

## Chapter 2 Appendices

## Appendix 2A. Distributional Results in the Central Case

## 2A.1. Distributions

Assume all distributions are central. We will take advantage of the well-known property of hierarchical (or nested) models that for  $k \geq k_*$  and  $L > 0$ ,

$$SSE_k - SSE_{k+L} \sim \sigma_*^2 \chi_L^2, \quad (2A.1)$$

$$SSE_k \sim \sigma_*^2 \chi_{n-k}^2 \quad (2A.2)$$

and

$$SSE_k - SSE_{k+L} \text{ is independent of } SSE_{k+L}. \quad (2A.3)$$

We will also need the distribution of linear combinations of  $SSE_k$  and  $SSE_{k+L}$ . Consider the linear combination  $aSSE_k - bSSE_{k+L}$  where  $a$  and  $b$  are scalars. It follows from Eq. (2A.2) that

$$E[aSSE_k - bSSE_{k+L}] = a(n-k)\sigma_*^2 - b(n-k-L)\sigma_*^2 \quad (2A.4)$$

and

$$\begin{aligned} \text{var}[aSSE_k - bSSE_{k+L}] &= \text{var}[aSSE_k - aSSE_{k+L} + aSSE_{k+L} - bSSE_{k+L}] \\ &= \text{var}[a(SSE_k - SSE_{k+L}) + (a-b)SSE_{k+L}]. \end{aligned}$$

Applying Eqs. (2A.1)–(2A.3), we have

$$\text{var}[aSSE_k - bSSE_{k+L}] = 2a^2 L \sigma_*^4 + 2(a-b)^2 (n-k-L) \sigma_*^4. \quad (2A.5)$$

Since many model selection criteria (such as AIC) use some function of  $\log(SSE_k)$ , we will also introduce some useful distributional results involving  $\log(SSE_k)$ . It can be shown (Gradshteyn, 1965, p. 576) that

$$\int_0^\infty z^{\nu-1} e^{-\mu z} \log(z) dz = \frac{1}{\mu^\nu} \Gamma(\nu) [\psi(\nu) - \log(\mu)], \quad \mu > 0, \nu > 0,$$

where

$$\psi(\nu) = -C - \sum_{j=0}^{\infty} \left( \frac{1}{j+\nu} - \frac{1}{j+1} \right),$$

$C = 0.577\,215\,664\,901$  is Euler's constant.  $Z \sim \chi_m^2$  with  $m$  degrees of freedom

$$E[\log(Z)] = \frac{1}{m} \sum_{j=1}^m \log(j) = \psi\left(\frac{m}{2}\right) + \log 2$$

which has no closed-form solution

$$E[\log(SSE_k)]$$

Although  $\psi$  has no closed-form expression, it is useful for computing exact values (see pp. 943–945):

$$\psi(v)$$

where

$$\psi\left(\frac{1}{2}\right) = -\gamma - \log 2$$

This recursion will be used in Section 2.3 as well as in Section 2.6.

The distribution of differences of  $SSE_k$  and  $SSE_{k+L}$  is involved. Since we have differences

$$\log(SSE_{k+L})$$

Let  $Q = SSE_{k+L}/SSE_k$ . According to Eq. (2A.3), we have

In nested models these two  $\chi^2$  variables are independent

$$Q \sim$$

Since  $Q = SSE_{k+L}/SSE_k$ , the

$$\log\left(\frac{SSE_{k+L}}{SSE_k}\right)$$

ral Case

$C = 0.577\ 215\ 664\ 901$  is Euler's constant, and  $\psi$  is Euler's psi function. For  $Z \sim \chi_m^2$  with  $m$  degrees of freedom,

$$E[\log(Z)] = \int_0^\infty \log(z) \frac{1}{2^{m/2} \Gamma(m/2)} z^{m/2-1} e^{-z/2} dz \\ = \log(2) + \psi\left(\frac{m}{2}\right),$$

(2A.1)

which has no closed-form solution. For  $SSE_k \sim \sigma_*^2 \chi_{n-k}^2$ ,

(2A.2)

$$E[\log(SSE_k)] = \log(\sigma_*^2) + \log(2) + \psi\left(\frac{n-k}{2}\right). \quad (2A.6)$$

(2A.3)

Although  $\psi$  has no closed-form solution, a simple recursion exists which is useful for computing exact expectations in small samples (Gradshteyn, 1965 pp. 943-945):

$$\psi(v+1) = \psi(v) + \frac{1}{v}, \quad v > 0, \quad (2A.7)$$

where

$$\psi\left(\frac{1}{2}\right) = -C - 2\log(2) \text{ and } \psi(1) = -C. \quad (2A.8)$$

This recursion will be used to check the accuracy of the Taylor expansion derived in Section 2.3 as well as in studying small-sample properties in Section 2.6.

The distribution of differences between  $\log(SSE_k)$  and  $\log(SSE_{k+L})$  is more involved. Since we have differences of logs,

$$\log(SSE_{k+L}) - \log(SSE_k) = \log\left(\frac{SSE_{k+L}}{SSE_k}\right).$$

Let  $Q = SSE_{k+L}/SSE_k$ . Assuming nested models and applying Eqs. (2A.1)-(2A.3), we have

$$Q \sim \frac{\chi_{n-k-L}^2}{\chi_{n-k-L}^2 + \chi_L^2}.$$

In nested models these two  $\chi^2$  are independent, and  $Q$  has the Beta distribution

$$Q \sim \text{Beta}\left(\frac{n-k-L}{2}, \frac{L}{2}\right).$$

Since  $Q = SSE_{k+L}/SSE_k$ , the log-distribution is

$$\log\left(\frac{SSE_{k+L}}{SSE_k}\right) \sim \log\text{-Beta}\left(\frac{n-k-L}{2}, \frac{L}{2}\right). \quad (2A.9)$$

It can be shown (Gradshteyn, 1965 p. 538 and p. 541) that

$$\int_0^1 t^{\mu-1}(1-t)^{\nu-1} \log(t) dt = \frac{1}{r^2} B\left(\frac{\mu}{r} + \nu, \nu\right) \left[ \psi\left(\frac{\mu}{r}\right) - \psi\left(\frac{\mu}{r} + \nu\right) \right]$$

and

$$\int_0^1 t^{\mu-1}(1-t)^{\nu-1} \log^2(t) dt = B\left(\frac{\mu}{r} + \nu, \nu\right) \left[ \left( \psi\left(\frac{\mu}{r}\right) - \psi\left(\frac{\mu}{r} + \nu\right) \right)^2 + \psi'(\mu) - \psi'(\mu + \nu) \right],$$

where

$$\psi'(v) = \sum_{j=0}^{\infty} \frac{1}{(j+v)^2}$$

and  $\mu > 0, \nu > 0$ . For  $Q \sim \text{Beta}(\frac{m}{2}, \frac{L}{2})$ ,

$$\begin{aligned} E[\log(Q)] &= \int_0^1 \log(t) B^{-1}\left(\frac{m}{2}, \frac{L}{2}\right) t^{\frac{m}{2}-1} (1-t)^{\frac{L}{2}-1} dt \\ &= \psi\left(\frac{m}{2}\right) - \psi\left(\frac{m}{2} + \frac{L}{2}\right) \end{aligned}$$

and

$$\begin{aligned} E[\log^2(Q)] &= \int_0^1 \log^2(t) B^{-1}\left(\frac{m}{2}, \frac{L}{2}\right) t^{\frac{m}{2}-1} (1-t)^{\frac{L}{2}-1} dt \\ &= \left( \psi\left(\frac{m}{2}\right) - \psi\left(\frac{m}{2} + \frac{L}{2}\right) \right)^2 + \psi'\left(\frac{m}{2}\right) - \psi'\left(\frac{m}{2} + \frac{L}{2}\right). \end{aligned}$$

Hence,

$$\text{var}[\log(Q)] = \psi'\left(\frac{m}{2}\right) - \psi'\left(\frac{m}{2} + \frac{L}{2}\right)$$

and

$$\text{var}\left[\log\left(\frac{\text{SSE}_{k+L}}{\text{SSE}_k}\right)\right] = \psi'\left(\frac{n-k-L}{2}\right) - \psi'\left(\frac{n-k}{2}\right). \quad (2A.10)$$

Eq. (2A.10) also has no closed-form, but again, convenient recursions for  $\psi'$  exist which are useful in small samples (Gradshteyn, 1965 pp. 945-946):

$$\psi'\left(\frac{m}{2}\right) = \psi'\left(\frac{m}{2} - 1\right) - \frac{4}{(m-2)^2}, \quad (2A.11)$$

where

## 2A.2. Appr

AIC, AI  
computing r  
of  $\log(\text{SSE}_k)$   
approximati  
Beginning w  
Taylor expan  
Suppose

and

Numerically  
paring the re  
good for  $E[\log(Q)]$   
mation with  
tion is used t  
in Section 2.3  
moments of l

## 2A.3. Appr

Since  $Q =$

However, ther  
(=  $\log(Q)$ ) in  
( $n-k-L$ )/( $n-k$ )

$\log(Q)$