# LARGE DEVIATION APPROXIMATIONS
# FOR MAXIMUM LIKELIHOOD ESTIMATORS*

BY

## I. M. SKOVGAARD (COPENHAGEN)

*Abstract.* A large deviation expansion of the density of a maximum likelihood estimator is derived in the case of replications from a multivariate curved subfamily of a continuous exponential family. Apart from an exponentially decreasing term, the approximation deviates only by a relative error of order $O(n^{-1})$ from the true density in a fixed neighbourhood of the true parameter value. An example is given which shows an excellent tail approximation even for small $n$. The results are specialized to the multidimensional nonlinear normal regression models and it is shown that, in these models, the approximation may be improved to deviate only by an exponentially decreasing error term.

**1. Introduction.** We shall derive an approximation to the density of the maximum likelihood estimator (MLE) of a vector parameter $\beta$ in the case of a smooth subfamily of an exponential family of continuous type. The expansion is a large deviation expansion in the sense that under simple replications the *relative error* of the approximation is $O(n^{-1})$ as $n \to \infty$, uniformly in a *fixed neighbourhood* of the true parameter value. This fact ensures a much better tail approximation to the distribution than that obtained by an Edgeworth expansion, where the density is only approximated up to a *fixed* (not relative) error over the whole range. This may be sufficient for large $n$ (or moderate $n$), but for small (or moderate) $n$ other approximations are needed.

The computational work required to derive the approximation is probably larger than that required to derive the first and may be the second term of the Edgeworth expansion, and also integration of the approximate density will usually have to be done numerically. However, in a very common class of models, namely the non-linear normal regression models, the result may be stated explicitly and is a very simple algebraic expression (see Section 6). Also, in any case, the complexity of the calculations is mainly

---

determined by the dimension of the parameter space and does not increase with increasing sample size. In fact, if $n$ is the number of replications, the approximation to the density $g(b)$ of $\hat{\beta}$ (the MLE) takes the form $\sqrt{n}\,g_1(b) \times$ $\times \exp\{-ng_2(b)\}$, where $g_1$ and $g_2$ are non-negative functions.

There are several ways of refining the approximation, some of which will be mentioned in the paper, but each, of course, at the prize of an increased amount of numerical work. We shall be concerned mainly with the simplest version.

The main idea of the approximations is Lemma 4.2, which gives an exact, although not directly computable, expression for the intensity of the process of local maxima of the likelihood function, together with an application of the saddlepoint approximation (see [5]) to this expression. The paper has been restricted to maximum likelihood estimators within curved exponential families; it will, however, be clear that the method may be applied to other estimators and other models.

The paper has been restricted to the derivation of the basic expansion with only a few remarks on the (rather obvious) applications. More efficient use of the expansion in the construction of critical regions and confidence regions is probably possible, but a discussion of these problems would be beyond the scope of this paper.

A few notational definitions, needed only for the multivariate algebra in Sections 4, 6 and 7, are given in Section 2. In Section 3 we review the basic method, deriving the expansion in the one-dimensional non-linear normal regression model without attention to mathematical rigour. In Section 4 we shall derive the approximation rigorously in the general (multivariate) curved exponential family model, including the basic proofs, but postponing technical proofs to the Appendix. Section 5 contains an example illustrating the behaviour of the approximation for small $n$. For large $n$ the behaviour can be deduced from the theorems on its asymptotic properties. In Section 6 we obtain the approximation for the important class of multivariate non-linear normal regression models, using the general results of Section 4. Finally, the Appendix contains the more technical proofs, whereas, as mentioned above, the conceptually important proofs are included in Section 4.

**2. Notation.** Most of the notation will be easily understood or explained, where it occurs. In Sections 4, 6 and 7 we shall, however, use a slightly generalized matrix notation. Vectors, matrices and 3-dimensional arrays of numbers are all regarded as matrices, e.g. $A = (a_{ijk})$, $i = 1, \ldots, n$; $j = 1, \ldots, m$; $k = 1, \ldots, m$, is an $(n \times m \times m)$-matrix. If $B = (b_{\varkappa\beta})$ is an $(m \times n)$-matrix, then $AB$ is the matrix product with respect to the last index of $A$ and the first index of $B$, i.e.

$$(AB)_{ij\beta} = \sum_{k=1}^{m} a_{ijk}\, b_{k\beta},$$

which is an $(n \times m \times n)$-matrix, etc. We shall sometimes emphasize the dimensions of a matrix by writing these as subscripts, e.g. $(a_{ijk})_{n \times m \times m}$ for $A$ or $(c_i)_n$ for a vector $c$ in $R^n$. To make the notation as conventional as possible, we shall still regard vectors as column-vectors (i.e. $c$ is an $(n \times 1)$-matrix) unless otherwise indicated, and write $c'$ for the transpose of $c$. Also $B'$ is the transpose of $B$.

If $f: R^m \to R^n$ is a differentiable function, we denote its differential by $Df$, i.e.

$$Df(x) = \left( \frac{df_i}{dx_j}(x) \right)_{n \times m}, \qquad x \in R^m,$$

and, similarly, $D^2 f(x)$ is the $(n \times m \times m)$-matrix of second partial derivatives at $x$, etc.

## 3. The one-dimensional non-linear normal regression model.

Let $X_i = \mu_i(\beta) + \varepsilon_i$, $\varepsilon_i \sim N(0, \sigma^2)$, $\beta \in B \subseteq R$, $B$ open, $\mu_i: B \to R$ twice continuously differentiable, $i = 1, \dots, k$ and $\varepsilon_1, \dots, \varepsilon_k$ mutually independent. $\beta$ is the unknown parameter; we shall consider $\sigma^2$ to be known, since this makes no difference in estimating $\beta$. We assume the existence of the maximum likelihood estimator $\hat{\beta}$ of $\beta$ and, to avoid technical details, we shall also assume that only one local maximum of the likelihood function can occur. Both of these assumptions will be relaxed in the next section.

Let $\beta_0 \in B$ be a fixed (true) value of the parameter, and let $g_0(b)$ denote the $\beta_0$-density of $\hat{\beta}$ at an arbitrary fixed point $b$. For the various functions of $\beta$, we shall use the convention that if the argument is omitted, $b$ is understood, whereas an index 0 means the function evaluated at $\beta_0$. Define

$$(3.1) \qquad D_j(\beta) = \frac{d^j}{d\beta^j} \log f(x; \beta), \quad j = 1, 2,$$

where $X = (X_1, \dots, X_k)$ and $f$ is the density of $X$, and let $Q_0$ be the $\beta_0$-distribution of $(D_1, D_2) = (D_1(b), D_2(b))$. Then a formal computation yields

$$(3.2) \quad g_0(b) = \lim_{\varepsilon \to 0} (2\varepsilon)^{-1} P_0 \{b - \varepsilon \leqslant \hat{\beta} \leqslant b + \varepsilon\}$$

$$= \lim_{\varepsilon \to 0} (2\varepsilon)^{-1} P_0 \{D_1(\beta) = 0 \text{ and } D_2(\beta) < 0 \text{ for some } \beta \text{ in } ]b-\varepsilon, b+\varepsilon[\}$$

$$= \lim_{\varepsilon \to 0} (2\varepsilon)^{-1} P_0 \{D_1 + D_2(\beta - b) = 0 \text{ and } D_2 < 0 \text{ for some } \beta$$
$$\text{in } ]b-\varepsilon, b+\varepsilon[\}$$

$$= \lim_{\varepsilon \to 0} (2\varepsilon)^{-1} \int_{-\infty}^{0} \int_{-\varepsilon|D_2|}^{\varepsilon|D_2|} dQ_0(d_1, d_2)$$

$$= h_0(0) E_0 \{|D_2| \cdot I_{\{D_2 < 0\}} | D_1 = 0\}$$

$$= h_0(0) e_0, \text{ say,}$$

where $h_0$ is the $\beta_0$-density of $D_1$, and $I_{\{\ldots\}}$ denotes the indicator function of the set $\{\ldots\}$. Notice that no approximations are involved in this computation; the general version will be given in Lemma 4.2 and its proof in the Appendix.

To compute (3.2) notice that the distribution of $(D_1, D_2)$ is bivariate normal with parameters

$$E_0\{D_1, D_2\} = (\gamma_1, \gamma_2), \quad V_0\{D_1\} = \sigma_{11}, \quad V_0\{D_2\} = \sigma_{22}, \quad V_0\{D_1, D_2\} = \sigma_{12}$$

given by

$$\gamma_1 = \big(\mu(\beta_0) - \mu(b)\big)'\left(\frac{d}{d\beta}\,\mu(b)\right)\big/\sigma^2,$$

$$\gamma_2 = \big(\mu(\beta_0) - \mu(b)\big)'\left(\frac{d^2}{d\beta^2}\,\mu(b)\right)\big/\sigma^2 - I(b),$$

(3.3)
$$\sigma_{11} = I(b) = \left(\frac{d}{d\beta}\,\mu(b)\right)'\left(\frac{d}{d\beta}\,\mu(b)\right)\big/\sigma^2,$$

$$\sigma_{22} = \left(\frac{d^2}{d\beta^2}\,\mu(b)\right)'\left(\frac{d^2}{d\beta^2}\,\mu(b)\right)\big/\sigma^2,$$

$$\sigma_{12} = \left(\frac{d}{d\beta}\,\mu(b)\right)'\left(\frac{d^2}{d\beta^2}\,\mu(b)\right)\big/\sigma^2,$$

where $\mu(\beta) = (\mu_1(\beta), \ldots, \mu_k(\beta))'$, and $I(\beta)$ is the Fisher information.

A direct computation now gives

(3.4)
$$h_0(0) = \big(2\pi I(b)\big)^{-1/2} \exp\left\{-\frac{1}{2}\gamma_1^2/I(b)\right\},$$

(3.5)
$$e_0 = \alpha\Phi(\alpha/\tau) + \tau\varphi(\alpha/\tau),$$

where $\alpha = -E_0\{D_2 | D_1 = 0\} = -\gamma_2 + \sigma_{12}\gamma_1/I(b)$, $\tau^2 = V_0\{D_2 | D_1\} = \sigma_{22} - \sigma_{12}^2/I(b)$, and $\Phi$, $\varphi$ are the standardized normal distribution and density function, respectively. Insertion in (3.2) now gives

(3.6)
$$g_0(b) = \big(2\pi I(b)\big)^{-1/2} \exp\left\{-\frac{1}{2}\gamma_1^2/I(b)\right\}\big(\alpha\Phi(\alpha/\tau) + \tau\varphi(\alpha/\tau)\big),$$

which is exact and hence solves our problem completely, since also its computation is feasible.

To illustrate the general approximation, we shall continue to approximate (3.6) by a simpler expression. We shall do this by expanding $e_0$ as a function of $\sigma^2$ as $\sigma^2 \to 0$, which is mathematically the same asymptotics as obtained by replications of the experiment. Using the expansion $\Phi(-x) \sim (\varphi(x)/x)(1-x^{-2})$ as $x \to +\infty$ we obtain $e_0 \sim \alpha(1 + o(\exp\{-c/\sigma^2\}))$ as $\sigma^2 \to 0$ if $\alpha > 0$ for some constant $c > 0$ depending on $\alpha$. Since $\alpha = I(\beta_0)$ if

$b = \beta_0$, $\alpha$ will by continuity be positive in a neighbourhood of $\beta_0$, which is independent of $\sigma^2$. Defining the approximation

(3.7)
$$\tilde{e}_0 = \begin{cases} \alpha, & \alpha > 0, \\ 0, & \text{otherwise,} \end{cases}$$

the relative error is $o(\exp\{-c/\sigma^2\})$ for some $c > 0$ uniformly in $b \in B_0$, where $B_0$ is some compact neighbourhood of $\beta_0$. Since the probability of $\hat{\beta}$ being outside $B_0$ is also decreasing at exponential rate as $\sigma^2 \to 0$, we shall not worry about that part of the approximation. As a final result we have

(3.8)
$$\tilde{g}_0(b) = h_0(0)\tilde{e}_0,$$

which satisfies

THEOREM 3.1. *There exist a constant $c > 0$ and a compact neighbourhood $B_0$ of $\beta_0$ such that*

(3.9)
$$\tilde{g}_0(b) = g_0(b)\big(1 + o(\exp\{-c/\sigma^2\})\big)$$

*uniformly in $b \in B_0$ as $\sigma^2 \to 0$, where $\tilde{g}_0$ is given by (3.4), (3.7) and (3.8).*

Proof. Follows from Theorem 4.8.

COROLLARY 3.2. *There exists a constant $c > 0$ such that*

(3.10)
$$0 \leqslant \int_A g_0(b)\,db - \int_A \tilde{g}_0(b)\,db = o(\exp\{-c/\sigma^2\})$$

*uniformly in all Borel sets $A \subseteq B$ as $\sigma^2 \to 0$, where $\tilde{g}_0$ is given by equalities (3.4), (3.7) and (3.8).*

Proof. The inequality follows from the trivial fact that $\tilde{e}_0 \leqslant e_0$. The second part follows from Theorem 3.1 and from the fact that $P_0\{\hat{\beta} \notin B_0\} = o(\exp\{-c_1/\sigma^2\})$ for some $c_1 > 0$, which is a consequence of the results in [4].

**4. Multivariate curved exponential families.** We shall now generalize the results of the previous section to multi-dimensional subfamilies of arbitrary exponential families of continuous type. In Section 3 we essentially used the normality *via* the normality of $(D_1, D_2)$ only. The idea in the general case is instead to use a saddlepoint approximation to the distribution of $(D_1, D_2)$ or rather a simple version of the mixed Edgeworth — saddlepoint approximation (see [3]). The results will, of course, be somewhat more complicated than in the normal case.

Let us first shortly review the saddlepoint approximations; for more thorough accounts on this type of approximations, see [3], [5] and [6], XVI.7.

Let $S = X_1 + \ldots + X_n$ be a sum of $n$ i.i.d. random vectors with density $f$ on $R^p$ and let $f^{*n}$ denote its $n$-th convolution with itself, i.e. the density of $S$.

Define

$$(4.1) \qquad g_t(x) \doteq f(x) e^{t'x}/\varphi(t), \qquad t \in \mathbf{R}^p,$$

where $\varphi(t) = \int e^{t'x} f(x)\, dx$ is the Laplace-transform of $X_1$; we shall not at the moment worry about its domain. Now,

$$(4.2) \qquad g_t^{*n}(s) = f^{*n}(s) e^{t's}/\varphi(t)^n$$

and, for a particular $s$, we may choose $\tilde{t}$ such that $g_{\tilde{t}}$ is "centered" at $s/n$, i.e.

$$(4.3) \qquad D \log \varphi(\tilde{t}) = \int x g_{\tilde{t}}(x)\, dx = s/n.$$

Applying the central limit theorem to $g_{\tilde{t}}^{*n}$, the saddlepoint approximation to $f^{*n}(s)$ follows from (4.2). Letting $h_n$ denote the density of $U = S/n$, this yields

$$(4.4) \quad h_n(u) = (n/2\pi)^{p/2} |\beta(\tilde{t})|^{-1/2} \exp\{-\tilde{t}'un\} \, \varphi(\tilde{t})^n (1 + O(n^{-1})) \quad \text{as } n \to \infty,$$

where $|\beta(t)|$ is the determinant of

$$\beta(t) = D^2 \log \varphi(t) = \left( \frac{d^2}{dt_i\, dt_j} \log \varphi(t) \right)_{p \times p}.$$

The remaining part of this section deals with the following setup. Let $X_1, \ldots, X_n$ be i.i.d. random vectors in $\mathbf{R}^k$, the density of $X_1$ with respect to some measure $\mu$ on $\mathbf{R}^k$ being

$$(4.5) \qquad f(x; \beta) = \exp\{x' \theta(\beta) - \psi(\theta(\beta))\},$$

where $\beta \in B \subseteq \mathbf{R}^p$, $B$ open; $\theta: B \to \mathbf{R}^k$ has a range satisfying

$$\theta(\beta) \in \operatorname{int} \Theta = \operatorname{int}\{\theta \in \mathbf{R}^k \mid \psi(\theta) = \int \exp\{x'\theta\} \mu(dx) < \infty\}$$

for all $\beta \in B$. Further, let $X = (X_1, \ldots, X_n)$, $f_n$ the density of $X$ and $\bar{X} = \sum X_i/n$. As in the previous section we define

$$D_1(\beta) = n^{-1}(D \log f_n(X; \beta))_p = (D\theta(\beta))'(\bar{X} - E_\beta\{X_1\}) \in \mathbf{R}^p,$$

$$D_2(\beta) = n^{-1}(D^2 \log f_n(X; \beta))_{p \times p}$$

$$= -I(\beta)/n + (\bar{X} - E_\beta\{X_1\})'_k (D^2 \theta(\beta))_{k \times p \times p},$$

where $I(\beta) = n(D\theta(\beta))'_{p \times k} (D^2 \psi(\theta(\beta)))_{k \times k} (D\theta(\beta))_{k \times p}$ is the Fisher information matrix. For later reference we shall need the following assumptions:

ASSUMPTIONS 4.1. (i) *$\bar{X}$ has a continuous density with respect to the Lebesgue measure on the closed convex support of $X_1$.*

(ii) *$\theta: B \to \mathbf{R}^k$ is one-to-one, bicontinuous and three times differentiable on $B$.*

(iii) *The Fisher information matrix $I(\beta)$ is regular for all $\beta \in B$.*

Let $\beta_0 \in B$ and $b \in B$ be fixed points, and let us again use the convention

that if an argument is omitted, $b$ is understood, whereas a subscript 0 means the value at $\beta_0$.

Due to the assumption in Section 3 that only one local maximum of the likelihood function could occur, the limit in (3.2) was equal to the density of $\hat\beta$ at $b$. In general this is not so, but we shall still consider the same quantity

$$(4.6) \qquad \lambda_0(b) = \lim_{\varepsilon \to 0} (\varepsilon^p A_p)^{-1} P_0 \{ f_n(X; \beta) \text{ has a local}$$
$$\text{maximum within } \|\beta - b\| < \varepsilon \},$$

where $A_p = \text{vol} \{ \beta \in R^p | \|\beta\| < 1 \}$. Thus $\lambda_0(b)$ is, when it exists, *the intensity of the point process of local maxima of the likelihood function*. Obviously, we have

$$(4.7) \qquad\qquad\qquad g_0(b) \leqslant \lambda_0(b), \quad b \in B,$$

when the MLE is formally defined as some external point, $\infty$ say, if the likelihood function has no maximum. The following lemma now generalizes equality (3.2):

LEMMA 4.2. *If Assumptions 4.1. are fulfilled, then*

$$(4.8) \qquad\qquad\qquad \lambda_0(b) = h_0(0) e_0,$$

*where $h_0$ is the $\beta_0$-density of $D_1$ and*

$$(4.9) \qquad\qquad e_0 = E_0 \{ |-D_2| \cdot I_{\{D_2 \text{ is neg. definite}\}} | D_1 = 0 \}.$$

Proof. See Appendix.

Remark 4.3. It is clear from the proof that Lemma 4.2 is not restricted to exponential families, but we have chosen this framework, since in these cases the results are fairly simple and almost no extra conditions are needed to prove their validity.

Our next step is to approximate $h_0(0)$ using the saddlepoint approximation. Since the expectation of $D_1$ is usually different from zero, the outcome $D_1 = 0$ is a "large deviation" and the usual normal approximation or the Edgeworth expansions could not be expected to give useful results.

LEMMA 4.4. *Let Assumptions 4.1 be fulfilled, and let $B_0 \subseteq B$ be any compact neighbourhood of $\beta_0$. Then*

$$(4.10) \qquad h_0(0) = (n/2\pi)^{p/2} |\beta(\tilde t)|^{-1/2} \times$$
$$\times \exp \{ -n [\psi(\theta_0) - \psi(\tilde\theta) + \tau(\theta)'(D\theta)\tilde t] \} (1 + O(n^{-1})) \quad \text{as } n \to \infty$$

*uniformly in $b \in B_0$, where $\tau(\theta) = D\psi(\theta)$, $\tilde t$ is the unique solution to*

$$(4.11) \qquad\qquad (\tau(\tilde\theta) - \tau(\theta))'(D\theta) = 0, \quad \tilde\theta = \theta_0 + (D\theta)\tilde t \in \Theta$$

*and $\beta(t)$ is given in (4.14).*

Proof. If $n = 1$, the Laplace transform of $D_1$ and its first two logarithmic derivatives are

$$(4.12) \quad \varphi(t) = E_0\{\exp(t' D_1)\}$$
$$= \exp\{\psi(\theta_0 + (D\theta)t) - \psi(\theta_0) - \tau(\theta)'(D\theta)t\}, \quad t \in \mathbf{R}^p,$$

$$(4.13) \quad D \log \varphi(t) = (\tau(\theta_0) + (D\theta)t - \tau(\theta))' D\theta,$$

$$(4.14) \quad D^2 \log \varphi(t) = (\beta(t))_{p \times p} = (D\theta)'(D^2\psi(\theta_0 + (D\theta)t)) D\theta.$$

Hence the equation defining the saddlepoint, cf. (4.3), becomes (4.11). $\tilde{t}$ may be recognized as the MLE in the model $\theta \in \{\theta_0 + (D\theta)t\}$, which is an affine hypothesis in the canonical parameter $\theta$. This fact ensures the existence and uniqueness of a solution to (4.11) satisfying $\tilde{\theta} \in \Theta$, since $\tau(\theta(b))$ belongs to the relative interior of the closed convex support of $X_1$. The result (4.10) now follows directly from (4.4).

The approximation to $e_0$, stated in the sequel, is derived from a simple kind of the mixed Edgeworth-saddlepoint approximation (see [3]), expanding the joint distribution of $D_1$ and $D_2$ around the same point $\tilde{\theta}$ as used to approximate $h_0(\theta)$.

LEMMA 4.5. *Let Assumptions 4.1 be fulfilled and define*

$$(4.15) \quad \gamma_2 = (\tau(\tilde{\theta}) - \tau(\theta))'(D^2\theta) - I(b)/n,$$

$$(4.16) \quad \tilde{e}_0 = \begin{cases} |-\gamma_2|, & \text{if neg. definite,} \\ 0 & \text{otherwise.} \end{cases}$$

*Then*

$$(4.17) \quad e_0 = \tilde{e}_0(1 + O(n^{-1})) \quad \text{as } n \to \infty$$

*uniformly on any compact set $B_0 \subseteq B$ on which $\gamma_2$ is negatively definite, where $e_0$ is defined in (4.9).*

Let $f_\theta$, $\theta \in \Theta$, denote the $\theta$-density of $\bar{X}$ induced by (4.5) with $\theta(\beta)$ replaced by $\theta$. For fixed $d_1$ let $t_1$ be the solution to

$$(4.18) \quad (\tau(\theta_1) - \tau(\theta(b)))' D\theta(b) = d_1, \quad \theta_1 = (D\theta(b))t_1 + \theta_0 \in \Theta,$$

i.e. $t_1$ is the saddlepoint corresponding to $D_1 = d_1$, when approximating the distribution of $D_1$. Then, by (4.2),

$$(4.19) \quad f_{\theta_0}(x) = f_{\theta_1}(x)\exp\{n[\psi(\theta_1) - \psi(\theta_0) - \tau(\theta_1)'(D\theta(b))t_1]\}.$$

Since $t_1$ depends only on $x$ through $d_1$, so does the entire exponential factor, and since $d_1$ is an affine function of $x$, it follows that the conditional $\theta_0$-density of $X$, given $D_1 = d_1$, is proportional to $f_{\theta_1}(x)$ on the affine support of $X$, given $D_1 = d_1$. In particular, we may approximate conditional moments, such as $e_0$, using a normal approximation to $f_{\theta_1}$. Since this

approach takes the "large deviation" event $D_1 = 0$ into account, it is preferable to a direct normal approximation. The result becomes

$$(4.20) \qquad E_0\{D_2|D_1 = 0\} = E_{\tilde{\partial}}\{D_2|D_1 = 0\} = \gamma_2 + O(n^{-1})$$

uniformly in $b$ in any compact set. Since the variance and higher cumulants of $D_2$, given $D_1 = 0$, are $O(n^{-1})$, while $\gamma_2$ is independent of $n$, (4.17) follows easily.

**Remark 4.6.** In some cases it may be possible to evaluate $\int(-d_2)q_{\tilde{\partial}}(0, d_2)d(d_2)$, which would provide a better approximation to $e_0$. Only in the one-dimensional case, however, would it be feasible to restrict the integration to the set, where $(-d_2)$ is positively definite. Other improvements are possible, e.g. by including variance-terms in the evaluation of the determinant rather than just computing the determinant of $(-\gamma_2)$.

On combining Lemma 4.2. with Lemma 4.4. and Lemma 4.5, we now have

**COROLLARY 4.7.** *Let Assumptions 4.1 be fulfilled, and let $B_0 \subseteq B$ be any compact set on which $\gamma_2$ is neg. definite; then*

$$(4.21) \qquad \lambda_0(b) = \tilde{\lambda}_0(b)\big(1 + O(n^{-1})\big) \quad \textit{uniformly in } b \in B_0,$$

*where*

$$(4.22) \quad \tilde{\lambda}_0(b)$$

$$= \begin{cases} (n/2\pi)^{p/2} |\beta(\tilde{t})|^{-1/2} \times \\ \quad \times \exp\{-n[\psi(\theta_0) - \psi(\tilde{\partial}) + \tau(\theta)'(D\theta)\tilde{t}]\}|-\gamma_2|, & \textit{if } \gamma_2 \textit{ is neg. definite,} \\ 0 & \textit{otherwise.} \end{cases}$$

**Proof.** Trivial.

It remains now only to be shown, how the approximation to $\lambda_0(b)$ provides an approximation to $g_0(b)$. The answer is simple. Since the probability of a local maximum in a neighbourhood of $\beta_0$ not being global tends exponentially to zero, the difference between $g_0(b)$ and $\lambda_0(b)$ tends rapidly to zero, such that $\tilde{\lambda}_0(b)$ also approximates $g_0(b)$.

**THEOREM 4.8.** *Let Assumptions 4.1 be fulfilled. Then, for some neighbourhood $B_0$ of $\beta_0$ and some constant $c > 0$,*

$$(4.23) \quad g_0(b) = \lambda_0(b)\big(1 + o(\exp\{-cn\})\big) = \tilde{\lambda}_0(b)\big(1 + O(n^{-1})\big) \quad \textit{as } n \to \infty$$

*uniformly in $b \in B_0$.*

**Proof.** See the Appendix.

Notice that since $B_0$ is independent of $n$, (4.23) is valid for large deviations of the type $\|\sqrt{n}(b - \beta_0)\| = o(\sqrt{n})$ in the normalized variable $\sqrt{n}(\hat{\beta} - \beta_0)$.

**COROLLARY 4.9.** *Let Assumptions 4.1 be fulfilled. Then there exists a*

*constant $c > 0$ such that*

$$(4.24) \quad \int_A g_0(b)\,db - \int_A \tilde{\lambda}_0(b)\,db$$

$$= O(n^{-1}) \int_A g_0(b)\,db + o(\exp\{-cn\}) \quad as\ n \to \infty$$

*uniformly in the class of Borel sets A.*

Proof. Follows easily from Theorem 4.8 and the fact that, for any neighbourhood $B_0$ of $\beta_0$, a constant $c > 0$ exists such that $P_0\{\hat{\beta} \notin B_0\} = o(\exp\{-cn\})$ as $n \to \infty$, which follows from the results in [4].

Remark 4.10. The advantage of this approximation, compared to the normal or Edgeworth approximations, is that apart from the exponentially decreasing term the *relative error* is $O(n^{-1})$ uniformly in all sets $A$. This makes the approximation particularly useful for calculating tail probabilities.

Remark 4.11. The relative error may be improved from $O(n^{-1})$ to $O(n^{-3/2})$ by a renormalization. There are several ways of doing this; the simplest is to divide $\tilde{\lambda}$ by its integral over $B$, but this may be infinite. Another method is to adjust $\tilde{\lambda}$ such that at $b = \beta_0$ it equals the value of the third order Edgeworth expansion, i.e. including the $O(n^{-1})$ terms; but the computational work is rather large. A simpler method, which is always valid, is to divide $\tilde{\lambda}_0$ by its integral as approximated by a Gauss–Hermite sum. In the one-dimensional case $(p = 1)$ only 4 terms are required, and the approximation becomes

$$g_0(b) \approx \tilde{\lambda}_0(b)/I,$$

$$(4.25) \qquad I = \sum_{i=1}^{4} w_i\, \tilde{\lambda}_0(b_i) \exp\{x_i^2\} \sqrt{2/I(\beta_0)},$$

where $b_i = \beta_0 + x_i\sqrt{2/I(\beta_0)}$ and $w_i, x_i$ $(i = 1, \ldots, 4)$ may be found in [1], Table 25.10. We shall not prove that this formula yields a valid asymptotic normalization improving the $O(n^{-1})$ error to $O(n^{-3/2})$. The proof relies on the fact that formula (4.25) yields the exact integral, if $\tilde{\lambda}_0$ is a third order Edgeworth expansion.

Remark 4.12. In calculating tail probabilities the best use of the approximation is probably to calculate the tail area rather than its complementary part. The reason is that essentially the relative error is bounded, such that small probabilities give smaller errors. However, one must convince oneself that the truncation introduced in (4.16) is of no great importance. This may be indicated by a small density at the boundary, where the truncation becomes effective.

**5. An example.** To illustrate the performance of the approximation for small *n*, we shall use an example chosen not because of practical relevance,

but rather as a case, where none of the steps in the approximation are "too accurate" as in the normal regression, where the saddlepoint approximation (4.10) is exact, and at the same time the computation of the exact density is feasible. The example is one-dimensional, since this makes pictures easier to look at.

Let $(Y, Z)$ be normally distributed with expectations zero and covariance matrix

$$\sum = \sigma^2 \begin{pmatrix} 1 & \varrho \\ \varrho & 1 \end{pmatrix}, \quad \sigma^2 > 0, \; -1 < \varrho < 1,$$

and consider the subfamily given by $\sigma^2 = 1$. Defining

$$X = (\tfrac{1}{2}(Y^2 + Z^2), \; YZ),$$

$$\theta(\varrho) = (\theta_1, \theta_2) = (-1/(1-\varrho^2), \; \varrho/(1-\varrho^2)),$$

$$\psi(\theta) = -\tfrac{1}{2} \log(\theta_1^2 - \theta_2^2),$$

the model is of the form of Section 4 with $\varrho$ being the unknown parameter corresponding to $\beta$. Consider $n$ independent replications; then

$$\bar{X} = \Big( \sum_{i=1}^{n} (Y_i^2 + Z_i^2)/2n, \; \sum_{i=1}^{n} Y_i Z_i/n \Big)$$

and

$$I(\varrho) = n(1+\varrho^2)/(1-\varrho^2)^2.$$

Let us first briefly sketch the derivation of the approximation (4.22), (4.23) to the density $g_\varrho(r)$ of $\hat{\varrho}$ at $r \in \,]-1, 1[$, when $\varrho$ is the true parameter.

The saddlepoint equation (4.11) becomes

(5.1) $$\tilde{t} - b = a^2 t - 2abt + c,$$

where

(5.2) $$a = r/(1-r^2), \quad b = (2r - \varrho(1+r^2))/(1-\varrho^2),$$
$$c = -r(1-r^2)/(1-\varrho^2).$$

If $r = 0$, then $\tilde{t} = b$, otherwise the saddlepoint is

(5.3) $$\tilde{t} = (1 + 2ab - (4a^2 b^2 + 1 - 4ac)^{1/2})/2a, \quad r \neq 0,$$

for the other solution to (5.1), $\tilde{\theta}$ is outside the range of $\theta$.

A straightforward computation now gives the *approximation* $\tilde{\lambda}_\varrho(r)$, defined in (4.22),

(5.4) $$\tilde{\lambda}_\varrho(0) = \begin{cases} (n/2\pi)^{1/2}(1-\varrho^2)^{(n-2)/2}(1-2\varrho^2), & |\varrho| \leq 1/\sqrt{2}, \\ 0, & |\varrho| > 1/\sqrt{2}, \end{cases}$$

(5.5) $$\tilde{\lambda}_\varrho(r) = (n/2\pi)^{1/2} g_1(\varrho, r)/\sqrt{g_2(\varrho, r)} \exp\{-ng_3(\varrho, r)\}, \quad r \neq 0,$$

where

$$g_1(\varrho, r) = \begin{cases} (4r^3 + \varrho(1 - 4r^2 - r^4))/(r(1 - r^2)(1 - \varrho^2)) + \\ \qquad\qquad\qquad\qquad + \tilde{t}/(r(1 - r^2)), \text{ if positive} \\ 0 \quad \text{otherwise,} \end{cases}$$

$$g_2(\varrho, r) = (\tilde{t} - b)(1 - r^2)/r + (2(1 + r^2)^2 + 8\varrho r(\varrho r - 1 - r^2))/(1 - \varrho^2)^2,$$

$$g_3(\varrho, r) = \tfrac{1}{2}\log\left[(b - \tilde{t})(1 - \varrho^2)/(r(1 - r^2))\right] - \tilde{t}r/(1 - r^2),$$

and $b$ is given in (5.2), $\tilde{t}$ in (5.3). Although the expression seems complicated, it is quite explicit and quickly calculated on a computer.

The *exact density*, $g_\varrho(r)$, is derived from the Wishart distribution of $\bar{X}$ on integration along the estimation lines (see e.g. [2], Example 1). The result is

$$(5.6) \quad f_\varrho(r) = \int_0^\infty (r^2 + (1 - r^2)s)(s^2 + qs + 1/4)^{(n-2)/2} \times$$

$$\times \exp\{-n(\alpha + \beta s)\}\, ds\, c_n(1 - r^2)^{n-2}/(1 - \varrho^2)^{n/2},$$

where

$$q = (1 + r^2)/(1 - r^2), \quad c_n = n^n/(2^{n-2}\,\Gamma(n/2)^2),$$

$$\alpha = \tfrac{1}{2}(1 - r^2)/(1 - \varrho^2), \quad \beta = (1 + r^2 - 2r\varrho)/(1 - \varrho^2).$$

If $n$ is even, the integral may be calculated explicitly. In particular, if $n = 2$, we obtain

$$g_\varrho(r) = (4/(1 - \varrho^2))\,[r^2/(2\beta) + (1 - r^2)/(2\beta)^2]\exp\{-2\alpha\}.$$

Note that here the computational work increases rapidly with $n$.

We shall compute the exact and approximate density in three cases and, for comparison, also give the usual *normal approximation* given by

$$(5.7) \quad \hat{g}_\varrho(r) = (n/2\pi)^{1/2}(1 + \varrho^2)^{1/2}/(1 - \varrho^2)\exp\left\{-\frac{n}{2}(r - \varrho)^2(1 + \varrho^2)/(1 - \varrho^2)^2\right\}.$$

We have also calculated the *renormalized approximation*

$$(5.8) \qquad\qquad\qquad \tilde{\lambda}_\varrho(r)/\int_{-1}^{1} \tilde{\lambda}_\varrho(r)\, dr.$$

The three cases are:

I. $\varrho = 0$, $n = 10$. Small $n$, symmetric distribution.

II. $\varrho = 0.9$, $n = 10$. Small $n$, skew distribution.

III. $\varrho = 0$, $n = 2$. Extremely small $n$, symmetric distribution.

The results are given in Figs. 1-3 and Tables 1-3 below. Rather than stating the densities of $\hat{\varrho}$ themselves, we have stated the densities of the normalized variable $\sqrt{n}(\hat{\varrho} - \varrho)$ as functions of $\hat{\varrho}$. These are obtained from (5.4), (5.6)-(5.8) on division by $\sqrt{n}$.
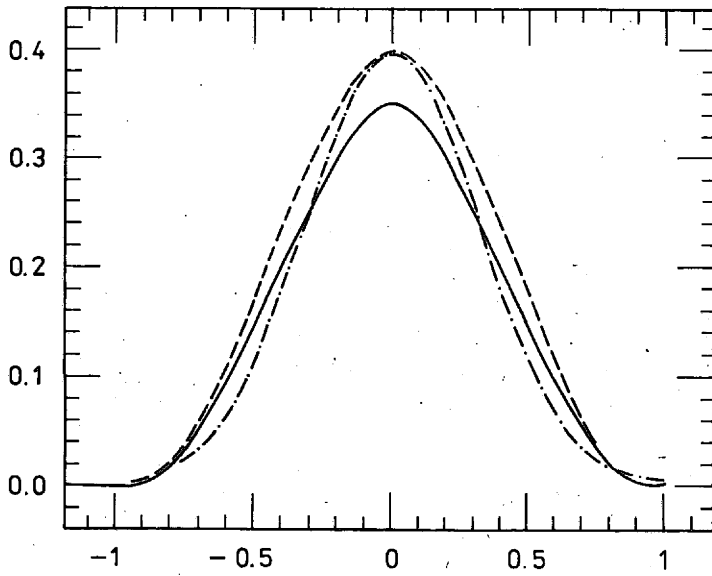
Fig. 1. Approximations to the density of $\sqrt{n}(\hat{\varrho} - \varrho)$ with $\varrho = 0$, $n = 10$.
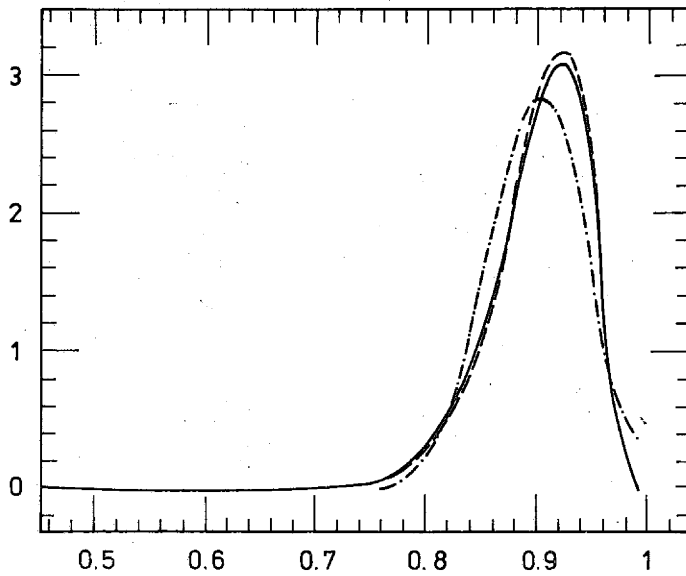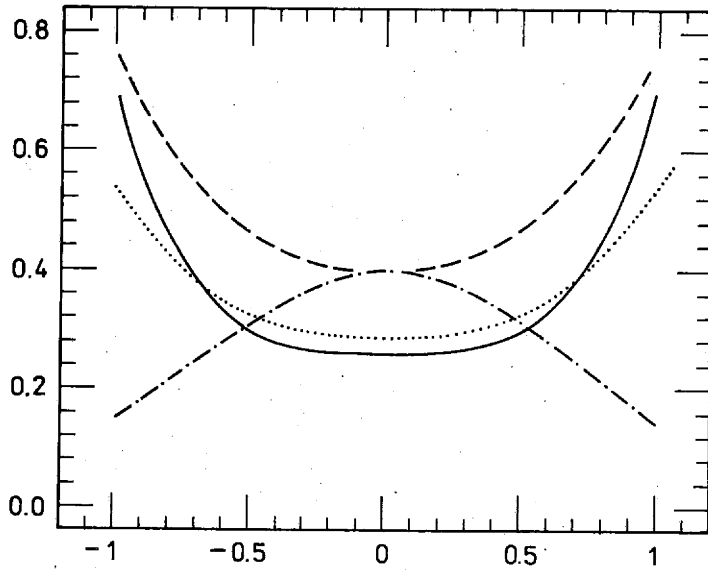Exact density (solid), approximation (5.4) (dashes) and normal approximation (dot-dash).



Fig. 2. Approximations to the density of $\sqrt{n}(\hat{\varrho} - \varrho)$ with $\varrho = 0.9$, $n = 10$.
Exact density (solid), approximation (5.4) (dashes) and normal approximation (dot-dash).

Fig. 3. Approximations to the density of $\sqrt{n}(\hat{\varrho}-\varrho)$ with $\varrho = 0$, $n = 2$.
Exact density (solid), approximation (5.4) (dashes), renormalized approximation (5.8) (dots) and normal approximation (dot-dash).

Table 1. Comparison of approximations to the density of $\sqrt{n}(\hat{\varrho}-\varrho)$ at $\hat{\varrho} = r$, when $\varrho = 0$, $n = 10$.

| $r$ | exact | normal | approx. | renorm. |
|------|-------|--------|---------|---------|
| 0.00 | 0.3505 | 0.3989 | 0.3989 | 0.3453 |
| 0.10 | 0.3366 | 0.3795 | 0.3853 | 0.3357 |
| 0.20 | 0.2996 | 0.3266 | 0.3477 | 0.3010 |
| 0.35 | 0.2244 | 0.2162 | 0.2645 | 0.2290 |
| 0.50 | 0.1476 | 0.1143 | 0.1720 | 0.1489 |
| 0.70 | $5.355 \cdot 10^{-2}$ | $3.443 \cdot 10^{-2}$ | $5.917 \cdot 10^{-2}$ | $5.122 \cdot 10^{-2}$ |
| 0.90 | $2.176 \cdot 10^{-3}$ | $6.951 \cdot 10^{-3}$ | $2.270 \cdot 10^{-3}$ | $1.965 \cdot 10^{-3}$ |
| 0.99 | $3.859 \cdot 10^{-7}$ | $2.969 \cdot 10^{-3}$ | $3.933 \cdot 10^{-7}$ | $3.405 \cdot 10^{-7}$ |

Table 2. Comparison of approximations to the density of $\sqrt{n}(\hat{\varrho}-\varrho)$ at $\hat{\varrho} = r$, when $\varrho = 0.9$, $n = 10$.

| $r$ | exact | normal | approx. | renorm. |
|------|-------|--------|---------|---------|
| 0.50 | $1.284 \cdot 10^{-5}$ | $1.074 \cdot 10^{-17}$ | $6.950 \cdot 10^{-6}$ | $6.893 \cdot 10^{-6}$ |
| 0.70 | $1.088 \cdot 10^{-2}$ | $1.247 \cdot 10^{-4}$ | $8.547 \cdot 10^{-3}$ | $8.477 \cdot 10^{-3}$ |
| 0.80 | 0.2983 | 0.2303 | 0.2761 | 0.2738 |
| 0.85 | 1.152 | 1.509 | 1.130 | 1.121 |
| 0.88 | 2.118 | 2.555 | 2.130 | 2.113 |
| 0.90 | 2.777 | 2.825 | 2.825 | 2.802 |
| 0.92 | 3.089 | 2.555 | 3.167 | 3.141 |
| 0.95 | 2.109 | 1.509 | 2.172 | 2.154 |
| 0.99 | $2.494 \cdot 10^{-2}$ | $3.708 \cdot 10^{-1}$ | $2.547 \cdot 10^{-2}$ | $2.526 \cdot 10^{-2}$ |

Table 3. Comparison of approximations to the density of $\sqrt{n}(\hat{\varrho}-\varrho)$ at $\hat{\varrho}=r$, when $\varrho=0$, $n=2$.

| $r$ | exact | normal | approx. | renorm. |
|------|--------|--------|---------|---------|
| 0.00 | 0.2601 | 0.3989 | 0.3989 | 0.2813 |
| 0.20 | 0.2611 | 0.3833 | 0.4069 | 0.2869 |
| 0.40 | 0.2748 | 0.3400 | 0.4340 | 0.3060 |
| 0.60 | 0.3264 | 0.2783 | 0.4950 | 0.3491 |
| 0.80 | 0.4511 | 0.2104 | 0.6041 | 0.4260 |
| 0.95 | 0.6258 | 0.1618 | 0.7208 | 0.5083 |

In case I ($\varrho=0$, $n=10$) it is seen (Fig. 1, Table 1) that, except for the extreme tail, the normal approximation does quite well. Approximation (5.4) is slightly worse in the main part of the distribution, but it keeps the shape better leading to an excellent renormalized approximation, which has not been drawn since it is hardly distinguishable from the exact density. Note that approximation (5.4), renormalized or not, keeps its degree of approximation throughout the range. No truncation (see (4.22) and (5.4)) is needed, when $\varrho=0$.

In case II ($\varrho=0.9$, $n=10$) the distribution is skew and the normal approximation is useless. Approximation (5.4), however, does quite well throughout the range. Here, the effect of renormalization is vanishing, since the integral of the approximation is 1.008. In this case the approximation is truncated at 0.35324 at a density of approx. $10^{-7}$.

In case III ($\varrho=0$, $n=2$) $n$ is so small that hardly any approximation can be expected to work. Both (5.4) and the normal approximation are numerically useless as direct approximations (Fig. 3, Table 3), but approximation (5.4) again has the right shape and its renormalized version does surprisingly well.

The comparison with the normal approximation has only been included to give an impression of the magnitude of the deviations. In a thorough investigation of the behaviour it would be more relevant to compare with the second-order Edgeworth expansion, which has an error of order $O(n^{-1})$ and which better approximates skew distributions. Its tail behaviour is, however, not in general better than that of the normal distribution.

It seems that approximation (5.4) behaves similar to the saddlepoint approximation on which it is based (see [5]); namely, keeping its relative error fairly constant throughout the range such that the renormalized approximation works extremely well. When calculating tail probabilities of magnitude 0.01, say, a relative error of 50°/₀ is often of no great importance and the approximation may safely be used directly, unless the truncation is of importance.

**6. The multi-dimensional normal regression model.** In this section we shall specialize the results of Section 4 to obtain a simple explicit approximation to the density of the MLE of the (multivariate) parameter in

the important class of nonlinear normal regression models. As in Section 3 we shall consider the asymptotics obtained by letting the variance tend to zero, which is equivalent (mathematically) to simple replications.

Let $X = (X_1, \ldots, X_k)$ be normally distributed with expectation vector $\mu(\beta) = (\mu_1(\beta), \ldots, \mu_k(\beta))$ and covariance matrix $\Sigma = \sigma^2 I_{k \times k}$, where $\beta \in B \subseteq R^p$ is the unknown parameter, $B$ is open, $\sigma^2$ is considered to be known, and $\mu: B \to R^k$ is a known function.

As previously, $\hat{\beta}$ is the MLE of $\beta$; $\beta_0 \in B$ and $b \in B$ are arbitrary fixed points, and

$$(6.1) \qquad I(\beta) = (D\mu(\beta))' D\mu(\beta)/\sigma^2$$

is the Fisher information matrix. Since the density of $X$ is of the form

$$f(X; \beta) = c(\beta) \{\exp \mu(\beta)' X/\sigma^2\},$$

this is a curved exponential family model as discussed in Section 4. Assumptions 4.1 are equivalent to

ASSUMPTIONS 6.1. (i) $\mu: B \to R^k$ is one-to-one, bicontinuous and three times differentiable.

(ii) The Fisher information matrix $I(\beta)$ is regular for all $\beta \in B$.

A direct computation, either by insertion in the results of Section 4 or using the normality of $(D_1, D_2)$ in combination with Lemma 4.2, now shows that approximation (4.22) becomes

$$(6.2) \quad \tilde{\lambda}_0(b) = (2\pi)^{-p/2} |I(b)|^{-1/2} \tilde{e}_0 \times$$

$$\times \exp\left\{-\tfrac{1}{2}(\mu(\beta_0) - \mu(b))' D\mu(b) I(b)^{-1} D\mu(b)' (\mu(\beta_0) - \mu(b))/\sigma^4\right\}$$

$$= (2\pi)^{-p/2} |I(b)|^{-1/2} \tilde{e}_0 \exp\left\{-\frac{1}{2\sigma^2} \|P_b(\mu(\beta_0) - \mu(b))\|^2\right\},$$

where $P_b$ is the projection matrix onto the subspace spanned by the columns of $D\mu(b)$, and

$$(6.3) \qquad \tilde{e}_0 = \begin{cases} |-\gamma_2| & \text{if } \gamma_2 \text{ is negatively definite,} \\ 0 & \text{otherwise,} \end{cases}$$

$$(6.4) \quad \gamma_2 = \frac{1}{\sigma^2}(\mu(\beta_0) - \mu(b))'(I_{k \times k} - D\mu(b)(\sigma^2 I(b))^{-1} D\mu(b)') D^2\mu(b) - I(b)$$

$$= \frac{1}{\sigma^2}(\mu(\beta_0) - \mu(b))'(I_{k \times k} - P_b) D^2\mu(b) - I(b).$$

Remark 6.2. The asymptotic behaviour of this approximation is stated in Theorem 4.8 and Corollary 4.9. There are, however, less approximations involved in this case, since approximation (4.10) to $h_0(0)$ here is exact, and also $\gamma_2$ is exactly equal to $E_0\{D_2|D_1 = 0\}$. Both of these approximations

contributed with a relative error of order $O(n^{-1})$ in the general case. Hence, it might be worthwhile also to remove the last $O(n^{-1})$-error by including the variance terms in the approximation to $e_0$. Even the higher-order terms (in powers of $n^{-1}$) may be removed in this way, leaving only exponentially decreasing errors but, if the dimension $p$ is large, this requires quite a lot of computation. We shall not state these refined expansions, which are easily written down for a specific $p$.

## 7. Appendix.

Proof of Lemma 4.2. The event $M_1(\varepsilon)$, say, that $f(X; \beta)$ has a local maximum at some $\beta$ in $\{\|\beta - b\| < \varepsilon\}$ is the same as the event that $D_1(\beta) = 0$ and $D_2(\beta)$ is negatively definite for some $\beta$ in the same set, except if $D_1(\beta) = 0$ and $D_2(\beta) = 0$, which may be disregarded because of Assumption 4.1(i). We shall show that, furthermore, $M_1(\varepsilon)$ may be replaced by the event

$$M_2(\varepsilon) = \{S_1(\beta) = 0 \text{ for some } \beta \text{ in } \{\|\beta - b\| < \varepsilon\},$$

$$\text{and } D_2(b) \text{ is neg. definite}\},$$

where $S_1(\beta) = D_1(b) + D_2(b)(\beta - b) = D_1 + D_2(\beta - b)$ is the linear approximation to $D_1(\beta)$ around $\beta = b$.

Observe that since the Laplace transform $E_0\{e^{s'X}\}$ exists in a neighbourhood of zero, there exists a $K_1 > 0$ such that

$$(7.1) \qquad P_0\{\|\bar{X} - \tau(\theta)\| \geqslant K_1 \log \varepsilon^{-1}\} = o(\varepsilon^p) \quad \text{as } \varepsilon \to 0$$

and for some positive constants $K_2, K_3$ we have

$$(7.2) \qquad \|D_1(\beta) - S_1(\beta)\| \leqslant K_2 \varepsilon^2 \log \varepsilon^{-1},$$

$$(7.3) \qquad \|D_2(b)\| \leqslant K_3 \log \varepsilon^{-1}$$

if $\|\bar{X} - \tau(\theta)\| < K_1 \log \varepsilon^{-1}$.

Now, let $0 < \delta < 1$ be fixed. Then

$$(7.4) \qquad (1 + \delta)^p \lim_{\varepsilon \to 0} (\varepsilon^p A_p)^{-1} P_0(M_1(\varepsilon))$$

$$= (1 + \delta)^p \lim_{\varepsilon \to 0} (\varepsilon^p (1 + \delta)^p A_p)^{-1} P_0(M_1(\varepsilon(1 + \delta)))$$

$$\geqslant \lim_{\varepsilon \to 0} (\varepsilon^p A_p)^{-1} [P_0(M_2(\varepsilon)) - P_0(M_2(\varepsilon) \backslash M_1(\varepsilon(1 + \delta)))].$$

But, if $S_1(\beta_1) = 0$, $\|\beta_1 - b\| < \varepsilon$, and $\|\bar{X} - \tau(\theta)\| < K_1 \log \varepsilon^{-1}$, then

$$\{S_1(\beta) | \|\beta - \beta_1\| < \delta\varepsilon\} \supseteq \{y \in \mathbf{R}^p | \|y\| < \lambda\delta\varepsilon\},$$

where $\lambda$ is the smallest eigenvalue of $D_2$ such that, by (7.2),

$$\{D_1(\beta) | \|\beta - \beta_1\| < \delta\varepsilon\} \supseteq \{y \in \mathbf{R}^p | \|y\| < \lambda\delta\varepsilon - K_2 \varepsilon^2 \log \varepsilon^{-1}\},$$

which contains zero if $\lambda > K_2 \varepsilon \log \varepsilon^{-1}/\delta$. Since also

$$\|D_1\| < \|D_2\| \|\beta_1 - b\| \leqslant K_3 \varepsilon \log \varepsilon^{-1},$$

we obtain

$$P_0\left(M_2(\varepsilon) \backslash M_1\left(\varepsilon(1+\delta)\right)\right) \leqslant P_0\left\{\|\bar{X} - \tau(\theta)\| \geqslant K_1 \log \varepsilon^{-1}\right\} +$$

$$+ P_0\left\{\|D_1\| \leqslant K_3 \varepsilon \log \varepsilon^{-1}\right\} \cap \left\{\lambda \leqslant K_2 \varepsilon \log \varepsilon^{-1}/\delta\right\}$$

$$= o(\varepsilon^p) + O\left((K_3 \varepsilon \log \varepsilon^{-1})^p (K_2 \varepsilon \log \varepsilon^{-1}/\delta)\right) = o(\varepsilon^p) \quad \text{as } \varepsilon \to 0,$$

proving that

$$(7.5) \qquad \lim_{\varepsilon \to 0} (\varepsilon^p A_p)^{-1} P_0\left(M_1(\varepsilon)\right) \geqslant \lim_{\varepsilon \to 0} (\varepsilon^p A_p)^{-1} P_0\left(M_2(\varepsilon)\right),$$

since $\delta$ was arbitrary. The other inequality follows similarly.

By Assumptions 4.1 we may write $D_2 = A(D_1) + Y - I(b)$, such that $A$ is a linear function and $(D_1, Y)$ has a continuous density $\zeta_0(d_1, y)$, say, on its closed convex support, which contains $(0, E_0\{Y | D_1 = 0\})$ as an interior point. Thus, by continuity, we have

$$\lambda_0(b) = \lim_{\varepsilon \to 0} (\varepsilon^p A_p)^{-1} P_0\left(M_2(\varepsilon)\right)$$

$$= \lim_{\varepsilon \to 0} (\varepsilon^p A_p)^{-1} \int_{d_2 \text{neg.def.}} \int_{d_1 \in d_2(B_\varepsilon)} \zeta_0(d_1, y) \, d(d_1) \, dy$$

$$= \int_{y - I(b) \text{neg.def.}} |I(b) - y| \, \zeta_0(0, y) \, dy = h_0(0) e_0,$$

where $B_\varepsilon = \{x \in \mathbf{R}^p \mid \|x\| < \varepsilon\}$, and $h_0$ and $e_0$ are defined in the Lemma.

Proof of Theorem 4.8. If $\bar{X} = \tau(\theta_0)$, then the likelihood function has a global maximum at $\beta_0$. Hence, by continuity, if $D_1(\beta_0) = 0$ and $\bar{X}$ lies within a certain neighbourhood of $\tau(\theta_0)$, the local maximum at $\beta_0$ will also be global. By another continuity argument, using Assumption 4.1 (ii), there is a neighbourhood $B_0'$ of $\beta_0$ and, for each $b \in B_0'$, a neighbourhood $T(b)$ such that if $D_1(b) = 0$ and $\bar{X} \in T(b)$, then the likelihood function has a global maximum at $b$. Now, for some constant $c(b)$,

$$P_0\left\{\bar{X} \notin T(b) \mid D_1(b) = 0\right\} = o\left(\exp\left\{-c(b) n\right\}\right) \quad \text{as } n \to \infty.$$

This proves the first equality of (4.23), since also the uniformity follows by continuity. The second equality follows by Corollary 4.7.

### REFERENCES

[1] M. Abramowitz and I. A. Stegun (ed.), *Handbook of Mathematical Functions*, National Bureau of Standards, Washington 1964.

[2] O. Barndorff-Nielsen, *Conditionality resolutions*, Biometrika 67 (1980), p. 293-310.

[3] — and D. R. Cox, *Edgeworth and saddlepoint approximations with statistical applications*, J. R. Statist. Soc. B. 41 (1979), p. 279-312.

[4] R. H. Berk, *Consistency and asymptotic normality of mle's for exponential models*, Ann. Math. Statist. 43 (1972), p. 193-204.

[5] H. E. Daniels, *Saddlepoint approximation in statistics*, ibidem 25 (1954), p. 631-650.

[6] W. Feller, *An Introduction to Probability Theory and its Applications*, vol. II, Wiley, New York 1971.

Royal Veterinary and Agricultural University
Department of Mathematics
Denmark