

## **NAG Fortran Library Chapter Introduction**

### **S – Approximations of Special Functions**

#### **Contents**

<b>1</b>	<b>Scope of the Chapter</b> .....	<b>2</b>
<b>2</b>	<b>Background to the Problems</b> .....	<b>2</b>
2.1	Functions of a Single Real Argument .....	2
2.2	Approximations to Elliptic Integrals .....	4
2.3	Bessel and Airy Functions of a Complex Argument .....	5
<b>3</b>	<b>Recommendations on Choice and Use of Available Routines</b> .....	<b>5</b>
3.1	Elliptic Integrals .....	5
3.2	Bessel and Airy Functions .....	6
<b>4</b>	<b>Index</b> .....	<b>6</b>
<b>5</b>	<b>Routines Withdrawn or Scheduled for Withdrawal</b> .....	<b>7</b>
<b>6</b>	<b>References</b> .....	<b>7</b>

## 1 Scope of the Chapter

This chapter is concerned with the provision of some commonly occurring physical and mathematical functions.

## 2 Background to the Problems

The majority of the routines in this chapter approximate real-valued functions of a single real argument, and the techniques involved are described in Section 2.1. In addition the chapter contains routines for elliptic integrals (see Section 2.2), Bessel and Airy functions of a complex argument (see Section 2.3), exponential of a complex argument, and complementary error function of a complex argument.

### 2.1 Functions of a Single Real Argument

Most of the routines for functions of a single real argument have been based on truncated Chebyshev expansions. This method of approximation was adopted as a compromise between the conflicting requirements of efficiency and ease of implementation on many different machine ranges. For details of the reasons behind this choice and the production and testing procedures followed in constructing this chapter see Schonfelder (1976).

Basically, if the function to be approximated is  $f(x)$ , then for  $x \in [a, b]$  an approximation of the form

$$f(x) = g(x) \sum'_{r=0} C_r T_r(t)$$

is used ( $\sum'$  denotes, according to the usual convention, a summation in which the first term is halved), where  $g(x)$  is some suitable auxiliary function which extracts any singularities, asymptotes and, if possible, zeros of the function in the range in question and  $t = t(x)$  is a mapping of the general range  $[a, b]$  to the specific range  $[-1, +1]$  required by the Chebyshev polynomials,  $T_r(t)$ . For a detailed description of the properties of the Chebyshev polynomials see Clenshaw (1962) and Fox and Parker (1968).

The essential property of these polynomials for the purposes of function approximation is that  $T_n(t)$  oscillates between  $\pm 1$  and it takes its extreme values  $n + 1$  times in the interval  $[-1, +1]$ . Therefore, provided the coefficients  $C_r$  decrease in magnitude sufficiently rapidly the error made by truncating the Chebyshev expansion after  $n$  terms is approximately given by

$$E(t) \simeq C_n T_n(t).$$

That is, the error oscillates between  $\pm C_n$  and takes its extreme value  $n + 1$  times in the interval in question. Now this is just the condition that the approximation be a mini-max representation, one which minimizes the maximum error. By suitable choice of the interval,  $[a, b]$ , the auxiliary function,  $g(x)$ , and the mapping of the independent variable,  $t(x)$ , it is almost always possible to obtain a Chebyshev expansion with rapid convergence and hence truncations that provide near mini-max polynomial approximations to the required function. The difference between the true mini-max polynomial and the truncated Chebyshev expansion is seldom sufficiently great enough to be of significance.

The evaluation of the Chebyshev expansions follows one of two methods. The first and most efficient, and hence the most commonly used, works with the equivalent simple polynomial. The second method, which is used on the few occasions when the first method proves to be unstable, is based directly on the truncated Chebyshev series, and uses backward recursion to evaluate the sum. For the first method, a suitably truncated Chebyshev expansion (truncation is chosen so that the error is less than the *machine precision*) is converted to the equivalent simple polynomial. That is, we evaluate the set of coefficients  $b_r$  such that

$$y(t) = \sum_{r=0}^{n-1} b_r t^r = \sum'_{r=0}^{n-1} C_r T_r(t).$$

The polynomial can then be evaluated by the efficient Horner's method of nested multiplications,

$$y(t) = (b_0 + t(b_1 + t(b_2 + \dots t(b_{n-2} + tb_{n-1})))) \dots).$$

This method of evaluation results in efficient routines but for some expansions there is considerable loss of accuracy due to cancellation effects. In these cases the second method is used. It is well known that if

$$\begin{aligned} b_{n-1} &= C_{n-1} \\ b_{n-2} &= 2tb_{n-1} + C_{n-2} \\ b_j &= 2tb_{j+1} - b_{j+2} + C_j, \quad j = n-3, n-4, \dots, 0 \end{aligned}$$

then

$$\sum_{r=0}^n C_r T_r(t) = \frac{1}{2}(b_0 - b_2)$$

and this is always stable. This method is most efficiently implemented by using three variables cyclically and explicitly constructing the recursion.

That is,

$$\begin{aligned} \alpha &= C_{n-1} \\ \beta &= 2t\alpha + C_{n-2} \\ \gamma &= 2t\beta - \alpha + C_{n-3} \\ \alpha &= 2t\gamma - \beta + C_{n-4} \\ \beta &= 2t\alpha - \gamma + C_{n-5} \\ &\vdots \\ \text{say } \alpha &= 2t\gamma - \beta + C_2 \\ \beta &= 2t\alpha - \gamma + C_1 \\ y(t) &= t\beta - \alpha + \frac{1}{2}C_0 \end{aligned}$$

The auxiliary functions used are normally functions compounded of simple polynomial (usually linear) factors extracting zeros, and the primary compiler-provided functions, sin, cos, ln, exp, sqrt, which extract singularities and/or asymptotes or in some cases basic oscillatory behaviour, leaving a smooth well-behaved function to be approximated by the Chebyshev expansion which can therefore be rapidly convergent.

The mappings of  $[a, b]$  to  $[-1, +1]$  used range from simple linear mappings to the case when  $b$  is infinite, and considerable improvement in convergence can be obtained by use of a bilinear form of mapping. Another common form of mapping is used when the function is even; that is, it involves only even powers in its expansion. In this case an approximation over the whole interval  $[-a, a]$  can be provided using a mapping  $t = 2(x/a)^2 - 1$ . This embodies the evenness property but the expansion in  $t$  involves all powers and hence removes the necessity of working with an expansion with half its coefficients zero.

For many of the routines an analysis of the error in principle is given, namely, if  $E$  and  $\nabla$  are the absolute errors in function and argument and  $\epsilon$  and  $\delta$  are the corresponding relative errors, then

$$E \simeq |f'(x)|\nabla$$

$$E \simeq |xf'(x)|\delta$$

$$\epsilon \simeq \left| \frac{xf'(x)}{f(x)} \right| \delta.$$

If we ignore errors that arise in the argument of the function by propagation of data errors, etc., and consider only those errors that result from the fact that a real number is being represented in the computer in floating-point form with finite precision, then  $\delta$  is bounded and this bound is independent of the magnitude of  $x$ . For example, on an 11-digit machine

$$|\delta| \leq 10^{-11}.$$

(This of course implies that the absolute error  $\nabla = x\delta$  is also bounded but the bound is now dependent on  $x$ .) However, because of this the last two relations above are probably of more interest. If possible the relative error propagation is discussed; that is, the behaviour of the error amplification factor  $|xf'(x)/f(x)|$  is described, but in some cases, such as near zeros of the function which cannot be extracted explicitly, absolute error in the result is the quantity of significance and here the factor  $|xf'(x)|$  is described. In general, testing of the functions has shown that their error behaviour follows fairly well these theoretical error behaviours. In regions where the error amplification factors are less than or of the order of one, the

errors are slightly larger than the above predictions. The errors are here limited largely by the finite precision of arithmetic in the machine, but  $\epsilon$  is normally no more than a few times greater than the bound on  $\delta$ . In regions where the amplification factors are large, of order ten or greater, the theoretical analysis gives a good measure of the accuracy obtainable.

It should be noted that the definitions and notations used for the functions in this chapter are all taken from Abramowitz and Stegun (1972). Users are strongly recommended to consult this book for details before using the routines in this chapter.

## 2.2 Approximations to Elliptic Integrals

Four functions provided here are symmetrised variants of the classic elliptic integrals. These alternative definitions have been suggested by Carlson (1965), Carlson (1977a) and Carlson (1977b) and he also developed the basic algorithms used in this chapter.

The standard integral of the first kind is represented by

$$R_F(x, y, z) = \frac{1}{2} \int_0^\infty \frac{dt}{\sqrt{(t+x)(t+y)(t+z)}},$$

where  $x, y, z \geq 0$  and at most one may be equal to zero.

The normalisation factor,  $\frac{1}{2}$ , is chosen so as to make

$$R_F(x, x, x) = 1/\sqrt{x}.$$

If any two of the variables are equal,  $R_F$  degenerates into the second function

$$R_C(x, y) = R_F(x, y, y) = \frac{1}{2} \int_0^\infty \frac{dt}{\sqrt{t+x}(t+y)},$$

where the argument restrictions are now  $x \geq 0$  and  $y \neq 0$ .

This function is related to the logarithm or inverse hyperbolic functions if  $0 < y < x$ , and to the inverse circular functions if  $0 \leq x \leq y$ .

The integrals of the second kind are defined by

$$R_D(x, y, z) = \frac{3}{2} \int_0^\infty \frac{dt}{\sqrt{(t+x)(t+y)(t+z)}^3}$$

with  $z > 0$ ,  $x \geq 0$  and  $y \geq 0$ , but only one of  $x$  or  $y$  may be zero.

The function is a degenerate special case of the integral of the third kind

$$R_J(x, y, z, \rho) = \frac{3}{2} \int_0^\infty \frac{dt}{\sqrt{(t+x)(t+y)(t+z)(t+\rho)}}$$

with  $\rho \neq 0$  and  $x, y, z \geq 0$  with at most one equality holding. Thus  $R_D(x, y, z) = R_J(x, y, z, z)$ . The normalisation of both these functions is chosen so that

$$R_D(x, x, x) = R_J(x, x, x, x) = 1/(x\sqrt{x}).$$

The algorithms used for all these functions are based on duplication theorems. These allow a recursion system to be established which constructs a new set of arguments from the old using a combination of arithmetic and geometric means. The value of the function at the original arguments can then be simply related to the value at the new arguments. These recursive reductions are used until the arguments differ from the mean by an amount small enough for a Taylor series about the mean to give sufficient accuracy when retaining terms of order less than six. Each step of the recurrences reduces the difference from the mean by a factor of four, and as the truncation error is of order six, the truncation error goes like  $(4096)^{-n}$ , where  $n$  is the number of iterations.

The above forms can be related to the more traditional canonical forms (see Section 17.2 in Abramowitz and Stegun (1972)).

If we write  $q = \cos^2 \phi$ ,  $r = 1 - m \cdot \sin^2 \phi$ ,  $s = 1 + n \cdot \sin^2 \phi$ , where  $0 < \phi \leq \frac{1}{2}\pi$ , we have

the elliptic integral of the first kind:

$$F(\phi|m) = \int_0^{\sin \phi} (1-t^2)^{-1/2}(1-mt^2)^{-1/2} dt = \sin \phi \cdot R_F(q, r, 1);$$

the elliptic integral of the second kind:

$$\begin{aligned} E(\phi|m) &= \int_0^{\sin \phi} (1-t^2)^{-1/2}(1-mt^2)^{1/2} dt \\ &= \sin \phi \cdot R_F(q, r, 1) - \frac{1}{3}m \cdot \sin^3 \phi \cdot R_D(q, r, 1) \end{aligned}$$

the elliptic integral of the third kind:

$$\begin{aligned} \Pi(n; \phi|m) &= \int_0^{\sin \phi} (1-t^2)^{-1/2}(1-mt^2)^{-1/2}(1+nt^2)^{-1} dt \\ &= \sin \phi \cdot R_F(q, r, 1) - \frac{1}{3}n \cdot \sin^3 \phi \cdot R_J(q, r, 1, s). \end{aligned}$$

Also the complete elliptic integral of the first kind:

$$K(m) = \int_0^{\pi/2} (1-m \cdot \sin^2 \theta)^{-1/2} d\theta = R_F(0, 1-m, 1);$$

the complete elliptic integral of the second kind:

$$E(m) = \int_0^{\pi/2} (1-m \cdot \sin^2 \theta)^{1/2} d\theta = R_F(0, 1-m, 1) - \frac{1}{3}m \cdot R_D(0, 1-m, 1).$$

### 2.3 Bessel and Airy Functions of a Complex Argument

The routines for Bessel and Airy functions of a real argument are based on Chebyshev expansions, as described in Section 2.1. The routines for functions of a complex argument, however, use different methods. These routines relate all functions to the modified Bessel functions  $I_\nu(z)$  and  $K_\nu(z)$  computed in the right-half complex plane, including their analytic continuations.  $I_\nu$  and  $K_\nu$  are computed by different methods according to the values of  $z$  and  $\nu$ . The methods include power series, asymptotic expansions and Wronskian evaluations. The relations between functions are based on well known formulae (see Abramowitz and Stegun (1972)).

## 3 Recommendations on Choice and Use of Available Routines

**Note:** refer to the Users' Note for your implementation to check that a routine is available.

### 3.1 Elliptic Integrals

**IMPORTANT ADVICE:** users who encounter elliptic integrals in the course of their work are strongly recommended to look at transforming their analysis directly to one of the Carlson forms, rather than to the traditional canonical Legendre forms. In general, the extra symmetry of the Carlson forms is likely to simplify the analysis, and these symmetric forms are much more stable to calculate.

The routine S21BAF for  $R_C$  is largely included as an auxiliary to the other routines for elliptic integrals. This integral essentially calculates elementary functions, e.g.,

$$\begin{aligned} \ln x &= (x-1) \cdot R_C\left(\left(\frac{1+x}{2}\right)^2, x\right), \quad x > 0; \\ \arcsin x &= x \cdot R_C(1-x^2, 1), \quad |x| \leq 1; \\ \operatorname{arcsinh} x &= x \cdot R_C(1+x^2, 1), \quad \text{etc.} \end{aligned}$$

In general this method of calculating these elementary functions is not recommended as there are usually much more efficient specific routines available in the Library. However, S21BAF may be used, for example, to compute  $\ln x/(x-1)$  when  $x$  is close to 1, without the loss of significant figures that occurs when  $\ln x$  and  $x-1$  are computed separately.

### 3.2 Bessel and Airy Functions

For computing the Bessel functions  $J_\nu(x)$ ,  $Y_\nu(x)$ ,  $I_\nu(x)$  and  $K_\nu(x)$  where  $x$  is real and  $\nu = 0$  or  $1$ , special routines are provided, which are much faster than the more general routines that allow a complex argument and arbitrary real  $\nu \geq 0$ . Similarly, special routines are provided for computing the Airy functions and their derivatives  $Ai(x)$ ,  $Bi(x)$ ,  $Ai'(x)$ ,  $Bi'(x)$  for a real argument which are much faster than the routines for complex arguments.

## 4 Index

Airy function, $Ai$ , real argument .....	S17AGF
Airy function, $Ai'$ , real argument .....	S17AJF
Airy function, $Ai$ or $Ai'$ , complex argument, optionally scaled .....	S17DGF
Airy function, $Bi$ , real argument .....	S17AHF
Airy function, $Bi'$ , real argument .....	S17AKF
Airy function, $Bi$ or $Bi'$ , complex argument, optionally scaled .....	S17DHF
Arccos, inverse circular cosine .....	S09ABF
Arccosh, inverse hyperbolic cosine .....	S11ACF
Arcsin, inverse circular sine .....	S09AAF
Arcsinh, inverse hyperbolic sine .....	S11ABF
Arctanh, inverse hyperbolic tangent .....	S11AAF
Bessel function, $J_0$ , real argument .....	S17AEF
Bessel function, $J_1$ , real argument .....	S17AFF
Bessel function, $J_\nu$ , complex argument, optionally scaled .....	S17DEF
Bessel function, $Y_0$ , real argument .....	S17ACF
Bessel function, $Y_1$ , real argument .....	S17ADF
Bessel function, $Y_\nu$ , complex argument, optionally scaled .....	S17DCF
Complement of the Cumulative Normal distribution .....	S15ACF
Complement of the Error function, real argument .....	S15ADF
Complement of the Error function, scaled, complex argument .....	S15DDF
Cosine, hyperbolic .....	S10ACF
Cosine Integral .....	S13ACF
Cumulative Normal distribution function .....	S15ABF
Dawson's Integral .....	S15AFF
Digamma function, scaled .....	S14ADF
Elliptic functions, Jacobian, sn, cn, dn .....	S21CAF
Elliptic integral, symmetrised, degenerate of 1st kind, $R_C$ .....	S21BAF
Elliptic integral, symmetrised, of 1st kind, $R_F$ .....	S21BBF
Elliptic integral, symmetrised, of 2nd kind, $R_D$ .....	S21BCF
Elliptic integral, symmetrised, of 3rd kind, $R_J$ .....	S21BDF
Elliptic integral, general, of 2nd kind, $F(z, k', a, b)$ .....	S21DAF
Erf, real argument .....	S15AEF
Erfc, real argument .....	S15ADF
Erfc, scaled, complex argument .....	S15DDF
Error function, real argument .....	S15AEF
Exponential, complex .....	S01EAF
Exponential Integral .....	S13AAF
Fresnel Integral, $C$ .....	S20ADF
Fresnel Integral, $S$ .....	S20ACF
Gamma function .....	S14AAF
Gamma function, incomplete .....	S14BAF
Generalized Factorial function .....	S14AAF
Hankel function $H_\nu^{(1)}$ or $H_\nu^{(2)}$ , complex argument, optionally scaled .....	S17DLF
Incomplete Gamma function .....	S14BAF
Jacobian elliptic functions, sn, cn, dn, real argument .....	S21CAF
Jacobian elliptic functions sn, cn, dn, complex argument .....	S21CBF
Jacobian theta functions $\theta_k(x, q)$ , real argument .....	S21CCF
Kelvin function, $bei x$ .....	S19ABF

Kelvin function, ber $x$ .....	S19AAF
Kelvin function, kei $x$ .....	S19ADF
Kelvin function, ker $x$ .....	S19ACF
Legendre functions of 1st kind $P_n^m(x)$ , $\overline{P}_n^m(x)$ .....	S22AAF
Logarithm of Gamma function .....	S14ABF
Logarithm of $1 + x$ .....	S01BAF
Modified Bessel function, $I_0$ , real argument .....	S18AEF
Modified Bessel function, $I_1$ , real argument .....	S18AFF
Modified Bessel function, $I_\nu$ , complex argument, optionally scaled .....	S18DEF
Modified Bessel function, $K_0$ , real argument .....	S18ACF
Modified Bessel function, $K_1$ , real argument .....	S18ADF
Modified Bessel function, $K_\nu$ , complex argument, optionally scaled .....	S18DCF
Polygamma function $\psi^{(n)}(x)$ , real $x$ .....	S14AEF
Polygamma function $\psi^{(n)}(z)$ , complex $z$ .....	S14AFF
Psi function .....	S14ACF
Psi function and derivatives, scaled .....	S14ADF
Scaled modified Bessel function, $e^{- x }I_0(x)$ , real argument .....	S18CEF
Scaled modified Bessel function, $e^{- x }I_1(x)$ , real argument .....	S18CFF
Scaled modified Bessel function, $e^x K_0(x)$ , real argument .....	S18CCF
Scaled modified Bessel function, $e^x K_1(x)$ , real argument .....	S18CDF
Sine, hyperbolic .....	S10ABF
Sine integral .....	S13ADF
Tangent, circular .....	S07AAF
Tangent, hyperbolic .....	S10AAF
Trigamma function, scaled .....	S14ADF
Zeros of Bessel functions $J_\alpha(x)$ , $J'_\alpha(x)$ , $Y_\alpha(x)$ , $Y'_\alpha(x)$ .....	S17ALF

## 5 Routines Withdrawn or Scheduled for Withdrawal

None.

## 6 References

- Abramowitz M and Stegun I A (1972) *Handbook of Mathematical Functions* (3rd Edition) Dover Publications
- Carlson B C (1965) On computing elliptic integrals and functions *J. Math. Phys.* **44** 36–51
- Carlson B C (1977a) Elliptic integrals of the first kind *SIAM J. Math. Anal.* **8** 231–242
- Carlson B C (1977b) *Special Functions of Applied Mathematics* Academic Press
- Clenshaw C W (1962) *Mathematical tables Chebyshev-series for Mathematical Functions* HMSO
- Fox L and Parker I B (1968) *Chebyshev Polynomials in Numerical Analysis* Oxford University Press
- Schonfelder J L (1976) The production of special function routines for a multi-machine library *Softw. Pract. Exper.* **6** (1)
-