

# SDS 387D: Statistical Modeling II

## Exercise 1

Travis Lilley

Jan 23, 2019

## Bayesian inference in simple conjugate families

(a)

The posterior is proportional to the product of the prior and the likelihood,  $f(x_1, \dots, x_N | w)$ . To find the likelihood, we write

$$\begin{aligned}
 f(x_1, \dots, x_N | w) &= \prod_{i=1}^N f(x_i | w) && \text{i.i.d. observations} \\
 &= \prod_{i=1}^N w^{x_i} (1-w)^{1-x_i} && \text{density of Bernoulli} \\
 &= w^{\sum_{i=1}^N x_i} (1-w)^{N - \sum_{i=1}^N x_i}. && \text{simplify product}
 \end{aligned}$$

We then obtain the posterior up to a constant of proportionality.

$$\begin{aligned}
 p(w | x_1, \dots, x_N) &\propto f(x_1, \dots, x_N | w) p(w) && \text{likelihood times prior} \\
 &= \left[ w^{\sum_{i=1}^N x_i} (1-w)^{N - \sum_{i=1}^N x_i} \right] \left[ \frac{\Gamma(a+b)}{\Gamma(a)\Gamma(b)} w^{a-1} (1-w)^{b-1} \right] && \text{substitute} \\
 &\propto w^{a + (\sum_{i=1}^N x_i) - 1} (1-w)^{b + N - (\sum_{i=1}^N x_i) - 1}. && \text{simplify}
 \end{aligned}$$

But if we view the last expression as a function of  $w$ , then we recognize that it is the kernel of a beta density with parameters  $a + \sum_{i=1}^N x_i$  and  $b + N - \sum_{i=1}^N x_i$ . Thus, the posterior  $p(w | x_1, \dots, x_N)$  is the density of a  $\text{Beta}\left(a + \sum_{i=1}^N x_i, b + N - \sum_{i=1}^N x_i\right)$  distribution.

(b)

We assume that  $X_1$  and  $X_2$  are independent, so that their joint distribution is the product of their marginal distributions. That is,

$$\begin{aligned}
 f_{X_1, X_2}(x_1, x_2) &= f_{X_1}(x_1)f_{X_2}(x_2) && \text{independence} \\
 &= \left[ \frac{1}{\Gamma(a_1)} x_1^{a_1-1} \exp(-x_1) \right] \left[ \frac{1}{\Gamma(a_2)} x_2^{a_2-1} \exp(-x_2) \right] && \text{gamma densities} \\
 &= \frac{1}{\Gamma(a_1)\Gamma(a_2)} x_1^{a_1-1} x_2^{a_2-1} \exp\{-(x_1 + x_2)\}. && \text{simplify}
 \end{aligned}$$

Since  $X_1$  and  $X_2$  are independent, the support of their joint distribution is  $\mathcal{A} = \{(x_1, x_2) \in \mathbb{R}^2 : x_1 > 0, x_2 > 0\}$ .

We now consider the transformation defined by  $Y_1 = \frac{X_1}{X_1 + X_2}$  and  $Y_2 = X_1 + X_2$ . We first need to determine the image of this transformation. To do this, we solve for  $X_1$  and  $X_2$  in terms of  $Y_1$  and  $Y_2$ . Since  $Y_2 = X_1 + X_2$ , we see that  $Y_1 = \frac{X_1}{Y_2}$ , so that  $X_1 = Y_1 Y_2$ . Substituting  $Y_1 Y_2$  in for  $X_1$  into the equation  $Y_2 = X_1 + X_2$  then shows that  $Y_2 = Y_1 Y_2 + X_2$ , so that  $X_2 = Y_2 - Y_1 Y_2$ .

Now, since  $X_2 > 0$  and  $X_2 = Y_2 - Y_1 Y_2$ , we see that  $Y_2 - Y_1 Y_2 > 0$ . But this implies that  $Y_2 > Y_1 Y_2 > 0$ , which in turn implies that  $0 < Y_1 < 1$ . Moreover, since  $Y_2$  is the sum of two positive real numbers, it may also equal any positive real number. Hence, the support of the joint distribution for  $Y_1$  and  $Y_2$  will be  $\mathcal{B} = \{(y_1, y_2) \in \mathbb{R}^2 : 0 < y_1 < 1, y_2 > 0\}$ .

We also need to confirm that the transformation defined by  $Y_1$  and  $Y_2$  is one-to-one. Suppose the points  $(x_1, x_2)$  and  $(x'_1, x'_2)$  in  $\mathcal{A}$  are mapped to the same point in  $\mathcal{B}$ . It suffices to show that this implies that  $(x_1, x_2)$  and  $(x'_1, x'_2)$  must actually be the same point in. In other words, suppose that  $\left(\frac{x_1}{x_1 + x_2}, x_1 + x_2\right)$  is the same as the point  $\left(\frac{x'_1}{x'_1 + x'_2}, x'_1 + x'_2\right)$ . This implies that  $\frac{x_1}{x_1 + x_2} = \frac{x'_1}{x'_1 + x'_2}$ , and that  $x_1 + x_2 = x'_1 + x'_2$ . Hence,  $\frac{x_1}{x_1 + x_2} = \frac{x'_1}{x_1 + x_2}$ , so that  $x_1 = x'_1$ . And if  $x_1 = x'_1$ , then  $x_1 + x_2 = x'_1 + x'_2$  implies that  $x_2 = x'_2$ . Thus,  $(x_1, x_2)$  and  $(x'_1, x'_2)$  are actually the same point, and the transformation is one-to-one.

The final piece we need before using the formula for bivariate transformations is the Jacobian of the transformation,  $J$ . Since  $X_1 = Y_1 Y_2$ , and  $X_2 = Y_2 - Y_1 Y_2$ , we have

$$\begin{aligned}
 J &= \begin{vmatrix} \frac{\partial x_1}{\partial y_1} & \frac{\partial x_1}{\partial y_2} \\ \frac{\partial x_2}{\partial y_1} & \frac{\partial x_2}{\partial y_2} \end{vmatrix} && \text{definition of Jacobian} \\
 &= \begin{vmatrix} y_2 & y_1 \\ -y_2 & 1 - y_1 \end{vmatrix} && \text{partial differentiation} \\
 &= y_2(1 - y_1) - (y_1)(-y_2) && \text{definition of determinant} \\
 &= y_2.
 \end{aligned}$$

We are now ready to find the joint density of  $Y_1$  and  $Y_2$ .

$$\begin{aligned}
 f_{Y_1, Y_2}(y_1, y_2) &= f_{X_1, X_2}(x_1, x_2)|J| && \text{bivariate transformation formula} \\
 &= f_{X_1, X_2}(y_1 y_2, y_2 - y_1 y_2)|y_2| && \text{inverse transformation and Jacobian} \\
 &= \frac{y_2}{\Gamma(a_1)\Gamma(a_2)} (y_1 y_2)^{a_1-1} (y_2 - y_1 y_2)^{a_2-1} \exp\{-(y_1 y_2 + y_2 - y_1 y_2)\} && \text{substitute and simplify} \\
 &= \frac{1}{\Gamma(a_1)\Gamma(a_2)} y_2 y_1^{a_1-1} y_2^{a_1-1} y_2^{a_2-1} (1 - y_1)^{a_2-1} \exp(-y_2) \\
 &= \frac{1}{\Gamma(a_1)\Gamma(a_2)} y_1^{a_1-1} (1 - y_1)^{a_2-1} y_2^{a_1+a_2-1} \exp(-y_2) \\
 &= \frac{\Gamma(a_1 + a_2)}{\Gamma(a_1)\Gamma(a_2)} \frac{1}{\Gamma(a_1 + a_2)} y_1^{a_1-1} (1 - y_1)^{a_2-1} y_2^{a_1+a_2-1} \exp(-y_2) \\
 &= \underbrace{\left[ \frac{\Gamma(a_1 + a_2)}{\Gamma(a_1)\Gamma(a_2)} y_1^{a_1-1} (1 - y_1)^{a_2-1} \right]}_{f(y_1)} \underbrace{\left[ \frac{1}{\Gamma(a_1 + a_2)} y_2^{a_1+a_2-1} \exp(-y_2) \right]}_{g(y_2)},
 \end{aligned}$$

which has support on  $\mathcal{B}$ , defined earlier. But notice that we have written the joint density of  $Y_1$  and  $Y_2$  as the product of two functions,  $f$  and  $g$ , the first of which depends only on  $y_1$ , and the second which depends only on  $y_2$ . This means that  $y_1$  and  $y_2$  must be independent. Moreover, the joint density must equal the product of the two marginal densities of  $Y_1$  and  $Y_2$ . We observe that since  $0 < y_1 < 1$ , and the function  $f(y_1)$  is the same as the density of a  $\text{Beta}(a_1, a_2)$  distribution. Likewise,  $y_2 > 0$ , and  $g(y_2)$  is the density of a  $\text{Gamma}(a_1 + a_2, 1)$  distribution. Therefore, we conclude that, marginally,

$$Y_1 \sim \text{Beta}(a_1, a_2), \text{ and} \\ Y_2 \sim \text{Gamma}(a_1 + a_2, 1).$$

Suppose we wish to generate a random sample from a  $\text{Beta}(a_1, a_2)$  distribution using only random samples from gamma distributions. Then the above results suggest the following procedure:

1. Sample some  $X_1$  from a  $\text{Gamma}(a_1, 1)$  distribution.
2. Sample some  $X_2$  from a  $\text{Gamma}(a_2, 1)$  distribution that is independent of the first gamma.
3. Compute  $Y^{(1)} = \frac{X_1}{X_1 + X_2}$ .
4. Repeat steps 1 through 3 above  $B$  times to generate  $Y^{(1)}, Y^{(2)}, \dots, Y^{(B)}$ .

These  $Y^{(1)}, Y^{(2)}, \dots, Y^{(B)}$  will then constitute an i.i.d. random sample of size  $B$  from a  $\text{Beta}(a_1, a_2)$  distribution.

(c)

As in part (a), we use the fact that the posterior is proportional to the product of the likelihood and the prior of the parameter(s) of interest. Now, since the prior on  $\theta$  is  $N(m, v)$ , we have  $p(\theta) = \frac{1}{\sqrt{2\pi v}} \exp(-\frac{1}{2v}(\theta - m)^2)$ . Next, we obtain the likelihood as follows:

$$\begin{aligned}
 f(x_1, \dots, x_N | \theta) &= \prod_{i=1}^N f(x_i | \theta) && \text{i.i.d. observations} \\
 &= \prod_{i=1}^N \left[ \frac{1}{\sigma\sqrt{2\pi}} \exp \left\{ -\frac{1}{2\sigma^2} (x_i - \theta)^2 \right\} \right] && \text{normal density} \\
 &= \sigma^{-N} (2\pi)^{-\frac{N}{2}} \exp \left\{ -\frac{1}{2\sigma^2} \sum_{i=1}^N (x_i - \theta)^2 \right\}. && \text{expand product and simplify}
 \end{aligned}$$

The posterior is then

$$\begin{aligned}
 p(\theta | x_1, \dots, x_N) &\propto f(x_1, \dots, x_N | \theta) p(\theta) && \text{Bayes' theorem} \\
 &= \left[ \sigma^{-N} (2\pi)^{-\frac{N}{2}} \exp \left\{ -\frac{1}{2\sigma^2} \sum_{i=1}^N (x_i - \theta)^2 \right\} \right] \left[ \frac{1}{\sqrt{2\pi v}} \exp \left\{ -\frac{1}{2v} (\theta - m)^2 \right\} \right] && \text{substitute} \\
 &\propto \exp \left\{ -\frac{1}{2\sigma^2} \sum_{i=1}^N (x_i - \theta)^2 - \frac{1}{2v} (\theta - m)^2 \right\} && \text{drop front constants} \\
 &= \exp \left\{ -\frac{1}{2} \left[ \frac{1}{\sigma^2} \sum_{i=1}^N (x_i - \theta)^2 + \frac{1}{v} (\theta - m)^2 \right] \right\} \\
 &= \exp \left\{ -\frac{1}{2} \left[ \frac{1}{\sigma^2} \sum_{i=1}^N x_i^2 - \frac{2\theta}{\sigma^2} \sum_{i=1}^N x_i + \frac{1}{\sigma^2} \sum_{i=1}^N \theta^2 + \frac{\theta^2}{v} - \frac{2\theta m}{v} + \frac{m^2}{v} \right] \right\} && \text{expand} \\
 &\propto \exp \left\{ -\frac{1}{2} \left[ -\frac{2\theta}{\sigma^2} \sum_{i=1}^N x_i + \frac{N\theta^2}{\sigma^2} + \frac{\theta^2}{v} - \frac{2\theta m}{v} \right] \right\} && \text{drop non-}\theta \text{ terms} \\
 &= \exp \left\{ -\frac{1}{2} \left[ \left( \frac{N}{\sigma^2} + \frac{1}{v} \right) \theta^2 - 2 \left( \frac{m}{v} + \frac{1}{\sigma^2} \sum_{i=1}^N x_i \right) \theta \right] \right\} && \text{group } \theta \text{ terms} \\
 &= \exp \left\{ -\frac{1}{2} \left( \frac{N}{\sigma^2} + \frac{1}{v} \right) \left[ \theta^2 - 2 \left( \frac{N}{\sigma^2} + \frac{1}{v} \right)^{-1} \left( \frac{m}{v} + \frac{1}{\sigma^2} \sum_{i=1}^N x_i \right) \theta \right] \right\} && \text{complete square} \\
 &\propto \exp \left\{ -\frac{1}{2} \left( \frac{N}{\sigma^2} + \frac{1}{v} \right) \left[ \theta - \left( \frac{N}{\sigma^2} + \frac{1}{v} \right)^{-1} \left( \frac{m}{v} + \frac{1}{\sigma^2} \sum_{i=1}^N x_i \right) \right]^2 \right\} && \text{complete square}
 \end{aligned}$$

But if we view this last expression as a function of  $\theta$ , we see that it is proportional to the density of a normal distribution whose mean is

$$\begin{aligned}
 m^* &= \left( \frac{N}{\sigma^2} + \frac{1}{v} \right)^{-1} \left( \frac{m}{v} + \frac{1}{\sigma^2} \sum_{i=1}^N x_i \right) \\
 &= \frac{\sigma^2 m + v \sum_{i=1}^N x_i}{\sigma^2 + nv},
 \end{aligned}$$

and whose variance is

$$\begin{aligned}
 v^* &= \left( \frac{N}{\sigma^2} + \frac{1}{v} \right)^{-1} \\
 &= \frac{\sigma^2 v}{\sigma^2 + nv}.
 \end{aligned}$$

Therefore, we conclude that the posterior of  $\theta$ , given the data, is a [Normal](#)  $\left( \frac{\sigma^2 m + v \sum_{i=1}^N x_i}{\sigma^2 + nv}, \frac{\sigma^2 v}{\sigma^2 + nv} \right)$  distribution.

(d)

As in parts (a) and (c), we can write the posterior, up to a constant of proportionality, as the product of the likelihood of the data and the prior of the parameter(s). The prior for  $\omega$  is  $p(\omega) = \frac{b^a}{\Gamma(a)} \omega^{a-1} \exp(-b\omega)$ . The likelihood is

$$\begin{aligned}
 f(x_1, \dots, x_N | \omega) &= \prod_{i=1}^N f(x_i | \omega) && \text{i.i.d.} \\
 &= \prod_{i=1}^N \left[ \left( \frac{\omega}{2\pi} \right)^{1/2} \exp \left\{ -\frac{\omega}{2} (x_i - \theta)^2 \right\} \right] && \text{gamma density} \\
 &= \omega^{\frac{N}{2}} (2\pi)^{\frac{N}{2}} \exp \left\{ -\frac{\omega}{2} \sum_{i=1}^N (x_i - \theta)^2 \right\}. && \text{expand product}
 \end{aligned}$$

We then can compute the posterior.

$$\begin{aligned}
 p(\omega | x_1, \dots, x_N) &\propto f(x_1, \dots, x_N | \omega) p(\omega) && \text{Bayes'} \\
 &= \left[ \omega^{\frac{N}{2}} (2\pi)^{\frac{N}{2}} \exp \left\{ -\frac{\omega}{2} \sum_{i=1}^N (x_i - \theta)^2 \right\} \right] \left[ \frac{b^a}{\Gamma(a)} \omega^{a-1} \exp(-b\omega) \right] \\
 &\propto \omega^{\frac{N}{2}} \omega^{a-1} \exp \left\{ -\frac{\omega}{2} \sum_{i=1}^N (x_i - \theta)^2 \right\} \exp(-b\omega) && \text{drop initial constants} \\
 &= \omega^{\frac{N}{2} + a - 1} \exp \left\{ -\omega \left[ b + \frac{1}{2} \sum_{i=1}^N (x_i - \theta)^2 \right] \right\}. && \text{simplify}
 \end{aligned}$$

Viewing this last expression as a function of  $\omega$ , we recognize it as the density of a gamma distribution whose parameters are

$$\begin{aligned}
 a^* &= a + \frac{N}{2}, \text{ and} \\
 b^* &= b + \frac{1}{2} \sum_{i=1}^N (x_i - \theta)^2.
 \end{aligned}$$

Thus we see that the posterior of  $\omega$ , given the data, is a **Gamma**( $a + \frac{N}{2}, b + \frac{1}{2} \sum_{i=1}^N (x_i - \theta)^2$ ) distribution. Moreover, since  $\sigma^2 = \frac{1}{\omega}$ , we can say that the posterior of  $\sigma^2$ , given the data, is **IG**( $a + \frac{N}{2}, b + \frac{1}{2} \sum_{i=1}^N (x_i - \theta)^2$ ).

(e)

As before, we will determine the posterior by first computing the product of the prior and the likelihood. The prior for  $\theta$  is  $p(\theta) = \frac{1}{\sqrt{2\pi v}} \exp \left\{ -\frac{1}{2v}(\theta - m)^2 \right\}$ . The likelihood is

$$\begin{aligned}
 f(x_1, \dots, x_N | \theta) &= \prod_{i=1}^N f(x_i | \theta) && \text{i.i.d.} \\
 &= \prod_{i=1}^N \left[ \frac{1}{\sigma_i \sqrt{2\pi}} \exp \left\{ -\frac{1}{2\sigma_i^2} (x_i - \theta)^2 \right\} \right] && \text{normal density} \\
 &= \left( \prod_{i=1}^N \sigma_i \right) (2\pi)^{\frac{N}{2}} \exp \left\{ -\frac{1}{2} \sum_{i=1}^N \left( \frac{x_i - \theta}{\sigma_i} \right)^2 \right\}. && \text{expand product}
 \end{aligned}$$

The posterior is then

$$\begin{aligned}
 p(\theta | x_1, \dots, x_N) &\propto f(x_1, \dots, x_N | \theta) p(\theta) && \text{Bayes'} \\
 &= \left[ \left( \prod_{i=1}^N \sigma_i \right) (2\pi)^{\frac{N}{2}} \exp \left\{ -\frac{1}{2} \sum_{i=1}^N \left( \frac{x_i - \theta}{\sigma_i} \right)^2 \right\} \right] \left[ \frac{1}{\sqrt{2\pi v}} \exp \left\{ -\frac{1}{2v}(\theta - m)^2 \right\} \right] \\
 &\propto \exp \left\{ -\frac{1}{2} \sum_{i=1}^N \left( \frac{x_i - \theta}{\sigma_i} \right)^2 \right\} \exp \left\{ -\frac{1}{2v}(\theta - m)^2 \right\} && \text{drop constants} \\
 &= \exp \left\{ -\frac{1}{2} \sum_{i=1}^N \frac{x_i^2}{\sigma_i^2} - \frac{1}{2} \sum_{i=1}^N \left( -\frac{2x_i\theta}{\sigma_i^2} \right) - \frac{1}{2} \sum_{i=1}^N \frac{\theta^2}{\sigma_i^2} - \frac{1}{2} \left( \frac{\theta^2}{v} \right) - \frac{1}{2} \left( -\frac{2\theta m}{v} \right) - \frac{1}{2} \left( \frac{m^2}{v} \right) \right\} \\
 &\propto \exp \left\{ -\frac{1}{2} \sum_{i=1}^N \left( -\frac{2x_i\theta}{\sigma_i^2} \right) - \frac{1}{2} \sum_{i=1}^N \frac{\theta^2}{\sigma_i^2} - \frac{1}{2} \left( \frac{\theta^2}{v} \right) - \frac{1}{2} \left( -\frac{2\theta m}{v} \right) \right\} && \text{drop non-}\theta \\
 &= \exp \left\{ -\frac{1}{2} \left[ \left( \frac{1}{v} + \sum_{i=1}^N \frac{1}{\sigma_i^2} \right) \theta^2 - 2 \left( \frac{m}{v} + \sum_{i=1}^N \frac{x_i}{\sigma_i^2} \right) \theta \right] \right\} && \text{group } \theta \\
 &= \exp \left\{ -\frac{1}{2} \left( \frac{1}{v} + \sum_{i=1}^N \frac{1}{\sigma_i^2} \right) \left[ \theta^2 - 2 \left( \frac{1}{v} + \sum_{i=1}^N \frac{1}{\sigma_i^2} \right)^{-1} \left( \frac{m}{v} + \sum_{i=1}^N \frac{x_i}{\sigma_i^2} \right) \theta \right] \right\} && \text{complete square} \\
 &\propto \exp \left\{ -\frac{1}{2} \left( \frac{1}{v} + \sum_{i=1}^N \frac{1}{\sigma_i^2} \right) \left[ \theta - \left( \frac{1}{v} + \sum_{i=1}^N \frac{1}{\sigma_i^2} \right)^{-1} \left( \frac{m}{v} + \sum_{i=1}^N \frac{x_i}{\sigma_i^2} \right) \right]^2 \right\}. && \text{complete square}
 \end{aligned}$$

As a function of theta, this last expression is proportional to the density of a normal distribution with mean and variance

$$\begin{aligned}
 m^* &= \left( \frac{1}{v} + \sum_{i=1}^N \frac{1}{\sigma_i^2} \right)^{-1} \left( \frac{m}{v} + \sum_{i=1}^N \frac{x_i}{\sigma_i^2} \right) \\
 &= \frac{m + v \sum_{i=1}^N x_i \sigma_i^{-2}}{1 + v \sum_{i=1}^N \sigma_i^{-2}}, \text{ and} \\
 v^* &= \left( \frac{1}{v} + \sum_{i=1}^N \frac{1}{\sigma_i^2} \right)^{-1} \\
 &= \frac{v}{1 + v \sum_{i=1}^N \sigma_i^{-2}}.
 \end{aligned}$$

Therefore, we conclude that the posterior of  $\theta$ , given the data is a [Normal](#)  $\left( \frac{m + v \sum_{i=1}^N x_i \sigma_i^{-2}}{1 + v \sum_{i=1}^N \sigma_i^{-2}}, \frac{v}{1 + v \sum_{i=1}^N \sigma_i^{-2}} \right)$  [distribution](#).

(f)

We can first write the joint distribution of  $X$  and  $\Omega$ , using definition of the conditional distribution of  $x|\omega$ , which states that

$$f_{X|\Omega}(x|\omega) = \frac{f_{X,\Omega}(x,\omega)}{f_{\Omega}(\omega)}.$$

Rearranging this last equation gives us

$$\begin{aligned} f_{X,\Omega}(x,\omega) &= f_{X|\Omega}(x|\omega)f_{\Omega}(\omega) \\ &= \left[ \frac{\sqrt{\omega}}{\sqrt{2\pi}} \exp \left\{ -\frac{\omega}{2}(x-m)^2 \right\} \right] \left[ \frac{(\frac{b}{2})^{\frac{a}{2}}}{\Gamma(\frac{a}{2})} \omega^{\frac{a}{2}-1} \exp \left( -\omega \frac{b}{2} \right) \right]. \end{aligned}$$

Now that we have the joint distribution of  $X$  and  $\Omega$ , we can integrate out  $\Omega$  to obtain the marginal distribution of  $X$ . That is,

$$\begin{aligned} f_X(x) &= \int_0^\infty f_{X,\Omega}(x,\omega) d\omega \\ &= \int_0^\infty \left[ \frac{\sqrt{\omega}}{\sqrt{2\pi}} \exp \left\{ -\frac{\omega}{2}(x-m)^2 \right\} \right] \frac{(\frac{b}{2})^{\frac{a}{2}}}{\Gamma(\frac{a}{2})} \omega^{\frac{a}{2}-1} \exp \left( -\omega \frac{b}{2} \right) d\omega \\ &= \frac{(\frac{b}{2})^{\frac{a}{2}}}{\sqrt{2\pi}\Gamma(\frac{a}{2})} \int_0^\infty \left[ \omega^{\frac{a}{2}+\frac{1}{2}-1} \exp \left\{ -\omega \left( \frac{b}{2} + \frac{(x-m)^2}{2} \right) \right\} \right] d\omega \\ &= \frac{(\frac{b}{2})^{\frac{a}{2}}}{\sqrt{2\pi}\Gamma(\frac{a}{2})} \frac{\Gamma(\frac{a}{2} + \frac{1}{2})}{(\frac{b}{2} + \frac{(x-m)^2}{2})^{\frac{a}{2}+\frac{1}{2}}} \underbrace{\int_0^\infty \left[ \frac{(\frac{b}{2} + \frac{(x-m)^2}{2})^{\frac{a}{2}+\frac{1}{2}}}{\Gamma(\frac{a}{2} + \frac{1}{2})} \omega^{\frac{a}{2}+\frac{1}{2}-1} \exp \left\{ -\omega \left( \frac{b}{2} + \frac{(x-m)^2}{2} \right) \right\} \right] d\omega}_{\text{Gamma}\left(\frac{a}{2} + \frac{1}{2}, \frac{b}{2} + \frac{(x-m)^2}{2}\right) \text{ density}} \\ &= \frac{(\frac{b}{2})^{\frac{a}{2}}}{\sqrt{2\pi}\Gamma(\frac{a}{2})} \frac{\Gamma(\frac{a}{2} + \frac{1}{2})}{(\frac{b}{2} + \frac{(x-m)^2}{2})^{\frac{a}{2}+\frac{1}{2}}}. \end{aligned}$$

This last expression represents the marginal distribution of  $X$ . We will show by some algebraic manipulations that this density is that of a  $t$  distribution with a shape and scale parameter. In general, such a  $t$  distribution has a pdf given by the form

$$f(x) = \frac{\Gamma(\frac{\nu+1}{2})}{\sigma\sqrt{\nu\pi}\Gamma(\frac{\nu}{2})} \left[ 1 + \left( \frac{x-\mu}{\sigma\sqrt{\nu}} \right)^2 \right]^{-\frac{\nu+1}{2}},$$

for all  $x \in \mathbb{R}$ , where  $\mu \in \mathbb{R}$  is called the location parameter,  $\sigma > 0$  is the scale parameter, and  $\nu > 0$  is the degrees of freedom.



We rearrange the marginal density of  $X$  as follows, highlighting changes in **green**.

$$\begin{aligned}
f_X(x) &= \frac{\left(\frac{b}{2}\right)^{\frac{a}{2}}}{\sqrt{2\pi}\Gamma\left(\frac{a}{2}\right)} \frac{\Gamma\left(\frac{a}{2} + \frac{1}{2}\right)}{\left(\frac{b}{2} + \frac{(x-m)^2}{2}\right)^{\frac{a}{2} + \frac{1}{2}}} \\
&= \frac{\left(\frac{b}{2}\right)^{\frac{a}{2}}}{\sqrt{2\pi}\Gamma\left(\frac{a}{2}\right)} \frac{\Gamma\left(\frac{a}{2} + \frac{1}{2}\right)}{\left(\frac{b}{2}\right)^{\frac{a}{2} + \frac{1}{2}} \left(1 + \frac{(x-m)^2}{b}\right)^{\frac{a}{2} + \frac{1}{2}}} \\
&= \frac{\left(\frac{b}{2}\right)^{-\frac{1}{2}}}{\sqrt{2\pi}\Gamma\left(\frac{a}{2}\right)} \frac{\Gamma\left(\frac{a}{2} + \frac{1}{2}\right)}{\left(1 + \frac{(x-m)^2}{b}\right)^{\frac{a}{2} + \frac{1}{2}}} \\
&= \frac{\left(\frac{b}{2}\right)^{-\frac{1}{2}} \Gamma\left(\frac{a+1}{2}\right)}{\sqrt{2\pi}\Gamma\left(\frac{a}{2}\right)} \left[1 + \frac{(x-m)^2}{b}\right]^{-\frac{a+1}{2}} \\
&= \frac{\left(\frac{b}{2}\right)^{-\frac{1}{2}} \Gamma\left(\frac{a+1}{2}\right)}{\sqrt{2\pi}\Gamma\left(\frac{a}{2}\right)} \left[1 + \left(\frac{x-m}{\sqrt{b}}\right)^2\right]^{-\frac{a+1}{2}} \\
&= \frac{\left(\frac{b}{2}\right)^{-\frac{1}{2}} \Gamma\left(\frac{a+1}{2}\right)}{\sqrt{2\pi}\Gamma\left(\frac{a}{2}\right)} \left[1 + \left(\frac{x-m}{\sqrt{\frac{b}{a}}\sqrt{a}}\right)^2\right]^{-\frac{a+1}{2}} \\
&= \frac{\Gamma\left(\frac{a+1}{2}\right)}{\left(\frac{b}{2}\right)^{\frac{1}{2}} \sqrt{2\pi}\Gamma\left(\frac{a}{2}\right)} \left[1 + \left(\frac{x-m}{\sqrt{\frac{b}{a}}\sqrt{a}}\right)^2\right]^{-\frac{a+1}{2}} \\
&= \frac{\Gamma\left(\frac{a+1}{2}\right)}{\sqrt{\frac{b}{a}}\sqrt{a\pi}\Gamma\left(\frac{a}{2}\right)} \left[1 + \left(\frac{x-m}{\sqrt{\frac{b}{a}}\sqrt{a}}\right)^2\right]^{-\frac{a+1}{2}}.
\end{aligned}$$

But by comparing this last expression to the general form of the  $t$  distribution's pdf, we see that  $X$  must follow a  $t$  distribution with location parameter  $m$ , scale parameter  $\sqrt{\frac{b}{a}}$ , and  $a$  degrees of freedom.

## The multivariate normal distribution

(a)

(i)

$$\begin{aligned}
 Cov(\mathbf{x}) &= E[(\mathbf{x} - \boldsymbol{\mu})(\mathbf{x} - \boldsymbol{\mu})^T] \\
 &= E[(\mathbf{x} - \boldsymbol{\mu})(\mathbf{x}^T - \boldsymbol{\mu}^T)] && \text{linearity of transposition} \\
 &= E[\mathbf{x}\mathbf{x}^T - \mathbf{x}\boldsymbol{\mu}^T - \boldsymbol{\mu}\mathbf{x}^T + \boldsymbol{\mu}\boldsymbol{\mu}^T] && \text{distributive property of matrix multiplication} \\
 &= E[\mathbf{x}\mathbf{x}^T] - E[\mathbf{x}\boldsymbol{\mu}^T] - E[\boldsymbol{\mu}\mathbf{x}^T] + E[\boldsymbol{\mu}\boldsymbol{\mu}^T] && \text{linearity of expectation of random vectors} \\
 &= E[\mathbf{x}\mathbf{x}^T] - \boldsymbol{\mu}\boldsymbol{\mu}^T - \boldsymbol{\mu}\boldsymbol{\mu}^T + E[\boldsymbol{\mu}\boldsymbol{\mu}^T] && \text{linearity of expectation of random vectors} \\
 &= E[\mathbf{x}\mathbf{x}^T] - \boldsymbol{\mu}\boldsymbol{\mu}^T - \boldsymbol{\mu}\boldsymbol{\mu}^T + \boldsymbol{\mu}\boldsymbol{\mu}^T && \text{expectation of constants} \\
 &= E[\mathbf{x}\mathbf{x}^T] - \boldsymbol{\mu}\boldsymbol{\mu}^T && \text{simplify}
 \end{aligned}$$

(ii)

$$\begin{aligned}
 Cov(\mathbf{Ax} + \mathbf{b}) &= E[(\mathbf{Ax} + \mathbf{b})(\mathbf{Ax} + \mathbf{b})^T] - E[\mathbf{Ax} + \mathbf{b}](E[\mathbf{Ax} + \mathbf{b}])^T && \text{result from (i)} \\
 &= E[(\mathbf{Ax} + \mathbf{b})(\mathbf{x}^T \mathbf{A}^T + \mathbf{b}^T)] - E[\mathbf{Ax} + \mathbf{b}](E[\mathbf{Ax} + \mathbf{b}])^T && \text{transposition} \\
 &= E[\mathbf{Axx}^T \mathbf{A}^T + \mathbf{Axb}^T + \mathbf{bx}^T \mathbf{A}^T + \mathbf{bb}^T] - E[\mathbf{Ax} + \mathbf{b}](E[\mathbf{Ax} + \mathbf{b}])^T && \text{distribute} \\
 &= E[\mathbf{Axx}^T \mathbf{A}^T + \mathbf{Axb}^T + \mathbf{bx}^T \mathbf{A}^T + \mathbf{bb}^T] - (\mathbf{A}\boldsymbol{\mu} + \mathbf{b})(\mathbf{A}\boldsymbol{\mu} + \mathbf{b})^T && \text{linearity of expectation} \\
 &= E[\mathbf{Axx}^T \mathbf{A}^T + \mathbf{Axb}^T + \mathbf{bx}^T \mathbf{A}^T + \mathbf{bb}^T] - (\mathbf{A}\boldsymbol{\mu} + \mathbf{b})(\boldsymbol{\mu}^T \mathbf{A}^T + \mathbf{b}^T) && \text{transposition} \\
 &= E[\mathbf{Axx}^T \mathbf{A}^T + \mathbf{Axb}^T + \mathbf{bx}^T \mathbf{A}^T + \mathbf{bb}^T] - \mathbf{A}\boldsymbol{\mu}\boldsymbol{\mu}^T \mathbf{A}^T - \mathbf{A}\boldsymbol{\mu}\mathbf{b}^T - \mathbf{b}\boldsymbol{\mu}^T \mathbf{A}^T - \mathbf{bb}^T && \text{distribute} \\
 &= E[\mathbf{Axx}^T \mathbf{A}^T] + E[\mathbf{Axb}^T] + E[\mathbf{bx}^T \mathbf{A}^T] + E[\mathbf{bb}^T] - \mathbf{A}\boldsymbol{\mu}\boldsymbol{\mu}^T \mathbf{A}^T - \mathbf{A}\boldsymbol{\mu}\mathbf{b}^T - \mathbf{b}\boldsymbol{\mu}^T \mathbf{A}^T - \mathbf{bb}^T && \text{distribute} \\
 &= \mathbf{A}E[\mathbf{xx}^T] \mathbf{A}^T + \mathbf{A}\boldsymbol{\mu}\mathbf{b}^T + \mathbf{b}\boldsymbol{\mu}^T \mathbf{A}^T + \mathbf{bb}^T - \mathbf{A}\boldsymbol{\mu}\boldsymbol{\mu}^T \mathbf{A}^T - \mathbf{A}\boldsymbol{\mu}\mathbf{b}^T - \mathbf{b}\boldsymbol{\mu}^T \mathbf{A}^T - \mathbf{bb}^T && \text{expectation properties} \\
 &= \mathbf{A}E[\mathbf{xx}^T] \mathbf{A}^T - \mathbf{A}\boldsymbol{\mu}\boldsymbol{\mu}^T \mathbf{A}^T && \text{simplify} \\
 &= \mathbf{A}(E[\mathbf{xx}^T] - \boldsymbol{\mu}\boldsymbol{\mu}^T) \mathbf{A}^T && \text{factor} \\
 &= \mathbf{A}Cov(\mathbf{x}) \mathbf{A}^T && \text{result from (i)}
 \end{aligned}$$

(b)

The density of  $\mathbf{z}$  is the joint distribution of its individual random variables. That is  $f(\mathbf{z}) = f_{Z_1, \dots, Z_p}(z_1, \dots, z_p)$ . But since the  $Z_i$  are independent, their joint distribution is simply the product of their marginal distributions. That is,  $f_{Z_1, \dots, Z_p} = \prod_{i=1}^p f_{Z_i}(z_i)$ . Moreover, since the  $Z_i$  are all identically normal with mean 0 and variance 1, we have  $f_{Z_i}(z_i) = \frac{1}{\sqrt{2\pi}} \exp\left\{-\frac{1}{2}z_i^2\right\}$  for all  $i$ , so that

$$\begin{aligned}
 f(\mathbf{z}) &= \prod_{i=1}^p f_{Z_i}(z_i) && \text{indp.} \\
 &= \prod_{i=1}^p \frac{1}{\sqrt{2\pi}} \exp\left\{-\frac{1}{2}z_i^2\right\} && \text{identical normals} \\
 &= \frac{1}{(2\pi)^{\frac{p}{2}}} \exp\left\{-\frac{1}{2}z_1^2 - \dots - \frac{1}{2}z_p^2\right\} && \text{property of exp()} \\
 &= \frac{1}{(2\pi)^{\frac{p}{2}}} \exp\left\{-\frac{1}{2}(z_1^2 + \dots + z_p^2)\right\} && \text{simplify} \\
 &= \frac{1}{(2\pi)^{\frac{p}{2}}} \exp\left\{-\frac{1}{2}\mathbf{z}^T \mathbf{z}\right\}, && \text{inner product}
 \end{aligned}$$

where  $\mathbf{z} \in \mathbb{R}^p$ .

Next, using the definition of the mgf of a random vector, we have

$$\begin{aligned}
 m_{\mathbf{z}}(\mathbf{t}) &= E_{\mathbf{z}}[\exp\{\mathbf{t}^T \mathbf{z}\}] && \text{definition} \\
 &= E_{\mathbf{z}}\left[\exp\left\{\sum_{i=1}^p t_i Z_i\right\}\right] && \text{inner product} \\
 &= E_{\mathbf{z}}\left[\prod_{i=1}^p \exp\{t_i Z_i\}\right] && \text{property of exp()} \\
 &= \prod_{i=1}^p E_{Z_i}\left[\exp\{t_i Z_i\}\right] && \text{expectation of ind. r.v.'s} \\
 &= \prod_{i=1}^p \left[ \int_{\mathbb{R}} \left( \exp\{t_i z_i\} \frac{1}{\sqrt{2\pi}} \exp\left\{-\frac{1}{2}z_i^2\right\} \right) dz_i \right] && \text{def. of univar expectation} \\
 &= \prod_{i=1}^p \left[ \int_{\mathbb{R}} \frac{1}{\sqrt{2\pi}} \exp\left\{-\frac{1}{2}z_i^2 + t_i z_i\right\} dz_i \right] && \text{combine exp() terms} \\
 &= \prod_{i=1}^p \left[ \int_{\mathbb{R}} \frac{1}{\sqrt{2\pi}} \exp\left\{-\frac{1}{2}(z_i^2 - 2t_i z_i + t_i^2 - t_i^2)\right\} dz_i \right] && \text{complete square} \\
 &= \prod_{i=1}^p \left[ \exp\left\{\frac{1}{2}t_i^2\right\} \int_{\mathbb{R}} \underbrace{\frac{1}{\sqrt{2\pi}} \exp\left\{-\frac{1}{2}(z_i - t_i)^2\right\}}_{\text{density of } N(t_i, 1) \text{ dist'n}} dz_i \right] && \text{complete square, pull out const} \\
 &= \prod_{i=1}^p \exp\left\{\frac{1}{2}t_i^2\right\} && \text{integral of density} = 1 \\
 &= \exp\left\{\sum_{i=1}^p \frac{1}{2}t_i^2\right\} && \text{property of exp()} \\
 &= \exp\left\{\frac{1}{2}\mathbf{t}^T \mathbf{t}\right\}, && \text{inner product}
 \end{aligned}$$

for all  $\mathbf{t} \in \mathbb{R}^p$ .

(c)

From the definition beginning at the problem, we see that  $Y = \mathbf{t}^T \mathbf{x}$ , as a linear combination of univariate normals, must be a univariate normal distribution itself. Moreover, using our results from earlier, we can obtain the mean and variance of  $Y$  as follows:

$$\begin{aligned} E[Y] &= E[\mathbf{t}^T \mathbf{x}] \\ &= \mathbf{t}^T E[\mathbf{x}] \\ &= \mathbf{t}^T \boldsymbol{\mu}. \end{aligned} \quad \text{part (a)}$$

$$\begin{aligned} \text{Var}(Y) &= \text{Var}(\mathbf{t}^T \mathbf{x}) \\ &= \mathbf{t}^T \text{Cov}(\mathbf{x}) \mathbf{t} \\ &= \mathbf{t}^T \boldsymbol{\Sigma} \mathbf{t}. \end{aligned} \quad \text{part (a)}$$

But from the exercise instructions, we know that the mgf of  $Y$  must be  $m_Y(s) = \exp\{\mu t + \frac{s^2 \sigma^2}{2}\}$ , where  $\mu$  and  $\sigma^2$  are its mean and variance. Therefore,

$$\begin{aligned} m_Y(s) &= E[\exp\{sY\}] \\ &= \exp\left\{\mu s + \frac{s^2 \sigma^2}{2}\right\} && \text{univariate normal mgf} \\ &= \exp\left\{\mathbf{t}^T \boldsymbol{\mu} s + \frac{s^2 \mathbf{t}^T \boldsymbol{\Sigma} \mathbf{t}}{2}\right\}. && \text{result from above} \end{aligned}$$

Now, suppose that we evaluate  $m_Y(s)$  at  $s = 1$ . This would give us

$$\begin{aligned} m_Y(1) &= E[\exp\{(1)(Y)\}] \\ &= E[Y] \\ &= E[\mathbf{t}^T \mathbf{x}] \\ &= \exp\left\{\mathbf{t}^T \boldsymbol{\mu}(1) + \frac{(1)^2 \mathbf{t}^T \boldsymbol{\Sigma} \mathbf{t}}{2}\right\} \\ &= \exp\left\{\mathbf{t}^T \boldsymbol{\mu} + \frac{\mathbf{t}^T \boldsymbol{\Sigma} \mathbf{t}}{2}\right\}. \end{aligned}$$

Thus,  $E[\mathbf{t}^T \mathbf{x}] = \exp\{\mathbf{t}^T \boldsymbol{\mu} + \frac{1}{2} \mathbf{t}^T \boldsymbol{\Sigma} \mathbf{t}\}$ , as desired, and we conclude that this is the mgf of a multivariate normal  $\mathbf{x}$ . Furthermore, since the mgf of any distribution uniquely characterizes it, we know that if any other distribution has this same mgf, then it, too, must be multivariate normal. Hence, we can define  $\mathbf{x}$  as being multivariate normal if and only if its mgf is  $\exp\{\mathbf{t}^T \boldsymbol{\mu} + \frac{1}{2} \mathbf{t}^T \boldsymbol{\Sigma} \mathbf{t}\}$ .

(d)

First, we show that if  $\mathbf{x}$  is a multivariate normal distribution, and  $\mathbf{c}$  is a (additively conformable) vector of constants, then  $\mathbf{x} + \mathbf{c}$  is also multivariate normal. To see why, we use the mgf developed in the previous problem. Suppose  $\mathbf{x}$  has mean  $\boldsymbol{\mu}$  and covariance matrix  $\boldsymbol{\Sigma}$ . From part (c), we know that its mgf must be

$$m_{\mathbf{x}}(\mathbf{t}) = \exp \left\{ \mathbf{t}^T \boldsymbol{\mu} + \frac{1}{2} \mathbf{t}^T \boldsymbol{\Sigma} \mathbf{t} \right\}.$$

The mgf of  $\mathbf{x} + \mathbf{c}$  is then

$$\begin{aligned} m_{\mathbf{x}+\mathbf{c}}(\mathbf{t}) &= E[\exp \{ \mathbf{t}^T (\mathbf{x} + \mathbf{c}) \}] && \text{definition of mgf} \\ &= E[\exp \{ \mathbf{t}^T \mathbf{x} \} \exp \{ \mathbf{t}^T \mathbf{c} \}] && \text{property of exp()} \\ &= \exp \{ \mathbf{t}^T \mathbf{c} \} E[\exp \{ \mathbf{t}^T \mathbf{x} \}] && \text{linearity of E[]} \\ &= \exp \{ \mathbf{t}^T \mathbf{c} \} \exp \left\{ \mathbf{t}^T \boldsymbol{\mu} + \frac{1}{2} \mathbf{t}^T \boldsymbol{\Sigma} \mathbf{t} \right\} && \text{mgf of multivar normal} \\ &= \exp \left\{ \mathbf{t}^T (\mathbf{c} + \boldsymbol{\mu}) + \frac{1}{2} \mathbf{t}^T \boldsymbol{\Sigma} \mathbf{t} \right\} && \text{property of exp().} \end{aligned}$$

But from the results of part (c), we know that this means  $\mathbf{x}$  must be multivariate normal with mean  $\boldsymbol{\mu} + \mathbf{c}$  and covariance  $\boldsymbol{\Sigma}$ .

Now, suppose that  $\mathbf{x} = \mathbf{L}\mathbf{z} + \boldsymbol{\mu}$ , where  $\mathbf{z}$  is standard multivariate normal, and  $\boldsymbol{\mu}$  is a vector of constants. If we can show that  $\mathbf{a}^T \mathbf{x}$  is univariate normal for any nonzero  $\mathbf{a}$ , then the definition at the beginning of part (c) suggests that  $\mathbf{x}$  must be multivariate normal. We observe that  $\mathbf{a}^T \mathbf{x} = \mathbf{a}^T [\mathbf{L}\mathbf{z} + \boldsymbol{\mu}] = \mathbf{a}^T \mathbf{L}\mathbf{z} + \mathbf{a}^T \boldsymbol{\mu}$ . But since  $\mathbf{z}$  is multivariate normal, we know that  $\mathbf{a}^T \mathbf{L}\mathbf{z}$ , as a linear combination of the random vectors in  $\mathbf{z}$ , must be univariate normal. Also, we showed earlier in this problem that adding a constant to a normal distribution results in another normal distribution. (We showed this in general for multivariate normal distributions, of which univariate normals are special case.) Therefore,  $\mathbf{a}^T \mathbf{L}\mathbf{z} + \mathbf{a}^T \boldsymbol{\mu}$  is univariate normal, and we see that  $\mathbf{x}$  must be multivariate normal.

To find the mean and variance of  $\mathbf{x}$ , we use results from part (a). Recall that the mean of the standard multivariate normal distribution is  $\mathbf{0}$ , and that its variance is  $\mathbf{I}$ , the identity matrix.

$$\begin{aligned} E[\mathbf{x}] &= E[\mathbf{L}\mathbf{z} + \boldsymbol{\mu}] \\ &= \mathbf{L}E[\mathbf{z}] + \boldsymbol{\mu} && \text{property from (a)} \\ &= \mathbf{L}(\mathbf{0}) + \boldsymbol{\mu} && \text{mean of } \mathbf{z} \\ &= \boldsymbol{\mu}. \end{aligned}$$

$$\begin{aligned} \text{Var}[\mathbf{x}] &= \text{Var}[\mathbf{L}\mathbf{z} + \boldsymbol{\mu}] \\ &= \text{Var}[\mathbf{L}\mathbf{z}] && \text{variance unchanged by addition of const} \\ &= \mathbf{L}\text{Var}[\mathbf{z}]\mathbf{L}^T && \text{property from (b)} \\ &= \mathbf{L}\mathbf{I}\mathbf{L}^T && \text{variance of } \mathbf{z} \\ &= \mathbf{L}\mathbf{L}^T. \end{aligned}$$

Hence,  $\mathbf{x}$  is multivariate normal with mean  $\boldsymbol{\mu}$  and variance  $\mathbf{L}\mathbf{L}^T$ .

(e)

Let  $\mathbf{x}$  be multivariate normal with mean  $\boldsymbol{\mu}$  and covariance matrix  $\boldsymbol{\Sigma}$ . Then  $\boldsymbol{\Sigma}$ , like every covariance matrix, is positive definite, so that its eigen-decomposition can be written as  $\boldsymbol{\Sigma} = \mathbf{Q}\boldsymbol{\Lambda}\mathbf{Q}^T$ , where  $\mathbf{Q}$  is an orthogonal matrix containing the eigenvectors of  $\boldsymbol{\Sigma}$ , and  $\boldsymbol{\Lambda}$  is a diagonal matrix containing the eigenvalues of  $\boldsymbol{\Sigma}$ . Moreover, since  $\boldsymbol{\Sigma}$  is positive definite, all of its eigenvalues are positive. Hence, we can write  $\boldsymbol{\Lambda}$  as  $\boldsymbol{\Lambda} = \boldsymbol{\Lambda}^{\frac{1}{2}}\boldsymbol{\Lambda}^{\frac{1}{2}}$ , where  $\boldsymbol{\Lambda}^{\frac{1}{2}}$  is a diagonal matrix containing the square roots of the eigenvalues of  $\boldsymbol{\Sigma}$ . Hence,  $\boldsymbol{\Sigma} = \mathbf{Q}\boldsymbol{\Lambda}^{\frac{1}{2}}\boldsymbol{\Lambda}^{\frac{1}{2}}\mathbf{Q}^T$ .

Now, suppose that  $\mathbf{z}$  is the standard multivariate normal distribution. From the results of part (d), we know that the affine transformation  $\mathbf{Q}\boldsymbol{\Lambda}^{\frac{1}{2}}\mathbf{z} + \boldsymbol{\mu}$  must result is a multivariate normal distribution. We will show that this transformation results in  $\mathbf{x}$ , i.e., that  $\mathbf{x} = \mathbf{Q}\boldsymbol{\Lambda}^{\frac{1}{2}}\mathbf{z} + \boldsymbol{\mu}$ . From our earlier results, it is sufficient to show that the mean and variance of  $\mathbf{Q}\boldsymbol{\Lambda}^{\frac{1}{2}}\mathbf{z} + \boldsymbol{\mu}$  are the same as that of  $\mathbf{x}$  in order to establish that  $\mathbf{x} = \mathbf{Q}\boldsymbol{\Lambda}^{\frac{1}{2}}\mathbf{z} + \boldsymbol{\mu}$ . We verify this now.

$$\begin{aligned} E[\mathbf{Q}\boldsymbol{\Lambda}^{\frac{1}{2}}\mathbf{z} + \boldsymbol{\mu}] &= \mathbf{Q}\boldsymbol{\Lambda}^{\frac{1}{2}}E[\mathbf{z}] + \boldsymbol{\mu} \\ &= \mathbf{Q}\boldsymbol{\Lambda}^{\frac{1}{2}}\mathbf{0} + \boldsymbol{\mu} \\ &= \boldsymbol{\mu}. \end{aligned}$$

$$\begin{aligned} Var[\mathbf{Q}\boldsymbol{\Lambda}^{\frac{1}{2}}\mathbf{z} + \boldsymbol{\mu}] &= Var[\mathbf{Q}\boldsymbol{\Lambda}^{\frac{1}{2}}\mathbf{z}] \\ &= \mathbf{Q}\boldsymbol{\Lambda}^{\frac{1}{2}}Var[\mathbf{z}](\mathbf{Q}\boldsymbol{\Lambda}^{\frac{1}{2}})^T \\ &= \mathbf{Q}\boldsymbol{\Lambda}^{\frac{1}{2}}\mathbf{I}(\mathbf{Q}\boldsymbol{\Lambda}^{\frac{1}{2}})^T \\ &= \mathbf{Q}\boldsymbol{\Lambda}^{\frac{1}{2}}(\boldsymbol{\Lambda}^{\frac{1}{2}})^T\mathbf{Q}^T \\ &= \mathbf{Q}\boldsymbol{\Lambda}\mathbf{Q}^T \\ &= \boldsymbol{\Sigma}. \end{aligned}$$

Therefore,  $\mathbf{x}$  does indeed equal  $\mathbf{Q}\boldsymbol{\Lambda}^{\frac{1}{2}}\mathbf{z} + \boldsymbol{\mu}$ , so that we can indeed express any multivariate normal distribution  $\mathbf{x}$  as an affine transformation of the standard multivariate normal distribution.

(f)

Again, suppose that  $\mathbf{x}$  is multivariate normal with mean  $\boldsymbol{\mu}$  and covariance  $\boldsymbol{\Sigma}$ . We consider the transformation  $\mathbf{x} = \mathbf{Q}\boldsymbol{\Lambda}^{\frac{1}{2}}\mathbf{z} + \boldsymbol{\mu}$ , where  $\boldsymbol{\Sigma} = \mathbf{Q}\boldsymbol{\Lambda}^{\frac{1}{2}}\boldsymbol{\Lambda}^{\frac{1}{2}}\mathbf{Q}^T$  is the eigendecomposition of  $\boldsymbol{\Sigma}$ , and  $\mathbf{z}$  is the standard multivariate normal distribution. We will use the formula for pdf of transformation random variables to derive the pdf of  $\mathbf{x}$ . Recall that the pdf of  $\mathbf{z}$  is

$$f_{\mathbf{z}}(\mathbf{z}) = \frac{1}{(2\pi)^{\frac{p}{2}}} \exp \left\{ -\frac{1}{2} \mathbf{z}^T \mathbf{z} \right\},$$

for all  $\mathbf{z} \in \mathbb{R}^p$ . Next, observe that the inverse of the transformation, in which we solve for  $\mathbf{z}$  in terms of  $\mathbf{x}$  is given by  $\mathbf{z} = \boldsymbol{\Lambda}^{-\frac{1}{2}}\mathbf{Q}^{-1}(\mathbf{x} - \boldsymbol{\mu})$ . From the instructions, we know that the Jacobian of this transformation is then  $\mathbf{J} = \boldsymbol{\Lambda}^{-\frac{1}{2}}\mathbf{Q}^{-1}$ . Using the formula for transformation of random variables, we then have

$$\begin{aligned} f_{\mathbf{x}}(\mathbf{x}) &= f_{\mathbf{z}}(\mathbf{z})|\mathbf{J}| \\ &= f_{\mathbf{z}}(\boldsymbol{\Lambda}^{-\frac{1}{2}}\mathbf{Q}^{-1}(\mathbf{x} - \boldsymbol{\mu}))|\boldsymbol{\Lambda}^{-\frac{1}{2}}\mathbf{Q}^{-1}| \\ &= |\boldsymbol{\Lambda}^{-\frac{1}{2}}\mathbf{Q}^{-1}|(2\pi)^{-\frac{p}{2}} \exp \left\{ -\frac{1}{2} [\boldsymbol{\Lambda}^{-\frac{1}{2}}\mathbf{Q}^{-1}(\mathbf{x} - \boldsymbol{\mu})]^T [\boldsymbol{\Lambda}^{-\frac{1}{2}}\mathbf{Q}^{-1}(\mathbf{x} - \boldsymbol{\mu})] \right\} \\ &= |\boldsymbol{\Lambda}^{-\frac{1}{2}}\mathbf{Q}^{-1}|(2\pi)^{-\frac{p}{2}} \exp \left\{ -\frac{1}{2} (\mathbf{x} - \boldsymbol{\mu})^T (\mathbf{Q}^{-1})^T (\boldsymbol{\Lambda}^{-\frac{1}{2}})^T \boldsymbol{\Lambda}^{-\frac{1}{2}} \mathbf{Q}^{-1} (\mathbf{x} - \boldsymbol{\mu}) \right\} \\ &= |\boldsymbol{\Lambda}^{-\frac{1}{2}}\mathbf{Q}^{-1}|(2\pi)^{-\frac{p}{2}} \exp \left\{ -\frac{1}{2} (\mathbf{x} - \boldsymbol{\mu})^T (\mathbf{Q}^T)^{-1} (\boldsymbol{\Lambda}^T)^{-\frac{1}{2}} \boldsymbol{\Lambda}^{-\frac{1}{2}} \mathbf{Q}^{-1} (\mathbf{x} - \boldsymbol{\mu}) \right\} \\ &= |\boldsymbol{\Lambda}^{-\frac{1}{2}}\mathbf{Q}^{-1}|(2\pi)^{-\frac{p}{2}} \exp \left\{ -\frac{1}{2} (\mathbf{x} - \boldsymbol{\mu})^T (\mathbf{Q}^T)^{-1} \boldsymbol{\Lambda}^{-\frac{1}{2}} \boldsymbol{\Lambda}^{-\frac{1}{2}} \mathbf{Q}^{-1} (\mathbf{x} - \boldsymbol{\mu}) \right\} \\ &= |\boldsymbol{\Lambda}^{-\frac{1}{2}}\mathbf{Q}^{-1}|(2\pi)^{-\frac{p}{2}} \exp \left\{ -\frac{1}{2} (\mathbf{x} - \boldsymbol{\mu})^T \boldsymbol{\Sigma}^{-1} (\mathbf{x} - \boldsymbol{\mu}) \right\} \\ &= |\boldsymbol{\Lambda}^{-\frac{1}{2}}\mathbf{Q}^{-1}|(2\pi)^{-\frac{p}{2}} \exp \left\{ -\frac{1}{2} (\mathbf{x} - \boldsymbol{\mu})^T \boldsymbol{\Sigma}^{-1} (\mathbf{x} - \boldsymbol{\mu}) \right\} \\ &= |\boldsymbol{\Sigma}^{-1}|^{\frac{1}{2}} (2\pi)^{-\frac{p}{2}} \exp \left\{ -\frac{1}{2} (\mathbf{x} - \boldsymbol{\mu})^T \boldsymbol{\Sigma}^{-1} (\mathbf{x} - \boldsymbol{\mu}) \right\} \\ &= |\boldsymbol{\Sigma}|^{-\frac{1}{2}} (2\pi)^{-\frac{p}{2}} \exp \left\{ -\frac{1}{2} (\mathbf{x} - \boldsymbol{\mu})^T \boldsymbol{\Sigma}^{-1} (\mathbf{x} - \boldsymbol{\mu}) \right\} \\ &= \frac{1}{(2\pi)^{\frac{p}{2}} |\boldsymbol{\Sigma}|^{\frac{1}{2}}} \exp \left\{ -\frac{1}{2} (\mathbf{x} - \boldsymbol{\mu})^T \boldsymbol{\Sigma}^{-1} (\mathbf{x} - \boldsymbol{\mu}) \right\}. \end{aligned}$$

Therefore, for a multivariate random vector  $\mathbf{x}$  with mean  $\boldsymbol{\mu}$  and covariance  $\boldsymbol{\Sigma}$ , its given is given by the last expression.

(g)

We proceed by finding the mgf of  $\mathbf{Ax}_1 + \mathbf{Bx}_2$ . We know from our earlier results that both  $\mathbf{Ax}_1$  and  $\mathbf{Bx}_2$  are both multivariate normal. Also, using other previous results, we have that the mean and covariance of  $\mathbf{Ax}_1$  are  $\mathbf{A}\boldsymbol{\mu}_1$  and  $\mathbf{A}\boldsymbol{\Sigma}_1\mathbf{A}^T$ , respectively. Similarly, the mean and covariance of  $\mathbf{Bx}_2$  are  $\mathbf{B}\boldsymbol{\mu}_2$  and  $\mathbf{B}\boldsymbol{\Sigma}_2\mathbf{B}^T$ .

$$\begin{aligned} m_{\mathbf{Ax}_1 + \mathbf{Bx}_2}(\mathbf{t}) &= E[\exp \{\mathbf{t}^T (\mathbf{Ax}_1 + \mathbf{Bx}_2)\}] \\ &= E[\exp \{\mathbf{t}^T \mathbf{Ax}_1\} \exp \{\mathbf{t}^T \mathbf{Bx}_2\}] \\ &= E[\exp \{\mathbf{t}^T \mathbf{Ax}_1\}] E[\exp \{\mathbf{t}^T \mathbf{Bx}_2\}] \\ &= m_{\mathbf{Ax}_1}(\mathbf{t}) \cdot m_{\mathbf{Bx}_2}(\mathbf{t}) \\ &= \exp \left\{ \mathbf{t}^T \mathbf{A}\boldsymbol{\mu}_1 + \frac{1}{2} \mathbf{t}^T \mathbf{A}\boldsymbol{\Sigma}_1 \mathbf{A}^T \mathbf{t} \right\} \cdot \exp \left\{ \mathbf{t}^T \mathbf{B}\boldsymbol{\mu}_2 + \frac{1}{2} \mathbf{t}^T \mathbf{B}\boldsymbol{\Sigma}_2 \mathbf{B}^T \mathbf{t} \right\} \\ &= \exp \left\{ \mathbf{t}^T \mathbf{A}\boldsymbol{\mu}_1 + \frac{1}{2} \mathbf{t}^T \mathbf{A}\boldsymbol{\Sigma}_1 \mathbf{A}^T \mathbf{t} + \mathbf{t}^T \mathbf{B}\boldsymbol{\mu}_2 + \frac{1}{2} \mathbf{t}^T \mathbf{B}\boldsymbol{\Sigma}_2 \mathbf{B}^T \mathbf{t} \right\} \\ &= \exp \left\{ \mathbf{t}^T (\mathbf{A}\boldsymbol{\mu}_1 + \mathbf{B}\boldsymbol{\mu}_2) + \frac{1}{2} \mathbf{t}^T (\mathbf{A}\boldsymbol{\Sigma}_1 \mathbf{A}^T + \mathbf{B}\boldsymbol{\Sigma}_2 \mathbf{B}^T) \mathbf{t} \right\}, \end{aligned}$$

which we know from our earlier results means that  $\mathbf{Ax}_1 + \mathbf{Bx}_2$  must be multivariate normal with mean  $\mathbf{A}\boldsymbol{\mu}_1 + \mathbf{B}\boldsymbol{\mu}_2$  and covariance  $\mathbf{A}\boldsymbol{\Sigma}_1\mathbf{A}^T + \mathbf{B}\boldsymbol{\Sigma}_2\mathbf{B}^T$ .



## Conditionals and Marginals

(A)

Let  $\mathbf{L} = (\mathbf{I}_k \quad \mathbf{0})$ , so that  $\mathbf{x}_1 = \mathbf{L} \begin{pmatrix} \mathbf{x}_1 \\ \mathbf{x}_2 \end{pmatrix} = (\mathbf{I}_k \quad \mathbf{0}_{k \times (p-k)}) \begin{pmatrix} \mathbf{x}_1 \\ \mathbf{x}_2 \end{pmatrix} = \mathbf{x}_1$ . Then from our earlier results,  $\mathbf{x}_1$  is multivariate normal. Now, using more previous results,

$$\begin{aligned} E[\mathbf{x}_1] &= E \left[ \mathbf{L} \begin{pmatrix} \mathbf{x}_1 \\ \mathbf{x}_2 \end{pmatrix} \right] \\ &= \mathbf{L} E \left[ \begin{pmatrix} \mathbf{x}_1 \\ \mathbf{x}_2 \end{pmatrix} \right] \\ &= \mathbf{L} \begin{pmatrix} \boldsymbol{\mu}_1 \\ \boldsymbol{\mu}_2 \end{pmatrix} \\ &= (\mathbf{I}_k \quad \mathbf{0}) \begin{pmatrix} \boldsymbol{\mu}_1 \\ \boldsymbol{\mu}_2 \end{pmatrix} \\ &= \boldsymbol{\mu}_1. \end{aligned}$$

Also, using earlier results,

$$\begin{aligned} \text{Var}[\mathbf{x}_1] &= \text{Var} \left[ \mathbf{L} \begin{pmatrix} \mathbf{x}_1 \\ \mathbf{x}_2 \end{pmatrix} \right] \\ &= \mathbf{L} \text{Var} \left[ \begin{pmatrix} \mathbf{x}_1 \\ \mathbf{x}_2 \end{pmatrix} \right] \mathbf{L}^T \\ &= \mathbf{L} \boldsymbol{\Sigma} \mathbf{L}^T \\ &= (\mathbf{I}_k \quad \mathbf{0}) \begin{pmatrix} \boldsymbol{\Sigma}_{11} & \boldsymbol{\Sigma}_{12} \\ \boldsymbol{\Sigma}_{21} & \boldsymbol{\Sigma}_{22} \end{pmatrix} (\mathbf{I}_k \quad \mathbf{0})^T \\ &= (\mathbf{I}_k \quad \mathbf{0}) \begin{pmatrix} \boldsymbol{\Sigma}_{11} & \boldsymbol{\Sigma}_{12} \\ \boldsymbol{\Sigma}_{21} & \boldsymbol{\Sigma}_{22} \end{pmatrix} \begin{pmatrix} \mathbf{I}_k^T \\ \mathbf{0}^T \end{pmatrix} \\ &= (\mathbf{I}_k \quad \mathbf{0}) \begin{pmatrix} \boldsymbol{\Sigma}_{11} & \boldsymbol{\Sigma}_{12} \\ \boldsymbol{\Sigma}_{21} & \boldsymbol{\Sigma}_{22} \end{pmatrix} \begin{pmatrix} \mathbf{I}_k \\ \mathbf{0} \end{pmatrix} \\ &= (\boldsymbol{\Sigma}_{11} \quad \boldsymbol{\Sigma}_{12}) \begin{pmatrix} \mathbf{I}_k \\ \mathbf{0} \end{pmatrix} \\ &= \boldsymbol{\Sigma}_{11}. \end{aligned}$$

Therefore, we conclude that  $\mathbf{x}_1$  is multivariate normal with mean  $\boldsymbol{\mu}_1$  and covariance  $\boldsymbol{\Sigma}_{11}$ .

(B)

First, we verify that if  $\mathbf{M}_{(m+n) \times (m+n)} = \begin{pmatrix} \mathbf{A}_{m \times m} & \mathbf{B}_{n \times m} \\ \mathbf{C}_{m \times n} & \mathbf{D}_{n \times n} \end{pmatrix}$  is an appropriately partitioned matrix whose blocks  $\mathbf{A}$  and  $\mathbf{D}$  are invertible, then

$$\mathbf{M}^{-1} = \begin{pmatrix} (\mathbf{A} - \mathbf{B}\mathbf{D}^{-1}\mathbf{C})^{-1} & -(\mathbf{A} - \mathbf{B}\mathbf{D}^{-1}\mathbf{C})^{-1}\mathbf{B}\mathbf{D}^{-1} \\ -(\mathbf{D} - \mathbf{C}\mathbf{A}^{-1}\mathbf{B})^{-1}\mathbf{C}\mathbf{A}^{-1} & (\mathbf{D} - \mathbf{C}\mathbf{A}^{-1}\mathbf{B})^{-1} \end{pmatrix}.$$

To verify this, we use the fact that for any two square matrices  $\mathbf{A}_{n \times n}$  and  $\mathbf{B}_{n \times n}$ , then  $\mathbf{A}\mathbf{B} = \mathbf{I}$  if and only if  $\mathbf{A} = \mathbf{B}^{-1}$  and  $\mathbf{B} = \mathbf{A}^{-1}$ .

$$\begin{aligned} & \begin{pmatrix} (\mathbf{A} - \mathbf{B}\mathbf{D}^{-1}\mathbf{C})^{-1} & -(\mathbf{A} - \mathbf{B}\mathbf{D}^{-1}\mathbf{C})^{-1}\mathbf{B}\mathbf{D}^{-1} \\ -(\mathbf{D} - \mathbf{C}\mathbf{A}^{-1}\mathbf{B})^{-1}\mathbf{C}\mathbf{A}^{-1} & (\mathbf{D} - \mathbf{C}\mathbf{A}^{-1}\mathbf{B})^{-1} \end{pmatrix} \begin{pmatrix} \mathbf{A} & \mathbf{B} \\ \mathbf{C} & \mathbf{D} \end{pmatrix} = \\ & \begin{pmatrix} (\mathbf{A} - \mathbf{B}\mathbf{D}^{-1}\mathbf{C})^{-1}\mathbf{A} - (\mathbf{A} - \mathbf{B}\mathbf{D}^{-1}\mathbf{C})^{-1}\mathbf{B}\mathbf{D}^{-1}\mathbf{C} & (\mathbf{A} - \mathbf{B}\mathbf{D}^{-1}\mathbf{C})^{-1}\mathbf{B} - (\mathbf{A} - \mathbf{B}\mathbf{D}^{-1}\mathbf{C})^{-1}\mathbf{B}\mathbf{D}^{-1}\mathbf{D} \\ -(\mathbf{D} - \mathbf{C}\mathbf{A}^{-1}\mathbf{B})^{-1}\mathbf{C}\mathbf{A}^{-1}\mathbf{A} + (\mathbf{D} - \mathbf{C}\mathbf{A}^{-1}\mathbf{B})^{-1}\mathbf{C} & -(\mathbf{D} - \mathbf{C}\mathbf{A}^{-1}\mathbf{B})^{-1}\mathbf{C}\mathbf{A}^{-1}\mathbf{B} + (\mathbf{D} - \mathbf{C}\mathbf{A}^{-1}\mathbf{B})^{-1}\mathbf{D} \end{pmatrix} \\ & = \begin{pmatrix} (\mathbf{A} - \mathbf{B}\mathbf{D}^{-1}\mathbf{C})^{-1}(\mathbf{A} - \mathbf{B}\mathbf{D}^{-1}\mathbf{C}) & (\mathbf{A} - \mathbf{B}\mathbf{D}^{-1}\mathbf{C})^{-1}\mathbf{B} - (\mathbf{A} - \mathbf{B}\mathbf{D}^{-1}\mathbf{C})^{-1}\mathbf{B} \\ -(\mathbf{D} - \mathbf{C}\mathbf{A}^{-1}\mathbf{B})^{-1}\mathbf{C} + (\mathbf{D} - \mathbf{C}\mathbf{A}^{-1}\mathbf{B})^{-1}\mathbf{C} & (\mathbf{D} - \mathbf{C}\mathbf{A}^{-1}\mathbf{B})^{-1}(-\mathbf{C}\mathbf{A}^{-1}\mathbf{B} + \mathbf{D}) \end{pmatrix} \\ & = \begin{pmatrix} \mathbf{I}_{m \times m} & \mathbf{0}_{n \times m} \\ \mathbf{0}_{m \times n} & \mathbf{I}_{n \times n} \end{pmatrix} \\ & = \mathbf{I}_{(m+n) \times (m+n)}. \end{aligned}$$

We now apply this formula to the matrix  $\mathbf{\Sigma} = \begin{pmatrix} \mathbf{\Sigma}_{11} & \mathbf{\Sigma}_{12} \\ \mathbf{\Sigma}_{21} & \mathbf{\Sigma}_{22} \end{pmatrix}$ .

$$\mathbf{\Sigma}^{-1} = \begin{pmatrix} (\mathbf{\Sigma}_{11} - \mathbf{\Sigma}_{12}\mathbf{\Sigma}_{22}^{-1}\mathbf{\Sigma}_{21})^{-1} & -(\mathbf{\Sigma}_{11} - \mathbf{\Sigma}_{12}\mathbf{\Sigma}_{22}^{-1}\mathbf{\Sigma}_{21})^{-1}\mathbf{\Sigma}_{12}\mathbf{\Sigma}_{22}^{-1} \\ -(\mathbf{\Sigma}_{22} - \mathbf{\Sigma}_{21}\mathbf{\Sigma}_{11}^{-1}\mathbf{\Sigma}_{12})^{-1}\mathbf{\Sigma}_{21}\mathbf{\Sigma}_{11}^{-1} & (\mathbf{\Sigma}_{22} - \mathbf{\Sigma}_{21}\mathbf{\Sigma}_{11}^{-1}\mathbf{\Sigma}_{12})^{-1} \end{pmatrix},$$

so that

$$\mathbf{\Omega} = \mathbf{\Sigma}^{-1} = \begin{pmatrix} \mathbf{\Omega}_{11} & \mathbf{\Omega}_{12} \\ \mathbf{\Omega}_{21} & \mathbf{\Omega}_{22} \end{pmatrix} = \begin{pmatrix} (\mathbf{\Sigma}_{11} - \mathbf{\Sigma}_{12}\mathbf{\Sigma}_{22}^{-1}\mathbf{\Sigma}_{21})^{-1} & -(\mathbf{\Sigma}_{11} - \mathbf{\Sigma}_{12}\mathbf{\Sigma}_{22}^{-1}\mathbf{\Sigma}_{21})^{-1}\mathbf{\Sigma}_{12}\mathbf{\Sigma}_{22}^{-1} \\ -(\mathbf{\Sigma}_{22} - \mathbf{\Sigma}_{21}\mathbf{\Sigma}_{11}^{-1}\mathbf{\Sigma}_{12})^{-1}\mathbf{\Sigma}_{21}\mathbf{\Sigma}_{11}^{-1} & (\mathbf{\Sigma}_{22} - \mathbf{\Sigma}_{21}\mathbf{\Sigma}_{11}^{-1}\mathbf{\Sigma}_{12})^{-1} \end{pmatrix}.$$

(C)

The multivariate pdf in terms of partitioned matrices is

$$f(\mathbf{x}) = f\left(\begin{pmatrix} \mathbf{x}_1 \\ \mathbf{x}_2 \end{pmatrix}\right) = \frac{1}{(2\pi)^{\frac{p}{2}}|\boldsymbol{\Sigma}|^{\frac{1}{2}}} \exp\left\{-\frac{1}{2}\left[\begin{pmatrix} \mathbf{x}_1 \\ \mathbf{x}_2 \end{pmatrix} - \begin{pmatrix} \boldsymbol{\mu}_1 \\ \boldsymbol{\mu}_2 \end{pmatrix}\right]^T \begin{pmatrix} \boldsymbol{\Omega}_{11} & \boldsymbol{\Omega}_{12} \\ \boldsymbol{\Omega}_{21} & \boldsymbol{\Omega}_{22} \end{pmatrix} \left[\begin{pmatrix} \mathbf{x}_1 \\ \mathbf{x}_2 \end{pmatrix} - \begin{pmatrix} \boldsymbol{\mu}_1 \\ \boldsymbol{\mu}_2 \end{pmatrix}\right]\right\}.$$

The conditional distribution is proportional to the joint, whereby we treat  $\mathbf{x}_2$  as known.

$$\begin{aligned} f(\mathbf{x}_1|\mathbf{x}_2) &\propto f\left(\begin{pmatrix} \mathbf{x}_1 \\ \mathbf{x}_2 \end{pmatrix}\right) \\ &\propto \exp\left\{-\frac{1}{2}\left[\begin{pmatrix} \mathbf{x}_1 \\ \mathbf{x}_2 \end{pmatrix} - \begin{pmatrix} \boldsymbol{\mu}_1 \\ \boldsymbol{\mu}_2 \end{pmatrix}\right]^T \begin{pmatrix} \boldsymbol{\Omega}_{11} & \boldsymbol{\Omega}_{12} \\ \boldsymbol{\Omega}_{21} & \boldsymbol{\Omega}_{22} \end{pmatrix} \left[\begin{pmatrix} \mathbf{x}_1 \\ \mathbf{x}_2 \end{pmatrix} - \begin{pmatrix} \boldsymbol{\mu}_1 \\ \boldsymbol{\mu}_2 \end{pmatrix}\right]\right\} \\ &= \exp\left\{-\frac{1}{2}(\mathbf{x}_1^T - \boldsymbol{\mu}_1^T \quad \mathbf{x}_2^T - \boldsymbol{\mu}_2^T) \begin{pmatrix} \boldsymbol{\Omega}_{11} & \boldsymbol{\Omega}_{12} \\ \boldsymbol{\Omega}_{21} & \boldsymbol{\Omega}_{22} \end{pmatrix} \begin{pmatrix} \mathbf{x}_1 - \boldsymbol{\mu}_1 \\ \mathbf{x}_2 - \boldsymbol{\mu}_2 \end{pmatrix}\right\} \\ &= \exp\left\{-\frac{1}{2}(\mathbf{x}_1^T \boldsymbol{\Omega}_{11} - \boldsymbol{\mu}_1^T \boldsymbol{\Omega}_{11} + \mathbf{x}_2^T \boldsymbol{\Omega}_{21} - \boldsymbol{\mu}_2^T \boldsymbol{\Omega}_{21} \quad \mathbf{x}_1^T \boldsymbol{\Omega}_{12} - \boldsymbol{\mu}_1^T \boldsymbol{\Omega}_{12} + \mathbf{x}_2^T \boldsymbol{\Omega}_{22} - \boldsymbol{\mu}_2^T \boldsymbol{\Omega}_{22}) \begin{pmatrix} \mathbf{x}_1 - \boldsymbol{\mu}_1 \\ \mathbf{x}_2 - \boldsymbol{\mu}_2 \end{pmatrix}\right\} \\ &= \exp\left\{-\frac{1}{2}[(\mathbf{x}_1^T \boldsymbol{\Omega}_{11} - \boldsymbol{\mu}_1^T \boldsymbol{\Omega}_{11} + \mathbf{x}_2^T \boldsymbol{\Omega}_{21} - \boldsymbol{\mu}_2^T \boldsymbol{\Omega}_{21})(\mathbf{x}_1 - \boldsymbol{\mu}_1) + (\mathbf{x}_1^T \boldsymbol{\Omega}_{12} - \boldsymbol{\mu}_1^T \boldsymbol{\Omega}_{12} + \mathbf{x}_2^T \boldsymbol{\Omega}_{22} - \boldsymbol{\mu}_2^T \boldsymbol{\Omega}_{22})(\mathbf{x}_2 - \boldsymbol{\mu}_2)]\right\} \\ &\propto \exp\left\{-\frac{1}{2}[\mathbf{x}_1^T \boldsymbol{\Omega}_{11} \mathbf{x}_1 - \boldsymbol{\mu}_1^T \boldsymbol{\Omega}_{11} \mathbf{x}_1 + \mathbf{x}_2^T \boldsymbol{\Omega}_{21} \mathbf{x}_1 - \boldsymbol{\mu}_2^T \boldsymbol{\Omega}_{21} \mathbf{x}_1 - \mathbf{x}_1^T \boldsymbol{\Omega}_{11} \boldsymbol{\mu}_1 + \mathbf{x}_1^T \boldsymbol{\Omega}_{12} \mathbf{x}_2 - \mathbf{x}_1^T \boldsymbol{\Omega}_{12} \boldsymbol{\mu}_2]\right\} \\ &= \exp\left\{-\frac{1}{2}[\mathbf{x}_1^T \boldsymbol{\Omega}_{11} \mathbf{x}_1 - \boldsymbol{\mu}_1^T \boldsymbol{\Omega}_{11} \mathbf{x}_1 + \mathbf{x}_2^T \boldsymbol{\Omega}_{21} \mathbf{x}_1 - \boldsymbol{\mu}_2^T \boldsymbol{\Omega}_{21} \mathbf{x}_1 - \boldsymbol{\mu}_1^T \boldsymbol{\Omega}_{11} \mathbf{x}_1 + \mathbf{x}_2^T \boldsymbol{\Omega}_{21} \mathbf{x}_1 - \boldsymbol{\mu}_2^T \boldsymbol{\Omega}_{21} \mathbf{x}_1]\right\} \\ &= \exp\left\{-\frac{1}{2}[\mathbf{x}_1^T \boldsymbol{\Omega}_{11} \mathbf{x}_1 - 2(\boldsymbol{\mu}_1^T \boldsymbol{\Omega}_{11} - \mathbf{x}_2^T \boldsymbol{\Omega}_{21} + \boldsymbol{\mu}_2^T \boldsymbol{\Omega}_{21})\mathbf{x}_1]\right\} \\ &\propto \exp\left\{-\frac{1}{2}[(\mathbf{x}_1 - \boldsymbol{\Omega}_{11}^{-1}(\boldsymbol{\mu}_1^T \boldsymbol{\Omega}_{11} - \mathbf{x}_2^T \boldsymbol{\Omega}_{21} + \boldsymbol{\mu}_2^T \boldsymbol{\Omega}_{21}))^T \boldsymbol{\Omega}_{11} (\mathbf{x}_1 - \boldsymbol{\Omega}_{11}^{-1}(\boldsymbol{\mu}_1^T \boldsymbol{\Omega}_{11} - \mathbf{x}_2^T \boldsymbol{\Omega}_{21} + \boldsymbol{\mu}_2^T \boldsymbol{\Omega}_{21}))^T]\right\} \end{aligned}$$

$$\propto \exp \left\{ -\frac{1}{2} \left[ \left( \mathbf{x}_1 - \boldsymbol{\Omega}_{11}^{-1}(\boldsymbol{\Omega}_{11}\boldsymbol{\mu}_1 - \boldsymbol{\Omega}_{12}\mathbf{x}_2 + \boldsymbol{\Omega}_{12}\boldsymbol{\mu}_2) \right)^T \boldsymbol{\Omega}_{11} \left( \mathbf{x}_1 - \boldsymbol{\Omega}_{11}^{-1}(\boldsymbol{\Omega}_{11}\boldsymbol{\mu}_1 - \boldsymbol{\Omega}_{12}\mathbf{x}_2 + \boldsymbol{\Omega}_{12}\boldsymbol{\mu}_2) \right) \right] \right\},$$

which we recognize as the kernel of a multivariate normal density whose mean is

$$\boldsymbol{\Omega}_{11}^{-1}(\boldsymbol{\Omega}_{11}\boldsymbol{\mu}_1 - \boldsymbol{\Omega}_{12}\mathbf{x}_2 + \boldsymbol{\Omega}_{12}\boldsymbol{\mu}_2),$$

and whose variance is

$$\boldsymbol{\Sigma}^* = \boldsymbol{\Omega}_{11}^{-1}.$$

Simplifying the mean a little bit gives us

$$\begin{aligned} \boldsymbol{\mu}^* &= \boldsymbol{\Omega}_{11}^{-1}(\boldsymbol{\Omega}_{11}\boldsymbol{\mu}_1 - \boldsymbol{\Omega}_{12}\mathbf{x}_2 + \boldsymbol{\Omega}_{12}\boldsymbol{\mu}_2) \\ &= \boldsymbol{\mu}_1 - \boldsymbol{\Omega}_{11}^{-1}\boldsymbol{\Omega}_{12}\mathbf{x}_2 + \boldsymbol{\Omega}_{11}^{-1}\boldsymbol{\Omega}_{12}\boldsymbol{\mu}_2 \\ &= \boldsymbol{\mu}_1 - \boldsymbol{\Omega}_{11}^{-1}\boldsymbol{\Omega}_{12}(\mathbf{x}_2 - \boldsymbol{\mu}_2). \end{aligned}$$

But using our result from part (B), we can then rewrite the mean and variance in terms of blocks of  $\boldsymbol{\Sigma}$ .

$$\begin{aligned} \boldsymbol{\mu}^* &= \boldsymbol{\mu}_1 - \boldsymbol{\Omega}_{11}^{-1}\boldsymbol{\Omega}_{12}(\mathbf{x}_2 - \boldsymbol{\mu}_2) \\ &= \boldsymbol{\mu}_1 - (\boldsymbol{\Sigma}_{11} - \boldsymbol{\Sigma}_{12}\boldsymbol{\Sigma}_{22}^{-1}\boldsymbol{\Sigma}_{21})(-(\boldsymbol{\Sigma}_{11} - \boldsymbol{\Sigma}_{12}\boldsymbol{\Sigma}_{22}^{-1}\boldsymbol{\Sigma}_{21})^{-1}\boldsymbol{\Sigma}_{12}\boldsymbol{\Sigma}_{22}^{-1})(\mathbf{x}_2 - \boldsymbol{\mu}_2) \\ &= \boldsymbol{\mu}_1 + \boldsymbol{\Sigma}_{12}\boldsymbol{\Sigma}_{22}^{-1}(\mathbf{x}_2 - \boldsymbol{\mu}_2) \end{aligned}$$

$$\begin{aligned} \boldsymbol{\Sigma}^* &= \boldsymbol{\Omega}_{11}^{-1} \\ &= [(\boldsymbol{\Sigma}_{11} - \boldsymbol{\Sigma}_{12}\boldsymbol{\Sigma}_{22}^{-1}\boldsymbol{\Sigma}_{21})^{-1}]^{-1} \\ &= \boldsymbol{\Sigma}_{11} - \boldsymbol{\Sigma}_{12}\boldsymbol{\Sigma}_{22}^{-1}\boldsymbol{\Sigma}_{21}. \end{aligned}$$

We conclude then that the conditional distribution of  $\mathbf{x}_1|\mathbf{x}_2$  is

$$\mathcal{N}_p(\boldsymbol{\mu}_1 + \boldsymbol{\Sigma}_{12}\boldsymbol{\Sigma}_{22}^{-1}(\mathbf{x}_2 - \boldsymbol{\mu}_2), \boldsymbol{\Sigma}_{11} - \boldsymbol{\Sigma}_{12}\boldsymbol{\Sigma}_{22}^{-1}\boldsymbol{\Sigma}_{21}).$$

We can view this conditional relationship between  $\mathbf{x}_1$  and  $\mathbf{x}_2$  as analogous to a regression problem. In regression with Gaussian errors, when given the value of the covariates, we say that the predictor is normal with a certain mean and variance. We have the same sort of relationship here. Given the value of  $\mathbf{x}_2$ , we know that  $\mathbf{x}_1$  has a multivariate normal distribution specified as above.

## Multiple regression: three classical principles for inference

(A)

*Least squares:*

$$\begin{aligned}\text{LSS}(\boldsymbol{\beta}) &= (\mathbf{y} - \mathbf{X}\boldsymbol{\beta})^T (\mathbf{y} - \mathbf{X}\boldsymbol{\beta}) \\ &= (\mathbf{y}^T - \boldsymbol{\beta}^T \mathbf{X}^T) (\mathbf{y} - \mathbf{X}\boldsymbol{\beta}) \\ &= \mathbf{y}^T \mathbf{y} - \mathbf{y}^T \mathbf{X}\boldsymbol{\beta} - \boldsymbol{\beta}^T \mathbf{X}^T \mathbf{y} + \boldsymbol{\beta}^T \mathbf{X}^T \mathbf{X}\boldsymbol{\beta}\end{aligned}$$

$$\begin{aligned}\frac{\partial \text{LSS}}{\partial \boldsymbol{\beta}} &= -\mathbf{X}^T \mathbf{y} - \mathbf{X}^T \mathbf{y} + 2\mathbf{X}^T \mathbf{X}\boldsymbol{\beta} \\ &= -2\mathbf{X}^T \mathbf{y} + 2\mathbf{X}^T \mathbf{X}\boldsymbol{\beta}\end{aligned}$$

Setting this last equation to the zero vector and solving gives us

$$\begin{aligned}2\mathbf{X}^T \mathbf{X}\boldsymbol{\beta} &= 2\mathbf{X}^T \mathbf{y} \\ \boldsymbol{\beta} &= (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y}\end{aligned}$$

Therefore,  $\hat{\boldsymbol{\beta}}_{\text{OLS}} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y}$ .

*Maximum Likelihood:*

We write the likelihood function as

$$\begin{aligned}L(\boldsymbol{\beta}) &= \prod_{i=1}^n p(y_i | \boldsymbol{\beta}, \sigma^2) \\ &= \prod_{i=1}^n \left[ \frac{1}{\sigma\sqrt{2\pi}} \exp \left\{ -\frac{1}{2\sigma^2} (y_i - \mathbf{x}_i^T \boldsymbol{\beta})^2 \right\} \right] \\ &= \sigma^{-n} (2\pi)^{-\frac{n}{2}} \exp \left\{ -\frac{1}{2\sigma^2} \sum_{i=1}^n (y_i - \mathbf{x}_i^T \boldsymbol{\beta})^2 \right\} \\ &= \sigma^{-n} (2\pi)^{-\frac{n}{2}} \exp \left\{ -\frac{1}{2\sigma^2} (\mathbf{y} - \mathbf{X}\boldsymbol{\beta})^T (\mathbf{y} - \mathbf{X}\boldsymbol{\beta}) \right\}.\end{aligned}$$

Maximizing  $L(\boldsymbol{\beta})$  is equivalent to maximizing the log-likelihood:

$$\ell(\boldsymbol{\beta}) = \log L(\boldsymbol{\beta}) = -n \log \sigma - \frac{n}{2} \log 2\pi - \frac{1}{2\sigma^2} (\mathbf{y} - \mathbf{X}\boldsymbol{\beta})^T (\mathbf{y} - \mathbf{X}\boldsymbol{\beta})$$

But notice that  $\ell(\boldsymbol{\beta}) = c\text{LSS}(\boldsymbol{\beta}) + d$ , where  $c$  and  $d$  are negative constants. Thus, the maximum of  $L(\boldsymbol{\beta})$  and the minimum of  $\text{LSS}(\boldsymbol{\beta})$  coincide, so that  $\hat{\boldsymbol{\beta}}_{\text{MLE}} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y}$ .

### *Method of Moments:*

We consider the sample covariance between the first predictor and the residuals. In general, the sample covariance between these could be expressed as  $\frac{1}{n}(\sum x_{1i}e_i) - \bar{x}_1\bar{e}$ , where  $e_i = y_i - \mathbf{x}_i^T \boldsymbol{\beta}$ ; and  $\bar{x}_1 = \frac{1}{n}\sum x_{i1}$ ; and  $\bar{e} = \frac{1}{n}\sum e_i$ . However, we assume that  $\bar{x}_i = 0$  for all  $i = 1, \dots, p$ , so that the sample covariance between the first predictor and the residuals is merely  $\frac{1}{n}\sum x_{1i}e_i$ , which we simplify and then set equal to zero.

$$\begin{aligned}\frac{1}{n}\sum x_{1i}e_i &= \frac{1}{n}\sum x_{1i}(y_i - \mathbf{x}_i^T \boldsymbol{\beta}) \\ &= \frac{1}{n}\sum (x_{1i}y_i - x_{1i}\mathbf{x}_i^T \boldsymbol{\beta}) \\ &= \frac{1}{n}\mathbf{x}_{1.}^T \mathbf{y} - \frac{1}{n}\sum x_{i1}\mathbf{x}_i^T \boldsymbol{\beta} \\ &= 0.\end{aligned}$$

Thus, the first equation we can use to solve for  $\boldsymbol{\beta}$  is  $\mathbf{x}_{1.}^T \mathbf{y} - \mathbf{x}_{1.}^T \mathbf{X} \boldsymbol{\beta} = 0$ . We obtain similar equations for the other predictors, which gives us the following system of equations:

$$\begin{aligned}\mathbf{x}_{1.}^T \mathbf{y} - \mathbf{x}_{1.}^T \mathbf{X} \boldsymbol{\beta} &= 0 \\ \vdots \\ \mathbf{x}_{p.}^T \mathbf{y} - \mathbf{x}_{p.}^T \mathbf{X} \boldsymbol{\beta} &= 0\end{aligned}$$

We can represent the above system more compactly with the single matrix equation;

$$\mathbf{X}^T \mathbf{y} - \mathbf{X}^T \mathbf{X} \boldsymbol{\beta} = \mathbf{0}$$

Rearranging this last equation and solving for  $\boldsymbol{\beta}$ .

$$\begin{aligned}\mathbf{X}^T \mathbf{y} &= \mathbf{X}^T \mathbf{X} \boldsymbol{\beta} \\ \hat{\boldsymbol{\beta}}_{\text{MOM}} &= (\mathbf{X}^T \mathbf{X})^{-1} (\mathbf{X}^T \mathbf{y})\end{aligned}$$

We thus conclude that

$$\hat{\boldsymbol{\beta}}_{\text{OLS}} = \hat{\boldsymbol{\beta}}_{\text{MLE}} = \hat{\boldsymbol{\beta}}_{\text{MOM}} = (\mathbf{X}^T \mathbf{X})^{-1} (\mathbf{X}^T \mathbf{y}).$$

(B)

First, we write

$$\begin{aligned}\text{WLSS} &= \sum_{i=1}^n w_i (y_i - \mathbf{x}_i^T \boldsymbol{\beta})^2 \\ &= \sum_{i=1}^n (y_i - \mathbf{x}_i^T \boldsymbol{\beta}) w_i (y_i - \mathbf{x}_i^T \boldsymbol{\beta}) \\ &= (\mathbf{y} - \mathbf{X}\boldsymbol{\beta})^T \mathbf{W} (\mathbf{y} - \mathbf{X}\boldsymbol{\beta}) \\ &= (\mathbf{y}^T - \boldsymbol{\beta}^T \mathbf{X}^T) \mathbf{W} (\mathbf{y} - \mathbf{X}\boldsymbol{\beta}) \\ &= \mathbf{y}^T \mathbf{W} \mathbf{y} - \boldsymbol{\beta}^T \mathbf{X}^T \mathbf{W} \mathbf{y} - \mathbf{y}^T \mathbf{W} \mathbf{X} \boldsymbol{\beta} + \boldsymbol{\beta}^T \mathbf{X}^T \mathbf{W} \mathbf{X} \boldsymbol{\beta}\end{aligned}$$

where  $\mathbf{W} = \begin{pmatrix} w_1 & 0 & 0 \\ 0 & \ddots & 0 \\ 0 & 0 & w_n \end{pmatrix}$ . We then have

$$\frac{\partial \text{WLSS}}{\partial \boldsymbol{\beta}} = -2\mathbf{X}^T \mathbf{W} \mathbf{y} + 2\mathbf{X}^T \mathbf{W} \mathbf{X} \boldsymbol{\beta} = \mathbf{0}$$

$$\hat{\boldsymbol{\beta}}_{\text{WLS}} = (\mathbf{X}^T \mathbf{W} \mathbf{X})^{-1} (\mathbf{X}^T \mathbf{W} \mathbf{y}).$$

Next, with heteroskedastic errors, we can write the likelihood as

$$\begin{aligned}L &= \prod_{i=1}^n p(y_i | \boldsymbol{\beta}, \sigma_i^2) \\ &= \prod_{i=1}^n \frac{1}{\sigma_i \sqrt{2\pi}} \exp \left\{ -\frac{1}{2\sigma_i^2} (y_i - \mathbf{x}_i^T \boldsymbol{\beta})^2 \right\} \\ &= \sigma_i^{-n} (2\pi)^{-\frac{n}{2}} \exp \left\{ -\frac{1}{2} \sum_{i=1}^n \left( \frac{y_i - \mathbf{x}_i^T \boldsymbol{\beta}}{\sigma_i} \right)^2 \right\},\end{aligned}$$

so that the log-likelihood becomes

$$l = -n \log \sigma_i - \frac{n}{2} \log(2\pi) - \frac{1}{2} \sum_{i=1}^n \left( \frac{y_i - \mathbf{x}_i^T \boldsymbol{\beta}}{\sigma_i} \right)^2.$$

Maximizing this last equation is equivalent to minimizing  $\sum_{i=1}^n \left( \frac{y_i - \mathbf{x}_i^T \boldsymbol{\beta}}{\sigma_i} \right)^2$ . But suppose that we let  $w_i = \frac{1}{\sigma_i^2}$ , so that  $\sum_{i=1}^n \left( \frac{y_i - \mathbf{x}_i^T \boldsymbol{\beta}}{\sigma_i} \right)^2 = \sum_{i=1}^n w_i (y_i - \mathbf{x}_i^T \boldsymbol{\beta})^2$ . Then our objective function to minimize becomes precisely the same one that we minimized when considering the weighted least squares problem. Hence, the weighted least squares problem is the same as considering a model with heteroskedastic variances, where we simply let  $w_i = \frac{1}{\sigma_i^2}$ . This corresponds to assigning smaller weights to the observations with larger variances, so that we minimize their impact on the parameter estimates.



## Quantifying uncertainty: some basic frequentist ideas

(A)

As we showed before, the OLS, MLE, and MOM estimator for  $\boldsymbol{\beta}$  is  $\widehat{\boldsymbol{\beta}} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y}$ . Conditional on  $\mathbf{X}$ , we see that  $\mathbf{y}$  must be multivariate normal with mean  $\mathbf{X}\boldsymbol{\beta}$  and covariance matrix  $\sigma^2 \mathbf{I}$ , since it is the sum of a constant vector  $\mathbf{X}\boldsymbol{\beta}$  and a multivariate normal vector  $\boldsymbol{\epsilon} \sim \mathcal{N}(\mathbf{0}, \sigma^2 \mathbf{I})$ . Hence,  $\widehat{\boldsymbol{\beta}}$ , which is the product of a multivariate normal vector and the constant matrix  $(\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T$ , will also be multivariate normal. We need only find its expected value and covariance matrix to completely characterize it.

$$\begin{aligned} E[\widehat{\boldsymbol{\beta}}] &= E[(\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y}] \\ &= (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T E[\mathbf{y}] \\ &= (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T (\mathbf{X}\boldsymbol{\beta}) \\ &= [(\mathbf{X}^T \mathbf{X})^{-1} (\mathbf{X}^T \mathbf{X})] \boldsymbol{\beta} \\ &= \mathbf{I} \boldsymbol{\beta} \\ &= \boldsymbol{\beta}. \end{aligned}$$

$$\begin{aligned} \text{Var}[\widehat{\boldsymbol{\beta}}] &= \text{Var}[(\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y}] \\ &= (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \text{Var}(\mathbf{y}) [(\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T]^T \\ &= (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T (\sigma^2 \mathbf{I}) \mathbf{X} ((\mathbf{X}^T \mathbf{X})^{-1})^T \\ &= \sigma^2 (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{X} ((\mathbf{X}^T \mathbf{X})^{-1})^T \\ &= \sigma^2 \mathbf{I} ((\mathbf{X}^T \mathbf{X})^T)^{-1} \\ &= \sigma^2 (\mathbf{X}^T \mathbf{X})^{-1} \end{aligned}$$

Therefore, we conclude that, conditional on the data,  $\widehat{\boldsymbol{\beta}} \sim \mathcal{N}(\mathbf{X}\boldsymbol{\beta}, \sigma^2 (\mathbf{X}^T \mathbf{X})^{-1})$ .

(B)

We will use the maximum likelihood estimate of  $\sigma^2$  based on the model. Recall that the log-likelihood of the model is given by

$$l = -\frac{n}{2} \log \sigma^2 - \frac{n}{2} \log(2\pi) - \frac{1}{2} \sum_{i=1}^n \frac{1}{\sigma^2} (y_i - \mathbf{x}_i^T \boldsymbol{\beta})^2,$$

so that

$$\frac{\partial l}{\partial \sigma^2} = -\frac{n}{2\sigma^2} + \frac{\sum_{i=1}^n (y_i - \mathbf{x}_i^T \boldsymbol{\beta})^2}{2(\sigma^2)^2}.$$

Setting this last equation equal to zero and solving for  $\sigma^2$  gives us

$$\hat{\sigma}^2_{\text{MLE}} = \frac{1}{n} \sum (y_i - \mathbf{x}_i^T \hat{\boldsymbol{\beta}})^2,$$

where  $\hat{\boldsymbol{\beta}}$  represents the MLE estimate of  $\boldsymbol{\beta}$ . Recall that the covariance of  $\hat{\boldsymbol{\beta}}$  is given by  $\sigma^2(\mathbf{X}^T \mathbf{X})^{-1}$ . Therefore, to estimate the covariance matrix of  $\hat{\boldsymbol{\beta}}$ , we can simply multiply this estimate of  $\sigma^2$  by  $(\mathbf{X}^T \mathbf{X})^{-1}$ . That is,

$$\widehat{\text{Cov}}(\hat{\boldsymbol{\beta}}) = \hat{\sigma}^2_{\text{MLE}} (\mathbf{X}^T \mathbf{X})^{-1},$$

so that the square roots of the diagonal elements of this matrix will be our estimates for the standard errors of the  $\hat{\beta}_i$ 's. We will use this formula for  $\hat{\sigma}^2_{\text{MLE}}$  and  $\widehat{\text{Cov}}(\hat{\boldsymbol{\beta}})$  on the *ozone* data set and then compare them with the estimates produced by the `lm()` function in R.

Std. Error	Our estimate	lm() estimate
$\hat{\beta}_1$	$7.07 \cdot 10^{-3}$	$7.25 \cdot 10^{-3}$
$\hat{\beta}_2$	0.170	0.174
$\hat{\beta}_3$	0.0232	0.0238
$\hat{\beta}_4$	0.0678	0.0693
$\hat{\beta}_5$	0.123	0.125
$\hat{\beta}_6$	$3.85 \cdot 10^{-4}$	$3.94 \cdot 10^{-4}$
$\hat{\beta}_7$	0.0144	0.0148
$\hat{\beta}_8$	0.116	0.119
$\hat{\beta}_9$	0.00477	0.00490

Our estimates of the standard errors of the coefficients are very similar to the ones produced by the `lm()` function in R. The reason for the minor discrepancy is that we used the MLE, to obtain our estimate of  $\sigma^2$ , and the MLE is biased in this case. The unbiased estimate, which `lm()` computes, is given below:

$$\hat{\sigma}^2_{\text{Unb.}} = \frac{1}{n - p - 1} \sum (y_i - \mathbf{x}_i^T \hat{\boldsymbol{\beta}})^2,$$

where  $p$  is the number of predictors in the model.

(A)

For two parameters

$$\begin{aligned}\widehat{\text{SE}}\left(f(\hat{\boldsymbol{\theta}})\right) &= \sqrt{\widehat{\text{Var}}\left(f(\hat{\boldsymbol{\theta}})\right)} \\ &= \sqrt{\widehat{\text{Var}}(\hat{\theta}_1 + \hat{\theta}_2)} \\ &= \sqrt{\widehat{\text{Var}}(\hat{\theta}_1) + \widehat{\text{Var}}(\hat{\theta}_2) + 2\widehat{\text{Cov}}(\hat{\theta}_1, \hat{\theta}_2)} \\ &= \sqrt{\hat{\boldsymbol{\Sigma}}_{11} + \hat{\boldsymbol{\Sigma}}_{22} + 2\hat{\boldsymbol{\Sigma}}_{12}}.\end{aligned}$$

If we generalize this to  $p$  parameters, we then have:

$$\begin{aligned}\widehat{\text{SE}}\left(f(\hat{\boldsymbol{\theta}})\right) &= \sqrt{\widehat{\text{Var}}\left(f(\hat{\boldsymbol{\theta}})\right)} \\ &= \sqrt{\widehat{\text{Var}}(\hat{\theta}_1 + \hat{\theta}_2 + \cdots + \hat{\theta}_p)} \\ &= \sqrt{\sum_{i=1}^p \widehat{\text{Var}}(\hat{\theta}_i) + 2 \sum_{i=1}^p \sum_{i < j}^p \widehat{\text{Cov}}(\hat{\theta}_i, \hat{\theta}_j)} \\ &= \sqrt{\sum_{i=1}^p \hat{\boldsymbol{\Sigma}}_{ii} + 2 \sum_{i=1}^p \sum_{i < j}^p \hat{\boldsymbol{\Sigma}}_{ij}}.\end{aligned}$$

(B)

The first-order Taylor polynomial centered around  $\boldsymbol{\theta}$  gives us

$$f(\hat{\boldsymbol{\theta}}) = f(\boldsymbol{\theta}) + \sum_{i=1}^p g'_i(\boldsymbol{\theta})(\hat{\theta}_i - \theta_i) + R_n,$$

where  $g'_i(\boldsymbol{\theta}) = \frac{d}{d\hat{\theta}_i} g(\hat{\boldsymbol{\theta}})|_{\hat{\theta}_1=\theta_1, \dots, \hat{\theta}_p=\theta_p}$ , and where  $R_n$  is some remainder term that goes to zero as  $n \rightarrow \infty$  and hence can be omitted, giving us the approximation:

$$f(\hat{\boldsymbol{\theta}}) \approx f(\boldsymbol{\theta}) + \sum_{i=1}^p g'_i(\boldsymbol{\theta})(\hat{\theta}_i - \theta_i).$$

Now, assuming that  $\hat{\boldsymbol{\theta}}$  is an unbiased estimator of  $\boldsymbol{\theta}$ , i.e., that  $E[\hat{\boldsymbol{\theta}}] = \boldsymbol{\theta}$ , and taking the expectation of both sides of the previous equation, we have

$$\begin{aligned} E[f(\hat{\boldsymbol{\theta}})] &= E[f(\boldsymbol{\theta})] + \sum_{i=1}^p E[g'_i(\boldsymbol{\theta})(\hat{\theta}_i - \theta_i)] \\ &= f(\boldsymbol{\theta}) + g'_i(\boldsymbol{\theta}) \sum_{i=1}^p E[\hat{\theta}_i - \theta_i] \\ &= f(\boldsymbol{\theta}) + 0 \\ &= f(\boldsymbol{\theta}). \end{aligned}$$

We use these results to approximate the variance of  $f(\hat{\boldsymbol{\theta}})$ .

$$\begin{aligned} \text{Var}(\hat{\boldsymbol{\theta}}) &= E \left[ (f(\hat{\boldsymbol{\theta}}) - E[f(\hat{\boldsymbol{\theta}})])^2 \right] \\ &\approx E \left[ (f(\hat{\boldsymbol{\theta}}) - f(\boldsymbol{\theta}))^2 \right] \\ &\approx E \left[ \left( \sum_{i=1}^p g'_i(\boldsymbol{\theta})(\hat{\theta}_i - \theta_i) \right)^2 \right] \\ &= E \left[ \sum_{i=1}^p (g'_i(\boldsymbol{\theta}))^2 (\hat{\theta}_i - \theta_i)^2 + 2 \sum_{i=1}^p \sum_{i < j} g'_i(\boldsymbol{\theta}) g'_j(\boldsymbol{\theta}) (\hat{\theta}_i - \theta_i) (\hat{\theta}_j - \theta_j) \right] \end{aligned}$$

$$\begin{aligned}
&= \sum_{i=1}^p E \left[ (g'_i(\boldsymbol{\theta}))^2 (\hat{\theta}_i - \theta_i)^2 \right] + 2 \sum_{i=1}^p \sum_{i < j} E [g'_i(\boldsymbol{\theta}) g'_j(\boldsymbol{\theta}) (\hat{\theta}_i - \theta_i) (\hat{\theta}_j - \theta_j)] \\
&= \sum_{i=1}^p (g'_i(\boldsymbol{\theta}))^2 E \left[ (\hat{\theta}_i - \theta_i)^2 \right] + 2 \sum_{i=1}^p \sum_{i < j} g'_i(\boldsymbol{\theta}) g'_j(\boldsymbol{\theta}) E [(\hat{\theta}_i - \theta_i) (\hat{\theta}_j - \theta_j)] \\
&= \sum_{i=1}^p (g'_i(\boldsymbol{\theta}))^2 \text{Var}(\hat{\theta}_i) + 2 \sum_{i=1}^p \sum_{i < j} g'_i(\boldsymbol{\theta}) g'_j(\boldsymbol{\theta}) \text{Cov}(\hat{\theta}_i, \hat{\theta}_j).
\end{aligned}$$

Hence, the estimated variance of  $f(\hat{\boldsymbol{\theta}})$  can be approximated by

$$\widehat{\text{Var}}(\hat{\boldsymbol{\theta}}) = \sum_{i=1}^p (g'_i(\hat{\boldsymbol{\theta}}))^2 \hat{\boldsymbol{\Sigma}}_{ii} + 2 \sum_{i=1}^p \sum_{i < j} g'_i(\hat{\boldsymbol{\theta}}) g'_j(\hat{\boldsymbol{\theta}}) \hat{\boldsymbol{\Sigma}}_{ij}.$$

Since this is just an approximation, it certainly won't be perfect, especially since it's only a first-order approximation. We have also assumed that  $f$  is differentiable (and hence continuous) everywhere. Probably the more poorly behaved  $f$  is, the worse our estimation will be. Discontinuities in particular might cause major problems. Finally, we have assumed unbiasedness of our estimators, which may not always be a reasonable assumption.