

3. SADDLEPOINT APPROXIMATIONS FOR THE MEAN

3.1. INTRODUCTION

The goal of this chapter is to introduce saddlepoint techniques for a simple problem, namely the approximation to the distribution of the mean of n iid random variables.

Although this case does not have much relevance from a practical point of view, the same basic idea is used in more complex models to derive saddlepoint approximations for very general statistics; cf. chapters 4 and 5. Thus, a good understanding of the technique in this simple case will allow a direct application to more complex and important situations. Historically, this was the first explicit statistical application of this method. It was developed by H.E. Daniels in a fundamental paper in 1954.

Basically, there are two ways to derive a saddlepoint approximation. The first one is presented in section 3.3 and is an application of the method of steepest descent (section 3.2). The second one is based on the idea of recentering by means of a conjugate or associate distribution (section 3.4) and shows the connection between saddlepoint techniques and Edgeworth expansions. Both ways lead to the same approximation and from a methodological point of view they both have their own merits. Finally, the examples in section 3.5 show the great accuracy of these approximations for very small sample sizes and far out in the tails.

3.2. THE METHOD OF STEEPEST DESCENT

We discuss here a general technique which allows us to compute asymptotic expansions of integrals of the form

$$\int_{\mathcal{P}} e^{v \cdot w(z)} \xi(z) dz \quad (3.1)$$

when the real parameter v is large and positive. Here w and ξ are analytic functions of z in a domain of the complex plane which contains the path of integration \mathcal{P} . This technique is called the *method of steepest descent* and will be used to derive saddlepoint approximations to the density of a mean (section 3.3) and later of a general statistic (chapter 4). In our exposition we follow Copson (1965). Other basic references are DeBruijn (1970), and Barndorff-Neilsen and Cox (1989).

Consider first the integral (3.1). In order to compute it we can deform arbitrarily the path of integration \mathcal{P} provided we remain in the domain where w and ξ are analytic. We deform \mathcal{P} such that

- (i) the new path of integration passes through a zero of the derivative $w'(z)$ of w ;
- (ii) the imaginary part of w , $\Im w(z)$ is constant on the new path.

Let us now look at the implications of (i) and (ii). If we write

$$z = x + iy, \quad z_0 = x_0 + iy_0,$$

$$w(z) = u(x, y) + iv(x, y), \quad w'(z_0) = 0,$$

and denote by S the surface $(x, y) \mapsto u(x, y)$, then by the Cauchy-Riemann differential equations

$$u_x = v_y, \quad u_y = -v_x,$$

it follows that the point (x_0, y_0) cannot be a maximum or a minimum but must be a *saddlepoint* on the surface S . Moreover, the orthogonal trajectories to the level curves $u(x, y) = \text{constant}$ are given (again by the Cauchy-Riemann differential equations) by the curves $\Im w(z) = v(x, y) = \text{constant}$. Since the paths on S corresponding to the orthogonal trajectories of the level curves are paths of steepest (ascent) descent, condition (ii) above means that the integration along a path where $\Im w(z)$ is constant implies that we are moving along the paths of steepest descent from the saddlepoint (x_0, y_0) on the surface S .

Therefore, on a steepest path through the saddlepoint we have

$$\begin{aligned} w(z) &= u(x, y) + iv(x_0, y_0) \\ &= u(x_0, y_0) + iv(x_0, y_0) - (u(x_0, y_0) - u(x, y)) = w(z_0) - \gamma(x, y), \end{aligned} \quad (3.2)$$

where γ is the real function

$$\gamma(x, y) = u(x_0, y_0) - u(x, y). \quad (3.3)$$

It follows directly that $d\gamma/ds = \pm|w'(z)|$, where s is the arc length of the path (on the plane). Thus, γ is monotonic on the steepest path from the saddlepoint and either increases to $+\infty$ or decreases to $-\infty$. Since by (3.1) and (3.2) $\gamma \rightarrow -\infty$ leads to a divergent integral, we choose the path where γ increases to $+\infty$. This is the path of steepest descent from the saddlepoint. Exhibit 3.1 shows this path for the function $w(z) = z^2$ and Exhibit 3.2 shows the surface $u = u(x, y)$ about the saddlepoint $(x_0, y_0) = (0.25, 0)$ and the path of steepest descent for the function $w(z) = K(z) - z \cdot t$, where $K(z) = -\beta \log(1 - z/\theta)$; $t > 0$, $\theta > 0$, $\beta > 0$ fixed. With this second choice of $w(z)$ and $v = n$, $\xi(z) \equiv n(2\pi i)^{-1}$, the integral (3.1) is just the density (evaluated at t) of the mean of n iid random variables from a Gamma distribution; cf. (3.6).

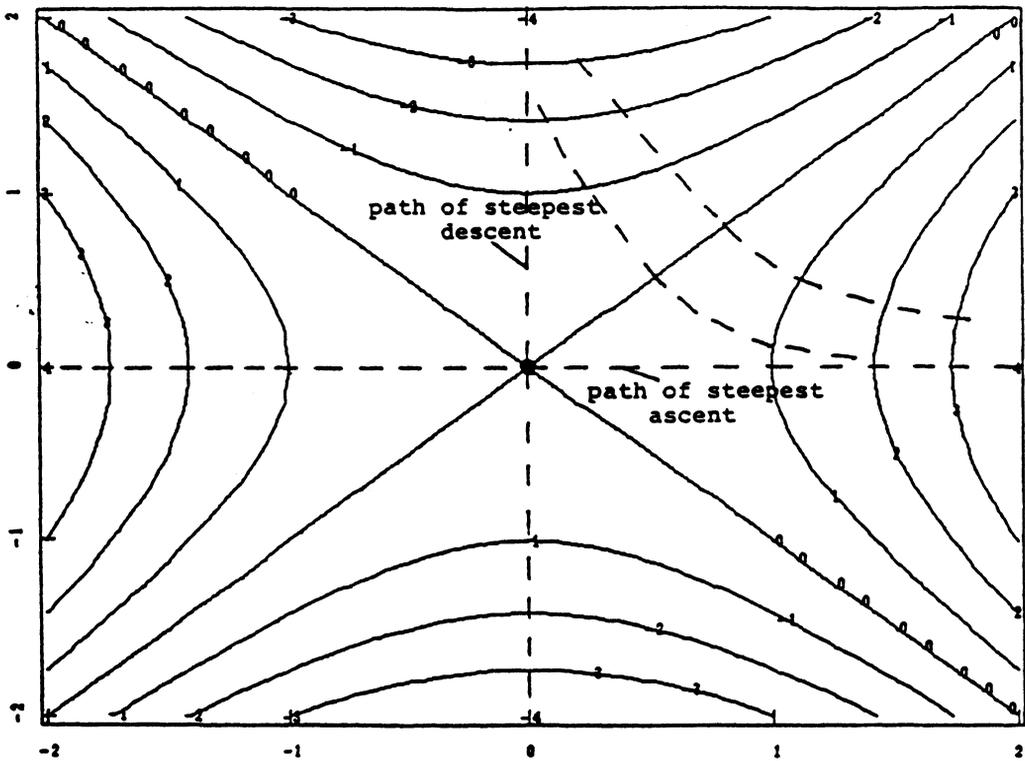


Exhibit 3.1

The path of steepest descent for the function

$$w(z) = z^2 = u(x, y) + iv(x, y)$$

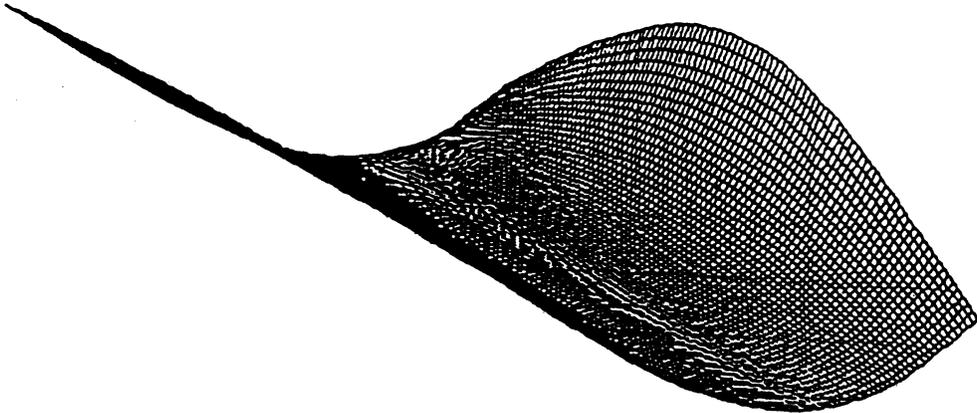
$$u(x, y) = x^2 - y^2; \quad v(x, y) = 2xy$$

$$\gamma(x, y) = 0 - u(x, y) = y^2 - x^2$$

$$\text{Saddlepoint} : (x_0, y_0) = (0, 0)$$

————— : $u(x, y) = \text{constant}$ (level curves)

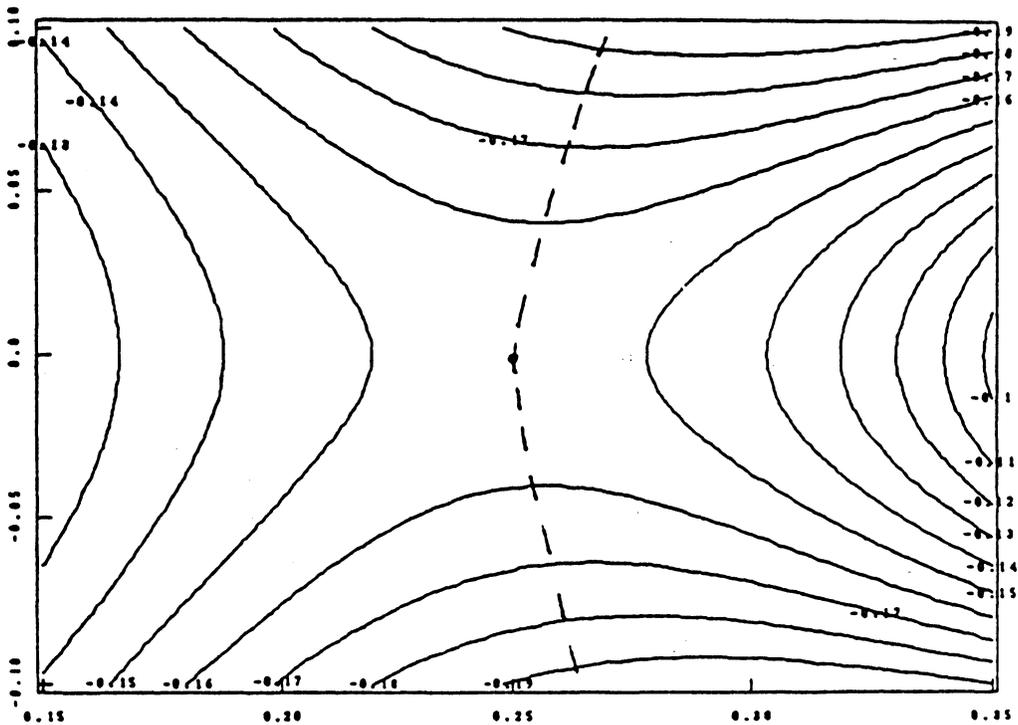
----- : $v(x, y) = \text{constant}$ (orthogonal trajectories)

**Exhibit 3.2a**

Surface $u = u(x, y)$ for the function

$$w(z) = u(x, y) + iv(x, y) = -\beta \log\left(1 - \frac{z}{\theta}\right) - z \cdot t,$$

$$t = 2, \theta = 0.5, \beta = 0.5.$$

**Exhibit 3.2b**

Level curves and paths of steepest descent from the saddlepoint $(x_0, y_0) = (0.25, 0)$ for the function of Exhibit 3.2a.

To summarize: if it is possible to deform the path of integration and express the integral as a sum of integrals along paths of steepest descent from saddlepoints, it follows from (3.1) and (3.2) that we have to consider only integrals of the form

$$\int_{\mathcal{P}} e^{v w(z)} \xi(z) dz = e^{v w(z_0)} \int_0^{\infty} e^{-v \cdot \gamma} \xi(z) \frac{dz}{d\gamma} d\gamma. \quad (3.4)$$

It can be seen from (3.4) that instead of approximating $w(z)$ in the exponential (where the error would be blown up), we approximate $dz/d\gamma$ which has been moved down from the exponent. This approximation can be obtained by expanding this expression into a series near the saddlepoint z_0 . A typical example is the application of this technique in statistics (see section 3.3). Further applications are given in chapter 7.

Remark 3.1

Historically, the method of steepest descent can be traced back to Riemann (1892) who found an asymptotic approximation to the hypergeometric function, that is a multiple of the integral (3.1) with $w(z) = \log[z(1-z)(1-sz)^{-1}]$ and $\xi(z) = z^a(1-z)^b(1-sz)^c$, where a, b, c, s are real parameters and \mathcal{P} is any curve which joins 0 and 1. Debye (1909) generalized the work of Riemann and obtained a complete asymptotic expansion for integrals of the form (3.1) by means of the idea presented in this section.

3.3. SADDLEPOINT APPROXIMATIONS FOR THE MEAN

This section serves two purposes. First, it is an application of the method of steepest descent. In particular, we will construct explicitly (i) and (ii) of section 3.2. Secondly, it shows a simple but conceptually important application of this technique in statistics, namely the approximation to the distribution of the mean of n iid random variables. In our exposition in this section we will follow the outline of Daniels' (1954) fundamental paper. Given n iid observations x_1, \dots, x_n with common known distribution $F(x)$ and density $f(x)$, we want to approximate the density $f_n(t)$ of the arithmetic mean $T_n(x_1, \dots, x_n) = n^{-1} \sum_{i=1}^n x_i$.

Denote by $M(\alpha) = \int_{-\infty}^{+\infty} e^{\alpha x} f(x) dx$ the moment generating function, by $K(\alpha) = \log M(\alpha)$ the cumulant generating function and suppose they exist for real values of α in some interval (c_1, c_2) containing the origin. Then, by Fourier inversion the density $f_n(t)$ can be written as

$$\begin{aligned} f_n(t) &= (n/2\pi) \int_{-\infty}^{+\infty} M^n(ir) e^{-inrt} dr \\ &= (n/2\pi i) \int_{\mathfrak{I}} M^n(z) e^{-niz} dz, \end{aligned} \quad (3.5)$$

where \mathfrak{I} is the imaginary axis in the complex plane. Since the contributions of the integral over the paths \mathcal{P}' and \mathcal{P}'' go to 0 as $\alpha \rightarrow \infty$ (see Exhibit 3.3), one can alternatively integrate over any straight line parallel to the imaginary axis. Therefore, (3.5) can be rewritten as

$$\begin{aligned}
 f_n(t) &= (n/2\pi i) \int_{\tau-i\infty}^{\tau+i\infty} M^n(z) e^{-nz t} dz \\
 &= (n/2\pi i) \int_{\tau-i\infty}^{\tau+i\infty} \exp\{n[K(z) - z t]\} dz,
 \end{aligned} \tag{3.6}$$

for any real number τ in the interval (c_1, c_2) .

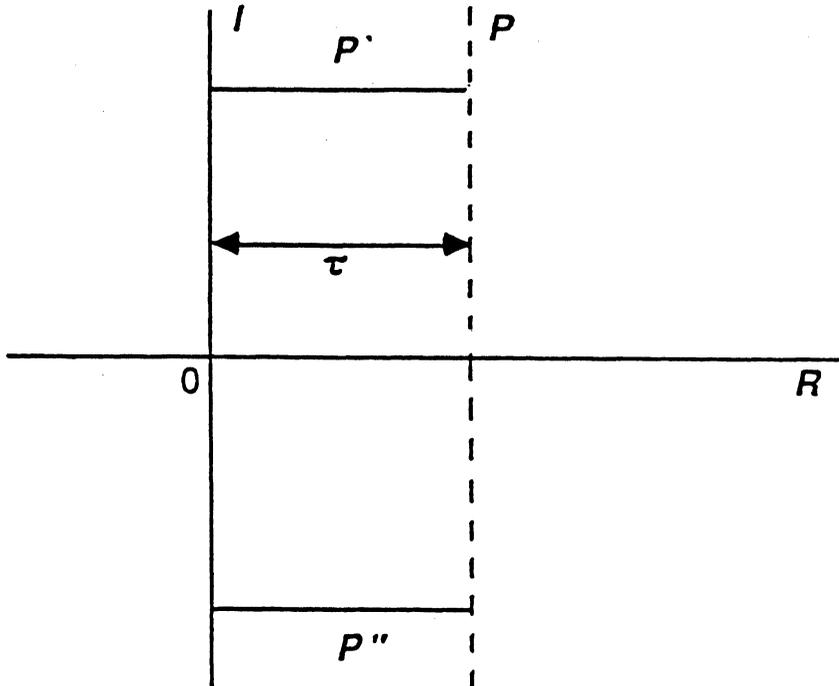


Exhibit 3.3

Shift of path of integration from \mathfrak{S} to \mathcal{P}

At this point the integral in (3.6) is of the form (3.1) where $v = n$, $w(z) = K(z) - z \cdot t$ with t fixed, $\xi(z) = n/2\pi i$, and the path \mathcal{P} is a straight line parallel to the imaginary axis going through the real point τ . Let us look at conditions (i) and (ii) of section 3.2. From (i) we see that the new path will have to go through a zero of $w'(z)$, that is

$$w'(z) = K'(z) - t = 0.$$

Thus, the new path will go through the saddlepoint z_0 defined as a solution to the equation

$$K'(z_0) = t.$$

Daniels (1954) shows in Theorems 6.1 and 6.2 under general conditions that the saddlepoint z_0 is *unique and real* on (c_1, c_2) , and that $K''(z_0) > 0$. Thus, from now on $z_0 = \alpha_0 \in \mathbb{R}$.

Condition (ii) requires that $\Im w(z) \equiv \text{constant}$ on the new path, that is $\Im w(z) \equiv \Im w(\alpha_0) = 0$ since $w(\alpha_0)$ is real. This allows us to deform \mathcal{P} as shown in Exhibit 3.4.

First, choose $\tau = \alpha_0$ and move the integration path on the straight line parallel to the imaginary axis which goes through the real point α_0 . Secondly, construct a small circle of radius ϵ around the saddlepoint α_0 and follow the path of steepest descent from the saddlepoint inside this circle (\mathcal{P}_0). On this path $\Im w(z) = 0$ by condition (ii). Then, continue on the curves orthogonal to (\mathcal{P}_0) at z_1 and z_2 . Since (\mathcal{P}_0) is a path of steepest descent, (\mathcal{P}_1) and (\mathcal{P}_2) are level curves defined by $\Re w(z) = \text{constant}$. From z_3 (respectively z_4) continue on the straight line ($\mathcal{P}_3, \mathcal{P}_4$). Therefore, the original integral (3.5) can be rewritten as

$$f_n(t) = I_0 + I_1, \tag{3.7}$$

where

$$I_0 = (n/2\pi i) \int_{\mathcal{P}_0} \exp\{n[K(z) - z \cdot t]\} dz \tag{3.8}$$

is the contribution to the integral inside the circle and

$$I_1 = (n/2\pi i) \int_{\tilde{\mathcal{P}} \setminus \mathcal{P}_0} \exp\{n[K(z) - z \cdot t]\} dz \tag{3.9}$$

is the contribution outside the circle.

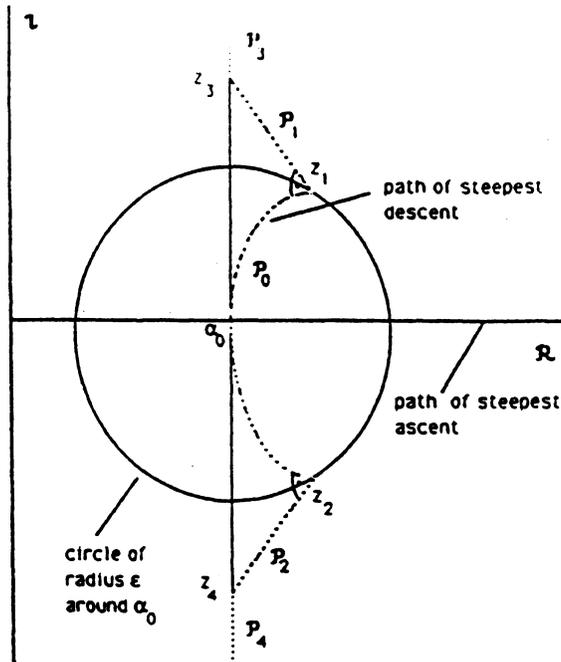


Exhibit 3.4

Deformation of the path of integration for the computation of the integral in (3.5).

Old path : imaginary axis

New path : $\tilde{\mathcal{P}} = \mathcal{P}_4 \cup \mathcal{P}_2 \cup \mathcal{P}_0 \cup \mathcal{P}_1 \cup \mathcal{P}_3$ (-----)

$\mathcal{P}_0 : \Im w(z) \equiv \Im w(\alpha_0) = 0$; this is a path of steepest descent which crosses \Re orthogonally at α_0 .

$\Re : \Im w(z) \equiv 0$; this is a path of steepest ascent ($K''(\alpha_0) > 0$)

$\mathcal{P}_1 : \Re w(z) \equiv \text{constant} = \Re w(z_3)$

$\mathcal{P}_2 : \Re w(z) \equiv \text{constant} = \Re w(z_4)$

We first look at I_1 . On the straight line $z = \alpha_0 + iy$ we have:

$$\begin{aligned} e^{w(z)} &= M(z)e^{-zt} = \int \exp\{(\alpha_0 + iy)x\} dF(x) \cdot e^{-zt} \\ &= \left\{ \int e^{iyx} e^{\alpha_0 x} dF(x) / M(\alpha_0) \right\} M(\alpha_0) e^{-zt} \\ &= \phi(y) M(\alpha_0) \exp\{-(\alpha_0 + iy)t\}, \end{aligned} \quad (3.10)$$

where $\phi(y)$ is the characteristic function of a random variable with density $e^{\alpha_0 x} f(x) / M(\alpha_0)$. Therefore,

$$|e^{w(z)}| = \exp\{\Re w(z)\} \leq \rho \cdot e^{w(\alpha_0)},$$

with $\rho < 1$, and the contribution to the integral outside the circle on \mathcal{P}_3 , and \mathcal{P}_4 is of order $O(\rho^n)$ and can be ignored. On \mathcal{P}_1 and \mathcal{P}_2 , $\Re w(z)$ is constant and $e^{w(z)}$ can be bounded as above. Hence the contribution on \mathcal{P}_1 and \mathcal{P}_2 can be ignored.

Let us now look at I_0 . By definition, $w(z)$ is real on \mathcal{P}_0 . Define $\gamma(x, y)$ as in (3.3)

$$\gamma = w(\alpha_0) - w(z) = K(\alpha_0) - \alpha_0 \cdot t - [K(z) - zt]$$

and expand the right hand side in a series around α_0

$$\begin{aligned} \gamma &= -(z - \alpha_0)[K'(\alpha_0) - t] - (z - \alpha_0)^2 K''(\alpha_0)/2 \\ &\quad - (z - \alpha_0)^3 K'''(\alpha_0)/6 - (z - \alpha_0)^4 K^{(iv)}(\alpha_0)/24 - \dots \end{aligned} \quad (3.11)$$

Since γ is real and steadily increasing from the saddlepoint, we rewrite it as $\gamma = \delta^2/2$. With the change of variable

$$\zeta = (z - \alpha_0)[K''(\alpha_0)]^{1/2}$$

and

$$\begin{aligned} \lambda_3(\alpha_0) &= K'''(\alpha_0)/[K''(\alpha_0)]^{3/2} \\ \lambda_4(\alpha_0) &= K^{(iv)}(\alpha_0)/[K''(\alpha_0)]^2, \end{aligned}$$

and recalling that $K'(\alpha_0) = t$ we can rewrite (3.11) as

$$-\delta^2/2 = \zeta^2/2 + \lambda_3(\alpha_0)\zeta^3/6 + \lambda_4(\alpha_0)\zeta^4/24 + \dots \quad (3.12)$$

The series (3.12) can be inverted in the neighborhood of $\zeta = 0$ ($z = \alpha_0$) that is ζ can be expressed as the following series of δ

$$\zeta = i\delta + \lambda_3(\alpha_0)\delta^2/6 + \{\lambda_4(\alpha_0)/24 - (5/72)\lambda_3^2\}i\delta^3 + \dots \quad (3.13)$$

At this point we are ready to rewrite I_0 according to (3.4). We obtain:

$$\begin{aligned} I_0 &= (n/2\pi i) \int_{\mathcal{P}_0} \exp\{n[K(z) - zt]\} dz \\ &= (n/2\pi i) \exp\{n[K(\alpha_0) - \alpha_0 t]\} \int_{\mathcal{P}_0} e^{-n\gamma} dz \\ &= (n/2\pi i) \frac{\exp\{n[K(\alpha_0) - \alpha_0 t]\}}{[K''(\alpha_0)]^{1/2}} \int_{-A}^B e^{-n\delta^2/2} \frac{d\zeta}{d\delta} d\delta, \end{aligned}$$

and from (3.13)

$$I_0 = (n/2\pi i) \frac{\exp\{n[K(\alpha_0) - \alpha_0 t]\}}{[K''(\alpha_0)]^{1/2}} \times \int_{-A}^B e^{-n\delta^2/2} \left\{ i + \lambda_3(\alpha_0)\delta/3 + i \left[\lambda_4(\alpha_0)/8 - \frac{5}{24}\lambda_3^2(\alpha_0) \right] \delta^2 + \dots \right\} d\delta, \quad (3.14)$$

where A and B are two positive numbers which correspond to the values of δ at z_2 and z_1 . By applying Watson's Lemma (see below) to (3.14), one finally obtains the following asymptotic expansion

$$f_n(t) = \left[\frac{n}{2\pi K''(\alpha_0)} \right]^{1/2} \exp\{n[K(\alpha_0 t) - \alpha_0 t]\} \times \left\{ 1 + \frac{1}{n} \left[\frac{1}{8}\lambda_4(\alpha_0) - \frac{5}{24}\lambda_3^2(\alpha_0) \right] + \dots \right\} \quad (3.15)$$

where α_0 is determined by the *saddlepoint equation*

$$K'(\alpha_0) = t, \quad (3.16)$$

and

$$\lambda_3(\alpha_0) = K'''(\alpha_0)/[K''(\alpha_0)]^{3/2} \quad (3.17)$$

$$\lambda_4(\alpha_0) = K^{(iv)}(\alpha_0)/[K''(\alpha_0)]^2 \quad (3.18)$$

are standardized measures of skewness and kurtosis respectively. The leading term of the expansion (3.15)

$$g_n(t) = \left[\frac{n}{2\pi K''(\alpha_0)} \right]^{1/2} \exp\{n[K(\alpha_0) - \alpha_0 t]\} \quad (3.19)$$

is called the *saddlepoint approximation*.

For the sake of completeness we give here a modification of Watson's lemma due to Jeffreys and Jeffreys (1950) which is used in the final step in the derivation of (3.15).

Lemma (Watson, 1948; Jeffreys and Jeffreys, 1950; Daniels, 1954).

If $\psi(\zeta)$ is analytic in a neighborhood of $\zeta = 0$ and bounded for real $\zeta = \delta$ in an interval $-A \leq \delta \leq B$ with $A > 0$ and $B > 0$, then

$$(n/2\pi)^{1/2} \int_{-A}^B e^{-n\delta^2/2} \psi(\delta) d\delta \sim \psi(0) + \frac{1}{2n} \psi''(0) + \dots + \frac{1}{(2n)^r} \frac{\psi^{(2r)}(0)}{r!} + \dots$$

is an asymptotic expansion in powers of n^{-1} .

In the following remarks we discuss some aspects of the saddlepoint approximation in some detail.

Remark 3.2: Error of the approximation.

From $f_n(t) = g_n(t)[1 + 0(1/n)]$ one can see that $g_n(t) \geq 0$ and that the relative error is of order n^{-1} . This is the most important property of this approximation and a major advantage with respect to Edgeworth expansions. Moreover, Daniels (1954), p. 640 ff. showed that for a wide class of underlying densities, the coefficient of the term of order n^{-1} doesn't depend on t . Thus, in such cases the relative error is of order n^{-1} *uniformly*. cf. Jensen (1988).

Since $g_n(t)$ doesn't necessarily integrate to 1, one can renormalize the approximation by dividing by $C_n = \int g_n(t)dt$. This operation comes out naturally by using an alternative derivation of the saddlepoint approximation proposed by Hampel (1973) which is based on the expansion of $f'_n(t)/f_n(t)$ rather than $f_n(t)$; see sections 4.2 and 5.2. By renormalization one actually improves the order of the approximation by getting a relative error of order $0(n^{-3/2})$ for values t in the range $t - \mu = 0(n^{-1/2})$, where $\mu = \int x dF(x)$. To see this write

$$f_n(t) = g_n(t)[1 + b(t)/n + 0(n^{-2})]$$

and

$$\begin{aligned} g_n(t) &= f_n(t)[1 - b(t)/n + 0(n^{-2})] \\ &= f_n(t)[1 - b(\mu)/n - (t - \mu)b'(\mu)/n + 0((t - \mu)^2/n) \\ &\quad + 0(n^{-2})] \\ &= f_n(t)[1 - b(\mu)/n - (t - \mu)b'(\mu)/n + 0(n^{-2})]. \end{aligned} \quad (3.20)$$

Therefore

$$C_n = \int g_n(t)dt = 1 - b(\mu)/n + 0(n^{-2}) \quad (3.21)$$

and from (3.20) and (3.21) by using $t - \mu = 0(n^{-1/2})$

$$g_n(t)/C_n = f_n(t)[1 + 0(n^{-3/2})].$$

If in addition the relative error is uniform of order n^{-1} , that is $b(t)$ does not depend on t , the relative error after renormalization is of order $0(n^{-2})$.

Remark 3.3: Exact saddlepoint approximations.

It turns out that in some cases the saddlepoint approximation $g_n(t)$ is exact or exact up to normalization. Daniels (1954,1980) proved that there are only three underlying densities $f(x)$ for which this is the case, namely the normal, the gamma, and the inverse normal distribution. In the case of the normal the leading term is exact and the higher order terms are zero. In the other two cases the leading term is exact up to a constant and the higher order terms are different from zero but independent of t and can therefore be included in the normalization constant. Moreover, Blaesild and Jensen (1985) showed that $f(x)$ has an exact saddlepoint approximation if and only if $f(x)$ is a reproductive exponential model.

Remark 3.4: Computational issue.

In order to compute the saddlepoint approximation $g_n(t)$, one has to solve the implicit saddlepoint equation (3.16) for each value t . Since $K'(\cdot)$ has an integral form, this can be computational intensive in multidimensional problems; see chapter 4. However, if the density has to be approximated on an interval $[t_1, t_2]$ one can find the saddlepoint $\alpha_0^{(1)}$ corresponding to t_1 and use this as a starting point for the next value t , and so on. Moreover, when α_0 as a function of t is monotone as in the case of the mean and the density $f_n(t)$ does not have

to be approximated at equally spaced points, one can just take a number of values α_0 and compute the corresponding values $t = K'(\alpha_0)$. cf. section 7.1.

Remark 3.5: Lattice underlying distribution.

A saddlepoint approximation for the mean can be derived when the underlying distribution is lattice; see Daniels (1983, 1987), Gamkrelidze (1980).

The numerical examples in section 3.5 show the great accuracy of the saddlepoint approximation for the mean. The same pattern can be found for more important and complex situations; cf. chapters 4, 6 and 7.

3.4. RELATIONSHIP WITH THE METHOD OF CONJUGATE DISTRIBUTIONS

Up to this point, the approximation to the density of the mean $f_n(t)$ has been derived using the method of steepest descent and the saddlepoint approximation. In this section we develop the approximation using the idea of conjugate densities and the normal approximation. Although both approaches lead to the same results, we can gain new insight into the approximation through the conjugate density.

Probably the simplest and most common approximation to the density of the mean is the normal approximation. This approximation works very well near the center of the distribution but breaks down in the tail. The idea here is to re-center our density at the point of interest t by means of a conjugate density and then to use a normal approximation in the re-centered problem. The approximation in the re-centered problem is then converted to an approximation for $f_n(t)$.

To be more specific, we introduce the conjugate density

$$h_t(x) = c(t) \exp\{\alpha(t)(x - t)\} f(x) \quad (3.22)$$

where $c(t)$ is chosen so that h_t is a density and $\alpha(t)$ is chosen so that

$$\int (x - t) \exp\{\alpha(t)(x - t)\} f(x) dx = 0 \quad \text{i.e. } E_{h_t} X = t. \quad (3.23)$$

The variance of X under h_t is denoted by $\sigma^2(t)$, i.e.

$$\sigma^2(t) = c(t) \int (x - t)^2 \exp\{\alpha(t)(x - t)\} f(x) dx.$$

Conjugate or associated densities are well known in probability (Khinchin, 1949; Feller, 1971, p. 518) and arise naturally in information theory (Kullback 1960).

To illustrate the way in which the conjugate density operates, it is informative to look at some graphs. Exhibit 3.5 shows the situation when $f(x)$ is uniform on $[-1, 1]$.

The plots give the conjugate densities centered at .3, .5, .7, .9. As we move to the right $\alpha(t)$ increases in order to put sufficient mass in the interval $[t, 1]$. The conjugate is plotted only for the interval $[0, 1]$. On $[-1, 0]$ all four conjugates are relatively flat and close together.

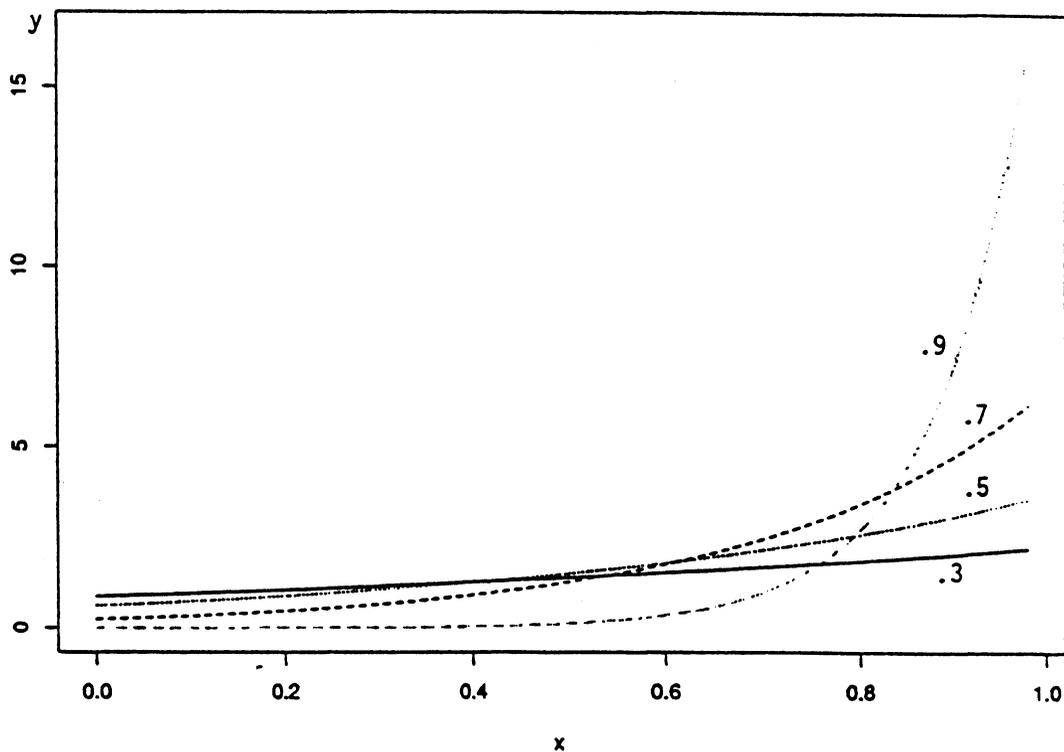


Exhibit 3.5

Uniform conjugate at $t = .3, .5, .7, .9$.

As a second example consider the extreme density with $f(x) = \exp\{x - \exp(x)\}$. Conjugate densities are plotted in Exhibit 3.6 for t values $-7, -3, 0, 0.5, 2$ along with f . As can be seen the shape of the density is changed substantially as the values of t are varied.

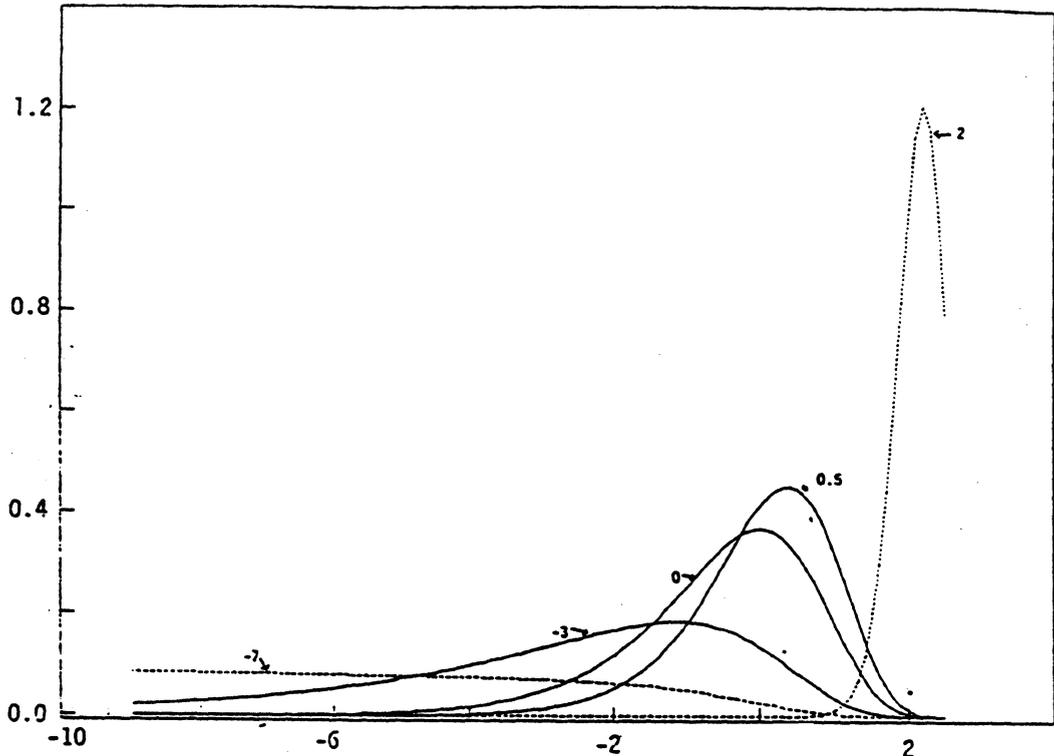


Exhibit 3.6

Extreme conjugate at $t = -7, -3, 0, 0.5, 2$.

Before proceeding with the development of the approximation, it is important to link the notation of the conjugate density with that of the cumulant generating functions that have been used up to now. Recall from section 3.3 that

$$K(\alpha) = \log \int e^{\alpha x} f(x) dx$$

and the saddlepoint (at a point t) is the solution of $K'(\alpha_0) = t$ (see (3.16))

$$\frac{\int x \exp(\alpha_0 x) f(x) dx}{\int \exp(\alpha_0 x) f(x) dx} = t$$

or

$$\frac{\int (x - t) \exp(\alpha_0 x) f(x) dx}{\int \exp(\alpha_0 x) f(x) dx} = 0.$$

Multiplying both numerator and denominator by $\exp(-\alpha_0 t)$ we have

$$\frac{\int (x - t) \exp\{\alpha_0(x - t)\} f(x) dx}{\int \exp\{\alpha_0(x - t)\} f(x) dx} = 0$$

or

$$\int (x - t) \exp\{\alpha_0(x - t)\} f(x) dx = 0.$$

By comparing this with (3.23), we see that $\alpha_0 = \alpha(t)$, i.e. $\alpha(t)$ is the saddlepoint at t . Moreover, $c^{-1}(t) = \int \exp\{\alpha(t)(x - t)\} f(x) dx = \exp\{K(\alpha(t)) - \alpha(t)t\}$, hence $-\log c(t) =$

$K(\alpha(t)) - \alpha(t)t = K(\alpha_0) - \alpha_0 t$. Similarly

$$\begin{aligned} K''(\alpha_0) &= K''(\alpha(t)) \\ &= \frac{\int \exp\{\alpha(t)x\} f(x) dx \int x^2 \exp\{\alpha(t)x\} f(x) dx - (\int x \exp\{\alpha(t)x\} f(x) dx)^2}{(\int \exp\{\alpha(t)x\} f(x) dx)^2} \end{aligned}$$

Multiplying numerator and denominator by $(e^{-\alpha(t)t})^2$ we have

$$K''(\alpha(t)) = E_{h_t} X^2 - (E_{h_t} X)^2 = \sigma^2(t).$$

To summarize, we have

$$\alpha(t) = \alpha_0, \quad -\log c(t) = K(\alpha_0) - \alpha_0 t, \quad \sigma^2(t) = K''(\alpha_0). \quad (3.24)$$

The next step is to consider the density of the mean under h_t , say $h_{t,n}$ and relate it to f_n , the density of the mean under f . We can write f_n as follows:

$$\begin{aligned} f_n(t) &= n \int \cdots \int f\left(nt - \sum_1^{n-1} x_i\right) \prod_1^{n-1} f(x_i) dx_1 \cdots dx_{n-1} \\ &= c^{-n}(t) n \int \cdots \int c(t) \exp\left\{\alpha(t)\left(nt - \sum_1^{n-1} x_i - t\right)\right\} f\left(nt - \sum_1^{n-1} x_i\right) \\ &\quad \prod_1^{n-1} c(t) \exp\{\alpha(t)(x_i - t)\} f(x_i) dx_1 \cdots dx_{n-1} \\ &= c^{-n}(t) n \int \cdots \int h_t\left(nt - \sum_1^{n-1} x_i\right) \prod_1^{n-1} h_t(x_i) dx_1 \cdots dx_{n-1} \\ &= c^{-n}(t) h_{t,n}(t). \end{aligned}$$

Hence we conclude

$$f_n(t) = c^{-n}(t) h_{t,n}(t) \quad (3.25)$$

This centering equation provides the link for relating the approximation to $h_{t,n}(t)$ to the desired approximation of $f_n(t)$.

The final step is to approximate $h_{t,n}(t)$. Now $h_{t,n}(t)$ is the density of \bar{X} under h_t where the X_i 's have mean t and variance $\sigma^2(t)$. Expression (2.10) gives the Edgeworth expansion at the origin for a standardized variable with

$$\begin{aligned} \lambda_3 &\equiv \lambda_3(t) = E_{h_t}(X - t)^3 / \sigma^3(t) \\ \lambda_4 &\equiv \lambda_4(t) = E_{h_t}(X - t)^4 / \sigma^4(t). \end{aligned}$$

The approximation to the density of $\sqrt{n}(\bar{X} - t)/\sigma(t)$ at 0 is

$$\frac{1}{\sqrt{2\pi}} \left[1 + \frac{1}{n} \left(\frac{\lambda_4(t)}{8} - \frac{5\lambda_3^2(t)}{24} \right) + o\left(\frac{1}{n^2}\right) \right].$$

From this it follows that

$$h_{t,n}(t) = \sqrt{\frac{n}{2\pi}} \frac{1}{\sigma(t)} \left(1 + o\left(\frac{1}{n}\right) \right)$$

and

$$f_n(t) = \sqrt{\frac{n}{2\pi}} \frac{c^{-n}(t)}{\sigma(t)} \left(1 + o\left(\frac{1}{n}\right) \right).$$

This leads to the small sample (or saddlepoint) approximation for the mean

$$g_n(t) = \sqrt{\frac{n}{2\pi}} c^{-n}(t)/\sigma(t). \quad (3.26)$$

Using (3.24), it can be seen that this is exactly (3.19).

It should be noted that (3.26) is obtained by shifting the underlying density to the point of interest, using a normal approximation at the mean, and then using the centering lemma. The process can be likened to using a low order Taylor's expansion at many points rather than one high order expansion at a single point as in an Edgeworth expansion. The approximation (3.26) can be thought of as a local normal approximation.

It is worth asking whether the form of the conjugate density is important for the argument above. If we look at the argument leading to the centering lemma, the exponential in the conjugate enabled us to go from f to h_t . It is hard to see how to obtain the necessary link between f_n and $h_{t,n}$ with any other form. From another point of view, start with the density f and ask for the density closest to f in Kullback-Liebler distance which is constrained to have mean t . Kullback (1960) shows that $h_t(x)$ is this density. By using the conjugate density, we have embedded our problem within an exponential family. This perspective becomes very useful in applying our techniques in multiparameter problems as we show in chapter 6.

3.5. EXAMPLES

We now consider numerical results from using approximation (3.26) with several underlying densities. In each case, the $\alpha(t)$ has been evaluated at a grid suitable for the underlying density. To solve this non-linear equation, we have used a secant style root finder, such as C05AJ5 in the NAG library. Given $\alpha(t)$, $c(t)$ and $\sigma(t)$ can be determined using numerical integration.

As a first example, let $f(x)$ be uniform on $[-1, 1]$ and consider the density of the mean for $n = 5$. The exact density is given by

$$f_n(t) = \frac{n^n}{2^n(n-1)!} \sum_{i=0}^n (-1)^i \binom{n}{i} \left\langle 1 - t - \frac{2i}{n} \right\rangle^{n-1} \quad |t| \leq 1$$

where $\langle z \rangle = z$ for $z \geq 0$ and $= 0$ for $z < 0$, cf. section 2.7. Exhibit 3.7 gives the exact and approximate density for some selected points. The error is measured by percent relative error = 100 (approximate-exact)/exact.

As can be seen the results are very accurate even for the extreme tail. They are clearly superior to the approximations obtained by Edgeworth expansions (see section 2.7). The maximum percent relative error is 1.65% for $t = .95$ and is actually smaller for values beyond .95. Such accuracy would certainly be more than adequate for all applications. The

property that the percent relative error stays bounded even in the extreme tail seems to be quite general with these approximations.

t	Exact	Approximate	% relative error
0.00	1.4974e0	1.4945e0	-0.19
0.05	1.4731e0	1.4700e0	-0.21
0.10	1.4022e0	1.3988e0	-0.24
0.15	1.2900e0	1.2872e0	-0.22
0.20	1.1458e0	1.1446e0	0.10
0.25	9.8216e-1	9.8277e-1	0.06
0.30	8.1217e-1	8.1361e-1	0.18
0.35	6.4687e-1	6.4836e-1	0.23
0.40	4.9479e-1	4.9622e-1	0.29
0.45	3.6204e-1	3.6369e-1	0.46
0.50	2.5228e-1	2.5428e-1	0.79
0.55	1.6673e-1	1.6872e-1	1.19
0.60	1.0417e-1	1.0548e-1	1.26
0.65	6.1061e-2	6.1494e-2	0.71
0.70	3.2959e-2	3.2902e-2	-0.17
0.75	1.5895e-2	1.5736e-2	-1.00
0.80	6.5104e-3	6.4131e-3	-1.49
0.85	2.0599e-3	2.0262e-3	-1.64
0.90	4.0690e-4	4.0085e-4	-1.49
0.95	2.5431e-5	2.5011e-5	-1.65
0.96	1.0417e-5	1.0245e-5	-1.65
0.97	3.2959e-6	3.2422e-6	-1.63
0.98	6.5105e-7	6.4780e-7	-0.50

Exhibit 3.7

Exact and approximate density for mean with uniform on $[-1, 1]$ and $n = 5$.

In many situations, it is not the density itself which is of interest but rather the tail area. In order to compute the tail area, it is possible to integrate the approximate density over a grid of points. However it is much easier and equally accurate to use a tail area approximation developed by Lugannani and Rice (1980) and discussed by Daniels (1987) and Tingley (1987). This approximation is discussed in some detail in section 6.1 but for completeness we give it here.

$$P[\bar{X} \geq t] \approx 1 - \Phi((2n \log c(t))^{1/2}) + c^{-n}(t) \left[\frac{1}{(\sigma(t)(\alpha(t)))} - \frac{1}{(2 \log c(t))^{1/2}} \right] / (n2\pi)^{1/2} \quad (3.27)$$

The approximation is very useful since to get the tail beyond t , we need only evaluate $\alpha(t)$, $\sigma(t)$ and $c(t)$. The numerical evidence is that the approximation gives an accuracy comparable to that obtained by numerical integration (cf. Daniels, 1987).

For the uniform case, Exhibit 3.8 gives the exact and approximate tail areas. Because

of symmetry, only upper tail areas are given.

t	Exact	Approximate	% relative error
0.05	4.2554e-1	4.2547e-1	-0.01
0.10	3.5347e-1	3.5338e-1	-0.02
0.15	2.8601e-1	2.8593e-1	-0.03
0.20	2.2500e-1	2.2493e-1	-0.03
0.25	1.7175e-1	1.7165e-1	-0.06
0.30	1.2689e-1	1.2675e-1	-0.11
0.35	9.0451e-2	9.0291e-2	-0.18
0.40	6.1979e-2	6.1837e-2	-0.23
0.45	4.0648e-2	4.0537e-2	-0.27
0.50	2.5391e-2	2.5298e-2	-0.36
0.55	1.5016e-2	1.4922e-2	-0.63
0.60	8.3333e-3	8.2371e-3	-1.15
0.65	4.2742e-3	4.1961e-3	-1.83
0.70	1.9775e-3	1.9307e-3	-1.37
0.75	7.9473e-4	7.7478e-4	-2.51
0.80	2.6042e-4	2.5486e-4	-2.13
0.85	6.1798e-5	6.0932e-5	-1.40
0.90	8.1380e-6	8.0784e-6	-0.73

Exhibit 3.8

Exact and approximate tail areas for mean with uniform on $[-1, 1]$ and $n = 5$.

As might be expected from the results with the density, we obtain the same order of accuracy for the tail areas. Again the relative error is much smaller than that of the Edgeworth approximation where it can reach 20%; cf section 2.7. These results indicate that the saddlepoint approximation gives very accurate results for small values of n .

To give some further numerical results, we now turn to the extreme case of $n = 1$. Since our expansion is asymptotic, there is no a priori reason to believe the approximation should work well in this case. For $n = 1$, the approximation (3.26) simply approximates the underlying density f . This point is discussed further in section 7.2. The convenient feature with $n = 1$ is that we can easily compute the exact results as a comparison. As a first example, we consider the case of the extreme density, $f(x) = \exp(x - \exp(x))$. Since f is asymmetric, we consider behavior in the upper and lower tails. The results are given in Exhibit 3.9.

t	Exact	Approximate	% relative error
-9.0	1.2339e-4	1.5693e-4	27.18
-8.0	3.3535e-4	4.4526e-4	32.78
-7.0	9.1105e-4	1.2270e-3	34.68
-6.0	2.4726e-3	3.2752e-3	32.46
-5.0	6.6927e-3	8.4597e-3	26.40
-4.0	1.7983e-2	2.1138e-2	17.54
-3.0	4.7369e-2	5.1054e-2	7.78
-2.5	7.5616e-2	7.8207e-2	3.43
-2.0	1.1821e-1	1.1818e-1	-0.02
-1.5	1.7851e-1	1.7451e-1	-2.24
-1.0	2.5465e-1	2.4639e-1	-3.24
-0.5	3.3070e-1	3.1970e-1	-3.33
0.0	3.6788e-1	3.6091e-1	-1.89
0.5	3.1704e-1	3.3015e-1	4.13
1.0	1.7937e-1	1.6912e-1	-5.71
1.5	5.0707e-2	4.0221e-2	-20.68
2.0	4.5663e-3	6.6205e-3	44.99
2.5	6.2366e-5	2.8859e-5	-53.73

Exhibit 3.9

Exact and approximate density for the extreme for $n = 1$.

Certainly the results are not as accurate as with $n = 5$. However the approximate density gives fairly reasonable results even as we go out into the tails. We can obtain another view by examining tail areas for the extreme with $n = 1$. The values for the negative t are lower tail areas and are upper tail errors for $t \geq 0$ (Exhibit 3.10).

t	Exact	Approximate	% relative error
-9.0	1.2338e-4	1.2382e-4	-0.35
-8.0	3.3540e-4	3.7148e-4	-9.71
-7.0	9.1147e-4	1.0723e-3	-15.00
-6.0	2.4757e-3	2.9806e-3	-16.94
-5.0	6.7153e-3	7.9918e-3	-15.97
-4.0	1.8149e-2	2.0710e-2	-12.37
-3.0	4.8568e-2	5.2058e-2	-6.70
-2.5	7.8806e-2	8.1733e-2	-3.58
-2.0	1.2658e-1	1.2751e-1	-0.73
-1.5	1.9999e-1	1.9662e-1	1.71
-1.0	3.0780e-1	2.8111e-1	9.49
0.0	3.6788e-1	3.7455e-1	-1.78
0.5	1.9230e-1	2.0463e-1	-6.03
1.0	6.5988e-2	5.9525e-2	10.86
1.5	1.1314e-2	9.7135e-3	16.48
2.0	6.1798e-4	7.7569e-4	-20.33
2.5	5.1193e-6	2.4462e-6	109.28

Exhibit 3.10

Exact and approximate tail area for extreme with $n = 1$.

The results here are remarkably good except in the extreme upper tail. To complete, the section we look at the results with $n = 1$ for three densities; the uniform on $[-1, 1]$ (Exhibit 3.11), a sum of exponentials where the number in the sum is Poisson with parameter 4 (Exhibit 3.12) and a density $f(x) = 1 + \cos 4\pi x$ on $[0, 1]$ (Exhibit 3.13). These were chosen to show the varying degrees of accuracy one can get.

t	Exact	Approximate	% relative error
0.05	0.475	0.4706	-0.92
0.10	0.450	0.4414	-1.91
0.15	0.425	0.4124	-2.96
0.20	0.400	0.3838	-4.05
0.25	0.375	0.3556	-5.16
0.30	0.350	0.3281	-6.27
0.35	0.325	0.3011	-7.34
0.40	0.300	0.2750	-8.34
0.45	0.275	0.2496	-9.22
0.50	0.250	0.2252	-9.91
0.55	0.225	0.2017	-10.35
0.60	0.200	0.1791	-10.44
0.65	0.175	0.1574	-10.07
0.70	0.150	0.1363	-9.11
0.75	0.125	0.1156	-7.48
0.80	0.100	0.09476	-5.24
0.85	0.075	0.07300	-2.66
0.90	0.050	0.05002	0.05
0.95	0.025	0.02562	2.47
0.96	0.020	0.02061	3.03
0.97	0.015	0.01555	3.65
0.98	0.010	0.01054	5.40

Exhibit 3.11

Exact and approximate tail area for uniform with $n = 1$.

t	Exact	Approximate	% relative error
5	0.3070	0.3079	-0.30
7	0.1425	0.1430	-0.32
9	0.05950	0.05961	-0.33
11	0.02277	0.02284	-0.34
12	0.01373	0.01378	-0.35
13	0.008151	0.008180	-0.36
14	0.004773	0.004790	-0.36
15	0.002759	0.002769	-0.38
16	0.001576	0.001582	-0.39
17	0.000890	0.000894	-0.41

Exhibit 3.12

Exact and approximate tail area for Poisson sum of exponentials with $n = 1$.

t	Exact	Approximate	% relative error
0.55	0.4032	0.4443	10.19
0.60	0.3243	0.3895	20.10
0.65	0.2743	0.3365	22.66
0.70	0.2532	0.2860	12.96
0.75	0.2500	0.2390	- 4.41
0.80	0.2468	0.1960	-20.57
0.85	0.2257	0.1581	-29.94
0.90	0.1757	0.1273	-27.51
0.95	0.09677	0.1134	17.15

Exhibit 3.13

Exact and approximate tail area for $f(x) = 1 + \cos 4\pi x$
on $[0, 1]$ with $n = 1$.

From the displays, we can see that the approximation is very accurate for the Poisson sum of exponentials, quite accurate for the uniform and less accurate for the cosine density. Since the approximation is based on a local normal approximation, the quality of the approximation is determined by the degree to which $(\bar{X} - t)/\sigma(t)$ under h_t is approximated by a normal. In some recent work, Field and Massam (1987) have developed a diagnostic function for the accuracy of the approximation.

We conclude from Exhibit 3.13 that even by taking multimodal f and choosing $n = 1$, we cannot make the approximation breakdown. Our evidence is that we have an asymptotic approximation which gives reasonable results for $n = 1$.