Department of Mathematics and Statistics

Preprint MPCS-2020-03

14 October 2020

Computation of the Complex Error Function using Modified Trapezoidal Rules

by

Mohammad Al Azah and Simon N Chandler-

for x > 0,

(6)
$$\operatorname{erfcz}(x) = \frac{hze^{-z^2}}{k_{\text{e}=1}} \frac{x^{\frac{1}{2}}}{z^2 + k^2h^2} + \frac{2H(-h-x)}{1 - e^{2z-h}} + E(h):$$

Here the $\,$ rst term is the trapezoidal rule approximation to (

\is very accurate, provided for given z and N the optimal step size h is selected. It is not easy, however, to determine this optimal a priori."

Our recommendations address this issue, detailing which approximation formula and what step size should be used for eath and z.

The bounds we obtain in carrying out ii) prove that the absolute error in our approximation fow(z) tends to zero exponentially with uniformly in the complex plane. This is a substantial improvement on the existing bound (7) which blows up when z = z and does not capture the additional truncation errors due to replacing in nite by nite sums in the approximations (6) and (9).

Concretely, our proposed approximation $\mathbf{x}(\mathbf{z})$, for $\mathbf{z} = \mathbf{x} + i\mathbf{y}$, with \mathbf{x} ; $\mathbf{y} = 0$, is

where is de ned by (8),N 2 $\mathbb{N}_0 := \mathbb{N}$ [f Og,

(11)
$$h := {q - (N + 1);}$$

W

Indeed, we have previously used, in the restricted casæ)a \neq g \in 4, an approximation resembling v_N^{MM} (z) when approximating Fresnel integrals [6], proving results in the spirit of Theorem 1.1.

Let us summarise the rest of the paper. In the largest/ve derive the above formulae and error bounds. Ix3 we review the existing, alternative approximate methods for computing era/cándw(z) for complex, for none of which has an error bound been proved, similar to Theorem 1.1. Ix4 we carry out numerical experiments that con rm the accuracy $\text{cuf}_N(z)$, showing that its absolute error is $2 \cdot 10^{15}$ throughout the complex plane wNth= 11, and that the same bound holds for the relative error in the upper half-plane. We also show that our new approximation is competitive in accuracy and computing times with the methods that we suckey in speci cally those of [24, 27, 26, 2].

We note that this paper is based, in signi cant part, on Chapter 3 of the rst author's thesis [4].

- 2. The proposed approximation and its error bounds. In this section we derive the approximation given by (10) based on modi ed trapezoidal rules. We also derive the error bounds of Theorem 1.1 that demonstrate that the absolute and relative errors in w_N (z) both decrease exponentially bisincreases.
- 2.1. The contour integral argument and its history. Given any f 2 $C(\mathbb{R})$ that decays su ciently rapidly at in nity, let

$$Z_{1}$$

$$I[f] := f(t) dt;$$

and, for h > 0 and 2 [Q 1), de ne the generalised trapezoidal rule approximation to I [f] by

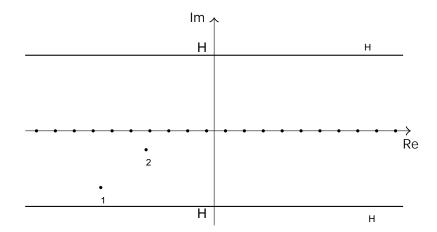
(17)
$$I_{h;} [f] := h X f ((k)h):$$

We note that h_i , $[f] = I_{h;0}[f]$, where f(t) := f(t + h) fort f(t) = f(t) and that f(t) = f(t) its composite midpoint rule approximation.

The approximation (17) fol [f] converges exponentially when the integrand is analytic in a strip surrounding the real axis and has su cient decay at The derivation of this result, using contour integration and Cauchy's residue theorem, dates back, for a particular case, at least to Turing [23], and has been analysed in more general cases by Goodwin [12], McNamee [16], Schwartz [20] and Stenger [21]. For a detailed history and discussion see Trefethen and Weideman [22].

The rate of exponential convergence depends on the width of the strip of analyticity around the real axis, and the accuracy of [f] deteriorates when has singularities close to the real line. But, in the case when these singularities are poles, the contour integral method for establishing the exponential convergence of that we will recall in Proposition 2.1 below, leads naturally to corrections for modifying the trapezoidal rule and recovering rapid convergence, these corrections expressed in terms of residues that these poles. This appears to have been observed explicitly rst by Chiarella and Reichel [8], in the context of evaluating (5) (and see Matta and Reichel [15], Hunter and Regan [13], and Mori [17]), and has been developed into a general theory by Bialecki [7] (and see La Porte [14]).

It is convenient to recall in a proposition the standard arguments ([8, 15, 13] and cf. [22, pp. 402(403]) that are made to prove exponential convergence, since we will



The representations (6) and (9

bounds hold also for = 0 and y = =h since the left hand sides of the bounds depend continuously on on [Q = h] (recall that b_{i} ; $[f_z]$ is an entire function of a and that b_{i} ; $[f_z]$ is bounded below on b0 by (32)).

Now suppose that y = h and takeH = y = h for some 2 (Q, y = h). Then y = h [f z] = y = h; [f z], and since $y = (y^2 + H^2)$ and y = h and y = h are decreasing as functions of on (H; 1], it follows from (38) and (39) whith y = h that

(42)
$$jw(z) I_{h;} [f_z]j = \frac{2^{0} \overline{2} e^{-2} = h^{2} + r^{2}}{1 - e^{-2} = 2^{2} = h^{2} + 2^{2}}$$

and

$$(43) \qquad \qquad \frac{jw(z) \quad I_{h;} \quad [f_z]j}{jw(z)j} \qquad \frac{2^p \, \overline{2} \, (h + ^p \overline{2} \,) \, e^{-^2 = h^2 + ^{"}^2}}{h \, 1 \, e^{-^2 \, 2^2 = h^2 + ^2 \, "=h} \quad "(2 \quad "h)};$$

If $=h > 1 = \frac{p}{2}$ we can again choose $1 = \frac{p}{2}$, obtaining the bounds (35) and (36) for y > =h; these bounds hold also for = =h since the left hand sides of the bounds depend continuously on on = =h; 1).

It follows immediately from the de nition (29) that xf 2r, y > 0,

(44)
$$jC_{h;} [f_z]j \frac{2e^{2y=h}}{1 e^{2y=h}} e^{y^2 x^2}.$$

Since $J_{h;}$ [f_z]j j $J_{h;}$ [f_z]j + jC_{h;} [f_z]j, the following corollary follows from the above proposition, (44), and (32).

Corollary 2.3. If z = x + iy with x = y 0 and $h < \frac{p}{2}$, then

(45)
$$w(z) I_{h;} [f_z] c_a \frac{e^{-2h^2}}{1 - e^{-2^2 - h^2 + \frac{p}{2} - h}}$$

and

(46)
$$\frac{w(z) - I_{h;} [f_z]}{jw(z)j} - \frac{c_r}{h} \frac{e^{-2} = h^2}{1 - e^{-2} = h^2 + \frac{p}{2} = h};$$

where

$$(47) \quad c_a := \frac{2(2e + {}^p \overline{}_{\underline{-}})}{p \overline{}_{\underline{-}}} \quad 4.934 \ \ \text{and} \quad c_r := \frac{2^p \overline{2} (1 + {}^p \overline{}_{\underline{-}})(2e + {}^p \overline{}_{\underline{-}})}{p \overline{}_{\underline{-}}} \quad 6077.$$

Proof. For 0 x = y =h these bounds follow immediately from the sharper bounds (33) and (34). Suppose now that 49 5.358 Td [(;)]TJ/F43 9.9626 Tf -276.101 -23.887 Td [(wher)51(e)]TJ/F8 9.9

Proposition 2.2 tells us that

for x 0, where

$$G_x(t) := e^{-t^2} j F_x(t + i = h) j =$$

But, if $h < \frac{p}{2}$, it follows from (31) and (38) applied whith 3 that, for some constant C > 0 independent of, $j E_h(z) j$ Cjzj if z = 2 with y = 2. Similarly, since $I_{h; ; H}$ $[f_z] = I_{h; } [f_z]$ if y > H and $I_{h; } [f_z] = I_{h; } [f_z] + C_{h; } [f_z]$, it follows from (38) applied with H = 1 and (44) that, for some constant C > 0 independent of C > 0 ind

The following corollary summarises and simplies, at the cost of a little sharpness, the results of Propositions 2.2 and 2.6 and of this subsection.

Corollary 2.7. Suppose that z=x+iy with x=Q,y=Q, and $h<\frac{p-1}{2}$. Then the bounds (45) and (46) hold with c_a and c_r given by (47) if $y=\max(x;=h)$. The same bounds hold as bounds of $w(z)=I_h$; $[f_z]j=jw(z)j$, respectively, with the same values of and c_r , if $y=\max(x;=h)$.

Proof. The rst claim of the corollary follows from Proposition 2.6 and (33) and (34), and the second follows from (35) and (36).

2.3. Truncating the in nite series. Propositions 2.2 and 2.6 together provide (ac)curate trapezoidal-rule-baeep38

on (0; 1) and noting (53), that

(57)
$$2h \sum_{k=M}^{1} e^{s_k^2} 2$$

Proof. From (25) and (55), for Oy x,

$$(63) \qquad jT_{h;}^{N} \ [f_{z}]j \qquad \frac{2hjzj}{k_{E-N+1}} \, \frac{x!}{jz^{2} - s_{k}^{2}j} \qquad \frac{2^{p}}{(x+s_{N+1})} \, \frac{x!}{k_{E-N+1}} \, \frac{e^{-s_{k}^{2}}}{jz - s_{k}j} :$$

Thus, and noting (32), the bounds (61) and (62) hold if O.

Choose with 0< < 1. Given x>0 letM be the smallest integer N + 1 such that $s_M>x$, so that, ifM > N + 1, $s_k=x$ and $jz=s_kj=(1-)x$ for k< M. If M > N + 1 it follows, using the bound (57), that

$$2hx \int_{k=\,N\,+1}^{N\!\!X} \frac{1}{jz} \frac{e^{-s_k^2}}{s_k j} - \frac{2h}{1} \int_{k=\,N\,+1}^{X\!\!1} e^{-s_k^2} - \frac{2h s_{N\,+1}\,+\,1}{(1}$$

Further, if x = 0 and

the nth convergent of the beautiful Laplace continued fraction representation forw(z) (speci cally suggesting n = 9). Gautschi notes that: i) by construction the nth convergent is asymptotically accurate, with error $O(jzj^{2n-1})$ as jzj!1, uniformly in the rst quadrant; ii) the nth convergent converges tow(z) as n!1 if and only if Im(z) > 0; iii) remarkably, for Im(z) > 0, the nth convergent coincides with the approximation obtained by approximating (3) by an n-point Gauss-Hermite rule. For smaller z Gautschi [11] proposed (rational) approximations that are truncated Taylor expansions with the coe cients approximated by convergents of continued fractions.

This methodology, carefully tuned, is the basis of TOMS Algorithm 680 (Poppe and Wijers [18]) which achieves a relative error of 10¹⁴ over nearly all the complex plane using, in the rst quadrant: i) Maclaurin polynomials of degree 55 for the odd function erf(iz) (substituted into (2)) in an ellipse around the origin; ii) the convergents (78) with n 18 outside a larger ellipse; iii) the more expensive mix of Taylor expansion and continued fraction calculation proposed by Gautschi [11] in between. This algorithm has been used as a benchmark by several later authors.

Weideman [24, 25] proposed (the derivation starts from (3)) the single rational approximation

(79)
$$w(z) = \frac{1}{(L-iz)} + \frac{2}{(L-iz)^2} \int_{z=0}^{N} a_{n+1} dz = \frac{L+iz}{L-iz}$$
; for Im(z) 0

where the size of N controls the accuracy of the approximation, L := 2^{-14} N $^{1-2}$, and the a_n are discrete Fourier coe cients that can be precomputed by the FFT. He argues, based on operation counts, that, for intermediate values g/zj, the work required to compute w(z) to 10^{-14} relative accuracy is much smaller for (79) than for the Poppe and Wijers algorithm [18].

Zagloul and Ali proposed a method (see TOMS Algorithm 916 [27] and the renements in [26], and cf. [19] and [3, (7.1.29)]) starting from

(80)
$$\operatorname{erf}(z) = \operatorname{erf}(x) + \frac{2e^{-x^2}}{p} \int_{0}^{x^2} e^{t^2} \sin(2xt) dt + \frac{2ie^{-x^2}}{p} \int_{0}^{x} e^{t^2} \cos(2xt) dt;$$

for z = x + i y. They approximate

(81)
$$w(z) u(x;y) + i v(x;y); x;y 0;$$

where

$$\begin{split} u(x;y) &:= e^{-x^2} \text{erfcx}(y) \cos(2xy) + \ \frac{2a \sin^2(xy)}{y} \, e^{-x^2} + \ \frac{ay}{} \, (\quad 2\cos(2xy) \, S_1 + \, S_2 + \, S_3) \, ; \\ v(x;y) &:= \quad e^{-x^2} \text{erfcx}(y) \sin(2xy) + \ \frac{a \sin(2xy)}{v} \, e^{-x^2} + \ \frac{a}{} \, (2y \sin(2xy) \, S_1 - \, S_4 + \, S_5) \, ; \end{split}$$

erfcx(y) := e^{y^2} erf(y), and S_j , j = 1; ...; 5, are the following summations reminiscent of the trapezoidal rule approximations (6):

$$S_{1} := \frac{X}{a^{2}k^{2} + y^{2}} \quad e^{-(a^{2}k^{2} + x^{2})}; \quad S_{2} := \frac{X}{a^{2}k^{2} + y^{2}} \quad e^{-(ak + x)^{2}};$$

$$(82) \quad S_{3} := \frac{X}{a^{2}k^{2} + y^{2}} \quad e^{-(ak - x)^{2}}; \quad S_{4} := \frac{X}{a^{2}k^{2} + y^{2}} \quad e^{-(ak + x)^{2}};$$

$$S_{5} := \frac{X}{a^{2}k^{2} + y^{2}} \quad e^{-(ak - x)^{2}};$$

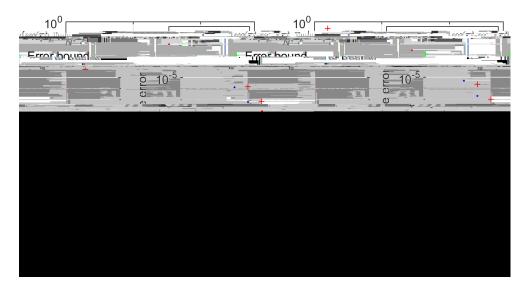


Fig. 2 . Maximum absolute and relative errors in the approximation (10) and the error bounds of Theorem 1.1, plotted against N .

The authors have supplied us with their Matlab implemental fixed beyon v2(z,M) [26], where the parame the rist the number of accurate significant gures required, and the code enforces 4M 13. In this code the choice 1=2 is made and the sums in (82) are truncated, the number of terms retained depending the target and Ali [27] (and see [26]) present numerical evidence that the approximation (81), with a = 1=2 and appropriate truncation of the in nite sums (82), is more accurate and faster than TOMS Algorithm 680 [18].

Abrarov et al. [2] (and see [1]) proposed recently another method for computing w(z) using modi ed rational approximations, namely

where $D_1 := fz = x + iy : jzj < 8$ and y > 0.05xg, $D_2 := fz = x + iy : jzj < 8$ and y = 0.05xg, y = 0.05

$$(84) \quad _{1}(z):= \begin{array}{c} x^{1} \\ \xrightarrow{m=1} \end{array} \frac{A_{m} + B_{m}(z+i=2)}{C_{m}^{2} \quad (z+i=2)^{2}}; \quad _{2}(z):= e^{-z^{2}} + z^{1} \xrightarrow{m=1} \frac{m \quad _{m}z^{2}}{m \quad _{m}z^{2} + z^{4}};$$

the coe cients A_m , B_m , C_m , m, m, and m are specified in [2], and m (z) is approximately (78) with = 10 (see [2, Equation (9)]). Abrarov et al. [2] present numerical evidence to show that (83) achieves an accuracy m of m = 2:75 and m = 23 in (84).

4. Numerical results. In this section we show calculations that illustrate and support Theorem 1.1, and that compare the accuracy and e ciency of our approximation w_N (z) given by (10

using the Matlab codeTrap(z,N) provided in Table 1 of the supplementary material to this paper [5]. The maximum values we plot are discrete maxima taken over the 16020801 points = $10^9 e^j$, with p = 6(00006)6 and = 0(=1600)=2, a superset of the ;4001 test values in Weideman [24, 25]. To compute errors we use as the exact values w(z) the independent approximation (79) wNth= 45, implemented through a cadef(z,45) to the Matlab code in [24, Table 1]. (The results in [24, 25, Figure 8], [6, Figure 2] suggestNthat5 in (79) is ample for accuracies close to machine precision, and we obtain almost identical results in Figure 2 and Table 1 below if we use, instead (z) given by (10) as the exact value.)

We observe in Figure 2 the rate of exponential convergence predicted by Theorem 1.1. The approximation $w_N(z)$ achieves, with $w_N(z)$ and $w_N(z)$ achieves, with $w_N(z)$ and $w_N(z)$ achieves, with $w_N(z)$ achieves, $w_N(z)$ and $w_N(z)$ achieves, $w_N(z)$ achieves, $w_N(z)$ and $w_N(z)$ achieves, $w_N(z)$ and $w_N(z)$ achieves, $w_N(z)$ achieves, $w_N(z)$ achieves, $w_N(z)$ achieves, $w_N(z)$ and $w_N(z)$ achieves, $w_N(z)$ and $w_N(z)$ achieves, $w_N(z)$ achieves, $w_N(z)$ achieves, $w_N(z)$ and $w_N(z)$ achieves, $w_N(z)$ achieves, $w_N(z)$ achieves, $w_N(z)$ achieves, $w_N(z)$ achieves, $w_N(z)$

Algorithm	Max	imum	Max	imum	Computing
Aigorttiili	abs.	error	rel.	error	time (seconds)
wTrap(z,11)	1:67	10 ¹⁵	1:89	10 ¹⁵	4:29 (0:08)
cef(z,40)	2:11	10 ¹⁵	2:15	10 ¹⁵	4:20 (0:02)
fadsamp(z)	3:86	10 ¹⁴	3:86	10 ¹⁴	5:74 (0:04)
Faddeyeva_v2(z,13)	4:07	10 ¹⁵	1:71	10 ¹³	11:00 (0:11)

Table 1

Maximum absolute and relative errors for the Matlab codes implementing the approximations (10), (79), (83), and (81). The computing times are mean and s.d. of 25 executions.

In Table 1 we compare the accuracy and e ciency of our approximation and Matlab code with Matlab implementations of the approximations (79), (81), and (83). Results are shown in Table 1 for:

- Our approximationw_N (z) withN = 11 implemented by the caWTrap(z,11) to the Matlab code provided in [5, Table 1];
- 2. Weideman's approximation (79) with= 40 (this choice of ensures high accuracble