

# Contents

<b>6</b>	<b>Conjugate families</b>	<b>1</b>
6.1	Binomial - beta prior . . . . .	2
6.1.1	Uninformative priors . . . . .	2
6.2	Gaussian (unknown mean, known variance) . . . . .	2
6.2.1	Uninformative prior . . . . .	4
6.3	Gaussian (known mean, unknown variance) . . . . .	4
6.3.1	Uninformative prior . . . . .	5
6.4	Gaussian (unknown mean, unknown variance) . . . . .	5
6.4.1	Completing the square . . . . .	6
6.4.2	Marginal posterior distributions . . . . .	7
6.4.3	Uninformative priors . . . . .	9
6.5	Multivariate Gaussian (unknown mean, known variance) . .	11
6.5.1	Completing the square . . . . .	11
6.5.2	Uninformative priors . . . . .	12
6.6	Multivariate Gaussian (unknown mean, unknown variance)	13
6.6.1	Completing the square . . . . .	14
6.6.2	Inverted-Wishart kernel . . . . .	14
6.6.3	Marginal posterior distributions . . . . .	15
6.6.4	Uninformative priors . . . . .	17
6.7	Bayesian linear regression . . . . .	18
6.7.1	Known variance . . . . .	19
6.7.2	Unknown variance . . . . .	22
6.7.3	Uninformative priors . . . . .	26
6.8	Bayesian linear regression with general error structure . . .	27

iv Contents

6.8.1	Known variance . . . . .	28
6.8.2	Unknown variance . . . . .	29
6.8.3	(Nearly) uninformative priors . . . . .	32
6.9	Appendix: summary of conjugacy . . . . .	34

# 6

## Conjugate families

Conjugate families arise when the likelihood times the prior produces a recognizable posterior kernel

$$p(\theta | y) \propto \ell(\theta | y) p(\theta)$$

where the *kernel* is the characteristic part of the distribution function that depends on the random variable(s) (the part excluding any normalizing constants). For example, the density function for a univariate Gaussian or normal is

$$\frac{1}{\sqrt{2\pi}\sigma} \exp\left[-\frac{1}{2\sigma^2} (x - \mu)^2\right]$$

and its kernel (for  $\sigma$  known) is

$$\exp\left[-\frac{1}{2\sigma^2} (x - \mu)^2\right]$$

as  $\frac{1}{\sqrt{2\pi}\sigma}$  is a normalizing constant. Now, we discuss a few common conjugate family results<sup>1</sup> and uninformative prior results to connect with classical results.

---

<sup>1</sup>A more complete set of conjugate families are summarized in chapter 7 of *Accounting and Causal Effects: Econometric Challenges* as well as tabulated in an appendix at the end of the chapter.

## 6.1 Binomial - beta prior

A binomial likelihood with unknown success probability,  $\theta$ ,

$$\ell(\theta | s; n) = \binom{n}{s} \theta^s (1 - \theta)^{n-s}$$

$$s = \sum_{i=1}^n y_i, \quad y_i = \{0, 1\}$$

combines with a beta( $\theta; a, b$ ) prior (i.e., with parameters  $a$  and  $b$ )

$$p(\theta) = \frac{\Gamma(a+b)}{\Gamma(a)\Gamma(b)} \theta^{a-1} (1-\theta)^{b-1}$$

to yield

$$\begin{aligned} p(\theta | y) &\propto \theta^s (1-\theta)^{n-s} \theta^{a-1} (1-\theta)^{b-1} \\ &\propto \theta^{s+a-1} (1-\theta)^{n-s+b-1} \end{aligned}$$

which is the kernel of a beta distribution with parameters  $(a+s)$  and  $(b+n-s)$ , beta( $\theta | y; a+s, b+n-s$ ).

### 6.1.1 Uninformative priors

Suppose priors for  $\theta$  are uniform over the interval zero to one or, equivalently, beta(1,1).<sup>2</sup> Then, the likelihood determines the posterior distribution for  $\theta$ .

$$p(\theta | y) \propto \theta^s (1-\theta)^{n-s}$$

which is beta( $\theta | y; 1+s, 1+n-s$ ).

## 6.2 Gaussian (unknown mean, known variance)

A single draw from a Gaussian likelihood with unknown mean,  $\theta$ , known standard deviation,  $\sigma$ ,

$$\ell(\theta | y, \sigma) \propto \exp \left[ -\frac{1}{2} \frac{(y - \theta)^2}{\sigma^2} \right]$$

combines with a Gaussian or normal prior for  $\theta$  given  $\sigma^2$  with prior mean  $\theta_0$  and prior variance  $\tau_0^2$

$$p(\theta | \sigma^2; \theta_0, \tau_0^2) \propto \exp \left[ -\frac{1}{2} \frac{(\theta - \theta_0)^2}{\tau_0^2} \right]$$

---

<sup>2</sup>Some would utilize Jeffreys' prior,  $p(\theta) \propto \text{beta}(\theta; \frac{1}{2}, \frac{1}{2})$ , which is invariant to transformation, as the uninformative prior.

or writing  $\tau_0^2 \equiv \sigma^2/\kappa_0$ , we have

$$p(\theta \mid \sigma^2; \theta_0, \sigma^2/\kappa_0) \propto \exp \left[ -\frac{1}{2} \frac{\kappa_0 (\theta - \theta_0)^2}{\sigma^2} \right]$$

to yield

$$p(\theta \mid y, \sigma, \theta_0, \sigma^2/\kappa_0) \propto \exp \left[ -\frac{1}{2} \left( \frac{(y - \theta)^2}{\sigma^2} + \frac{\kappa_0 (\theta - \theta_0)^2}{\sigma^2} \right) \right]$$

Expansion and rearrangement gives

$$p(\theta \mid y, \sigma, \theta_0, \sigma^2/\kappa_0) \propto \exp \left[ -\frac{1}{2\sigma^2} (y^2 + \kappa_0 \theta_0^2 - 2y\theta + \theta^2 + \kappa_0 (\theta^2 - 2\theta\theta_0)) \right]$$

Any terms not involving  $\theta$  are constants and can be discarded as they are absorbed on normalization of the posterior

$$p(\theta \mid y, \sigma, \theta_0, \sigma^2/\kappa_0) \propto \exp \left[ -\frac{1}{2\sigma^2} (\theta^2 (\kappa_0 + 1) - 2\theta (\kappa_0 \theta_0 + y)) \right]$$

Completing the square (add and subtract  $\frac{(\kappa_0 \theta_0 + y)^2}{\kappa_0 + 1}$ ), dropping the term subtracted (as it's all constants), and factoring out  $(\kappa_0 + 1)$  gives

$$p(\theta \mid y, \sigma, \theta_0, \sigma^2/\kappa_0) \propto \exp \left[ -\frac{\kappa_0 + 1}{2\sigma^2} \left( \theta - \frac{\kappa_0 \theta_0 + y}{\kappa_0 + 1} \right)^2 \right]$$

Finally, we have

$$p(\theta \mid y, \sigma, \theta_0, \sigma^2/\kappa_0) \propto \exp \left[ -\frac{1}{2} \frac{(\theta - \theta_1)^2}{\tau_1^2} \right]$$

where  $\theta_1 = \frac{\kappa_0 \theta_0 + y}{\kappa_0 + 1} = \frac{\frac{1}{\tau_0} \theta_0 + \frac{1}{\sigma^2} y}{\frac{1}{\tau_0} + \frac{1}{\sigma^2}}$  and  $\tau_1^2 = \frac{\sigma^2}{\kappa_0 + 1} = \frac{1}{\frac{1}{\tau_0} + \frac{1}{\sigma^2}}$ , or the posterior distribution of the mean given the data and priors is Gaussian or normal. Notice, the posterior mean,  $\theta_1$ , weights the data and prior beliefs by their relative precisions.

For a sample of  $n$  exchangeable draws, the likelihood is

$$\ell(\theta \mid y, \sigma) \propto \prod_{i=1}^n \exp \left[ -\frac{1}{2} \frac{(y_i - \theta)^2}{\sigma^2} \right]$$

combined with the above prior yields

$$p(\theta \mid y, \sigma, \theta_0, \sigma^2/\kappa_0) \propto \exp \left[ -\frac{1}{2} \frac{(\theta - \theta_n)^2}{\tau_n^2} \right]$$

where  $\theta_n = \frac{\kappa_0 \theta_0 + n \bar{y}}{\kappa_0 + n} = \frac{\frac{1}{\tau_0} \theta_0 + \frac{n}{\sigma^2} \bar{y}}{\frac{1}{\tau_0} + \frac{n}{\sigma^2}}$ ,  $\bar{y}$  is the sample mean, and  $\tau_n^2 = \frac{\sigma^2}{\kappa_0 + n} = \frac{1}{\frac{1}{\tau_0} + \frac{n}{\sigma^2}}$ , or the posterior distribution of the mean,  $\theta$ , given the data and priors is again Gaussian or normal and the posterior mean,  $\theta_n$ , weights the data and priors by their relative precisions.

### 6.2.1 Uninformative prior

An uninformative prior for the mean,  $\theta$ , is the (improper) uniform,  $p(\theta | \sigma^2) = 1$ . Hence, the likelihood

$$\begin{aligned} \ell(\theta | y, \sigma) &\propto \prod_{i=1}^n \exp\left[-\frac{1}{2} \frac{(y_i - \theta)^2}{\sigma^2}\right] \\ &\propto \exp\left[-\frac{1}{2\sigma^2} \left\{ \sum_{i=1}^n y_i^2 - 2n\bar{y}\theta + n\theta^2 \right\}\right] \\ &\propto \exp\left[-\frac{1}{2\sigma^2} \left\{ \sum_{i=1}^n y_i^2 - n\bar{y}^2 + n(\theta - \bar{y})^2 \right\}\right] \\ &\propto \exp\left[-\frac{1}{2\sigma^2} n(\theta - \bar{y})^2\right] \end{aligned}$$

determines the posterior

$$p(\theta | \sigma^2, y) \propto \exp\left[-\frac{n}{2} \frac{(\theta - \bar{y})^2}{\sigma^2}\right]$$

which is the kernel for a Gaussian or  $N(\theta | \sigma^2, y; \bar{y}, \frac{\sigma^2}{n})$ , the classical result.

## 6.3 Gaussian (known mean, unknown variance)

For a sample of  $n$  exchangeable draws with known mean,  $\mu$ , and unknown variance,  $\theta$ , a Gaussian or normal likelihood is

$$\ell(\theta | y, \mu) \propto \prod_{i=1}^n \theta^{-\frac{1}{2}} \exp\left[-\frac{1}{2} \frac{(y_i - \mu)^2}{\theta}\right]$$

combines with an inverted-gamma( $a, b$ )

$$p(\theta; a, b) \propto \theta^{-(a+1)} \exp\left[-\frac{b}{\theta}\right]$$

to yield an inverted-gamma( $\frac{n+2a}{2}, b + \frac{1}{2}t$ ) posterior distribution where

$$t = \sum_{i=1}^n (y_i - \mu)^2$$

Alternatively and conveniently (but equivalently), we could parameterize the prior as an inverted-chi square<sup>3</sup> $(\nu_0, \sigma_0^2)$

$$p(\theta; \nu_0, \sigma_0^2) \propto (\theta)^{-\left(\frac{\nu_0}{2}+1\right)} \exp\left[-\frac{\nu_0 \sigma_0^2}{2\theta}\right]$$

and combine with the above likelihood to yield

$$p(\theta | y) \propto \theta^{-\left(\frac{n+\nu_0}{2}+1\right)} \exp\left[-\frac{1}{2\theta}(\nu_0 \sigma_0^2 + t)\right]$$

an inverted chi-square $\left(\nu_0 + n, \frac{\nu_0 \sigma_0^2 + t}{\nu_0 + n}\right)$ .

### 6.3.1 Uninformative prior

An uninformative prior for scale is

$$p(\theta) \propto \theta^{-1}$$

Hence, the posterior distribution for scale is

$$p(\theta | y) \propto \theta^{-\left(\frac{n}{2}+1\right)} \exp\left[-\frac{t}{2\theta}\right]$$

which is the kernel of an inverted-chi square $\left(\theta; n, \frac{t}{n}\right)$ .

## 6.4 Gaussian (unknown mean, unknown variance)

For a sample of  $n$  exchangeable draws, a normal likelihood with unknown mean,  $\theta$ , and unknown (but constant) variance,  $\sigma^2$ , is

$$\ell(\theta, \sigma^2 | y) \propto \prod_{i=1}^n \sigma^{-1} \exp\left[-\frac{1}{2} \frac{(y_i - \theta)^2}{\sigma^2}\right]$$

Expanding and rewriting the likelihood gives

$$\ell(\theta, \sigma^2 | y) \propto \sigma^{-n} \exp\left[\sum_{i=1}^n -\frac{1}{2} \frac{y_i^2 - 2y_i\theta + \theta^2}{\sigma^2}\right]$$

Adding and subtracting  $\sum_{i=1}^n 2y_i\bar{y} = 2n\bar{y}^2$ , we write

$$\ell(\theta, \sigma^2 | y) \propto (\sigma^2)^{-\frac{n}{2}} \exp\left[-\frac{1}{2\sigma^2} \sum_{i=1}^n \{(y_i^2 - 2y_i\bar{y} + \bar{y}^2) + (\bar{y}^2 - 2\bar{y}\theta + \theta^2)\}\right]$$

---

<sup>3</sup>  $\frac{\sigma_0^2 \nu_0}{X}$  is a scaled, inverted-chi square $(\nu_0, \sigma_0^2)$  with scale  $\sigma_0^2$  where  $X$  is a chi square $(\nu_0)$  random variable.

or

$$\ell(\theta, \sigma^2 | y) \propto (\sigma^2)^{-\frac{n}{2}} \exp \left[ -\frac{1}{2\sigma^2} \sum_{i=1}^n \left\{ (y_i - \bar{y})^2 + (\bar{y} - \theta)^2 \right\} \right]$$

which can be rewritten as

$$\ell(\theta, \sigma^2 | y) \propto (\sigma^2)^{-\frac{n}{2}} \exp \left[ -\frac{1}{2\sigma^2} \left\{ (n-1)s^2 + n(\bar{y} - \theta)^2 \right\} \right]$$

where  $s^2 = \frac{1}{n-1} \sum_{i=1}^n (y_i - \bar{y})^2$ . The above likelihood combines with a Gaussian or normal( $\theta | \sigma^2; \theta_0, \sigma^2/\kappa_0$ )  $\times$  inverted-chi square( $\sigma^2; \nu_0, \sigma_0^2$ ) prior<sup>4</sup>

$$\begin{aligned} p(\theta | \sigma^2; \theta_0, \sigma^2/\kappa_0) \times p(\sigma^2; \nu_0, \sigma_0^2) &\propto \frac{1}{\sigma} \exp \left[ -\frac{\kappa_0 (\theta - \theta_0)^2}{2\sigma^2} \right] \\ &\times (\sigma^2)^{-(\nu_0/2+1)} \exp \left[ -\frac{\nu_0 \sigma_0^2}{2\sigma^2} \right] \\ &\propto (\sigma^2)^{-\left(\frac{\nu_0+3}{2}\right)} \\ &\times \exp \left[ -\frac{\nu_0 \sigma_0^2 + \kappa_0 (\theta - \theta_0)^2}{2\sigma^2} \right] \end{aligned}$$

to yield a normal( $\theta | \sigma^2; \theta_n, \sigma_n^2/\kappa_n$ )  $\times$  inverted-chi square( $\sigma^2; \nu_n, \sigma_n^2$ ) joint posterior distribution<sup>5</sup> where

$$\begin{aligned} \nu_n &= \nu_0 + n \\ \kappa_n &= \kappa_0 + n \\ \nu_n \sigma_n^2 &= \nu_0 \sigma_0^2 + (n-1)s^2 + \frac{\kappa_0 n}{\kappa_0 + n} (\theta_0 - \bar{y})^2 \end{aligned}$$

That is, the joint posterior is

$$\begin{aligned} p(\theta, \sigma^2 | y; \theta_0, \sigma^2/\kappa_0, \nu_0, \sigma_0^2) &\propto (\sigma^2)^{-\frac{n+\nu_0+3}{2}} \\ &\times \exp \left[ -\frac{1}{2\sigma^2} \left\{ \begin{array}{l} \nu_0 \sigma_0^2 + (n-1)s^2 \\ + \kappa_0 (\theta - \theta_0)^2 \\ + n(\theta - \bar{y})^2 \end{array} \right\} \right] \end{aligned}$$

#### 6.4.1 Completing the square

The expression for the joint posterior is written by completing the square. Completing the weighted square for  $\theta$  centered around

$$\theta_n = \frac{1}{\kappa_0 + n} (\kappa_0 \theta_0 + n\bar{y})$$

<sup>4</sup>The prior for the mean,  $\theta$ , is conditional on the scale of the data,  $\sigma^2$ .

<sup>5</sup>The product of normal or Gaussian kernels produces a Gaussian kernel.



where  $\bar{y} = \frac{1}{n} \sum_{i=1}^n y_i$  gives

$$\begin{aligned} (\kappa_0 + n) (\theta - \theta_n)^2 &= (\kappa_0 + n) \theta^2 - 2\theta (\kappa_0 + n) \theta_n + (\kappa_0 + n) \theta_n^2 \\ &= (\kappa_0 + n) \theta^2 - 2\theta (\kappa_0 \theta_0 + n\bar{y}) + (\kappa_0 + n) \theta_n^2 \end{aligned}$$

While expanding the exponent includes the square plus additional terms as follows

$$\begin{aligned} \kappa_0 (\theta - \theta_0)^2 + n (\theta - \bar{y})^2 &= \kappa_0 (\theta^2 - 2\theta\theta_0 + \theta_0^2) + n (\theta^2 - 2\theta\bar{y} + \bar{y}^2) \\ &= (\kappa_0 + n) \theta^2 - 2\theta (\kappa_0 \theta_0 + n\bar{y}) + \kappa_0 \theta_0^2 + n\bar{y}^2 \end{aligned}$$

Add and subtract  $(\kappa_0 + n) \theta_n^2$  and simplify.

$$\begin{aligned} \kappa_0 (\theta - \theta_0)^2 + n (\theta - \bar{y})^2 &= (\kappa_0 + n) \theta^2 - 2\theta (\kappa_0 + n) \theta_n + (\kappa_0 + n) \theta_n^2 \\ &\quad - (\kappa_0 + n) \theta_n^2 + \kappa_0 \theta_0^2 + n\bar{y}^2 \\ &= (\kappa_0 + n) (\theta - \theta_n)^2 \\ &\quad \frac{1}{(\kappa_0 + n)} \left\{ \begin{array}{l} (\kappa_0 + n) (\kappa_0 \theta_0^2 + n\bar{y}^2) \\ - (\kappa_0 \theta_0 + n\bar{y})^2 \end{array} \right\} \end{aligned}$$

Expand and simplify the last term.

$$\kappa_0 (\theta - \theta_0)^2 + n (\theta - \bar{y})^2 = (\kappa_0 + n) (\theta - \theta_n)^2 + \frac{\kappa_0 n}{\kappa_0 + n} (\theta_0 - \bar{y})^2$$

Now, the joint posterior can be rewritten as

$$\begin{aligned} p(\theta, \sigma^2 | y; \theta_0, \sigma^2 / \kappa_0, \nu_0, \sigma_0^2) &\propto (\sigma^2)^{-\frac{n+\nu_0+3}{2}} \\ &\quad \times \exp \left[ -\frac{1}{2\sigma^2} \left\{ \begin{array}{l} \nu_0 \sigma_0^2 + (n-1) s^2 \\ + \frac{\kappa_0 n}{\kappa_0 + n} (\theta_0 - \bar{y})^2 \\ + (\kappa_0 + n) (\theta - \theta_n)^2 \end{array} \right\} \right] \end{aligned}$$

or

$$\begin{aligned} p(\theta, \sigma^2 | y; \theta_0, \sigma^2 / \kappa_0, \nu_0, \sigma_0^2) &\propto (\sigma^2)^{-\frac{n+\nu_0}{2}-1} \exp \left[ -\frac{1}{2\sigma^2} \nu_n \sigma_n^2 \right] \\ &\quad \times \sigma^{-1} \exp \left[ -\frac{1}{2\sigma^2} (\kappa_0 + n) (\theta - \theta_n)^2 \right] \end{aligned}$$

Hence, the conditional posterior distribution for the mean,  $\theta$ , given  $\sigma^2$  is Gaussian or normal  $\left( \theta | \sigma^2; \theta_n, \frac{\sigma^2}{\kappa_0 + n} \right)$ .

### 6.4.2 Marginal posterior distributions

We're often interested in the marginal posterior distributions which are derived by integrating out the other parameter from the joint posterior. The

marginal posterior for the mean,  $\theta$ , on integrating out  $\sigma^2$  is a noncentral, scaled-Student  $t\left(\theta; \theta_n, \frac{\sigma_n^2}{\kappa_n}, \nu_n\right)$ <sup>6</sup> for the mean

$$p(\theta; \theta_n, \sigma_n^2, \kappa_n, \nu_n) \propto \left( \frac{\nu_n}{\nu_n + \frac{\kappa_n(\theta - \theta_n)^2}{\sigma_n^2}} \right)^{\frac{\nu_n+1}{2}}$$

or

$$p\left(\theta; \theta_n, \frac{\nu_n \sigma_n^2}{\kappa_n}, \nu_n\right) \propto \left( 1 + \frac{\kappa_n(\theta - \theta_n)^2}{\nu_n \sigma_n^2} \right)^{-\frac{\nu_n+1}{2}}$$

and the marginal posterior for the variance,  $\sigma^2$ , is an inverted-chi square  $(\sigma^2; \nu_n, \sigma_n^2)$  on integrating out  $\theta$ .

$$p(\sigma^2; \nu_n, \sigma_n^2) \propto (\sigma^2)^{-(\nu_n/2+1)} \exp\left[-\frac{\nu_n \sigma_n^2}{2\sigma^2}\right]$$

Derivation of the marginal posterior for the mean,  $\theta$ , is as follows. Let  $z = \frac{A}{2\sigma^2}$  where

$$\begin{aligned} A &= \nu_0 \sigma_0^2 + (n-1)s^2 + \frac{\kappa_0 n}{\kappa_0 + n} (\theta_0 - \bar{y})^2 + (\kappa_0 + n) (\theta - \theta_n)^2 \\ &= \nu_n \sigma_n^2 + (\kappa_0 + n) (\theta - \theta_n)^2 \end{aligned}$$

The marginal posterior for the mean,  $\theta$ , integrates out  $\sigma^2$  from the joint posterior

$$\begin{aligned} p(\theta | y) &= \int_0^\infty p(\theta, \sigma^2 | y) d\sigma^2 \\ &= \int_0^\infty (\sigma^2)^{-\frac{n+\nu_0+3}{2}} \exp\left[-\frac{A}{2\sigma^2}\right] d\sigma^2 \end{aligned}$$

Utilizing  $\sigma^2 = \frac{A}{2z}$  and  $dz = -\frac{2z^2}{A} d\sigma^2$  or  $d\sigma^2 = -\frac{A}{2z^2} dz$ ,

$$\begin{aligned} p(\theta | y) &\propto \int_0^\infty \left(\frac{A}{2z}\right)^{-\frac{n+\nu_0+3}{2}} \frac{A}{2z^2} \exp[-z] dz \\ &\propto \int_0^\infty \left(\frac{A}{2z}\right)^{-\frac{n+\nu_0+1}{2}} z^{-1} \exp[-z] dz \\ &\propto A^{-\frac{n+\nu_0+1}{2}} \int_0^\infty z^{\frac{n+\nu_0+1}{2}-1} \exp[-z] dz \end{aligned}$$

---

<sup>6</sup>The noncentral, scaled-Student  $t(\theta; \theta_n, \sigma_n^2/\kappa_n, \nu_n)$  implies  $\frac{\theta - \theta_n}{\sigma_n/\sqrt{\kappa_n}}$  has a standard Student- $t(\nu_n)$  distribution  $p(\theta | y) \propto \left[ 1 + \frac{\left(\frac{\theta - \theta_n}{\sigma_n/\sqrt{\kappa_n}}\right)^2}{\nu_n} \right]^{-\frac{\nu_n+1}{2}}$ .

The integral  $\int_0^\infty z^{\frac{n+\nu_0+1}{2}-1} \exp[-z] dz$  is a constant since it is the kernel of a gamma density and therefore can be ignored when deriving the kernel of the marginal posterior for the mean

$$\begin{aligned} p(\theta | y) &\propto A^{-\frac{n+\nu_0+1}{2}} \\ &\propto \left[ \nu_n \sigma_n^2 + (\kappa_0 + n) (\theta - \theta_n)^2 \right]^{-\frac{n+\nu_0+1}{2}} \\ &\propto \left[ 1 + \frac{(\kappa_0 + n) (\theta - \theta_n)^2}{\nu_n \sigma_n^2} \right]^{-\frac{n+\nu_0+1}{2}} \end{aligned}$$

which is the kernel for a noncentral, scaled Student  $t\left(\theta; \theta_n, \frac{\sigma_n^2}{\kappa_0+n}, n + \nu_0\right)$ .

Derivation of the marginal posterior for  $\sigma^2$  is somewhat simpler. Write the joint posterior in terms of the conditional posterior for the mean multiplied by the marginal posterior for  $\sigma^2$ .

$$p(\theta, \sigma^2 | y) = p(\theta | \sigma^2, y) p(\sigma^2 | y)$$

Marginalization of  $\sigma^2$  is achieved by integrating out  $\theta$ .

$$p(\sigma^2 | y) = \int_{-\infty}^{\infty} p(\sigma^2 | y) p(\theta | \sigma^2, y) d\theta$$

Since only the conditional posterior involves  $\theta$  the marginal posterior for  $\sigma^2$  is immediate.

$$\begin{aligned} p(\theta, \sigma^2 | y) &\propto (\sigma^2)^{-\frac{n+\nu_0+3}{2}} \exp\left[-\frac{A}{2\sigma^2}\right] \\ &\propto (\sigma^2)^{-\frac{n+\nu_0+2}{2}} \exp\left[-\frac{\nu_n \sigma_n^2}{2\sigma^2}\right] \sigma^{-1} \exp\left[-\frac{(\kappa_0 + n) (\theta - \theta_n)^2}{2\sigma^2}\right] \end{aligned}$$

Integrating out  $\theta$  yields

$$\begin{aligned} p(\sigma^2 | y) &\propto (\sigma^2)^{-\frac{n+\nu_0+2}{2}} \exp\left[-\frac{\nu_n \sigma_n^2}{2\sigma^2}\right] \\ &\quad \times \int_{-\infty}^{\infty} \sigma^{-1} \exp\left[-\frac{(\kappa_0 + n) (\theta - \theta_n)^2}{2\sigma^2}\right] d\theta \\ &\propto (\sigma^2)^{-\left(\frac{\nu_n}{2}+1\right)} \exp\left[-\frac{\nu_n \sigma_n^2}{2\sigma^2}\right] \end{aligned}$$

which we recognize as the kernel of an inverted-chi square  $(\sigma^2; \nu_n, \sigma_n^2)$ .

### 6.4.3 Uninformative priors

The case of uninformative priors is relatively straightforward. Since priors convey no information, the prior for the mean is uniform (proportional to

a constant,  $\kappa_0 \rightarrow 0$ ) and an uninformative prior for  $\sigma^2$  has  $\nu_0 \rightarrow 0$  degrees of freedom so that the joint prior is

$$p(\theta, \sigma^2) \propto (\sigma^2)^{-1}$$

The joint posterior is

$$\begin{aligned} p(\theta, \sigma^2 | y) &\propto (\sigma^2)^{-(n/2+1)} \exp\left[-\frac{1}{2\sigma^2} \left\{ (n-1)s^2 + n(\theta - \bar{y})^2 \right\}\right] \\ &\propto (\sigma^2)^{-[(n-1)/2+1]} \exp\left[-\frac{\sigma_n^2}{2\sigma^2}\right] \\ &\quad \times \sigma^{-1} \exp\left[-\frac{n}{2\sigma^2} (\theta - \bar{y})^2\right] \end{aligned}$$

where

$$\sigma_n^2 = (n-1)s^2$$

The conditional posterior for  $\theta$  given  $\sigma^2$  is Gaussian  $\left(\bar{y}, \frac{\sigma^2}{n}\right)$ . And, the marginal posterior for  $\theta$  is noncentral, scaled Student  $t\left(\bar{y}, \frac{s^2}{n}, n-1\right)$ , the classical estimator.

Derivation of the marginal posterior proceeds as above. The joint posterior is

$$p(\theta, \sigma^2 | y) \propto (\sigma^2)^{-(n/2+1)} \exp\left[-\frac{1}{2\sigma^2} \left\{ (n-1)s^2 + n(\theta - \bar{y})^2 \right\}\right]$$

Let  $z = \frac{A}{2\sigma^2}$  where  $A = (n-1)s^2 + n(\theta - \bar{y})^2$ . Now integrate  $\sigma^2$  out of the joint posterior following the transformation of variables.

$$\begin{aligned} p(\theta | y) &\propto \int_0^\infty (\sigma^2)^{-(n/2+1)} \exp\left[-\frac{A}{2\sigma^2}\right] d\sigma^2 \\ &\propto A^{-n/2} \int_0^\infty z^{n/2-1} e^{-z} dz \end{aligned}$$

As before, the integral involves the kernel of a gamma density and therefore is a constant which can be ignored. Hence,

$$\begin{aligned} p(\theta | y) &\propto A^{-n/2} \\ &\propto \left[ (n-1)s^2 + n(\theta - \bar{y})^2 \right]^{-\frac{n}{2}} \\ &\propto \left[ 1 + \frac{n(\theta - \bar{y})^2}{(n-1)s^2} \right]^{-\frac{n-1+1}{2}} \end{aligned}$$

which we recognize as the kernel of a noncentral, scaled Student  $t\left(\theta; \bar{y}, \frac{s^2}{n}, n-1\right)$ .

## 6.5 Multivariate Gaussian (unknown mean, known variance)

More than one random variable (the multivariate case) with joint Gaussian or normal likelihood is analogous to the univariate case with Gaussian conjugate prior. Consider a vector of  $k$  random variables (the sample is comprised of  $n$  draws for each random variable) with unknown mean,  $\theta$ , and known variance,  $\Sigma$ . For  $n$  exchangeable draws of the random vector (containing each of the  $m$  random variable), the multivariate Gaussian likelihood is

$$\ell(\theta | y, \Sigma) \propto \prod_{i=1}^n \exp \left[ -\frac{1}{2} (y_i - \theta)^T \Sigma^{-1} (y_i - \theta) \right]$$

where superscript  $T$  refers to transpose,  $y_i$  and  $\theta$  are  $k$  length vectors and  $\Sigma$  is a  $k \times k$  variance-covariance matrix. A Gaussian prior for the mean vector,  $\theta$ , with prior mean,  $\theta_0$ , and prior variance,  $\Upsilon_0$ , is

$$p(\theta | \Sigma; \theta_0, \Upsilon_0) \propto \exp \left[ -\frac{1}{2} (\theta - \theta_0)^T \Upsilon_0^{-1} (\theta - \theta_0) \right]$$

The product of the likelihood and prior yields the kernel of a multivariate posterior Gaussian distribution for the mean

$$\begin{aligned} p(\theta | \Sigma, y; \theta_0, \Upsilon_0) &\propto \exp \left[ -\frac{1}{2} (\theta - \theta_0)^T \Upsilon_0^{-1} (\theta - \theta_0) \right] \\ &\quad \times \exp \left[ \sum_{i=1}^n -\frac{1}{2} (y_i - \theta)^T \Sigma^{-1} (y_i - \theta) \right] \end{aligned}$$

### 6.5.1 Completing the square

Expanding terms in the exponent leads to

$$\begin{aligned} &(\theta - \theta_0)^T \Upsilon_0^{-1} (\theta - \theta_0) + \sum_{i=1}^n (y_i - \theta)^T \Sigma^{-1} (y_i - \theta) \\ = &\theta^T (\Upsilon_0^{-1} + n\Sigma^{-1}) \theta - 2\theta^T (\Upsilon_0^{-1}\theta_0 + n\Sigma^{-1}\bar{y}) \\ &+ \theta_0^T \Upsilon_0^{-1} \theta_0 + \sum_{i=1}^n y_i^T \Sigma^{-1} y_i \end{aligned}$$

where  $\bar{y}$  is the sample average. While completing the (weighted) square centered around

$$\bar{\theta} = (\Upsilon_0^{-1} + n\Sigma^{-1})^{-1} (\Upsilon_0^{-1}\theta_0 + n\Sigma^{-1}\bar{y})$$

leads to

$$\begin{aligned}
(\theta - \bar{\theta})^T (\Upsilon_0^{-1} + n\Sigma^{-1}) (\theta - \bar{\theta}) &= \theta^T (\Upsilon_0^{-1} + n\Sigma^{-1}) \theta \\
&\quad - 2\theta^T (\Upsilon_0^{-1} + n\Sigma^{-1}) \bar{\theta} \\
&\quad + \bar{\theta}^T (\Upsilon_0^{-1} + n\Sigma^{-1}) \bar{\theta}
\end{aligned}$$

Thus, adding and subtracting  $\bar{\theta}^T (\Upsilon_0^{-1} + n\Sigma^{-1}) \bar{\theta}$  in the exponent completes the square (with three extra terms).

$$\begin{aligned}
&(\theta - \theta_0)^T \Upsilon_0^{-1} (\theta - \theta_0) + \sum_{i=1}^n (y_i - \theta)^T \Sigma^{-1} (y_i - \theta) \\
&= \theta^T (\Upsilon_0^{-1} + n\Sigma^{-1}) \theta - 2\theta^T (\Upsilon_0^{-1} + n\Sigma^{-1}) \bar{\theta} + \bar{\theta}^T (\Upsilon_0^{-1} + n\Sigma^{-1}) \bar{\theta} \\
&\quad - \bar{\theta}^T (\Upsilon_0^{-1} + n\Sigma^{-1}) \bar{\theta} + \theta_0^T \Upsilon_0^{-1} \theta_0 + \sum_{i=1}^n y_i^T \Sigma^{-1} y_i \\
&= (\theta - \bar{\theta})^T (\Upsilon_0^{-1} + n\Sigma^{-1}) (\theta - \bar{\theta}) \\
&\quad - \bar{\theta}^T (\Upsilon_0^{-1} + n\Sigma^{-1}) \bar{\theta} + \theta_0^T \Upsilon_0^{-1} \theta_0 + \sum_{i=1}^n y_i^T \Sigma^{-1} y_i
\end{aligned}$$

Dropping constants (the last three extra terms unrelated to  $\theta$ ) gives

$$p(\theta \mid \Sigma, y; \theta_0, \Upsilon_0) \propto \exp \left[ -\frac{1}{2} (\theta - \bar{\theta})^T (\Upsilon_0^{-1} + n\Sigma^{-1}) (\theta - \bar{\theta}) \right]$$

Hence, the posterior for the mean  $\theta$  has expected value  $\bar{\theta}$  and variance

$$\text{Var}[\theta \mid y, \Sigma, \theta_0, \Upsilon_0] = (\Upsilon_0^{-1} + n\Sigma^{-1})^{-1}$$

As in the univariate case, the data and prior beliefs are weighted by their relative precisions.

### 6.5.2 Uninformative priors

Uninformative priors for  $\theta$  are proportional to a constant. Hence, the likelihood determines the posterior

$$\ell(\theta \mid \Sigma, y) \propto \exp \left[ -\frac{1}{2} \sum_{i=1}^n (y_i - \theta)^T \Sigma^{-1} (y_i - \theta) \right]$$

Expanding the exponent and adding and subtracting  $n\bar{y}^T \Sigma^{-1} \bar{y}$  (to complete the square) gives

$$\begin{aligned} \sum_{i=1}^n (y_i - \theta)^T \Sigma^{-1} (y_i - \theta) &= \sum_{i=1}^n y_i^T \Sigma^{-1} y_i - 2n\theta^T \Sigma^{-1} \bar{y} + n\theta^T \Sigma^{-1} \theta \\ &\quad + n\bar{y}^T \Sigma^{-1} \bar{y} - n\bar{y}^T \Sigma^{-1} \bar{y} \\ &= n(\bar{y} - \theta)^T \Sigma^{-1} (\bar{y} - \theta) \\ &\quad + \sum_{i=1}^n y_i^T \Sigma^{-1} y_i - n\bar{y}^T \Sigma^{-1} \bar{y} \end{aligned}$$

The latter two terms are constants, hence, the posterior kernel is

$$p(\theta \mid \Sigma, y) \propto \exp \left[ -\frac{n}{2} (\bar{y} - \theta)^T \Sigma^{-1} (\bar{y} - \theta) \right]$$

which is Gaussian or  $N(\theta; \bar{y}, \frac{1}{n}\Sigma)$ , the classical result.

## 6.6 Multivariate Gaussian (unknown mean, unknown variance)

When both the mean,  $\theta$ , and variance,  $\Sigma$ , are unknown, the multivariate Gaussian cases remains analogous to the univariate case. Specifically, a Gaussian likelihood

$$\begin{aligned} \ell(\theta, \Sigma \mid y) &\propto \prod_{i=1}^n |\Sigma|^{-\frac{1}{2}} \exp \left[ -\frac{1}{2} (y_i - \theta)^T \Sigma^{-1} (y_i - \theta) \right] \\ &\propto |\Sigma|^{-\frac{n}{2}} \exp \left[ -\frac{1}{2} \left\{ \sum_{i=1}^n (y_i - \bar{y})^T \Sigma^{-1} (y_i - \bar{y}) + n(\bar{y} - \theta)^T \Sigma^{-1} (\bar{y} - \theta) \right\} \right] \\ &\propto |\Sigma|^{-\frac{n}{2}} \exp \left[ -\frac{1}{2} \left\{ (n-1)s^2 + n(\bar{y} - \theta)^T \Sigma^{-1} (\bar{y} - \theta) \right\} \right] \end{aligned}$$

where  $s^2 = \frac{1}{n-1} \sum_{i=1}^n (y_i - \bar{y})^T \Sigma^{-1} (y_i - \bar{y})$  combines with a Gaussian-inverted Wishart prior

$$\begin{aligned} p\left(\theta \mid \Sigma; \theta_0, \frac{\Sigma}{\kappa_0}\right) \times p(\Sigma^{-1}; \nu, \Psi) &\propto |\Sigma|^{-\frac{1}{2}} \exp \left[ -\frac{1}{2} (\theta - \theta_0)^T \kappa_0 \Sigma^{-1} (\theta - \theta_0) \right] \\ &\quad \times |\Psi|^{\frac{\nu}{2}} |\Sigma|^{-\frac{\nu+k+1}{2}} \exp \left[ -\frac{\text{tr}(\Psi \Sigma^{-1})}{2} \right] \end{aligned}$$

where  $tr(\cdot)$  is the trace of the matrix and  $\nu$  is degrees of freedom, to produce

$$p(\theta, \Sigma | y) \propto |\Psi|^{\frac{\nu}{2}} |\Sigma|^{-\frac{\nu+n+k+1}{2}} \exp \left[ -\frac{tr(\Psi \Sigma^{-1})}{2} \right] \\ \times |\Sigma|^{-\frac{1}{2}} \exp \left[ -\frac{1}{2} \left\{ \begin{array}{l} (n-1) s^2 + n(\bar{y} - \theta)^T \Sigma^{-1} (\bar{y} - \theta) \\ + \kappa_0 (\theta - \theta_0)^T \Sigma^{-1} (\theta - \theta_0) \end{array} \right\} \right]$$

### 6.6.1 Completing the square

Completing the square involves the matrix analog to the univariate unknown mean and variance case. Consider the exponent (in braces)

$$\begin{aligned} & (n-1) s^2 + n(\bar{y} - \theta)^T \Sigma^{-1} (\bar{y} - \theta) + \kappa_0 (\theta - \theta_0)^T \Sigma^{-1} (\theta - \theta_0) \\ = & (n-1) s^2 + n\bar{y}^T \Sigma^{-1} \bar{y} - 2n\theta^T \Sigma^{-1} \bar{y} + n\theta^T \Sigma^{-1} \theta \\ & + \kappa_0 \theta^T \Sigma^{-1} \theta - 2\kappa_0 \theta^T \Sigma^{-1} \theta_0 + \kappa_0 \theta_0^T \Sigma^{-1} \theta_0 \\ = & (n-1) s^2 + (\kappa_0 + n) \theta^T \Sigma^{-1} \theta - 2\theta^T \Sigma^{-1} (\kappa_0 \theta_0 + n\bar{y}) + (\kappa_0 + n) \theta_n^T \Sigma^{-1} \theta_n \\ & - (\kappa_0 + n) \theta_n^T \Sigma^{-1} \theta_n + \kappa_0 \theta_0^T \Sigma^{-1} \theta_0 + n\bar{y}^T \Sigma^{-1} \bar{y} \\ = & (n-1) s^2 + (\kappa_0 + n) (\theta - \theta_n)^T \Sigma^{-1} (\theta - \theta_n) \\ & + \frac{\kappa_0 n}{\kappa_0 + n} (\theta_0 - \bar{y})^T \Sigma^{-1} (\theta_0 - \bar{y}) \end{aligned}$$

Hence, the joint posterior can be rewritten as

$$p(\theta, \Sigma | y) \propto |\Psi|^{\frac{\nu}{2}} |\Sigma|^{-\frac{\nu+n+k+1}{2}} \exp \left[ -\frac{tr(\Psi \Sigma^{-1})}{2} \right] \\ \times |\Sigma|^{-\frac{1}{2}} \exp \left[ -\frac{1}{2} \left\{ \begin{array}{l} (\kappa_0 + n) (\theta - \theta_n)^T \Sigma^{-1} (\theta - \theta_n) \\ + (n-1) s^2 \\ + \frac{\kappa_0 n}{\kappa_0 + n} (\theta_0 - \bar{y})^T \Sigma^{-1} (\theta_0 - \bar{y}) \end{array} \right\} \right] \\ \propto |\Psi|^{\frac{\nu}{2}} |\Sigma|^{-\frac{\nu+n+k+1}{2}} \exp \left[ -\frac{1}{2} \left\{ \begin{array}{l} tr(\Psi \Sigma^{-1}) + (n-1) s^2 \\ + \frac{\kappa_0 n}{\kappa_0 + n} (\theta_0 - \bar{y})^T \Sigma^{-1} (\theta_0 - \bar{y}) \end{array} \right\} \right] \\ \times |\Sigma|^{-\frac{1}{2}} \exp \left[ -\frac{1}{2} \left\{ (\kappa_0 + n) (\theta - \theta_n)^T \Sigma^{-1} (\theta - \theta_n) \right\} \right]$$

### 6.6.2 Inverted-Wishart kernel

We wish to identify the exponent with Gaussian by inverted-Wishart kernels where the inverted-Wishart involves the trace of a square, symmetric matrix, call it  $\Psi_n$ , multiplied by  $\Sigma^{-1}$ .

To make this connection we utilize the following general results. Since a quadratic form, say  $x^T \Sigma^{-1} x$ , is a scalar, it's equal to its trace,

$$x^T \Sigma^{-1} x = tr(x^T \Sigma^{-1} x)$$



Further, for conformable matrices  $A, B$  and  $C, D$ ,

$$\text{tr}(A) + \text{tr}(B) = \text{tr}(A + B)$$

and

$$\text{tr}(CD) = \text{tr}(DC)$$

We immediately have the results

$$\text{tr}(x^T x) = \text{tr}(xx^T)$$

and

$$\text{tr}(x^T \Sigma^{-1} x) = \text{tr}(\Sigma^{-1} x x^T) = \text{tr}(x x^T \Sigma^{-1})$$

Therefore, the above joint posterior can be rewritten as a  $\text{N}(\theta; \theta_n, (\kappa_0 + n)^{-1} \Sigma) \times$  inverted-Wishart( $\Sigma^{-1}; \nu + n, \Psi_n$ )

$$\begin{aligned} p(\theta, \Sigma | y) &\propto |\Psi_n|^{\frac{\nu+n}{2}} |\Sigma|^{-\frac{\nu+n+k+1}{2}} \exp\left[-\frac{1}{2} \text{tr}(\Psi_n \Sigma^{-1})\right] \\ &\times |\Sigma|^{-\frac{1}{2}} \exp\left[-\frac{\kappa_0 + n}{2} (\theta - \theta_n)^T \Sigma^{-1} (\theta - \theta_n)\right] \end{aligned}$$

where

$$\theta_n = \frac{1}{\kappa_0 + n} (\kappa_0 \theta_0 + n \bar{y})$$

and

$$\Psi_n = \Psi + \sum_{i=1}^n (y_i - \bar{y})(y_i - \bar{y})^T + \frac{\kappa_0 n}{\kappa_0 + n} (\bar{y} - \theta_0)(\bar{y} - \theta_0)^T$$

Now, it's apparent the conditional posterior for  $\theta$  given  $\Sigma$  is  $\text{N}(\theta_n, (\kappa_0 + n)^{-1} \Sigma)$

$$p(\theta | \Sigma, y) \propto \exp\left[-\frac{\kappa_0 + n}{2} (\theta - \theta_n)^T \Sigma^{-1} (\theta - \theta_n)\right]$$

### 6.6.3 Marginal posterior distributions

Integrating out the other parameter gives the marginal posteriors, a multivariate Student  $t$  for the mean,

$$\text{Student } t_k(\theta; \theta_n, \Gamma, \nu + n - k + 1)$$

and an inverted-Wishart for the variance,

$$\text{I-W}(\Sigma^{-1}; \nu + n, \Psi_n)$$

where

$$\Gamma = (\kappa_0 + n)^{-1} (\nu + n - k + 1)^{-1} \Psi_n$$

Marginalization of the mean derives from the following identities (see Box and Tiao [1973], p. 427, 441). Let  $Z$  be a  $m \times m$  positive definite symmetric matrix consisting of  $\frac{1}{2}m(m+1)$  distinct random variables  $z_{ij}$  ( $i, j = 1, \dots, m; i \geq j$ ). And let  $q > 0$  and  $B$  be an  $m \times m$  positive definite symmetric matrix. Then, the distribution of  $z_{ij}$ ,

$$p(Z) \propto |Z|^{\frac{1}{2}q-1} \exp\left[-\frac{1}{2}\text{tr}(ZB)\right], \quad Z > 0$$

is a multivariate generalization of the  $\chi^2$  distribution obtained by Wishart [1928]. Integrating out the distinct  $z_{ij}$  produces the first identity.

$$\int_{Z>0} |Z|^{\frac{1}{2}q-1} \exp\left[-\frac{1}{2}\text{tr}(ZB)\right] dZ = |B|^{-\frac{1}{2}(q+m-1)} \times 2^{\frac{1}{2}(q+m-1)} \Gamma_m\left(\frac{q+m-1}{2}\right) \quad (\text{I.1})$$

where  $\Gamma_p(b)$  is the generalized gamma function (Siegel [1935])

$$\Gamma_p(b) = \left[\Gamma\left(\frac{1}{2}\right)\right]^{\frac{1}{2}p(p-1)} \prod_{\alpha=1}^p \Gamma\left(b + \frac{\alpha-p}{2}\right), \quad b > \frac{p-1}{2}$$

and

$$\Gamma(z) = \int_0^\infty t^{z-1} e^{-t} dt$$

or for integer  $n$ ,

$$\Gamma(n) = (n-1)!$$

The second identity involves the relationship between determinants that allows us to express the above as a quadratic form. The identity is

$$|I_k - PQ| = |I_l - QP| \quad (\text{I.2})$$

for  $P$  a  $k \times l$  matrix and  $Q$  a  $l \times k$  matrix.

If we transform the joint posterior to  $p(\theta, \Sigma^{-1} | y)$ , the above identities can be applied to marginalize the joint posterior. The key to transformation is

$$p(\theta, \Sigma^{-1} | y) = p(\theta, \Sigma | y) \left| \frac{\partial \Sigma}{\partial \Sigma^{-1}} \right|$$

where  $\left| \frac{\partial \Sigma}{\partial \Sigma^{-1}} \right|$  is the (absolute value of the) determinant of the Jacobian or

$$\begin{aligned} \left| \frac{\partial \Sigma}{\partial \Sigma^{-1}} \right| &= \left| \frac{\partial(\sigma_{11}, \sigma_{12}, \dots, \sigma_{kk})}{\partial(\sigma^{11}, \sigma^{12}, \dots, \sigma^{kk})} \right| \\ &= |\Sigma|^{k+1} \end{aligned}$$

with  $\sigma_{ij}$  the elements of  $\Sigma$  and  $\sigma^{ij}$  the elements of  $\Sigma^{-1}$ . Hence,

$$\begin{aligned} p(\theta, \Sigma | y) &\propto |\Sigma|^{-\frac{\nu+n+k+1}{2}} \exp \left[ -\frac{1}{2} \text{tr} (\Psi_n \Sigma^{-1}) \right] \\ &\quad \times |\Sigma|^{-\frac{1}{2}} \exp \left[ -\frac{\kappa_0 + n}{2} (\theta - \theta_n)^T \Sigma^{-1} (\theta - \theta_n) \right] \\ &\propto |\Sigma|^{-\frac{\nu+n+k}{2}-1} \exp \left[ -\frac{1}{2} \text{tr} (S(\theta) \Sigma^{-1}) \right] \end{aligned}$$

where  $S(\theta) = \Psi_n + (\kappa_0 + n) (\theta - \theta_n) (\theta - \theta_n)^T$ , can be rewritten

$$\begin{aligned} p(\theta, \Sigma^{-1} | y) &\propto |\Sigma|^{-\frac{\nu+n+k+2}{2}} \exp \left[ -\frac{1}{2} \text{tr} (S(\theta) \Sigma^{-1}) \right] |\Sigma|^{\frac{2k+2}{2}} \\ &\propto |\Sigma^{-1}|^{\frac{\nu+n-k}{2}} \exp \left[ -\frac{1}{2} \text{tr} (S(\theta) \Sigma^{-1}) \right] \end{aligned}$$

Now, applying the first identity yields

$$\begin{aligned} \int_{\Sigma^{-1} > 0} p(\theta, \Sigma^{-1} | y) d\Sigma^{-1} &\propto |S(\theta)|^{-\frac{1}{2}(\nu+n+1)} \\ &\propto \left| \Psi_n + (\kappa_0 + n) (\theta - \theta_n) (\theta - \theta_n)^T \right|^{-\frac{1}{2}(\nu+n+1)} \\ &\propto \left| I + (\kappa_0 + n) \Psi_n^{-1} (\theta - \theta_n) (\theta - \theta_n)^T \right|^{-\frac{1}{2}(\nu+n+1)} \end{aligned}$$

And the second identity gives

$$p(\theta | y) \propto \left[ 1 + (\kappa_0 + n) (\theta - \theta_n)^T \Psi_n^{-1} (\theta - \theta_n) \right]^{-\frac{1}{2}(\nu+n+1)}$$

We recognize this is the kernel of a multivariate Student  $t_k(\theta; \theta_n, \Gamma, \nu + n - k + 1)$  distribution.

#### 6.6.4 Uninformative priors

The joint uninformative prior (with a locally uniform prior for  $\theta$ ) is

$$p(\theta, \Sigma) \propto |\Sigma|^{-\frac{k+1}{2}}$$

and the joint posterior is

$$\begin{aligned} p(\theta, \Sigma | y) &\propto |\Sigma|^{-\frac{k+1}{2}} |\Sigma|^{-\frac{n}{2}} \exp \left[ -\frac{1}{2} \left\{ (n-1) s^2 + n (\bar{y} - \theta)^T \Sigma^{-1} (\bar{y} - \theta) \right\} \right] \\ &\propto |\Sigma|^{-\frac{n+k+1}{2}} \exp \left[ -\frac{1}{2} \left\{ (n-1) s^2 + n (\bar{y} - \theta)^T \Sigma^{-1} (\bar{y} - \theta) \right\} \right] \\ &\propto |\Sigma|^{-\frac{n+k+1}{2}} \exp \left[ -\frac{1}{2} \text{tr} (S(\theta) \Sigma^{-1}) \right] \end{aligned}$$

where now  $S(\theta) = \sum_{i=1}^n (\bar{y} - y_i)(\bar{y} - y_i)^T + n(\bar{y} - \theta)(\bar{y} - \theta)^T$ . Then, the conditional posterior for  $\theta$  given  $\Sigma$  is  $N(\bar{y}, n^{-1}\Sigma)$

$$p(\theta | \Sigma, y) \propto \exp \left[ -\frac{n}{2} (\theta - \bar{y})^T \Sigma^{-1} (\theta - \bar{y}) \right]$$

The marginal posterior for  $\theta$  is derived analogous to the above informed conjugate prior case. Rewriting the posterior in terms of  $\Sigma^{-1}$  yields

$$\begin{aligned} p(\theta, \Sigma^{-1} | y) &\propto |\Sigma|^{-\frac{n+k+1}{2}} \exp \left[ -\frac{1}{2} \text{tr} (S(\theta) \Sigma^{-1}) \right] |\Sigma|^{\frac{2k+2}{2}} \\ &\propto |\Sigma^{-1}|^{\frac{n-k-1}{2}} \exp \left[ -\frac{1}{2} \text{tr} (S(\theta) \Sigma^{-1}) \right] \end{aligned}$$

$$\begin{aligned} p(\theta | y) &\propto \int_{\Sigma^{-1} > 0} p(\theta, \Sigma^{-1} | y) d\Sigma^{-1} \\ &\propto \int_{\Sigma^{-1} > 0} |\Sigma^{-1}|^{\frac{n-k-1}{2}} \exp \left[ -\frac{1}{2} \text{tr} (S(\theta) \Sigma^{-1}) \right] d\Sigma^{-1} \end{aligned}$$

The first identity (I.1) produces

$$\begin{aligned} p(\theta | y) &\propto |S(\theta)|^{-\frac{n}{2}} \\ &\propto \left| \sum_{i=1}^n (\bar{y} - y_i)(\bar{y} - y_i)^T + n(\bar{y} - \theta)(\bar{y} - \theta)^T \right|^{\frac{n}{2}} \\ &\propto \left| I + n \left[ \sum_{i=1}^n (\bar{y} - y_i)(\bar{y} - y_i)^T \right]^{-1} (\bar{y} - \theta)(\bar{y} - \theta)^T \right|^{\frac{n}{2}} \end{aligned}$$

The second identity (I.2) identifies the marginal posterior for  $\theta$  as (multivariate) Student  $t_k(\theta; \bar{y}, \frac{1}{n}s^2, n-k)$

$$p(\theta | y) \propto \left[ 1 + \frac{n}{(n-k)s^2} (\bar{y} - \theta)^T (\bar{y} - \theta)^T \right]^{-\frac{n}{2}}$$

where  $(n-k)s^2 = \sum_{i=1}^n (\bar{y} - y_i)^T (\bar{y} - y_i)$ . The marginal posterior for the variance is I-W( $\Sigma^{-1}; n, \Psi_n$ ) where now  $\Psi_n = \sum_{i=1}^n (\bar{y} - y_i)(\bar{y} - y_i)^T$ .

## 6.7 Bayesian linear regression

Linear regression is the starting point for more general data modeling strategies, including nonlinear models. Hence, Bayesian linear regression is foundational. Suppose the data are generated by

$$y = X\beta + \varepsilon$$

where  $X$  is a  $n \times p$  full column rank matrix of (weakly exogenous) regressors and  $\varepsilon \sim N(0, \sigma^2 I_n)$  and  $E[\varepsilon | X] = 0$ . Then, the sample conditional density is  $(y | X, \beta, \sigma^2) \sim N(X\beta, \sigma^2 I_n)$ .

### 6.7.1 Known variance

If the error variance,  $\sigma^2 I_n$ , is known and we have informed Gaussian priors for  $\beta$  conditional on  $\sigma^2$ ,

$$p(\beta | \sigma^2) \sim N(\beta_0, \sigma^2 V_0)$$

where we can think of  $V_0 = (X_0^T X_0)^{-1}$  as if we had a prior sample  $(y_0, X_0)$  such that

$$\beta_0 = (X_0^T X_0)^{-1} X_0^T y_0$$

then the conditional posterior for  $\beta$  is

$$p(\beta | \sigma^2, y, X; \beta_0, V_0) \sim N(\bar{\beta}, V_\beta)$$

where

$$\begin{aligned} \bar{\beta} &= (X_0^T X_0 + X^T X)^{-1} (X_0^T X_0 \beta_0 + X^T X \hat{\beta}) \\ \hat{\beta} &= (X^T X)^{-1} X^T y \end{aligned}$$

and

$$V_\beta = \sigma^2 (X_0^T X_0 + X^T X)^{-1}$$

The variance expression follows from rewriting the estimator

$$\begin{aligned} \bar{\beta} &= (X_0^T X_0 + X^T X)^{-1} (X_0^T X_0 \beta_0 + X^T X \hat{\beta}) \\ &= (X_0^T X_0 + X^T X)^{-1} (X_0^T X_0 (X_0^T X_0)^{-1} X_0^T y_0 + X^T X (X^T X)^{-1} X^T y) \\ &= (X_0^T X_0 + X^T X)^{-1} (X_0^T y_0 + X^T y) \end{aligned}$$

Since the DGP is

$$\begin{aligned} y_0 &= X_0 \beta + \varepsilon_0, & \varepsilon_0 &\sim N(0, \sigma^2 I_{n_0}) \\ y &= X \beta + \varepsilon, & \varepsilon &\sim N(0, \sigma^2 I_n) \end{aligned}$$

then

$$\bar{\beta} = (X_0^T X_0 + X^T X)^{-1} (X_0^T X_0 \beta + X_0^T \varepsilon_0 + X^T X \beta + X^T \varepsilon)$$

The conditional (and by iterated expectations, unconditional) expected value of the estimator is

$$E[\bar{\beta} | X, X_0] = (X_0^T X_0 + X^T X)^{-1} (X_0^T X_0 + X^T X) \beta = \beta$$

Hence,

$$\begin{aligned}\bar{\beta} - E[\bar{\beta} | X, X_0] &= \bar{\beta} - \beta \\ &= (X_0^T X_0 + X^T X)^{-1} (X_0^T \varepsilon_0 + X^T \varepsilon)\end{aligned}$$

so that

$$\begin{aligned}V_\beta &\equiv \text{Var}[\bar{\beta} | X, X_0] \\ &= E\left[(\bar{\beta} - \beta)(\bar{\beta} - \beta)^T | X, X_0\right] \\ &= E\left[\begin{aligned} &(X_0^T X_0 + X^T X)^{-1} (X_0^T \varepsilon_0 + X^T \varepsilon) (X_0^T \varepsilon_0 + X^T \varepsilon)^T \\ &\times (X_0^T X_0 + X^T X)^{-1} | X, X_0 \end{aligned}\right] \\ &= E\left[\begin{aligned} &(X_0^T X_0 + X^T X)^{-1} \begin{pmatrix} X_0^T \varepsilon_0 \varepsilon_0^T X_0 + X^T \varepsilon \varepsilon^T X_0 \\ + X_0^T \varepsilon_0 \varepsilon^T X + X^T \varepsilon^T \varepsilon^T X \end{pmatrix} \\ &\times (X_0^T X_0 + X^T X)^{-1} | X, X_0 \end{aligned}\right] \\ &= (X_0^T X_0 + X^T X)^{-1} (X_0^T \sigma^2 I X_0 + X^T \sigma^2 I X) (X_0^T X_0 + X^T X)^{-1} \\ &= \sigma^2 (X_0^T X_0 + X^T X)^{-1} (X_0^T X_0 + X^T X) (X_0^T X_0 + X^T X)^{-1} \\ &= \sigma^2 (X_0^T X_0 + X^T X)^{-1}\end{aligned}$$

Now, let's backtrack and derive the conditional posterior as the product of conditional priors and the likelihood function. The likelihood function for known variance is

$$\ell(\beta | \sigma^2, y, X) \propto \exp\left[-\frac{1}{2\sigma^2} (y - X\beta)^T (y - X\beta)\right]$$

Conditional Gaussian priors are

$$p(\beta | \sigma^2) \propto \exp\left[-\frac{1}{2\sigma^2} (\beta - \beta_0)^T V_0^{-1} (\beta - \beta_0)\right]$$

The conditional posterior is the product of the prior and likelihood

$$\begin{aligned}p(\beta | \sigma^2, y, X) &\propto \exp\left[-\frac{1}{2\sigma^2} \left\{ \begin{aligned} &(y - X\beta)^T (y - X\beta) \\ &+ (\beta - \beta_0)^T V_0^{-1} (\beta - \beta_0) \end{aligned} \right\}\right] \\ &= \exp\left[-\frac{1}{2\sigma^2} \left\{ \begin{aligned} &y^T y - 2y^T X\beta + \beta^T X^T X\beta \\ &+ \beta^T X_0^T X_0 \beta - 2\beta_0^T X_0^T X_0 \beta \\ &+ \beta_0^T X_0^T X_0 \beta_0 \end{aligned} \right\}\right]\end{aligned}$$

The first and last terms in the exponent do not involve  $\beta$  (are constants) and can be ignored as they are absorbed through normalization. This leaves

$$\begin{aligned}p(\beta | \sigma^2, y, X) &\propto \exp\left[-\frac{1}{2\sigma^2} \left\{ \begin{aligned} &-2y^T X\beta + \beta^T X^T X\beta + \beta^T X_0^T X_0 \beta \\ &-2\beta_0^T X_0^T X_0 \beta \end{aligned} \right\}\right] \\ &= \exp\left[-\frac{1}{2\sigma^2} \left\{ \begin{aligned} &\beta^T (X_0^T X_0 + X^T X) \beta \\ &-2(y^T X + \beta_0^T X_0^T X_0) \beta \end{aligned} \right\}\right]\end{aligned}$$

which can be recognized as the expansion of the conditional posterior claimed above.

$$\begin{aligned}
p(\beta \mid \sigma^2, y, X) &\sim N(\bar{\beta}, V_\beta) \\
&\propto \exp\left[-\frac{1}{2}(\beta - \bar{\beta})^T V_\beta^{-1}(\beta - \bar{\beta})\right] \\
&= \exp\left[-\frac{1}{2\sigma^2}(\beta - \bar{\beta})^T (X_0^T X_0 + X^T X)(\beta - \bar{\beta})\right] \\
&= \exp\left[-\frac{1}{2\sigma^2} \begin{pmatrix} \beta^T (X_0^T X_0 + X^T X) \beta \\ -2\bar{\beta}^T (X_0^T X_0 + X^T X) \beta \\ +\bar{\beta}^T (X_0^T X_0 + X^T X) \bar{\beta} \end{pmatrix}\right] \\
&= \exp\left[-\frac{1}{2\sigma^2} \begin{pmatrix} \beta^T (X_0^T X_0 + X^T X) \beta \\ -2(X_0^T X_0 \beta_0 + X^T y)^T \beta \\ +\bar{\beta}^T (X_0^T X_0 + X^T X) \bar{\beta} \end{pmatrix}\right]
\end{aligned}$$

The last term in the exponent is all constants (does not involve  $\beta$ ) so its absorbed through normalization and disregarded for comparison of kernels. Hence,

$$\begin{aligned}
p(\beta \mid \sigma^2, y, X) &\propto \exp\left[-\frac{1}{2}(\beta - \bar{\beta})^T V_\beta^{-1}(\beta - \bar{\beta})\right] \\
&\propto \exp\left[-\frac{1}{2\sigma^2} \begin{pmatrix} \beta^T (X_0^T X_0 + X^T X) \beta \\ -2(y^T X + \beta_0^T X_0^T X_0) \beta \end{pmatrix}\right]
\end{aligned}$$

as claimed.

#### *Uninformative priors*

If the prior for  $\beta$  is uniformly distributed conditional on known variance,  $\sigma^2$ ,  $p(\beta \mid \sigma^2) \propto 1$ , then it's as if  $X_0^T X_0 \rightarrow 0$  (the information matrix for the prior is null) and the posterior for  $\beta$  is

$$p(\beta \mid \sigma^2, y, X) \sim N(\hat{\beta}, \sigma^2 (X^T X)^{-1})$$

equivalent to the classical parameter estimators.

To see this intuition holds, recognize combining the likelihood with the uninformative prior indicates the posterior is proportional to the likelihood.

$$p(\beta \mid \sigma^2, y, X) \propto \exp\left[-\frac{1}{2\sigma^2}(y - X\beta)^T (y - X\beta)\right]$$

Expanding this expression yields

$$p(\beta \mid \sigma^2, y, X) \propto \exp\left[-\frac{1}{2\sigma^2}(y^T y - 2y^T X\beta + \beta^T X^T X\beta)\right]$$

The first term in the exponent doesn't depend on  $\beta$  and can be dropped as it's absorbed via normalization. This leaves

$$p(\beta \mid \sigma^2, y, X) \propto \exp \left[ -\frac{1}{2\sigma^2} \left( -2y^T X\beta + \beta^T X^T X\beta \right) \right]$$

Now, write  $p(\beta \mid \sigma^2, y, X) \sim N(\hat{\beta}, \sigma^2 (X^T X)^{-1})$

$$p(\beta \mid \sigma^2, y, X) \propto \exp \left[ -\frac{1}{2\sigma^2} \left( \beta - \hat{\beta} \right)^T X^T X \left( \beta - \hat{\beta} \right) \right]$$

and expand

$$p(\beta \mid \sigma^2, y, X) \propto \exp \left[ -\frac{1}{2\sigma^2} \left( \beta^T X^T X\beta - 2\beta^T X^T X\hat{\beta} + \hat{\beta}^T X^T X\hat{\beta} \right) \right]$$

The last term in the exponent doesn't depend on  $\beta$  and is absorbed via normalization. This leaves

$$\begin{aligned} p(\beta \mid \sigma^2, y, X) &\propto \exp \left[ -\frac{1}{2\sigma^2} \left( \beta^T X^T X\beta - 2\beta^T X^T X\hat{\beta} \right) \right] \\ &\propto \exp \left[ -\frac{1}{2\sigma^2} \left( \beta^T X^T X\beta - 2\beta^T X^T X (X^T X)^{-1} X^T y \right) \right] \\ &\propto \exp \left[ -\frac{1}{2\sigma^2} \left( \beta^T X^T X\beta - 2\beta^T X^T y \right) \right] \end{aligned}$$

As this latter expression matches the simplified likelihood expression, the demonstration is complete,  $p(\beta \mid \sigma^2, y, X) \sim N(\hat{\beta}, \sigma^2 (X^T X)^{-1})$ .

### 6.7.2 Unknown variance

In the usual case where the variance as well as the regression coefficients,  $\beta$ , are unknown, the likelihood function can be expressed as

$$\ell(\beta, \sigma^2 \mid y, X) \propto \sigma^{-n} \exp \left[ -\frac{1}{2\sigma^2} (y - X\beta)^T (y - X\beta) \right]$$

Rewriting gives

$$\ell(\beta, \sigma^2 \mid y, X) \propto \sigma^{-n} \exp \left[ -\frac{1}{2\sigma^2} \varepsilon^T \varepsilon \right]$$

since  $\varepsilon = y - X\beta$ . The estimated model is  $y = Xb + e$ , therefore  $X\beta + \varepsilon = Xb + e$  where  $b = (X^T X)^{-1} X^T y$  and  $e = y - Xb$  are estimates of  $\beta$  and  $\varepsilon$ , respectively. This implies  $\varepsilon = e - X(\beta - b)$  and

$$\ell(\beta, \sigma^2 \mid y, X) \propto \sigma^{-n} \exp \left[ -\frac{1}{2\sigma^2} \left\{ \begin{array}{l} e^T e - 2(\beta - b)^T X^T e \\ + (\beta - b)^T X^T X (\beta - b) \end{array} \right\} \right]$$



Since,  $X^T e = 0$  by construction, this simplifies as

$$\ell(\beta, \sigma^2 | y, X) \propto \sigma^{-n} \exp \left[ -\frac{1}{2\sigma^2} \left\{ e^T e + (\beta - b)^T X^T X (\beta - b) \right\} \right]$$

or

$$\ell(\beta, \sigma^2 | y, X) \propto \sigma^{-n} \exp \left[ -\frac{1}{2\sigma^2} \left\{ (n-p) s^2 + (\beta - b)^T X^T X (\beta - b) \right\} \right]$$

where  $s^2 = \frac{1}{n-p} e^T e$ .<sup>7</sup>

The conjugate prior for linear regression is the Gaussian( $\beta | \sigma^2; \beta_0, \sigma^2 \Omega_0^{-1}$ )-inverse chi square( $\sigma^2; \nu_0, \sigma_0^2$ )

$$\begin{aligned} p(\beta | \sigma^2; \beta_0, \sigma^2 \Omega_0^{-1}) \times p(\sigma^2; \nu_0, \sigma_0^2) &\propto \sigma^{-p} \exp \left[ -\frac{(\beta - \beta_0)^T \Omega_0 (\beta - \beta_0)}{2\sigma^2} \right] \\ &\times \sigma^{-(\nu_0/2+1)} \exp \left[ -\frac{\nu_0 \sigma_0^2}{2\sigma^2} \right] \end{aligned}$$

Combining the prior with the likelihood gives a joint Gaussian( $\bar{\beta}, \sigma^2 \Omega_n^{-1}$ )-inverse chi square( $\nu_0 + n, \sigma_n^2$ ) posterior

$$\begin{aligned} p(\beta, \sigma^2 | y, X; \beta_0, \sigma^2 \Omega_0^{-1}, \nu_0, \sigma_0^2) &\propto \sigma^{-n} \exp \left[ -\frac{(n-p) s^2}{2\sigma^2} \right] \\ &\times \exp \left[ -\frac{(\beta - b)^T X^T X (\beta - b)}{2\sigma^2} \right] \\ &\times \sigma^{-p} \exp \left[ -\frac{(\beta - \beta_0)^T \Omega_0 (\beta - \beta_0)}{2\sigma^2} \right] \\ &\times (\sigma^2)^{-(\nu_0/2+1)} \exp \left[ -\frac{\nu_0 \sigma_0^2}{2\sigma^2} \right] \end{aligned}$$

---

<sup>7</sup>Notice, the univariate Gaussian case is subsumed by linear regression where  $X = \iota$  (a vector of ones). Then, the likelihood as described earlier,

$$\ell(\beta, \sigma^2 | y, X) \propto \sigma^{-n} \exp \left[ -\frac{1}{2\sigma^2} \left\{ (n-p) s^2 + (\beta - b)^T X^T X (\beta - b) \right\} \right]$$

becomes

$$\ell(\beta = \theta, \sigma^2 | y, X = \iota) \propto \sigma^{-n} \exp \left[ -\frac{1}{2\sigma^2} \left\{ (n-1) s^2 + n(\theta - \bar{y})^2 \right\} \right]$$

where  $\theta = \beta$ ,  $b = (X^T X)^{-1} X^T y = \bar{y}$ ,  $p = 1$ , and  $X^T X = n$ .

Collecting terms and rewriting, we have

$$p(\beta, \sigma^2 \mid y, X; \beta_0, \sigma^2 \Omega_0^{-1}, \nu_0, \sigma_0^2) \propto (\sigma^2)^{-[(\nu_0+n)/2+1]} \exp\left[-\frac{\sigma_n^2}{2\sigma^2}\right] \\ \times \sigma^{-p} \exp\left[-\frac{1}{2\sigma^2} (\beta - \bar{\beta})^T \Omega_n (\beta - \bar{\beta})\right]$$

where

$$\bar{\beta} = (\Omega_0 + X^T X)^{-1} (\Omega_0 \beta_0 + X^T X b) \\ \Omega_n = (\Omega_0 + X^T X)$$

and

$$\nu_n \sigma_n^2 = \nu_0 \sigma_0^2 + (n-p) s^2 + (\beta_0 - \bar{\beta})^T \Omega_0 (\beta_0 - \bar{\beta}) + (\hat{\beta} - \bar{\beta})^T X^T X (\hat{\beta} - \bar{\beta})$$

where  $\nu_n = \nu_0 + n$ . The conditional posterior of  $\beta$  given  $\sigma^2$  is Gaussian( $\bar{\beta}, \sigma^2 \Omega_n^{-1}$ ).

*Completing the square*

The derivation of the above joint posterior follows from the matrix version of completing the square where  $\Omega_0$  and  $X^T X$  are square, symmetric, full rank  $p \times p$  matrices. The exponents from the prior for the mean and likelihood are

$$(\beta - \beta_0)^T \Omega_0 (\beta - \beta_0) + (\beta - \hat{\beta})^T X^T X (\beta - \hat{\beta})$$

Expanding and rearranging gives

$$\beta^T (\Omega_0 + X^T X) \beta - 2 (\Omega_0 \beta_0 + X^T X \hat{\beta})^T \beta + \beta_0^T \Omega_0 \beta_0 + \hat{\beta}^T X^T X \hat{\beta} \quad (6.1)$$

The latter two terms are constants not involving  $\beta$  (and can be ignored when writing the kernel for the conditional posterior) which we'll add to when we complete the square. Now, write out the square centered around  $\bar{\beta}$

$$(\beta - \bar{\beta})^T (\Omega_0 + X^T X) (\beta - \bar{\beta}) = \beta^T (\Omega_0 + X^T X) \beta \\ - 2\bar{\beta}^T (\Omega_0 + X^T X) \beta + \bar{\beta}^T (\Omega_0 + X^T X) \bar{\beta}$$

Substitute for  $\bar{\beta}$  in the second term on the right hand side and the first two terms are identical to the two terms in equation (6.1). Hence, the exponents from the prior for the mean and likelihood in (6.1) are equal to

$$(\beta - \bar{\beta})^T (\Omega_0 + X^T X) (\beta - \bar{\beta}) \\ - \bar{\beta}^T (\Omega_0 + X^T X) \bar{\beta} + \beta_0^T \Omega_0 \beta_0 + \hat{\beta}^T X^T X \hat{\beta}$$

which can be rewritten as

$$\begin{aligned} & (\beta - \bar{\beta})^T (\Omega_0 + X^T X) (\beta - \bar{\beta}) \\ & + (\beta_0 - \bar{\beta})^T \Omega_0 (\beta_0 - \bar{\beta}) + (\hat{\beta} - \bar{\beta})^T X^T X (\hat{\beta} - \bar{\beta}) \end{aligned}$$

or (in the form analogous to the univariate Gaussian case)

$$\begin{aligned} & (\beta - \bar{\beta})^T (\Omega_0 + X^T X) (\beta - \bar{\beta}) \\ & + (\beta_0 - \hat{\beta})^T (\Omega_1 \Omega_n^{-1} \Omega_0 \Omega_n^{-1} \Omega_1 + \Omega_0 \Omega_n^{-1} \Omega_1 \Omega_n^{-1} \Omega_0) (\beta_0 - \hat{\beta}) \end{aligned}$$

where  $\Omega_1 = X^T X$ .

### Stacked regression

Bayesian linear regression with conjugate priors works as if we have a prior sample  $\{X_0, y_0\}$ ,  $\Omega_0 = X_0^T X_0$ , and initial estimates

$$\beta_0 = (X_0^T X_0)^{-1} X_0^T y_0$$

Then, we combine this initial "evidence" with new evidence to update our beliefs in the form of the posterior. Not surprisingly, the posterior mean is a weighted average of the two "samples" where the weights are based on the relative precision of the two "samples".

### Marginal posterior distributions

The marginal posterior for  $\beta$  on integrating out  $\sigma^2$  is noncentral, scaled multivariate Student  $t_p(\bar{\beta}, \sigma_n^2 \Omega_n^{-1}, \nu_0 + n)$

$$\begin{aligned} p(\beta | y, X) & \propto \left[ \nu_n \sigma_n^2 + (\beta - \bar{\beta})^T \Omega_n (\beta - \bar{\beta}) \right]^{-\frac{\nu_0 + n + p}{2}} \\ & \propto \left[ 1 + \frac{1}{\nu_n \sigma_n^2} (\beta - \bar{\beta})^T \Omega_n (\beta - \bar{\beta}) \right]^{-\frac{\nu_0 + n + p}{2}} \end{aligned}$$

where  $\Omega_n = \Omega_0 + X^T X$ . This result corresponds with the univariate Gaussian case and is derived analogously by transformation of variables where  $z = \frac{A}{2\sigma^2}$  where  $A = \sigma_n^2 + (\beta - \bar{\beta})^T \Omega_n (\beta - \bar{\beta})$ . The marginal posterior for  $\sigma^2$  is inverted-chi square( $\sigma^2; \nu_n, \sigma_n^2$ ).

Derivation of the marginal posterior for  $\beta$  is as follows.

$$\begin{aligned} p(\beta | y) & = \int_0^\infty p(\beta, \sigma^2 | y) d\sigma^2 \\ & = \int_0^\infty (\sigma^2)^{-\frac{n + \nu_0 + p + 2}{2}} \exp\left[-\frac{A}{2\sigma^2}\right] d\sigma^2 \end{aligned}$$

Utilizing  $\sigma^2 = \frac{A}{2z}$  and  $dz = -\frac{2z^2}{A} d\sigma^2$  or  $d\sigma^2 = -\frac{A}{2z^2} dz$ , ( $-1$  and  $2$  are constants and can be ignored when deriving the kernel)

$$\begin{aligned} p(\beta | y) &\propto \int_0^\infty \left(\frac{A}{2z}\right)^{-\frac{n+\nu_0+p+2}{2}} \frac{A}{2z^2} \exp[-z] dz \\ &\propto A^{-\frac{n+\nu_0+p}{2}} \int_0^\infty z^{\frac{n+\nu_0+p}{2}-1} \exp[-z] dz \end{aligned}$$

The integral  $\int_0^\infty z^{\frac{n+\nu_0+k}{2}-1} \exp[-z] dz$  is a constant since it is the kernel of a gamma density and therefore can be ignored when deriving the kernel of the marginal posterior for beta

$$\begin{aligned} p(\beta | y) &\propto A^{-\frac{n+\nu_0+p}{2}} \\ &\propto \left[ \nu_n \sigma_n^2 + (\beta - \bar{\beta})^T \Omega_n (\beta - \bar{\beta}) \right]^{-\frac{n+\nu_0+p}{2}} \\ &\propto \left[ 1 + \frac{(\beta - \bar{\beta})^T \Omega_n (\beta - \bar{\beta})}{\nu_n \sigma_n^2} \right]^{-\frac{n+\nu_0+p}{2}} \end{aligned}$$

the kernel for a noncentral, scaled (multivariate) Student  $t_p(\beta; \bar{\beta}, \sigma_n^2 \Omega_n^{-1}, n + \nu_0)$ .

### 6.7.3 Uninformative priors

Again, the case of uninformative priors is relatively straightforward. Since priors convey no information, the prior for the mean is uniform (proportional to a constant,  $\Omega_0 \rightarrow 0$ ) and the prior for  $\sigma^2$  has  $\nu_0 \rightarrow 0$  degrees of freedom so that the joint prior is  $p(\beta, \sigma^2) \propto (\sigma^2)^{-1}$ .

The joint posterior is

$$p(\beta, \sigma^2 | y) \propto (\sigma^2)^{-[n/2+1]} \exp\left[-\frac{1}{2\sigma^2} (y - X\beta)^T (y - X\beta)\right]$$

Since  $y = Xb + e$  where  $b = (X^T X)^{-1} X^T y$ , the joint posterior can be written

$$p(\beta, \sigma^2 | y) \propto (\sigma^2)^{-[n/2+1]} \exp\left[-\frac{1}{2\sigma^2} \left\{ (n-p)s^2 + (\beta - b)^T X^T X (\beta - b) \right\}\right]$$

Or, factoring into the conditional posterior for  $\beta$  and marginal for  $\sigma^2$ , we have

$$\begin{aligned} p(\beta, \sigma^2 | y) &\propto p(\sigma^2 | y) p(\beta | \sigma^2, y) \\ &\propto (\sigma^2)^{-[(n-p)/2+1]} \exp\left[-\frac{\sigma_n^2}{2\sigma^2}\right] \\ &\quad \times \sigma^{-p} \exp\left[-\frac{1}{2\sigma^2} (\beta - b)^T X^T X (\beta - b)\right] \end{aligned}$$

where

$$\sigma_n^2 = (n - p) s^2$$

Hence, the conditional posterior for  $\beta$  given  $\sigma^2$  is Gaussian  $\left(b, \sigma^2 (X^T X)^{-1}\right)$ .

The marginal posterior for  $\beta$  is multivariate Student  $t_p \left(\beta; b, s^2 (X^T X)^{-1}, n - p\right)$ , the classical estimator. Derivation of the marginal posterior for  $\beta$  is analogous to that above. Let  $z = \frac{A}{2\sigma^2}$  where  $A = (n - p) s^2 + (\beta - b)^T X^T X (\beta - b)$ . Integrating  $\sigma^2$  out of the joint posterior produces the marginal posterior for  $\beta$ .

$$\begin{aligned} p(\beta | y) &\propto \int p(\beta, \sigma^2 | y) d\sigma^2 \\ &\propto \int (\sigma^2)^{-\frac{n+2}{2}} \exp\left[-\frac{A}{2\sigma^2}\right] d\sigma^2 \end{aligned}$$

Substitution yields

$$\begin{aligned} p(\beta | y) &\propto \int \left(\frac{A}{2z}\right)^{-\frac{n+2}{2}} \frac{A}{2z^2} \exp[-z] dz \\ &\propto A^{-\frac{n}{2}} \int z^{\frac{n}{2}-1} \exp[-z] dz \end{aligned}$$

As before, the integral involves the kernel of a gamma distribution, a constant which can be ignored. Therefore, we have

$$\begin{aligned} p(\beta | y) &\propto A^{-\frac{n}{2}} \\ &\propto \left[(n - p) s^2 + (\beta - b)^T X^T X (\beta - b)\right]^{-\frac{n}{2}} \\ &\propto \left[1 + \frac{(\beta - b)^T X^T X (\beta - b)}{(n - p) s^2}\right]^{-\frac{n}{2}} \end{aligned}$$

which is multivariate Student  $t_p \left(\beta; b, s^2 (X^T X)^{-1}, n - p\right)$ .

## 6.8 Bayesian linear regression with general error structure

Now, we consider Bayesian regression with a more general error structure. That is, the *DGP* is

$$y = X\beta + \varepsilon, \quad (\varepsilon | X) \sim N(0, \Sigma)$$

First, we consider the known variance case, then take up the unknown variance case.

### 6.8.1 Known variance

If the error variance,  $\Sigma$ , is known, we simply repeat the Bayesian linear regression approach discussed above for the known variance case after transforming all variables via the Cholesky decomposition of  $\Sigma$ . Let

$$\Sigma = \Gamma \Gamma^T$$

and

$$\Sigma^{-1} = (\Gamma^T)^{-1} \Gamma^{-1}$$

Then, the *DGP* is

$$\Gamma^{-1} \mathbf{y} = \Gamma^{-1} \mathbf{X} \beta + \Gamma^{-1} \varepsilon$$

where

$$\Gamma^{-1} \varepsilon \sim N(0, I_n)$$

With informed priors for  $\beta$ ,  $p(\beta | \Sigma) \sim N(\beta_0, \Sigma_\beta)$  where it is as if  $\Sigma_\beta = (X_0^T \Sigma_0^{-1} X_0)^{-1}$ , the posterior distribution for  $\beta$  conditional on  $\Sigma$  is

$$p(\beta | \Sigma, \mathbf{y}, \mathbf{X}; \beta_0, \Sigma_\beta) \sim N(\bar{\beta}, V_\beta)$$

where

$$\begin{aligned} \bar{\beta} &= (X_0^T \Sigma_0^{-1} X_0 + X^T \Sigma^{-1} X)^{-1} (X_0^T \Sigma_0^{-1} X_0 \beta_0 + X^T \Sigma^{-1} X \hat{\beta}) \\ &= (\Sigma_\beta^{-1} + X^T \Sigma^{-1} X)^{-1} (\Sigma_\beta^{-1} \beta_0 + X^T \Sigma^{-1} X \hat{\beta}) \\ \hat{\beta} &= (X^T \Sigma^{-1} X)^{-1} X^T \Sigma^{-1} \mathbf{y} \end{aligned}$$

and

$$\begin{aligned} V_\beta &= (X_0^T \Sigma_0^{-1} X_0 + X^T \Sigma^{-1} X)^{-1} \\ &= (\Sigma_\beta^{-1} + X^T \Sigma^{-1} X)^{-1} \end{aligned}$$

It is instructive to once again backtrack to develop the conditional posterior distribution. The likelihood function for known variance is

$$\ell(\beta | \Sigma, \mathbf{y}, \mathbf{X}) \propto \exp \left[ -\frac{1}{2} (\mathbf{y} - \mathbf{X} \beta)^T \Sigma^{-1} (\mathbf{y} - \mathbf{X} \beta) \right]$$

Conditional Gaussian priors are

$$p(\beta | \Sigma) \propto \exp \left[ -\frac{1}{2\sigma^2} (\beta - \beta_0)^T V_\beta^{-1} (\beta - \beta_0) \right]$$

The conditional posterior is the product of the prior and likelihood

$$\begin{aligned} p(\beta | \sigma^2, \mathbf{y}, \mathbf{X}) &\propto \exp \left[ -\frac{1}{2\sigma^2} \left\{ \begin{array}{l} (\mathbf{y} - \mathbf{X} \beta)^T \Sigma^{-1} (\mathbf{y} - \mathbf{X} \beta) \\ + (\beta - \beta_0)^T V_\beta^{-1} (\beta - \beta_0) \end{array} \right\} \right] \\ &= \exp \left[ -\frac{1}{2\sigma^2} \left\{ \begin{array}{l} \mathbf{y}^T \Sigma^{-1} \mathbf{y} - 2\mathbf{y}^T \Sigma^{-1} \mathbf{X} \beta + \beta^T \mathbf{X}^T \Sigma^{-1} \mathbf{X} \beta \\ + \beta^T V_\beta^{-1} \beta - 2\beta_0^T V_\beta^{-1} \beta + \beta_0^T V_\beta^{-1} \beta_0 \end{array} \right\} \right] \end{aligned}$$

The first and last terms in the exponent do not involve  $\beta$  (are constants) and can be ignored as they are absorbed through normalization. This leaves

$$\begin{aligned} p(\beta \mid \sigma^2, y, X) &\propto \exp \left[ -\frac{1}{2\sigma^2} \left\{ \begin{array}{l} -2y^T \Sigma^{-1} X \beta + \beta^T X^T \Sigma^{-1} X \beta \\ + \beta^T V_\beta^{-1} \beta - 2\beta_0^T V_\beta^{-1} \beta \end{array} \right\} \right] \\ &= \exp \left[ -\frac{1}{2\sigma^2} \left\{ \begin{array}{l} \beta^T (V_\beta^{-1} + X^T \Sigma^{-1} X) \beta \\ -2(y^T \Sigma^{-1} X + \beta_0^T V_\beta^{-1}) \beta \end{array} \right\} \right] \end{aligned}$$

which can be recognized as the expansion of the conditional posterior claimed above.

$$\begin{aligned} p(\beta \mid \Sigma, y, X) &\sim N(\bar{\beta}, V_\beta) \\ &\propto \exp \left[ -\frac{1}{2} (\beta - \bar{\beta})^T V_\beta^{-1} (\beta - \bar{\beta}) \right] \\ &= \exp \left[ -\frac{1}{2} (\beta - \bar{\beta})^T (V_\beta^{-1} + X^T \Sigma^{-1} X) (\beta - \bar{\beta}) \right] \\ &= \exp \left[ -\frac{1}{2} \left\{ \begin{array}{l} \beta^T (V_\beta^{-1} + X^T \Sigma^{-1} X) \beta \\ -2\bar{\beta}^T (V_\beta^{-1} + X^T \Sigma^{-1} X) \beta \\ + \bar{\beta}^T (V_\beta^{-1} + X^T \Sigma^{-1} X) \bar{\beta} \end{array} \right\} \right] \\ &= \exp \left[ -\frac{1}{2} \left\{ \begin{array}{l} \beta^T (X_0^T X_0 + X^T X) \beta \\ -2(y^T \Sigma^{-1} X + \beta_0^T V_\beta^{-1})^T \beta \\ + \bar{\beta}^T (X_0^T X_0 + X^T X) \bar{\beta} \end{array} \right\} \right] \end{aligned}$$

The last term in the exponent is all constants (does not involve  $\beta$ ) so its absorbed through normalization and disregarded for comparison of kernels. Hence,

$$\begin{aligned} p(\beta \mid \sigma^2, y, X) &\propto \exp \left[ -\frac{1}{2} (\beta - \bar{\beta})^T V_\beta^{-1} (\beta - \bar{\beta}) \right] \\ &\propto \exp \left[ -\frac{1}{2\sigma^2} \left\{ \begin{array}{l} \beta^T (X_0^T X_0 + X^T X) \beta \\ -2(y^T \Sigma^{-1} X + \beta_0^T V_\beta^{-1})^T \beta \end{array} \right\} \right] \end{aligned}$$

as claimed.

### 6.8.2 Unknown variance

Bayesian linear regression with unknown general error structure,  $\Sigma$ , is something of a composite of ideas developed for exchangeable ( $\sigma^2 I_n$  error structure) Bayesian regression and the multivariate Gaussian case with mean

and variance unknown where each draw is an element of the  $y$  vector and  $X$  is an  $n \times p$  matrix of regressors. A Gaussian likelihood is

$$\begin{aligned} \ell(\beta, \Sigma | y, X) &\propto |\Sigma|^{-\frac{n}{2}} \exp \left[ -\frac{1}{2} (y - \beta)^T \Sigma^{-1} (y - \beta) \right] \\ &\propto |\Sigma|^{-\frac{n}{2}} \exp \left[ -\frac{1}{2} \left\{ \begin{array}{l} (y - Xb)^T \Sigma^{-1} (y - Xb) \\ + (b - \beta)^T X^T \Sigma^{-1} X (b - \beta) \end{array} \right\} \right] \\ &\propto |\Sigma|^{-\frac{n}{2}} \exp \left[ -\frac{1}{2} \left\{ (n-p) s^2 + (b - \beta)^T X^T \Sigma^{-1} X (b - \beta) \right\} \right] \end{aligned}$$

where  $b = (X^T \Sigma^{-1} X)^{-1} X^T \Sigma^{-1} y$  and  $s^2 = \frac{1}{n-p} (y - Xb)^T \Sigma^{-1} (y - Xb)$ . Combine the likelihood with a Gaussian-inverted Wishart prior

$$\begin{aligned} p(\beta | \Sigma; \beta_0, \Sigma_\beta) \times p(\Sigma^{-1}; \nu, \Psi) &\propto \exp \left[ -\frac{1}{2} (\beta - \beta_0)^T \Sigma_\beta^{-1} (\beta - \beta_0) \right] \\ &\quad \times |\Psi|^{\frac{\nu}{2}} |\Sigma|^{-\frac{\nu+p+1}{2}} \exp \left[ -\frac{\text{tr}(\Psi \Sigma^{-1})}{2} \right] \end{aligned}$$

where  $\text{tr}(\cdot)$  is the trace of the matrix, it is as if  $\Sigma_\beta = (X_0^T \Sigma_0^{-1} X_0)^{-1}$ , and  $\nu$  is degrees of freedom to produce the joint posterior

$$\begin{aligned} p(\beta, \Sigma | y, X) &\propto |\Psi|^{\frac{\nu}{2}} |\Sigma|^{-\frac{\nu+n+p+1}{2}} \exp \left[ -\frac{\text{tr}(\Psi \Sigma^{-1})}{2} \right] \\ &\quad \times \exp \left[ -\frac{1}{2} \left\{ \begin{array}{l} (n-p) s^2 \\ + (b - \beta)^T X^T \Sigma^{-1} X (b - \beta) \\ + (\beta - \beta_0)^T \Sigma_\beta^{-1} (\beta - \beta_0) \end{array} \right\} \right] \end{aligned}$$

*Completing the square*

Completing the square involves the matrix analog to the univariate unknown mean and variance case. Consider the exponent (in braces)

$$\begin{aligned} &(n-p) s^2 + (b - \beta)^T X^T \Sigma^{-1} X (b - \beta) + (\beta - \beta_0)^T \Sigma_\beta^{-1} (\beta - \beta_0) \\ &= (n-p) s^2 + b^T X^T \Sigma^{-1} X b - 2\beta^T X^T \Sigma^{-1} X b + \beta^T X^T \Sigma^{-1} X \beta \\ &\quad + \beta^T \Sigma_\beta^{-1} \beta - 2\beta^T \Sigma_\beta^{-1} \beta_0 + \beta_0^T \Sigma_\beta^{-1} \beta_0 \\ &= (n-p) s^2 + \beta^T \left( \Sigma_\beta^{-1} + X^T \Sigma^{-1} X \right) \beta \\ &\quad - 2\beta^T V_\beta^{-1} \bar{\beta} + b^T X^T \Sigma^{-1} X b + \beta_0^T \Sigma_\beta^{-1} \beta_0 \\ &= (n-p) s^2 + \beta^T V_\beta^{-1} \beta - 2\beta^T V_\beta^{-1} \bar{\beta} + b^T X^T \Sigma^{-1} X b + \beta_0^T \Sigma_\beta^{-1} \beta_0 \end{aligned}$$

where

$$\begin{aligned} \bar{\beta} &= \left( \Sigma_\beta^{-1} + X^T \Sigma^{-1} X \right)^{-1} \left( \Sigma_\beta^{-1} \beta_0 + X^T \Sigma^{-1} X b \right) \\ &= V_\beta \left( \Sigma_\beta^{-1} \beta_0 + X^T \Sigma^{-1} X b \right) \end{aligned}$$



and  $V_\beta = \left( \Sigma_\beta^{-1} + X^T \Sigma^{-1} X \right)^{-1}$ .

Variation in  $\beta$  around  $\bar{\beta}$  is

$$(\beta - \bar{\beta})^T V_\beta^{-1} (\beta - \bar{\beta}) = \beta^T V_\beta^{-1} \beta - 2\bar{\beta}^T V_\beta^{-1} \beta + \bar{\beta}^T V_\beta^{-1} \bar{\beta}$$

The first two terms are identical to two terms in the posterior involving  $\beta$  and there is apparently no recognizable kernel from these expressions. The joint posterior is

$$\begin{aligned} & p(\beta, \Sigma \mid y, X) \\ & \propto |\Psi|^{\frac{\nu}{2}} |\Sigma|^{-\frac{\nu+n+p+1}{2}} \exp \left[ -\frac{\text{tr}(\Psi \Sigma^{-1})}{2} \right] \\ & \quad \times \exp \left[ -\frac{1}{2} \left\{ \begin{array}{l} (\beta - \bar{\beta})^T V_\beta^{-1} (\beta - \bar{\beta}) \\ + (n-p) s^2 - \bar{\beta}^T V_\beta^{-1} \bar{\beta} \\ + b^T X^T \Sigma^{-1} X b + \beta_0^T \Sigma_\beta^{-1} \beta_0 \end{array} \right\} \right] \\ & \propto |\Psi|^{\frac{\nu}{2}} |\Sigma|^{-\frac{\nu+n+p+1}{2}} \exp \left[ -\frac{1}{2} \left\{ \begin{array}{l} \text{tr}(\Psi \Sigma^{-1}) + (n-p) s^2 \\ - \bar{\beta}^T V_\beta^{-1} \bar{\beta} \\ + b^T X^T \Sigma^{-1} X b + \beta_0^T \Sigma_\beta^{-1} \beta_0 \end{array} \right\} \right] \\ & \quad \times \exp \left[ -\frac{1}{2} \left\{ (\beta - \bar{\beta})^T V_\beta^{-1} (\beta - \bar{\beta}) \right\} \right] \end{aligned}$$

Therefore, we write the conditional posteriors for the parameters of interest. First, we focus on  $\beta$  then we take up  $\Sigma$ .

The conditional posterior for  $\beta$  conditional on  $\Sigma$  involves collecting all terms involving  $\beta$ . Hence, the conditional posterior for  $\beta$  is  $(\beta \mid \Sigma) \sim N(\bar{\beta}, V_\beta)$  or

$$p(\beta \mid \Sigma, y, X) \propto \exp \left[ -\frac{1}{2} \left\{ (\beta - \bar{\beta})^T V_\beta^{-1} (\beta - \bar{\beta}) \right\} \right]$$

*Inverted-Wishart kernel*

Now, we gather all terms involving  $\Sigma$  and write the conditional posterior for  $\Sigma$ .

$$\begin{aligned} & p(\Sigma \mid \beta, y, X) \\ & \propto |\Psi|^{\frac{\nu}{2}} |\Sigma|^{-\frac{\nu+n+p+1}{2}} \exp \left[ -\frac{1}{2} \left\{ \begin{array}{l} \text{tr}(\Psi \Sigma^{-1}) + (n-p) s^2 \\ + (b - \beta)^T X^T \Sigma^{-1} X (b - \beta) \end{array} \right\} \right] \\ & \propto |\Psi|^{\frac{\nu}{2}} |\Sigma|^{-\frac{\nu+n+p+1}{2}} \exp \left[ -\frac{1}{2} \left\{ \begin{array}{l} \text{tr}(\Psi \Sigma^{-1}) + \\ (y - Xb)^T \Sigma^{-1} (y - Xb) \\ + (b - \beta)^T X^T \Sigma^{-1} X (b - \beta) \end{array} \right\} \right] \\ & \propto |\Psi|^{\frac{\nu}{2}} |\Sigma|^{-\frac{\nu+n+p+1}{2}} \exp \left[ -\frac{1}{2} \left\{ \text{tr} \left( \left[ \begin{array}{l} \Psi + (y - Xb)^T (y - Xb) \\ + (b - \beta)^T X^T X (b - \beta) \end{array} \right] \Sigma^{-1} \right) \right\} \right] \end{aligned}$$

We can identify the kernel as an inverted-Wishart involving the trace of a square, symmetric matrix, call it  $\Psi_n$ , multiplied by  $\Sigma^{-1}$ .

The above joint posterior can be rewritten as an inverted-Wishart ( $\Sigma^{-1}; \nu + n, \Psi_n$ )

$$p(\theta, \Sigma | y) \propto |\Psi_n|^{\frac{\nu+n}{2}} |\Sigma|^{-\frac{\nu+n+p+1}{2}} \exp \left[ -\frac{1}{2} \text{tr} (\Psi_n \Sigma^{-1}) \right]$$

where

$$\Psi_n = \Psi + (y - Xb)^T (y - Xb) + (b - \beta)^T X^T X (b - \beta)$$

With conditional posteriors in hand, we can employ *McMC* strategies (namely, a Gibbs sampler) to draw inferences around the parameters of interest,  $\beta$  and  $\Sigma$ . That is, we sequentially draw  $\beta$  conditional on  $\Sigma$  and  $\Sigma$ , in turn, conditional on  $\beta$ . We discuss *McMC* strategies (both the Gibbs sampler and its generalization, the Metropolis-Hastings algorithm) later.

### 6.8.3 (Nearly) uninformative priors

As discussed by Gelman, et al [2004] uninformative priors for this case is awkward, at best. What does it mean to posit uninformative priors for a regression with general error structure? Consistent probability assignment suggests that either we have some priors about the correlation structure or heteroskedastic nature of the errors (informative priors) or we know nothing about the error structure (uninformative priors). If priors are uninformative, then maximum entropy probability assignment suggests we assign independent and unknown homoskedastic errors. Hence, we discuss nearly uninformative priors for this general error structure regression.

The joint uninformative prior (with a locally uniform prior for  $\beta$ ) is

$$p(\beta, \Sigma) \propto |\Sigma|^{-\frac{1}{2}}$$

and the joint posterior is

$$\begin{aligned} p(\beta, \Sigma | y, X) &\propto |\Sigma|^{-\frac{1}{2}} |\Sigma|^{-\frac{n}{2}} \exp \left[ -\frac{1}{2} \left\{ (n-p) s^2 + (b - \beta)^T X^T \Sigma^{-1} X (b - \beta) \right\} \right] \\ &\propto |\Sigma|^{-\frac{n+1}{2}} \exp \left[ -\frac{1}{2} \left\{ (n-p) s^2 + (b - \beta)^T X^T \Sigma^{-1} X (b - \beta) \right\} \right] \\ &\propto |\Sigma|^{-\frac{n+1}{2}} \exp \left[ -\frac{1}{2} \text{tr} (S(\beta) \Sigma^{-1}) \right] \end{aligned}$$

where now  $S(\beta) = (y - Xb)^T (y - Xb) + (b - \beta)^T X^T X (b - \beta)$ . Then, the conditional posterior for  $\beta$  given  $\Sigma$  is  $N \left( b, (X^T \Sigma^{-1} X)^{-1} | \Sigma \right)$

$$p(\beta | \Sigma, y, X) \propto \exp \left[ -\frac{n}{2} (b - \beta)^T X^T \Sigma^{-1} X (b - \beta) \right]$$

The conditional posterior for  $\Sigma$  given  $\beta$  is inverted-Wishart( $\Sigma^{-1}; n, \Psi_n$ )

$$p(\beta, \Sigma | y) \propto |\Psi_n|^{\frac{n}{2}} |\Sigma|^{-\frac{n+1}{2}} \exp \left[ -\frac{1}{2} \text{tr}(\Psi_n \Sigma^{-1}) \right]$$

where

$$\Psi_n = (y - Xb)^T (y - Xb) + (b - \beta)^T X^T X (b - \beta)$$

As with informed priors, a Gibbs sampler (sequential draws from the conditional posteriors) can be employed to draw inferences for the uninformative prior case.

Next, we discuss posterior simulation, a convenient and flexible strategy for drawing inference from the evidence and (conjugate) priors.

## 6.9 Appendix: summary of conjugacy

focal parameter(s) $\theta$	<i>prior</i> $\pi(\theta)$	<i>likelihood</i> $\ell(\theta   y)$	<i>posterior</i> $\pi(\theta   y)$
discrete data:			
beta-binomial $p$	beta $\propto p^{a-1}(1-p)^{b-1}$	binomial $\propto p^s(1-p)^{n-s}$	beta $\propto p^{a+s-1}(1-p)^{b+n-s-1}$
gamma-poisson $\lambda$	gamma $\propto \lambda^{a-1}e^{-b\lambda}$	poisson $\propto \lambda^s e^{-n\lambda}$	gamma $\propto \lambda^{a+s-1}e^{-(b+n)\lambda}$
gamma-exponential $\theta$	gamma $\propto \theta^{a-1}e^{-b\theta}$	exponential $\propto \theta^n e^{-s\theta}$	gamma $\propto \theta^{a+n-1}e^{-(b+s)\theta}$
beta-negative binomial $p$	beta $\propto p^{a-1}(1-p)^{b-1}$	negative binomial $\propto p^{nr}(1-p)^s$	beta $\propto p^{a+nr-1}(1-p)^{b+s-1}$
beta-binomial- hypergeometric <sup>8</sup> $k$ $k$ unknown population success $N$ known population size $n$ known sample size $x$ known sample success	beta-binomial $\binom{n}{x}$ $\frac{\Gamma(a+x)\Gamma(b+n-x)\Gamma(a+b)}{\Gamma(a)\Gamma(b)\Gamma(a+b+n)}$ , $x = 0, 1, 2, \dots, n$	hypergeometric $\frac{\binom{k}{x}\binom{N-k}{n-x}}{\binom{N}{n}}$ sampling without replacement	beta-binomial $\binom{N-n}{k-x}$ $\frac{\Gamma(a+k)\Gamma(b+N-k)\Gamma(a+b+n)}{\Gamma(a+x)\Gamma(b+n-x)\Gamma(a+b+N)}$ , $k = x, x+1, \dots,$ $x+N-n$
multinomial- Dirichlet $\theta$ (vector)	Dirichlet $\propto \prod_{i=1}^K \theta_i^{a_i-1}$	multinomial $\propto \theta_1^{s_1} \dots \theta_K^{s_K}$	Dirichlet $\propto \prod_{i=1}^K \theta_i^{a_i+s_i-1}$

$$s = \sum_{i=1}^n y_i, \quad \binom{n}{x} = \frac{n!}{x!(n-x)!},$$

$$\Gamma(z) = \int_0^\infty e^{-t} t^{z-1} dt, \quad \Gamma(n) = (n-1)! \text{ for } n \text{ a positive integer, } B(a, b) = \frac{\Gamma(a)\Gamma(b)}{\Gamma(a+b)}$$

<sup>8</sup>See Dyer and Pierce [1993].

focal parameter(s) $\theta$	prior $\pi(\theta)$	likelihood $\ell(\theta   y)$	posterior $\pi(\theta   y)$
continuous data:			
Pareto-uniform $w$ $w$ unknown upper bound, 0 known lower bound	Pareto $\frac{ab^a}{w^{a+1}}$	uniform $\frac{1}{w^n}$ , $w > \max(x_i)$	Pareto $\propto \frac{(a+n) \max[b, x_i]^{a+n}}{w^{a+n+1}}$
Pareto-Pareto $\beta$ $\beta$ unknown precision, $\alpha$ known shape	Pareto $\frac{ab^{\alpha a}}{\beta^{a+\alpha}}$ , $\beta > b$	Pareto $\propto \beta^{n\alpha}$ , $0 < \beta < \min(x_i)$	Pareto $\propto \frac{(a-\alpha n)b^{(a-\alpha n)\alpha}}{\beta^{a-\alpha n+1}}$ , $a > \alpha n, \beta > b$
gamma-Pareto $\alpha$ $\alpha$ unknown shape, $\beta$ known precision	gamma $\propto \frac{\alpha^{a-1} e^{-\alpha/b}}{b^a \Gamma(a)}$ , $\alpha > 0$	Pareto $\propto \frac{\alpha^n \beta^{n\alpha}}{m^{a+1}}$ , $m = \prod_{i=1}^n x_i$ , $0 < \beta < \min(x_i)$	gamma $\propto \frac{\alpha^{a+n-1} e^{-\alpha/b'}}{(b')^{a+n} \Gamma(a+n)}$ , $b' = \frac{1}{\frac{1}{b} + \log m - n \log \beta} > 0$
gamma-exponential $\theta$	gamma $\propto \frac{\theta^{a-1} e^{-\theta/b}}{\Gamma(a)b^a}$	exponential $\propto \theta^n e^{-s\theta}$ , $s = \sum_{i=1}^n x_i$	gamma $\propto \frac{\theta^{a+n-1} e^{-\theta/b'}}{\Gamma(a+n)(b')^{a+n}}$ , $b' = \frac{b}{1+bs}$
inverse gamma-gamma $\beta$ $\beta$ unknown rate, $\alpha$ known shape	inverse gamma $\propto \frac{\beta^{1-a} e^{-1/\beta b}}{\Gamma(a)b^a}$	gamma $\propto \frac{e^{-s/\beta}}{\beta^{\alpha n}}$ , $s = \sum_{i=1}^n x_i$	inverse gamma $\propto \frac{\beta^{1-a-\alpha n} e^{-1/\beta b'}}{\Gamma(a+\alpha n)(b')^{a+\alpha n}}$ , $b' = \frac{b}{1+bs}$
conjugate prior-gamma $\alpha$ $\alpha$ unknown shape, $\beta$ known rate	nonstandard $\propto \frac{\alpha^{\alpha-1} \beta^{\alpha c}}{\Gamma(\alpha)^b}$ , $a, b, c > 0$ $\alpha > 0$	gamma $\propto \frac{m^{\alpha-1}}{\beta^{-\alpha n} \Gamma(\alpha)^n}$ , $m = \prod_{i=1}^n x_i$ , $x_i > 0$	nonstandard $\propto \frac{(am)^{\alpha-1} \beta^{\alpha(c+n)}}{\Gamma(\alpha)^{b+n}}$

focal parameter(s) $\theta$	prior $\pi(\theta)$	likelihood $\ell(\theta   y)$	posterior $\pi(\theta   y)$
continuous data:			
normal-normal $\mu$	normal $\propto \exp\left[-\frac{(\mu-\mu_0)^2}{2\sigma_0^2}\right],$ $\sigma_0^2 = \frac{\sigma^2}{\kappa_0}$	normal $\propto \prod_{i=1}^n \exp\left[-\frac{(y_i-\mu)^2}{2\sigma^2}\right]$ $= \exp\left[-\frac{ss}{2\sigma^2}\right]$	normal $\propto \exp\left[-\frac{(\mu-\mu_n)^2}{2\sigma_n^2}\right],$ $\mu_n = \frac{\kappa_0\mu_0+n\bar{y}}{\kappa_0+n},$ $\sigma_n^2 = \frac{\sigma^2}{\kappa_0+n}$
inverse gamma -normal $\sigma^2$	inverse gamma $\propto (\sigma^2)^{-(a+1)}$ $\exp\left[-\frac{b}{\sigma^2}\right]$	normal $\propto \frac{1}{(\sigma^2)^{n/2}}$ $\exp\left[-\frac{ss}{2\sigma^2}\right]$	inverse gamma $\propto (\sigma^2)^{-\left(\frac{n+2a}{2}+1\right)}$ $\exp\left[-\frac{b+\frac{1}{2}ss}{\sigma^2}\right]$
(normal   $\sigma^2$ ) $\times$ inverse gamma- normal $\mu, \sigma^2$	(normal   $\sigma^2$ ) $\times$ inverse gamma $\propto \sigma_0^{-1} \exp\left[-\frac{(\mu-\mu_0)^2}{2\sigma_0^2}\right]$ $\times (\sigma^2)^{-(a+1)}$ $\exp\left[-\frac{b}{\sigma^2}\right],$ $\sigma_0^2 = \frac{\sigma^2}{\kappa_0}$	normal $\propto \frac{1}{(\sigma^2)^{n/2}}$ $\exp\left[-\frac{ss}{2\sigma^2}\right]$	joint posterior: (normal   $\sigma^2$ ) $\times$ inverse gamma $\propto \sqrt{\frac{\kappa'_0}{\sigma^2}} \exp\left[-\frac{\kappa'_0(\mu-\mu'_0)^2}{2\sigma^2}\right]$ $\times (\sigma^2)^{-(a'+1)}$ $\exp\left[-\frac{b'}{\sigma^2}\right];$
			Student t marginal posterior for $\mu$ : $\propto \left(1 + \frac{\kappa_0 b' (\mu - \mu'_0)^2}{2}\right)^{-\frac{2a'+1}{2}};$ inverse gamma marginal posterior for $\sigma^2$ : $\propto (\sigma^2)^{-(a'+1)}$ $\exp\left[-\frac{b'}{\sigma^2}\right],$ $a' = a + \frac{n}{2}, \kappa'_0 = \kappa_0 + n,$ $b' = \frac{\frac{1}{b} + \frac{ss}{2}}{+ \frac{\kappa_0 n (\bar{y} - \mu_0)^2}{2(\kappa_0 + n)}},$ $\mu'_0 = \frac{\kappa_0 \mu_0 + n \bar{y}}{\kappa_0 + n}$

$$ss = \sum_{i=1}^n (y_i - \mu)^2$$

focal parameter(s) $\theta$	<i>prior</i> $\pi(\theta)$	<i>likelihood</i> $\ell(\theta   y)$	<i>posterior</i> $\pi(\theta   y)$
continuous data:			
bilateral bivariate Pareto- uniform $l, u$	bilateral bivariate Pareto $\frac{a(a+1)(r_2-r_1)^a}{(u-l)^{a+2}}$ , $l < r_1, u > r_2$	uniform $\left(\frac{1}{u-l}\right)^n$	bilateral bivariate Pareto $(a+n)(a+n+1)$ $\frac{(r_2-r_1)^{a+n}}{(u-l)^{a+n+2}}$ , $r_1' = \min(r_1, x_i)$ , $r_2' = \max(r_2, x_i)$
normal- lognormal $\mu$	normal $\propto \exp\left[-\frac{(\mu-\mu_0)^2}{2\sigma_0^2}\right]$ , $\sigma_0^2 = \frac{\sigma^2}{\kappa_0}$	log normal $\propto \prod_{i=1}^n \exp\left[-\frac{(\log y_i - \mu)^2}{2\sigma^2}\right]$ $= \exp\left[-\frac{lss}{2\sigma^2}\right]$	normal $\propto \exp\left[-\frac{(\mu-\mu_n)^2}{2\sigma_n^2}\right]$ , $\mu_n = \frac{\kappa_0\mu_0 + n\log y}{\kappa_0 + n}$ , $\sigma_n^2 = \frac{\sigma^2}{\kappa_0 + n}$
inverse gamma- lognormal $\sigma^2$	inverse gamma $\propto (\sigma^2)^{-(a+1)}$ $\exp\left[-\frac{b}{\sigma^2}\right]$	normal $\propto \frac{1}{(\sigma^2)^{n/2}}$ $\exp\left[-\frac{lss}{2\sigma^2}\right]$	inverse gamma $\propto (\sigma^2)^{-\left(\frac{n+2a}{2}+1\right)}$ $\exp\left[-\frac{b+\frac{1}{2}lss}{\sigma^2}\right]$

$$lss = \sum_{i=1}^n (\log y_i - \mu)^2$$

---

continuous data: multivariate normal × inverted Wishart- multivariate normal $\mu, \Sigma$ prior $\pi(\mu, \Sigma)$	
multivariate normal $\pi(\mu   \Sigma)$	$\propto  \Sigma ^{-\frac{1}{2}} \exp \left[ -\frac{\kappa_0}{2} (\mu - \mu_0)^T \Sigma^{-1} (\mu - \mu_0) \right]$
× inverted Wishart $\pi(\Sigma)$	$\propto  \Psi ^{\frac{k}{2}}  \Sigma ^{-\frac{\nu+k+1}{2}} \exp \left[ -\frac{\text{tr}(\Psi \Sigma^{-1})}{2} \right]$
likelihood $\ell(\mu, \Sigma   y)$ multivariate normal	$\propto  \Sigma ^{-\frac{n}{2}} \exp \left[ -\frac{1}{2} \left\{ (n-1) s^2 + (\bar{y} - \mu)^T \Sigma^{-1} (\bar{y} - \mu) \right\} \right]$
joint posterior $\pi(\mu, \Sigma   y)$ multivariate normal $\pi(\mu   \Sigma, y)$	$\propto \exp \left[ -\frac{\kappa_0+n}{2} (\mu - \mu_n)^T \Sigma^{-1} (\mu - \mu_n) \right]$
× inverted Wishart $\pi(\Sigma   y)$	$\propto  \Psi ^{\frac{\nu+n}{2}}  \Sigma ^{-\frac{\nu+n+k+1}{2}} \exp \left[ -\frac{\text{tr}(\Psi_n \Sigma^{-1})}{2} \right]$
marginal posterior multivariate Student t $\pi(\mu   y)$	$\propto \left[ I + (\kappa_0 + n) \left( (\mu - \mu_n)^T \Psi_n^{-1} (\mu - \mu_n) \right) \right]^{-\frac{1}{2}(\nu+n+1)}$
inverted Wishart $\pi(\Sigma   y)$	$\propto  \Psi ^{\frac{\nu+n}{2}}  \Sigma ^{-\frac{\nu+n+k+1}{2}} \exp \left[ -\frac{\text{tr}(\Psi_n \Sigma^{-1})}{2} \right]$

where

$\text{tr}(\cdot)$  is the trace of a matrix,

$$\mu_n = \frac{\kappa_0 \mu_0 + n \bar{y}}{\kappa_0 + n},$$

$$s^2 = \frac{1}{n-1} \sum_{i=1}^n (y_i - \bar{y})^T \Sigma^{-1} (y_i - \bar{y}),$$

$$\Psi_n = \Psi + \sum_{i=1}^n (y_i - \bar{y})^T (y_i - \bar{y}) + \frac{\kappa_0 n}{\kappa_0 + n} (\mu_0 - \bar{y})^T (\mu_0 - \bar{y})$$


---



---

continuous data:	
linear regression	
normal $\times$ inverse chi square-normal	
$\beta, \sigma^2$	
prior $\pi(\beta, \sigma^2)$	
normal $\pi(\beta   \sigma^2)$	$\propto \sigma^{-p} \exp \left[ -\frac{1}{2\sigma^2} (\beta - \beta_0)^T \Omega_0 (\beta - \beta_0) \right]$
$\times$ inverse chi square $\pi(\sigma^2)$	$\times \sigma^{-(\nu_0/2+1)} \exp \left[ -\frac{\nu_0 \sigma_0^2}{2\sigma^2} \right]$
normal likelihood $\ell(\beta, \sigma^2   y, X)$	
normal	$\propto \sigma^{-n} \exp \left[ -\frac{1}{2\sigma^2} \left\{ e^T e + (\beta - b)^T X^T X (\beta - b) \right\} \right]$
joint posterior $p(\beta, \sigma^2   y, X)$	
normal $p(\beta   \sigma^2, y, X)$	$\propto \sigma^{-p} \exp \left[ -\frac{1}{2\sigma^2} (\beta - \bar{\beta})^T \Omega_n (\beta - \bar{\beta}) \right]$
$\times$ inverse chi square $\pi(\sigma^2   y, X)$	$\times (\sigma^2)^{-[(\nu_0+n)/2+1]} \exp \left[ -\frac{\sigma_n^2}{2\sigma^2} \right]$
marginal posterior	
Student t $\pi(\beta   \sigma^2, y, X)$	$\propto \left[ 1 + \frac{1}{\nu_n \sigma_n^2} (\beta - \bar{\beta})^T \Omega_n (\beta - \bar{\beta}) \right]^{-\frac{\nu_0+n+p}{2}}$
$\times$ inverse chi square $\pi(\sigma^2   y, X)$	$\propto (\sigma^2)^{-[(\nu_0+n)/2+1]} \exp \left[ -\frac{\sigma_n^2}{2\sigma^2} \right]$

where

$$e = y - Xb,$$

$$b = (X^T X)^{-1} X^T y,$$

$$\bar{\beta} = (\Omega_0 + X^T X)^{-1} (\Omega_0 \beta_0 + X^T X b),$$

$$\Omega_n = \Omega_0 + X^T X,$$

$$\nu_n \sigma_n^2 = \nu_0 \sigma_0^2 + e^T e + (\beta_0 - \bar{\beta})^T \Omega_0 (\beta_0 - \bar{\beta}) + (\hat{\beta} - \bar{\beta})^T X^T X (\hat{\beta} - \bar{\beta}),$$

$$\text{and } \nu_n = \nu_0 + n$$


---

---

continuous data: linear regression with general variance	
normal $\times$ inverted Wishart-normal $\beta, \Sigma$ prior $\pi(\beta, \Sigma)$ normal $\pi(\beta   \Sigma)$	$\propto \exp \left[ -\frac{1}{2} (\beta - \beta_0)^T \Sigma_\beta^{-1} (\beta - \beta_0) \right]$
$\times$ inverted Wishart $\pi(\Sigma)$	$\times  \Psi ^{\frac{p}{2}}  \Sigma ^{-\frac{\nu+p+1}{2}} \exp \left[ -\frac{\text{tr}(\Psi \Sigma^{-1})}{2} \right]$
normal likelihood $\ell(\beta, \Sigma   y, X)$ normal	$\propto  \Sigma ^{-\frac{n}{2}} \exp \left[ -\frac{1}{2} \left\{ (n-p) s^2 + (\beta - b)^T X^T \Sigma^{-1} X (\beta - b) \right\} \right]$
conditional posterior normal $p(\beta   \Sigma, y, X)$	$\propto \exp \left[ -\frac{1}{2} (\beta - \bar{\beta})^T V_\beta^{-1} (\beta - \bar{\beta}) \right]$
inverted Wishart $\pi(\Sigma   \beta, y, X)$	$\propto  \Psi ^{\frac{\nu+n}{2}}  \Sigma ^{-\frac{\nu+n+p+1}{2}} \exp \left[ -\frac{\text{tr}(\Psi_n \Sigma^{-1})}{2} \right]$
where	$\text{tr}(\cdot)$ is the trace of a matrix,
$s^2 = \frac{1}{n-p} (y - Xb)^T \Sigma^{-1} (y - Xb),$	$b = (X^T \Sigma^{-1} X)^{-1} X^T \Sigma^{-1} y,$
$V_\beta = \left( \Sigma_\beta^{-1} + X^T \Sigma^{-1} X \right)^{-1},$	$\bar{\beta} = \left( \Sigma_\beta^{-1} + X^T \Sigma^{-1} X \right)^{-1} \left( \Sigma_\beta^{-1} \beta_0 + X^T \Sigma^{-1} X b \right),$
$\Psi_n = \Psi + (y - Xb)^T (y - Xb) + (b - \bar{\beta})^T X^T X (b - \bar{\beta}),$	
$\text{and } \Sigma_\beta = (X_0^T \Sigma_0^{-1} X_0)^{-1}$	

---