b}' \mathbf{X}'\mathbf{X} \mathbf{V}_1 \mathbf{X}'\mathbf{X} \mathbf{b} \right) = \\
& \boldsymbol{\mu}_0' \mathbf{V}_0^{-1} \boldsymbol{\mu}_0 - \boldsymbol{\mu}_0' \mathbf{V}_0^{-1} \mathbf{V}_1 \mathbf{V}_0^{-1} \boldsymbol{\mu}_0 + \boldsymbol{\mu}_0' \mathbf{X}'\mathbf{X} \mathbf{V}_1 \mathbf{V}_0^{-1} \boldsymbol{\mu}_0 - \boldsymbol{\mu}_0' \mathbf{X}'\mathbf{X} \mathbf{V}_1 \mathbf{V}_0^{-1} \boldsymbol{\mu}_0 + \\
& \mathbf{b}' \mathbf{X}'\mathbf{X} \mathbf{b} - \mathbf{b}' \mathbf{X}'\mathbf{X} \mathbf{V}_1 \mathbf{X}'\mathbf{X} \mathbf{b} + \mathbf{b}' \mathbf{X}'\mathbf{X} \mathbf{V}_1 \mathbf{V}_0^{-1} \mathbf{b} - \mathbf{b}' \mathbf{X}'\mathbf{X} \mathbf{V}_1 \mathbf{V}_0^{-1} \mathbf{b} - \\
& \mathbf{b}' \mathbf{X}'\mathbf{X} \mathbf{V}_1 \mathbf{V}_0^{-1} \boldsymbol{\mu}_0 - \boldsymbol{\mu}_0' \mathbf{V}_0^{-1} \mathbf{V}_1 \mathbf{X}'\mathbf{X} \mathbf{b} = \\
& (\mathbf{b} - \boldsymbol{\mu}_0)' \mathbf{X}'\mathbf{X} \mathbf{V}_1 \mathbf{V}_0^{-1} (\mathbf{b} - \boldsymbol{\mu}_0) \quad \text{since}
\end{aligned} \tag{12}$$

$$\mathbf{b}' \mathbf{X}'\mathbf{X} \mathbf{b} - \mathbf{b}' \mathbf{X}'\mathbf{X} \mathbf{V}_1 \mathbf{X}'\mathbf{X} \mathbf{b} - \mathbf{b}' \mathbf{X}'\mathbf{X} \mathbf{V}_1 \mathbf{V}_0^{-1} \mathbf{b} = \mathbf{b}' \left( \mathbf{X}'\mathbf{X} (\mathbf{I} - \mathbf{V}_1 (\mathbf{X}'\mathbf{X} + \mathbf{V}_0^{-1})) \right) \mathbf{b} = 0$$

$$\boldsymbol{\mu}_0' \mathbf{V}_0^{-1} \boldsymbol{\mu}_0 - \boldsymbol{\mu}_0' \mathbf{V}_0^{-1} \mathbf{V}_1 \mathbf{V}_0^{-1} \boldsymbol{\mu}_0 - \boldsymbol{\mu}_0' \mathbf{X}'\mathbf{X} \mathbf{V}_1 \mathbf{V}_0^{-1} \boldsymbol{\mu}_0 = \boldsymbol{\mu}_0' \left( \mathbf{V}_0^{-1} (\mathbf{I} - \mathbf{V}_1 (\mathbf{V}_0^{-1} + \mathbf{X}'\mathbf{X})) \right) \boldsymbol{\mu}_0 = 0$$

Almost done – one final simplification:

$$\begin{aligned}
& (\mathbf{b} - \boldsymbol{\mu}_0)' \mathbf{X}' \mathbf{X} \mathbf{V}_1 \mathbf{V}_0^{-1} (\mathbf{b} - \boldsymbol{\mu}_0) = \\
& (\mathbf{b} - \boldsymbol{\mu}_0)' \mathbf{X}' \mathbf{X} (\mathbf{V}_0^{-1} + \mathbf{X}' \mathbf{X})^{-1} \mathbf{V}_0^{-1} (\mathbf{b} - \boldsymbol{\mu}_0) = \\
& (\mathbf{b} - \boldsymbol{\mu}_0)' \mathbf{X}' \mathbf{X} (\mathbf{X}' \mathbf{X})^{-1} \left( \mathbf{V}_0 (\mathbf{V}_0^{-1} + \mathbf{X}' \mathbf{X}) (\mathbf{X}' \mathbf{X})^{-1} \right)^{-1} \mathbf{V}_0^{-1} \mathbf{V}_0 (\mathbf{b} - \boldsymbol{\mu}_0) = \\
& (\mathbf{b} - \boldsymbol{\mu}_0)' \left( \mathbf{V}_0 + (\mathbf{X}' \mathbf{X})^{-1} \right)^{-1} (\mathbf{b} - \boldsymbol{\mu}_0).
\end{aligned} \tag{13}$$

After all this tedious work we can now write the posterior kernel as:

$$\begin{aligned}
& p(\boldsymbol{\beta}, h | \mathbf{y}, \mathbf{X}) \propto \\
& h^{\frac{k+n\mathbf{V}_0^{-2}+n}{2}} \exp\left(-\frac{h}{2} \left( v_0 \tau_0^{-1} + SSE + (\boldsymbol{\beta} - \boldsymbol{\mu}_1)' \mathbf{V}_1^{-1} (\boldsymbol{\beta} - \boldsymbol{\mu}_1) + (\mathbf{b} - \boldsymbol{\mu}_0)' \left( \mathbf{V}_0 + (\mathbf{X}' \mathbf{X})^{-1} \right)^{-1} (\mathbf{b} - \boldsymbol{\mu}_0) \right)\right) = \\
& h^{\frac{k+n\mathbf{V}_0^{-2}+n}{2}} \exp\left(-\frac{h}{2} \left( v_0 \tau_0^{-1} + SSE + (\mathbf{b} - \boldsymbol{\mu}_0)' \left( \mathbf{V}_0 + (\mathbf{X}' \mathbf{X})^{-1} \right)^{-1} (\mathbf{b} - \boldsymbol{\mu}_0) \right)\right) * \exp\left(-\frac{h}{2} \left( (\boldsymbol{\beta} - \boldsymbol{\mu}_1)' \mathbf{V}_1^{-1} (\boldsymbol{\beta} - \boldsymbol{\mu}_1) \right)\right).
\end{aligned} \tag{14}$$

Step 5:

It is customary in Bayesian posterior analysis to break the posterior kernel into conditional components that have a well-understood density. Here, we want to first look at the posterior for  $\boldsymbol{\beta}$ , conditional on  $h$ . This means we can ignore all terms in (14) that are multiplicatively unrelated to  $\boldsymbol{\beta}$ . We get:

$$p(\boldsymbol{\beta} | h, \mathbf{y}, \mathbf{X}) \propto \exp\left(-\frac{h}{2} \left( (\boldsymbol{\beta} - \boldsymbol{\mu}_1)' \mathbf{V}_1^{-1} (\boldsymbol{\beta} - \boldsymbol{\mu}_1) \right)\right) = \left( -\frac{1}{2} \left( (\boldsymbol{\beta} - \boldsymbol{\mu}_1)' (h^{-1} \mathbf{V}_1)^{-1} (\boldsymbol{\beta} - \boldsymbol{\mu}_1) \right) \right). \tag{15}$$

This is *exactly* the kernel of a multivariate normal distribution. Thus we have established that

$$\boldsymbol{\beta} | h, \mathbf{y}, \mathbf{X} \sim n(\boldsymbol{\mu}_1, h^{-1} \mathbf{V}_1) \tag{16}$$

While this is not our end-product, it is a useful finding which we will also exploit in future models.

Step 6:

To derive the marginal posterior density for  $h$ , we can use a "reversed" version of Bayes' Rule:

$$p(h | \mathbf{y}, \mathbf{X}) = \frac{p(\boldsymbol{\beta}, h | \mathbf{y}, \mathbf{X})}{p(\boldsymbol{\beta} | h, \mathbf{y}, \mathbf{X})} \propto \frac{p(\boldsymbol{\beta} | h) p(h) p(\mathbf{y} | \boldsymbol{\beta}, h, \mathbf{X})}{p(\boldsymbol{\beta} | h, \mathbf{y}, \mathbf{X})} \tag{17}$$

Thus, we should get the kernel for  $p(h | \mathbf{y}, \mathbf{X})$  by dividing the joint posterior kernel in (14) by the full conditional density of  $\boldsymbol{\beta}$ . Dropping terms that are multiplicatively unrelated to  $h$ , this yields:

$$\begin{aligned}
p(h | \mathbf{y}, \mathbf{X}) &\propto h^{\frac{k+v_0-2+n}{2}} \exp\left(-\frac{h}{2}\left(v_0\tau_0^{-1} + SSE + (\mathbf{b} - \boldsymbol{\mu}_0)' \left(\mathbf{V}_0 + (\mathbf{X}'\mathbf{X})^{-1}\right)^{-1} (\mathbf{b} - \boldsymbol{\mu}_0)\right)\right) * \\
&\exp\left(-\frac{h}{2}\left((\boldsymbol{\beta} - \boldsymbol{\mu}_1)' \mathbf{V}_1^{-1} (\boldsymbol{\beta} - \boldsymbol{\mu}_1)\right)\right) * \\
&\left((2\pi)^{-k/2} h^{k/2} |\mathbf{V}_1|^{-1/2} \exp\left(-\frac{h}{2}(\boldsymbol{\beta} - \boldsymbol{\mu}_1)' \mathbf{V}_1^{-1} (\boldsymbol{\beta} - \boldsymbol{\mu}_1)\right)\right)^{-1} \propto \\
&h^{\frac{v_0-2+n}{2}} \exp\left(-\frac{h}{2}\left(v_0\tau_0^{-1} + SSE + (\mathbf{b} - \boldsymbol{\mu}_0)' \left(\mathbf{V}_0 + (\mathbf{X}'\mathbf{X})^{-1}\right)^{-1} (\mathbf{b} - \boldsymbol{\mu}_0)\right)\right) = \\
&h^{\frac{v_1-2}{2}} \exp\left(-\frac{h}{2}\left(v_0\tau_0^{-1} + v_1\tau_1^{-1}\right)\right) \quad \text{where} \\
v_1 = v_0 + n \quad \text{and} \quad \tau_1^{-1} &= \frac{1}{v_1}\left(v_0\tau_0^{-1} + SSE + (\mathbf{b} - \boldsymbol{\mu}_0)' \left(\mathbf{V}_0 + (\mathbf{X}'\mathbf{X})^{-1}\right)^{-1} (\mathbf{b} - \boldsymbol{\mu}_0)\right)
\end{aligned} \tag{18}$$

Comparing the final expression to the prior of  $h$  in (5), we can see that we have, again, the kernel of a gamma distribution. Thus we can state that

$$h | \mathbf{y}, \mathbf{X} \sim g(\tau_1, v_1) \tag{19}$$

The marginal posterior for  $\boldsymbol{\beta}$  can now be derived by marginalizing the joint posterior over  $h$ , i.e

$$p(\boldsymbol{\beta} | \mathbf{y}, \mathbf{X}) = \int p(\boldsymbol{\beta}, h | \mathbf{y}, \mathbf{X}) dh = \int p(\boldsymbol{\beta} | h, \mathbf{y}, \mathbf{X}) p(h | \mathbf{y}, \mathbf{X}) dh \tag{20}$$

All functions in the integrand are fully known, so this is straightforward, but tedious. We jump directly to the final result. It turns out that the marginal posterior of  $\boldsymbol{\beta}$  follows a  $t$ -distribution with mean  $\boldsymbol{\mu}_1$ , variance  $\tau_1^{-1}\mathbf{V}_1$ , and DoF  $v_1$ , i.e.

$$\boldsymbol{\beta} | \mathbf{y}, \mathbf{X} \sim t(\boldsymbol{\mu}_1, \tau_1^{-1}\mathbf{V}_1, v_1). \tag{21}$$

#### Step 7:

The final construct of interest is the marginal likelihood  $p(\mathbf{y})$ . It is useful for model comparison and hypothesis testing. For the conjugate normal model it has an analytical solution derived via

$$p(\mathbf{y}) = \int p(\mathbf{y} | \boldsymbol{\theta}) p(\boldsymbol{\theta}) d\boldsymbol{\theta} \quad \boldsymbol{\theta} = [\boldsymbol{\beta}' \quad h] \tag{22}$$

We will again skip the details of this derivation. The final result is

$$p(\mathbf{y}) = \frac{\Gamma(v_1/2)(v_0\tau_0^{-1})^{v_0/2}}{\Gamma(v_0/2)\pi^{n/2}} \left(\frac{|\mathbf{V}_1|}{|\mathbf{V}_0|}\right)^{1/2} (v_1\tau_1^{-1})^{-v_1/2} \tag{23}$$

## R Implementation

`R` script `mod6s2b` follows these lecture notes exactly. Since with conjugate priors we know the exact form of marginal posteriors for all parameters, there is no “posterior simulation”, such as a Gibbs Sampler, required for this model.

We are usually interested in the posterior mean and standard deviation for our parameters. From the lecture notes we know that

$$\boldsymbol{\beta} | \mathbf{y}, \mathbf{X} \sim t(\boldsymbol{\mu}_1, \tau_1^{-1} \mathbf{V}_1, \nu_1) \quad \text{and} \quad h | \mathbf{y}, \mathbf{X} \sim g(\tau_1, \nu_1)$$

Thus, focusing on a specific element of  $\boldsymbol{\beta}$ , say  $\beta_j$ , we immediately know  $E(\beta_j) = \mu_{1,j}$  and  $std(\beta_j) = \sqrt{\tau_1^{-1} \mathbf{V}_{1,jj}}$ , where  $\mu_{1,j}$  is the  $j^{\text{th}}$  element of  $\boldsymbol{\mu}_1$ , and  $\mathbf{V}_{1,jj}$  is the  $(j,j)$  element of  $\mathbf{V}_1$  (i.e. the  $j^{\text{th}}$  diagonal element). Similarly, we know  $E(h) = \tau_1$  and, by the properties of the gamma density we can directly compute  $std(h) = \sqrt{\frac{2\tau_1^2}{\nu_1}}$ . These posterior moments are sent to the log file.

Script `mod6s2b` also takes draws from the prior and posterior distributions for all parameters and compares the resulting densities in a set of nonparametric plots. It is clear that the posterior densities are much tighter than the priors in all cases. Thus, the sample data have indeed contributed to “learning” about the parameters.

## Model Comparison

We compare models with different sets of regressors using their (logged) marginal likelihood and Bayes Factors. Specifically, we estimate a known-to-be- false model that drops the last regressor in  $\mathbf{X}$  in script `mod6s2c`.

Scripts `mod6s2d` loads the marginal likelihood from both models and computes the logged Bayes Factor of model 1 (“main”) over model 2 (“constrained”).

Kass and Raftery (1995) propose the following decision rules for one model over another based on logged Bayes Factors:

<u>Log BF</u>	<u>Strength of evidence in favor of model with higher marg. LH</u>
0 - 6.9	“positive”
6.9 – 11.5	“strong”
>11.5	“decisive”

In this case, the first model receives “decisive” support from the sample data. In other words, it is much more probable for the first model to have generated the sample data than for the second model.

### References:

Kass, R. E. and A. E. Raftery, 1995, Bayes Factors. *Journal of the American Statistical Association*, 90, 773-795.