

1. Adjusted profile log-likelihoods

As usual, we assume  $y = (y_1, \dots, y_n)$  have a distribution with density  $f(y; \theta)$ . The maximum likelihood estimator of  $\theta$  is asymptotically normal, and its asymptotic variance achieves the ‘information bound’; any other consistent estimator has an asymptotic variance greater than or equal to the inverse Fisher information. (In the case of vector parameters, this means the difference between the two matrices is non-negative definite.)

However, if we are particularly interested in a small number of components of  $\theta = (\psi, \lambda)$ , with the remainder treated as nuisance parameters, then the maximum likelihood estimator of these components, while ‘asymptotically optimal’ may have poor finite sample properties. We have

$$\ell_p(\hat{\psi}) = \ell(\hat{\psi}, \hat{\lambda}_{\hat{\psi}}) = \ell(\hat{\psi}, \hat{\lambda}),$$

i.e. we can maximize the full log likelihood by first finding the maximum at each fixed  $\psi$  and then maximizing over  $\psi$ . (It is worth double-checking this starting with the profile score equation  $\ell'_p(\psi) = 0$ . While you are at it, you can check that  $-\ell''_p(\hat{\psi}) = \{j^{\psi\psi}(\hat{\theta})\}^{-1}$ . But the profile log-likelihood  $\ell_p(\psi)$  does not make any adjustment for the estimation of  $\lambda$ , so can lead to a point estimator of  $\psi$  with poor finite sample properties, especially if the number of nuisance parameters is large relative to the sample size. If the number of nuisance parameters increases with the sample size, then the usual asymptotic theory does not apply, and  $\hat{\psi}$  may not be consistent, or it may be consistent, but not asymptotically efficient.

This observation, as well as the detailed formulas for higher order approximation, suggest that an adjustment to the profile likelihood might give better finite sample performance. The generic form for an adjusted profile log-likelihood is

$$\ell_A(\psi) = \ell_p(\psi) + A(\psi) = \ell(\psi, \hat{\lambda}_{\psi}) + A(\psi), \tag{1}$$

where  $A(\psi)$  is assumed to be  $O_p(1)$  (recall that  $\ell_p(\psi)$  is  $O_p(n)$ , or at least we are assuming so).

The most accurate version of adjusted profile likelihood uses  $A(\psi)$  in Fraser (2003):

$$A_{FR}(\psi) = +\frac{1}{2} \log |j_{\lambda\lambda}(\psi, \hat{\lambda}_{\psi})| - \log \left| \frac{d(\lambda)}{d\hat{\lambda}_{\psi}} \right|,$$

but the precise interpretation of  $d(\lambda)/d\hat{\lambda}$  needs some background on tangent exponential models. A closely related version given in Barndorff-Nielsen (1983)

is

$$A_{BN}(\psi) = -\frac{1}{2} \log |j_{\lambda\lambda}(\psi, \hat{\lambda}_\psi)| + \left| \log \frac{d\hat{\lambda}}{d\hat{\lambda}_\psi} \right|,$$

and you can find the derivation sketched in Davison (2003, Ch.12.4). However, again the last term is difficult to calculate in general. Furthermore, if we have set up our parameterization so that  $\lambda$  is orthogonal to  $\psi$  with respect to expected Fisher information, then  $\hat{\lambda}_\psi/\hat{\lambda} = 1 + O_p(n^{-1})$ , as you showed in Problems 2.2, so we might think of dropping the last term to get

$$A_{CR}(\psi) = -\frac{1}{2} \log |j_{\lambda\lambda}(\psi, \hat{\lambda}_\psi)|.$$

This last was suggested in Cox & Reid (1987) as an approximation to conditional likelihood, although ‘approximation to marginal likelihood’ is probably more accurate. A drawback of this version of adjusted log-likelihood is that it is not invariant to changes in parameterization of  $\lambda$ , so, for example, the adjusted log-likelihood for  $\psi$ , with nuisance parameter  $\sqrt{\lambda}$  is different than the adjusted log-likelihood for  $\psi$ , with nuisance parameter  $\lambda$ . Also, an orthogonal nuisance parameterization can only be found, in general, if the parameter of interest is scalar.

We can use any of these adjusted log-likelihoods in the same fashion as we use the profile log-likelihood, i.e. constructing log-likelihood ratio statistics, standardized maximum ‘likelihood’ estimators, and so on. So, for example, if we use

$$\ell_{CR}(\psi) = \ell_p(\psi) - \frac{1}{2} \log |j_{\lambda\lambda}(\psi, \hat{\lambda}_\psi)|,$$

assuming  $\psi \in \mathbb{R}$  and  $\lambda \perp \psi$ , then we would have

$$\begin{aligned} (\hat{\psi}_{CR} - \psi) \{-\ell''_{CR}(\hat{\psi}_{CR})\}^{-1/2} &\xrightarrow{d} N(0, 1), \\ \pm \sqrt{\{\ell_{CR}(\hat{\psi}_{CR}) - \ell_{CR}(\psi)\}} &\xrightarrow{d} N(0, 1). \end{aligned}$$

We don’t automatically get improved inference from the adjusted log-likelihoods, at least in the asymptotic theory, but in finite samples the inferences do seem to be better. Detailed discussion is given in Fraser (2003); see also Sartori (2003) for a discussion of Neyman-Scott problems, and DiCiccio et al. (1996).

## 2. Transformation models

### (a) The models

*Location models* A one parameter family of distributions on  $R$  is said to be a location family if the density function  $f(y; \theta)$  takes the form  $f_0(y - \theta)$ , for  $-\infty < \theta < \infty$ . The density  $f_0(y)$  is the standard form of the density and  $\theta$  is the location parameter.

*Examples* The normal distribution with mean  $\theta$  and known variance is a location family. The  $t_\nu$  distribution with density function given by

$$f(y; \theta) = c \left\{ 1 + \frac{(y - \theta)^2}{\nu} \right\}^{-(\nu+1)/2}$$

is also a location family. The standard form of the density is just the usual  $t_\nu$  density. The Cauchy distribution is a special case of this. The exponential location density is

$$f(y; \theta) = e^{-(y-\theta)}, \quad y - \theta \geq 0;$$

note that the support of the density is the interval  $(\theta, \infty)$ , although  $\theta$  can be any real value. Members of the same location family are simply shifted along the axis, relative to each other, and all have the same shape.

*Scale models* A one parameter family of distributions on  $R$  is said to be a scale family if the density function  $f(y; \theta)$  takes the form  $\theta^{-1} f_0(y/\theta)$ , for  $0 < \theta < \infty$ . The density  $f_0(y)$  is the standard form of the density and  $\theta$  is the scale parameter.

*Examples* The normal distribution with known mean and unknown variance is a scale family. The gamma distribution with known shape parameter is a scale family. A special case of this is the simple exponential distribution:

$$f(y; \theta) = \theta^{-1} \exp(-y/\theta); \quad y > 0.$$

Note that  $Z = \log(Y)$  has the density function

$$g(z; \eta) = \exp\{z - \eta - e^{(z-\eta)}\}; \quad -\infty < z < \infty$$

where  $\eta = \log \theta$ ; this is a location family.

*Location-scale models* A one parameter family of distributions on  $R$  is said to be a location-scale family if the density function  $f(y; \theta)$  takes the form  $\theta_2^{-1} f_0((y - \theta_1)/\theta_2)$ , for  $-\infty < \theta_1 < \infty, 0 < \theta_2 < \infty$ . The density  $f_0(y)$  is the standard form of the density,  $\theta_1$  is the location parameter and  $\theta_2$  is the scale parameter.

*Examples* The normal distribution with mean  $\theta_1$  and variance  $\theta_2^2$  is a location scale family, and the standard normal is the standard form. The  $t_\nu(\theta_1, \theta_2)$  is

$$f(t; \theta) = c \left\{ 1 + (y - \theta_1)^2 / \nu \theta_2^2 \right\}^{-(\nu+1)/2}$$

and the logistic( $\theta_1, \theta_2$ ) density is

$$f(y; \theta) = \frac{e^{-(y-\theta_1)/\theta_2} \{1 + e^{-(y-\theta_1)/\theta_2}\}^2}{\theta_2}$$

It is more conventional to use  $\mu$  and  $\sigma$  for the location and scale parameter, although they do not always correspond to the mean and variance of the distribution. Any continuous density on  $R$  can be embedded in a location-scale family, and in fact most location-scale families are constructed this way. It is easily proved that if the distribution of the random variable  $Y$  is a member of the location-scale family, then it can be expressed as

$$Y = \theta_2 Z + \theta_1$$

where  $Z$  has the standard distribution with density function  $f_0(z)$ .

Note that discrete distributions are not members of the location-scale family, essentially because the parameter space and the variable must take values on the same space.

*Transformation families* The location, scale and location-scale families are examples of transformation families. The basic idea is that a transformation on the sample space has a corresponding transformation on the parameter space that leaves the density function unchanged. For example, if  $Y$  has the density  $f_0(y - \theta)$ , then  $Z = Y + a$  has the density  $f_0(z - a - \theta) = f_0(z - \theta')$  so the family of densities for  $Y, \{f_0(y - \theta); \theta \in R\}$ , is the same as that for  $Z = Y + a$ . Similarly, the family of location-scale densities is unchanged under location and scale transformations: if  $f(y; \theta) = \theta_2^{-1} f_0((y - \theta_1)/\theta_2)$  then  $Z = cY + a$  has density  $f(z; \theta') = \theta_2'^{-1} f_0((z - \theta_1')/\theta_2')$  where  $\theta_1' = c\theta_1 + a$  and  $\theta_2' = c\theta_2$ .

In general, we denote by  $g$  a transformation on the sample space  $\mathcal{Y}$ , and by  $g^*$  the induced transformation on the parameter space. Then the formalization of the above is the statement that

$$\text{pr}(gY \in A; \theta) = \text{pr}(Y \in A; g^*\theta).$$

If  $g^*\Theta = \Theta$ , the family of densities indexed by  $\theta \in \Theta$  is invariant under the transformation  $g$  on  $\mathcal{Y}$ . Let  $\mathcal{C}$  be a class of transformations on  $\mathcal{Y}$  satisfying this condition, and  $\mathcal{G}$  the smallest class containing  $\mathcal{C}$  that is a group. Then  $\mathcal{G}^* = \{g^* \text{ induced by } g \in \mathcal{G}\}$  is a group on  $\Theta$ .

*Example: linear regression* A generalization of the location-scale model is the regression model

$$y = X\beta + \sigma\epsilon$$

where  $y$  is a vector of length  $n$ ,  $X$  is a known  $n \times p$  matrix,  $\beta$  is a vector of unknown parameters of length  $p$ ,  $\epsilon$  is a vector of length  $n$  that follows a known distribution  $f_0(\cdot)$ . If we let  $y^* = Xb + cy$ , where  $b \in R^p$  and  $c > 0$ , then we can write

$$y^* = X(b + c\beta) + c\sigma\epsilon$$

which is a member of the same family, as long as the parameter space is  $R^p \times R^+$ .

(b) Dimension reduction in transformation models

A key feature of transformation models is that the model for a sample of size  $n$  permits a reduction in dimension of the sufficient statistic. This reduction is obtained by conditioning. Thus, there exist functions of the data, say  $s(Y)$  and  $a(Y)$ , for which we can write

$$f(y; \theta) \propto f_1(s(y)|a(y); \theta) f_2(a(y))$$

where the marginal density of  $a(Y)$  does not depend on the parameter  $\theta$ . More importantly, the conditional density  $f_1$  is itself a transformation family density, with parameter  $\theta$  and sample space variable  $s(Y)$ .

*Location family* Suppose  $Y_1, \dots, Y_n$  are i.i.d. observations from the location family  $f_0(y - \theta)$ . Letting  $s(y) = y_n$  and  $a_i = y_i - y_n, i = 1, \dots, n - 1$ , we can write

$$f(y; \theta) = \prod f_0(y_i - \theta) = f_1(y_n|a; \theta) f_2(a) \quad (2)$$

where

$$f_2(a) = \int f_0(a_1 + t) \cdots f_0(a_{n-1} + t) f_0(t) dt \quad (3)$$

and

$$f_1(y_n|a; \theta) = \frac{f_0(y_n + a_1 - \theta) \cdots f_0(y_n + a_{n-1} - \theta) f_0(y_n - \theta)}{f_2(a)}. \quad (4)$$

The numerator is just a rewriting of  $\prod f(y_i; \theta)$ , and an expression equivalent to (3) is

$$f_1(y_n|a; \theta) = \frac{f_0(y_1 - \theta) \cdots f_0(y_n - \theta)}{\int f_0(y_1 - \theta) \cdots f_0(y_n - \theta) d\theta}. \quad (5)$$

The density of  $A$  does not depend on  $\theta$ , and the conditional density of  $Y_n$  given  $A$  is a location family density on  $R$ .

The functions  $s(Y)$  and  $a(Y)$  are not uniquely determined, but they are uniquely determined up to a location transformation. We could for example let  $S = \bar{Y}$  and  $A_i = Y_i - \bar{Y}$ . The vector  $A$  has  $n$  components but lies in  $R^{n-1}$  (as all its components must sum to 0, or in other words it is orthogonal to the 1-vector). The marginal density for  $a$  is again given by (2), and

$$f(y; \theta) \propto f(\bar{y}, a; \theta) = f_1(\bar{y}|a; \theta) f_2(a).$$

We might choose instead to let  $S$  be  $\hat{\theta}$ , the maximum likelihood estimate of  $\theta$ , and define  $A_i = Y_i - \hat{\theta}$ .

*Location-scale model* A version of  $S$  and  $A$  that can be used for the location-scale model is  $s(Y) = (\bar{Y}, s_Y)$ , where  $s_Y^2 = (n - 1)^{-1} \sum (Y_i - \bar{Y})^2$ ,

and  $a_i(Y) = (Y_i - \bar{Y})/s_Y$ ,  $i = 1, \dots, n$ . The vector  $A$  has  $n$  components, but is restricted to lie in  $R^{n-2}$ . To prove that the distribution of  $A$  is indeed free of  $\theta$ , we write

$$\begin{aligned} f(y; \theta) dy &= \theta_2^{-n} \prod \left\{ f_0 \left( \frac{y_i - \theta_1}{\theta_2} \right) \right\} dy_1 \dots dy_n \\ &= \theta_2^{-n} \prod f_0 \left( \frac{a_i s + \bar{y} - \theta_1}{\theta_2} \right) |J| da_1 \dots da_n d\bar{y} ds \end{aligned} \quad (6)$$

where  $|J|$  is the Jacobian of the transformation from  $y$  to  $(a, \bar{y}, s)$ . To compute  $f(a)$  we need to integrate out  $\bar{y}$  and  $s$  from this expression, so we need to figure out the dependence of  $|J|$  on  $\bar{y}$  and  $s$ . The computation is a little bit tricky, but by writing  $y_i = a_i s + \bar{y}$  we can see that  $dy_i = s da_i$ , and since  $a$  has  $n - 2$  free dimensions, the factor  $s^{n-2}$  will be part of the Jacobian. It turns out that this is the only part that depends on  $\bar{y}$  and  $s$ . The details are presented in the next section. The result is

$$\begin{aligned} f(a) da &= \int \int \frac{s^{n-2}}{\theta_2^n} f_0 \left( \frac{a_1 s + \bar{y} - \theta_1}{\theta_2} \right) \dots f_0 \left( \frac{a_n s + \bar{y} - \theta_1}{\theta_2} \right) d\bar{y} ds \\ &= \int \int \frac{(\theta_2 v)^{n-2}}{\theta_2^n} f_0 \left( \frac{a_1 \theta_2 v + \bar{y} - \theta_1}{\theta_2} \right) \dots f_0 \left( \frac{a_n \theta_2 v + \bar{y} - \theta_1}{\theta_2} \right) d\bar{y} \theta_2 dv \\ &= \int \int \frac{v^{n-2}}{\theta_2} f_0 \left( a_1 v + \frac{\bar{y} - \theta_1}{\theta_2} \right) \dots f_0 \left( a_n v + \frac{\bar{y} - \theta_1}{\theta_2} \right) d\bar{y} dv \\ &= \int \int v^{n-2} f_0(a_1 v + t) \dots f_0(a_n v + t) dt dv \end{aligned} \quad (7)$$

which shows that the marginal distribution of  $A$  does not depend on  $\theta$ . The conditional distribution of  $s(Y)$ , given  $A$ , is simply the ratio of the joint density to this marginal density. Again, the ancillary is not uniquely determined, but it is unique up to choice of location and scale variable. We could use  $(a', x_{(1)}, x_{(n)} - x_{(1)})$ , where  $a'_i = (x_{(i)} - x_{(1)}) / (x_{(n)} - x_{(1)})$ , instead of  $(a, \bar{x}, s)$ , or many other equivalent formulations.

(c) Details on the Jacobian:

As mentioned above, it is necessary to compute the Jacobian in the transformation from  $(y_1, \dots, y_n)$  to  $(\bar{y}, s, a)$ , where  $a = (a_1, \dots, a_n)$  and  $a_i = (y_i - \bar{y})/s$ . This will be done below both algebraically and geometrically. The algebraic derivation is due to Angelo Canty, and the geometric one is taken from Chapter 2 of *Inference and Linear Models*.

Since  $s^2 = \sum (y_i - \bar{y})^2$ , we can see that the vector  $a$ , although it has  $n$  components, in fact lies in  $R^{n-2}$ , because  $\sum a_i = a \cdot 1 = 0$  and  $\sum a_i^2 = \|a\|^2 = 1$ . To compute the Jacobian we make the transformation one-to-one by letting

$$t_1 = \bar{y}, \quad t_2 = s, \quad t_i = a_i \quad i = 3, \dots, n;$$

note that we are explicitly using only  $n - 2$  components of  $a$ . To find the inverse transformation we have

$$y_i = t_1 + t_2 t_i \quad i = 3, \dots, n$$

and using the restrictions on  $a$  we have

$$\begin{aligned} a_1 + a_2 &= -\sum_3^n t_i \\ 1 - a_1^2 - a_2^2 &= \sum_3^n t_i^2 \end{aligned}$$

from which we can write

$$\begin{aligned} a_1 &= f_1(t_3, \dots, t_n) = f_1(t_{(2)}), \text{ say} \\ a_2 &= f_2(t_3, \dots, t_n) = f_2(t_{(2)}), \text{ say} \\ y_1 &= \bar{y} + s \cdot f_1(t_3, \dots, t_n) = t_1 + t_2 f_1(t_{(2)}) \\ y_2 &= \bar{y} + s \cdot f_2(t_3, \dots, t_n) = t_1 + t_2 f_2(t_{(2)}) \end{aligned}$$

thus completing the transformation. We now have

$$\left| \frac{\partial y}{\partial t} \right| = \begin{vmatrix} 1 & f_1(t_{(2)}) & s g_{13}(t_{(2)}) & \dots & s g_{1n}(t_{(2)}) \\ 1 & f_2(t_{(2)}) & s g_{23}(t_{(2)}) & \dots & s g_{2n}(t_{(2)}) \\ 1 & t_3 & s & \dots & 0 \\ 1 & t_4 & 0 & s & \dots \\ \vdots & \vdots & \vdots & & \vdots \\ 1 & t_n & 0 & \dots & s \end{vmatrix}$$

where  $g_{ij}(t_{(2)}) = \partial f_i(t_{(2)}) / \partial t_j$ ,  $i = 1, 2$ ;  $j = 3, \dots, n$ . A well-known formula for the determinant of a partitioned matrix shows that the Jacobian is of the form

$$s^{n-2} h(t_{(2)}) = s^{n-2} h(a)$$

which is the result we were looking for.

In the above derivation the location and scale estimates were  $\bar{y}$  and  $s$ , with the result that  $a$  is orthogonal to the 1-vector, and has length 1. However, the same derivation applies for a variety of other location and scale estimates. For example, suppose we wanted to use the maximum likelihood estimates of  $\mu$  and  $\sigma$ , which are defined as the solutions to the equations  $\partial \log f(y; \hat{\mu}, \hat{\sigma}) / \partial \mu = 0$ ,  $\partial \log f(y; \hat{\mu}, \hat{\sigma}) / \partial \sigma = 0$  i.e.

$$\begin{aligned} \frac{-1}{\hat{\sigma}} \sum g' \left( \frac{y_i - \hat{\mu}}{\hat{\sigma}} \right) &= 0 \\ \frac{-n}{\hat{\sigma}} + \frac{y_i - \hat{\mu}}{\hat{\sigma}^2} \sum g' \left( \frac{y_i - \hat{\mu}}{\hat{\sigma}} \right) &= 0 \end{aligned}$$

where  $g(y_i) = \log f(y_i)$ . The ancillary statistic  $a$  is defined by  $a_i = (y_i - \hat{\mu})/\hat{\sigma}$ , so that  $y_i = \hat{\sigma}a_i + \hat{\mu}$ , and these two equations can be reexpressed as

$$\begin{aligned}\sum g'(a_i) &= 0 \\ \sum a_i g'(a_i) &= n\end{aligned}$$

which gives two restrictions on the  $a_i$ . Thus we can proceed as above and express  $(\hat{\mu}, \hat{\sigma}, a_1, \dots, a_n)$  as a one-to-one function of  $\bar{y}$ , and find the inverse transformation. It is of exactly the same form (with  $s$  replaced by  $\hat{\sigma}$ ), but the functions called  $f_1$  and  $f_2$  in the above derivation are different.

The geometric derivation of the result is actually very similar, but a little more elegant. It again uses  $\bar{y}$  and  $s$  as coordinates to get the result, and then argues that this choice of coordinates is arbitrary. Although it's not necessary, it's a little bit easier to first define  $z_i = (y_i - \mu)/\sigma$ ; we want to construct the conditional distribution of  $\bar{z}, s(\mathbf{z})$ , where  $\bar{z} = n^{-1} \sum z_i$ , and  $s^2(\mathbf{z}) = \sum (z_i - \bar{z})^2$ , given  $\mathbf{d}(\mathbf{z}) = (\mathbf{z} - \bar{z}\mathbf{1})/\|\mathbf{z} - \bar{z}\mathbf{1}\|$ . Note that in terms of the original variables  $\bar{z} = (\bar{y} - \mu)/\sigma$ ,  $s(\mathbf{z}) = s/\sigma$ , and  $\mathbf{d}(\mathbf{z}) = \mathbf{a}$ . (Bold font is used for vectors here to try to clarify the geometric argument.)

Now we compute the Jacobian of the transformation from  $\mathbf{z}$  to  $(\bar{z}, s(\mathbf{z}), \mathbf{d}(\mathbf{z}))$  by figuring out what the differential element  $d\mathbf{z}$  is in the new coordinates. That is, in the joint density of  $\mathbf{z}$ ,

$$f(\mathbf{z})d\mathbf{z} = \prod f(z_i)dz_i$$

we consider the differential element  $\prod dz_i$  as giving the volume of a small box at the point  $\mathbf{z}$ . We want to express this volume in the new coordinates. The coordinates  $(\bar{z}, s(\mathbf{z}), \mathbf{d}(\mathbf{z}))$  provide locally orthogonal coordinates at the point  $\mathbf{z} = \bar{z}\mathbf{1} + s(\mathbf{z})\mathbf{d}(\mathbf{z})$ , and we want to know how they change as we change to point  $\mathbf{z}$  to  $\mathbf{z} + d\mathbf{z}$ . The coordinate specified by  $\bar{z}$  lies on the  $\mathbf{1}$ -vector, so a small change in the coordinates  $\mathbf{z}$  cause a change in  $\bar{z}$  of  $\sqrt{n}d\bar{z}$ . Since  $s(\mathbf{z})$  measures the length of  $\mathbf{z} - \bar{z}\mathbf{1}$ , its rate of change is simply  $ds$ . (Think of the picture in  $R^2$ .) Now  $\mathbf{d}(\mathbf{z})$  is orthogonal to the  $\mathbf{1}$ -vector, and lies on a unit sphere in the  $n - 1$ -dimensional subspace of  $R^n$  that is orthogonal to the  $\mathbf{1}$ -vector. Thus the volume element is the surface volume on the sphere defined by  $s(\mathbf{z})\mathbf{d}(\mathbf{z})$ , i.e. the sphere of radius  $s(\mathbf{z})$ . This volume is  $s(\mathbf{z})^{n-2}du$ , where  $du$  is surface volume on the unit sphere in  $R^{n-1}$  (which is  $n - 2$ -dimensional). (In fact an explicit expression for the surface area of the unit sphere in  $R^d$  is  $(2\pi)^{d/2}/\Gamma(d/2)$ .) The coordinates in this development are orthogonal, so the volume element is the product of the three pieces. For this reason these coordinates are a convenient choice for computing the differential. However, if we choose to coordinatize the point using other location and scale functions,



the result is unchanged. The only thing we need to make sure of is that the location coordinate, say  $\tilde{\mu}(\mathbf{z})$ , satisfies the property  $\tilde{\mu}(a\mathbf{z}+b\mathbf{1}) = a\tilde{\mu}(\mathbf{z})+b$ , the scale coordinate, say  $\tilde{\sigma}(\mathbf{z})$ , satisfies  $\tilde{\sigma}(a\mathbf{z} + b\mathbf{1}) = a\tilde{\sigma}(\mathbf{z})$ , and  $\mathbf{d}$  is appropriately defined in terms of these two coordinates. We can show that such location and scale coordinates must themselves be location scale transformations of  $\bar{z}$  and  $s(\mathbf{z})$ , so that we can convert the above result to a more general one. (Although  $\hat{\mu}\mathbf{1}$  and  $\hat{\sigma}$  will not give orthogonal coordinates, so that in these two dimensions the ‘box’ is a parallelogram, we can still figure out the volume by multiplying the base by the height!) Using the more general coordinates will not provide an explicit expression for the normalizing constant in terms of the surface area on the sphere, because the new vector  $\mathbf{d}$  isn’t forced to lie in the orthogonal complement of the  $\mathbf{1}$ -vector.

This latter part of the argument is formalized in Chapter 2 of Fraser’s *The Structure of Inference*.

## References

### *Transformation models*

- [SM] Davison, A.C. (2003). *Statistical Models*. Cambridge University Press, Cambridge. Ch 5.3
- [BNC94] Barndorff-Nielsen, O.E. & Cox, D.R. (1994). *Inference and Asymptotics*. Chapman & Hall. Ch 2.8
- [PS] Pace, L. and Salvan, A. (1997). *Principles of Statistical Inference*. World Scientific, Singapore. Ch. 7.

### *Adjusted profile likelihoods*

- SM Ch 12.4, BNC94 Ch 8, PS Ch 4.3.
- Fraser, D.A.S. (2003). Likelihood for component parameters. *Biometrika* **90**, 327–339.
- Barndorff-Nielsen, O.E. (1983). On a formula for the distribution of the maximum likelihood estimator. *Biometrika* **70**, 343–65.
- DiCiccio T. J., Martin, M. A., Stern, S. E. & Young, G. A. (1996). Information bias and adjusted profile likelihoods. *J. R. Statist. Soc. B* **58**, 189–203.
- Cox, D.R. and Reid, N. (1987). Parameter orthogonality and approximate conditional inference. *J. R. Statist. Soc. B* **49**, 1–39.