# Fast Evaluation of Special Functions by the Modified Trapezium Rule

## MOHAMMAD AL AZAH

Department of Mathematics and Statistics University of Reading

This dissertation is submitted for the degree of Doctor of Philosophy

School of Mathematical, Physical and Computational Sciences

April 2017

I would like to dedicate this thesis to my loving parents, especially to my mother who sadly passed away on 2015 ...

## Declaration

I confirm that this is my own work and that the use of all material from other sources has been properly and fully acknowledged. Chapter 2 is based on the paper [3], joint work with Chandler-Wilde and La Porte, for which I was the principal contributor.

> MOHAMMAD AL AZAH April 2017

Acknowledgements

# Table of contents

1	Introduction 1				
	1.1	Special functions	1		
	1.2	The trapezium rule approximation	2		
	1.3	Numerical Examples	8		
		1.3.1 Example 1	8		
		1.3.2 Example 2	9		
	1.4	The contributions of this Thesis	10		
2	Fres	snel integrals 15			
	2.1		15		
	2.2	Summary of the main Results	19		
	2.3	The proposed approximation and its error bounds	20		
		2.3.1 Extensions of the error bounds	27		
	2.4	The approximations $\mathbf{G}(x)$ and $\mathbf{S}(x)$	)		
	2.5	Numerical results	32		
3	The	Faddeeva function 37			
	3.1		37		
	3.2	Summary of the main results	41		
	3.3	The proposed approximation and its error bounds	43		
	3.4	Numerical results	61		
4	The 2D impedance half-space Green's function for the Helmholtz equation 67				
	4.1		67		
	4.2	Summary of the main results	74		
	4.3	The proposed approximation and its error bounds	77		
		4.3.1 Bounding the discretisation error	80		
		4.3.2 Bounding the truncation error	83		

	4.4	4.3.3Choices of the step-size h	85 88 90				
5	Con	cluding remarks and further work	97				
	5.1	Concluding Remarks	97				
	5.2	Further work	98				
References							
Appendix A Matlab codes 105							
	A.1	Matlab codes to compute Fresnel integrals	105				
	A.2	Matlab code to compute Faddeeva function	107				
	A.3	Matlab code to compute $P_b$	109				

## Chapter 1

# Introduction

## 1.1 Special functions

Special functions arise in the mathematical sciences as non-elementary solutions of differential equations, and these solutions can be represented in different ways. Computing these special functions efficiently is of major interest for scientific applications and we can find formulas for approximating many of them in Abramowitz and Stegun [2] and Luke [42]. (1.1) and evaluated effectively using the trapezium rule (1.6): this method of approximation has been proposed for the incomplete gamma function in [4]; for Bessel functions in [20, 29, 56], for the Airy function in [23], for the gamma function in [53]; and for the error function in [15, 43, 31, 45].

It is well-known [18] that integrals of the form (1.1) with f is given by (1.2) can be approximated by the Hermite-Gaussian quadrature rule, denoted by  $J_N$ , which is given by

$$J_{N} := \frac{1}{p_{\overline{r}}} \mathop{a}\limits^{N}_{i=1} w_{i} F(x_{i} = p_{\overline{r}}); \qquad (1.3)$$

where  $w_1$ ;:::; $w_N$  and  $x_1$ ;:::; $x_N$  are the weights and abscissae, respectively. The Hermite-Gaussian quadrature rule is very accurate, and sometimes outperforms the trapezium rule, when the function F is smooth; but the accuracy deteriorate when F is meromorphic with simple poles near the real axis. For example, approximating the integral

$$Z_{\text{¥}} e^{t^2} \cos(t) dt \qquad (1.4)$$

using  $J_N$  with N = 12 (see http://www.chebfun.org/examples/quad/HermiteQuad)htmls

before Propositions 1.2.3 and 1.2.4. We assume in the following results that the function F in (1.2) satisfies the following assumption.

Assumption 1.2.1. For H > 0 and  $S_H = f z 2 C : jIm(z)j < Hg$ , we have that

- (i) F is meromorphic with simple poles at  $2 S_H$ , Im( $z_j$ ) 6 0 and j = 1; ...; m;
- (ii) F is continuous or  $\overline{S}_H$  nf  $z_1$ ;  $z_2$ ;  $z_3$ ; ...;  $z_m g$ ;
- (iii)  $F(z) = O(1) \operatorname{asjRe}(z)j! + \operatorname{uniformly forjIm}(z)j + H.$

Given h > 0 and a 2 [0; 1), define the function g(z) by

$$g(z) := i \cot p \frac{z}{h} + a ;$$
 (1.7)

which is a meromorphic function with simple poles at z = (k a)h, k 2 Z, which has the properties that, for z = x + iy with y > 0,

j1 g(z)j 
$$\frac{2e^{-2py=h}}{1 - e^{-2py=h}};$$
 (1.8)

and for z = x + iy with y < 0,

j1+g(z)j 
$$\frac{2e^{2py=h}}{1 e^{2py=h}}$$
: (1.9)

We will make use in the following results of the signum function, sign(t), which is defined by sign(t) = 1 for t > 0, sign(0) = 0 and sign(t) = 1 for t < 0. We will make use also of the paths  $G_H$  and  $G_H^0$  in the complex plane which are defined as the lines Im(z) = H and Im(z) = H, respectively, traversed in the direction of increasing Re(z).

**Proposition 1.2.1.** If Assumption 1.2.1 holds, the (h; a) as de ned in(1.6) exists as the limit

$$\lim_{n;j!} h \overset{n}{\overset{a}{a}} f((k a)h);$$

and has the value

$$I(h;a) = \frac{1}{2} \int_{G_{H}}^{Z} f(z)g(z) dz = \int_{G_{H}^{0}}^{Z} f(z)g(z) dz + pi \overset{m}{a}_{k=1}^{m} g(z_{k}) R_{k}: \quad (1.10)$$

where  $\mathbf{R} = \operatorname{Res}(f; \mathbf{z}_k)$ .

**Proof.** Let  $A_k = (k + \frac{1}{2})h$  for  $k \ge N$  and define  $C_H$  as the positively oriented rectangular contour with vertices at  $A_j$  iH and  $A_n$  iH. Using Cauchy's residue theorem for  $C_H$  (which encloses j + n + 1 simple poles of the integrand) we find that

$$Z = f(z)g(z) dz = 2pi \qquad \overset{n+j+1}{\overset{a}{a}} \operatorname{Res}(fg;(k a)h) + \overset{m}{\overset{a}{a}} \operatorname{Res}(fg;z_k)$$

The following proposition is well-known from many papers. It is in Goodwin [24] for the case when a = 0 and the integrand is analytic in  $S_H$ , in Chiarella and Reichel [15,1.956(Reichef)1(1.02)]

integrand. For example, in Hunter [29, 30] we find this result for the case where the integrand is even and analytic in  $S_H$  and a = 0; in Hunter and Regan [31] for a = 0 and a = 1=2 with  $F(t) = 1=(t^2 + a^2)$ , for some a 2 C; in Theorem 2.2 of Bialecki [5] for a = 0 when the integrand is meromorphic with poles of arbitrary order, in Theorem 2.3.2 of La Porte [38] for a = 0, and recently in Theorem 5.1 of [60] for the case where a = 0 and the integrand is analytic in  $S_H$ .

Proposition 1.2.4. For h > 0 and a 2 [0;1) let E (h;a) := I I (h;a). If Assumption (1.2.1) holds, then \_pan [

**Definition 1.2.2.** For h > 0 and a = 0 or 1=2, we denote by  $y_N(h; a)$  the truncated trapezium rule de ned by

$$I_N(h;0) := h f(0) + 2h a_{k=1}^N f(kh)$$
 and  $I_N(h;1=2) := 2h a_{k=0}^N f((k+1=2)h)$ : (1.22)

We denote also by, (h; a) the truncated modified trapezium rule de ned by

$$I_N(h;a) := I_N(h;a) + C(h;a):$$
 (1.23)

Note that the truncation of I(h; a) induces the additional error

$$T_{N}(h;a) := 2h \mathop{a}\limits_{k=N+1}^{4} f((k+a)h); \qquad (1.24)$$

which will be considered in the coming chapters. The total error in approximating the integral I (1.1) by  $I_N(h;a)$  will be denoted by  $E_N(h;a)$  where

$$E_N(h;a) = E(h;a) + T_N(h;a)$$
: (1.25)

#### 1.3 Numerical Examples

To give a flavour and preview of the extraordinary efficiency of the modified trapezium rule we present here two examples that demonstrate the convergence rate of the rule (1.18). In the first example the integrand is an entire function; and in the second example the integrand is a meromorphic function. In both examples, we approximate the integral by  $I_N(h;a)$  with a = 0.

#### 1.3.1 Example 1

The following integral is a famous example (see Goodwin [24]):

$$I = \sum_{i=1}^{Z} e^{t^{2}} dt = p^{2} \overline{p} = 1:7724538509055160273:...$$
(1.26)

The integrand here is an entire function and hence we have that C(h; 0) = 0 so that

$$I(h;0) = I(h;0) = h \overset{i}{\overset{k}{a}}_{k=} e^{k^2h^2}$$
 and  $I_N(h;0) := I_N(h;0) = 1 + 2h \overset{N}{\overset{k}{a}}_{k=1} e^{k^2h^2}$ :

Table 1.1 shows the computed values of  ${\sf I}_{\sf N}(h;0)$ 

N	$h = 0.7(N + 1)^{2=3}$	I <sub>N</sub> (h;0)
10	0:142	0:910749
20	0:092	<b>0:88</b> 9598
40	0:059	<b>0.8875</b> 7706
80	0:037	0.88753706862

 To derive completely rigorous and explicit bounds on both the absolute and relative errors when approximating particular special functions by the truncated modified trapezium rule. The bounds we obtain justify theoretically the choices that we recommend for the parameters a, H, h and N, and prove exponential (or near exponential) convergence as N ! ¥. These theoretical predictions are supported by systematic and comprehensive numerical experiments.

The largest part of this thesis is concerned with the application of the truncated modified trapezium rule (1.23) (with a = 0 or a = 1=2) to the computation of the complex error function  $w(z) = e^{-z^2} \operatorname{erfc}(-iz)$  (Chapter 3), and with the related problem of computing Fresnel integrals (Chapter 2). The application of the modified trapezium rule (1.18) with a = 0 to compute the complementary error function, denoted by  $\operatorname{erfc}(z)$  with z = x + iy, starting from the integral representation

erfc(z) = 
$$\frac{z e^{z^2}}{p} \sum_{i=1}^{z^2} \frac{e^{t^2}}{z^2 + t^2} dt; x > 0;$$
 (1.28)

was proposed by Chiarella and Reichel [15] and Matta and Reichel [43] who proposed to use I (h; 0) given by (1.18) with H = p=h, i.e.

erfc(z) 
$$\frac{he^{z^2}}{pz} + \frac{2hze^{z^2}}{p} \overset{k}{\overset{a}{a}}_{k=1} \frac{e^{k^2h^2}}{z^2 + k^2h^2} + \frac{2H(H-x)}{1 e^{2pz=h}};$$
 (1.29)

where **H** is the Heaviside step function. This proposal was refined later by Hunter and Regan [31]. In particular, Hunter and Regan [31] noted that (1.29) blows up if the simple poles of the integrand at t = iz coincide with any quadrature point at kh. They proposed to use the approximation I (h; 1=2) with H = p=h, i.e.

erfc(z) 
$$\frac{2hz e^{z^2}}{p} \stackrel{\text{¥}}{\overset{\text{a}}{a}}_{k=1} \frac{e^{(k-1=2)^2 h^2}}{z^2 + (k-1=2)^2 h^2} + \frac{2H(H-x)}{1 + e^{2pz=h}};$$
 (1.30)

when (1.29) fails or suffers from numerical instability. They proposed precisely the approximation 8

erfc(z) 
$$(h; 0);$$
 if 1=4 f (y=h) 3=4 (1.31)  
: I (h; 1=2); otherwise;

where f (t) denotes the fractional part of t, i.e. f (t) = t [t]. They also proved, essentially applying Proposition 1.2.4 with H = p=h, and noting for

$$F(t) = \frac{ze^{z^2}}{p(z^2 + t^2)}$$

it holds that

$$M_H(F) = \frac{jze^{z^2}j}{pjx^2 p^2 = h^2j};$$

that the error in this approximation is

$$p \frac{jz e^{z^{2}} j e^{p^{2} = h^{2}}}{\overline{p} j x^{2} p^{2} = h^{2} j (1 e^{2p^{2} = h^{2}})};$$
(1.32)

Clearly this error bound blows up when x = p = h, and so is inadequate as a bound for x = p = h. This can be fixed by finding an improved version for jx p = hj e, for some e > 0, by taking H = p = h e in Proposition 1.2.4, but the bounds obtained with this modification are still unsatisfactory as they don't imply small absolute and relative errors as h ! = 0 uniformly in z = x + iy.

Mori [45] studied the approximation I (h;0) in (1.29) specifically for z = x > 0. He bounded the error in this approximation by (1.32) and by another bound obtained from Proposition 1.2.4 with  $H = p = h + 1 = \frac{p}{2}$ , namely that the error is

$$\frac{x e^{x^2} e^{1=2} e^{p^2=h^2}}{p \overline{p} j x^2 (p=h+1=\overline{2})^2 j (1 e^{2p=h(p=h+1=\overline{2})})}$$
: (1.33)

Mori [45] used the minimum of the bounds (1.32) and (1.33), i.e. he used (1.32) for x > b, (1.33) for 0 < x = b, where b is the value (given by (2.8) and (2.9) in [45], but here we correct a calculation error in [45])

$$b := \frac{1}{1+1} + \frac{p}{h} + \frac{p^{1}}{2} + 1 + \frac{p}{h} + \frac{p^{2}}{2}; \qquad (1.34)$$

with

$$I := \frac{1 e^{2p^2 = h^2} e^{1=2}}{1 e^{2p = h(p=h+1=\frac{p}{2})}};$$
 (1.35)

for this value of b the two bounds (1.32) and (1.33) coincide. Mori [45] also bounded the relative error, using that

erfc(x) 
$$\frac{2e^{x^2}}{\overline{p}(x+\overline{x^2+2})}; x = 0:$$
 (1.36)

Mori [45] showed further that the relative error in (1.29) is

$$\frac{b(b + p^{2} + 2)}{(b^{2} - p^{2} + h^{2})(1 - e^{-2p^{2} + h^{2}})} e^{-p^{2} + h^{2}}; \qquad (1.37)$$

for all z = x = 0.

The work in this thesis extends and improves significantly, by more sophisticated and delicate analysis, the previous works. In Chapter 2 we propose methods for computing Fresnel integrals based on the truncated modified trapezium rule in (1.23) where a = 1=2. We construct approximations in Sections §2.3 and §2.4 which we prove are exponentially convergent as a function of N, the number of quadrature points, obtaining completely explicit error bounds in Theorems 2.3.3 and 2.3.5 which show that accuracies of 10<sup>15</sup> uniformly on the real line are achieved with N = 12, this confirmed by computations in Section §2.5. The approximations we obtain are attractive in that they maintain small relative errors for small and large argument, are analytic on the real axis (echoing the analyticity of the Fresnel integrals), and are straightforward to implement.

In Chapter 3 we propose a method for computing the complex error function w z)

Matlab codes are provided (see Listings A.1, A.2, A.3 and A.4) for computing all these functions, and these codes are easily adaptable to other programming languages.

# Chapter 2

# **Fresnel integrals**

## 2.1 Introduction

Let C(x),

It also depends on the integral representation [2, (7.1.4)] that

$$w(z) = \frac{i}{p} \sum_{i=1}^{Z} \frac{e^{-t^{2}}}{z - t} dt = \frac{iz}{p} \sum_{i=1}^{Z} \frac{e^{-t^{2}}}{z^{2} - t^{2}} dt; \quad Im(z) > 0: \emptyset$$
()

#### 2.1 Introduction

where the size of N controls the accuracy of the approximation  $\neq 2^{1=4}N^{1=2}$  and the coef cients are computed as

1 a<sub>n</sub> :=

### 2.2 Summary of the main Results

Based on the truncated modified trapezium rule (1.23) with a = 1=2 and  $H = A_N$  (given by (2.13)), the approximation to F(x) we propose is

$$F_{N}(x) := \frac{1}{2} + \frac{i}{2} \tan A_{N} x e^{ip=4} + \frac{x}{A_{N}} e^{i(x^{2}+p=4)} \mathop{a}\limits_{k=1}^{N} \frac{e^{-t_{k}^{2}}}{x^{2}+it_{k}^{2}}$$
(2.11)  
$$= \frac{1}{\exp 2A_{N} x e^{-ip=4}}$$

in modified form into this strip. This implies exponentially convergent error estimates, presented in §2.3.1 and §2.4, for the difference between the coefficients in the Maclaurin series of F, C, and S and those in the corresponding series for  $F_N$ ,  $C_N$  and  $S_N$ . In turn (see §2.4), this implies that the approximations all retain small relative error for jxj small, and the computations in §2.5 demonstrate this.

These approximations inherit symmetries of the Fresnel integrals. In particular, our

< 10 <sup>15</sup>. From (2.6) we have that, for > 0,

$$F(x) := \sum_{i=1}^{Z} f(t) dt; \text{ where } f(t) := e^{i(x^2 + p = 4)} x$$

Here

$$d_{1}(x) := p \frac{x e^{p^{2} = h^{2}}}{\overline{p} j p^{2} = h^{2} x^{2} = 2j \ 1 \ e^{-2p^{2} = h^{2}}};$$
(2.29)

$$d_{2}(x) := p \frac{4hx e^{p^{2} = h^{2}}}{p \overline{p} p j p = h + x} \frac{2j}{2j} \frac{1}{1} e^{-2p^{2} = h^{2}}}{1 + 2^{p} \overline{p} e^{-bp^{2} = h^{2}}}; \qquad (2.30)$$

with b = 
$$\frac{15 \ 10^{p} \overline{2}}{16}$$
 0:0536, and  
d<sub>3</sub>(x) := d<sub>1</sub>(x) +  $\frac{e^{p \overline{2}px=h}}{1 \ e^{p \overline{2}px=h}}$ : (2.31)

**Proof.** Applying Proposition 1.2.4, for  $0 < x < p \overline{2}p = h$ , with H = p = h, and noting for

$$F(t) = \frac{x e^{j(x^2 + p = 4)}}{2p(x^2 + ith_{\frac{1}{2}} + ith_$$

Thus, and applying (1.8), similarly to (2.29) we deduce that

$$Z_{g^{0}}f(z)(1+g(z)) dz = \frac{x e^{p^{2}=h^{2}}}{p \overline{p} e j p = h + x = \overline{2}j} \frac{1}{1 e^{-2p^{2}=h^{2}}} (2.35)$$

To bound the integral over we note that, for  $z = X + iY = z_0 + ee^{iq} 2$  g, (2.34) is true and Y H. Further, je  $z^2 j = e^{P}$ , where

$$P = Y^2 \quad X^2 = 2xe sin(q p=4) \quad e^2 cos(2q) < 2xe + e^2 \quad 2^p \overline{2}He + (2^p \overline{2} + 1)e^2;$$

since  $x={}^{p}\overline{2}$  H < e. From these bounds an( $\mathfrak{A}$ .8), de ning a = e=H 2 (0; 1), we deduce that

$$\sum_{g}^{Z} f(z)(1+g(z)) dz = \frac{2x \exp((2^{p} \overline{2}a + (2^{p} \overline{2} + 1)a^{2} - 2)p^{2} = h^{2})}{ejp = h + x = 2^{p} \overline{2}j - 1 - e^{-2p^{2} = h^{2}}}$$
(2.36)

For  $x=p \overline{2}$  H < e we can bound using (2.33), (2.35), (2.36), and the triangle inequality, to get that

$$je_{h}j \quad d_{2}(x) := p \frac{4hxe^{p^{2}=h^{2}}}{\overline{p}pjp=h+x=\overline{2}j \quad 1 \quad e^{-2p^{2}=h^{2}}} \quad 1+2^{p}\overline{p}e^{-bp^{2}=h^{2}};$$
 (2.37)

where

$$b = 1 \quad 2^{p} \overline{2}a \quad (2^{p} \overline{2} + 1)a^{2}$$
: (2.38)

**Proposition 2.3.1.** For x > 0,

$$jT_N(h; 1=2)j = \frac{(2ht_{N+1} + 1)x}{2pt_{N+1}} e^{t_{N+1}^2}$$

Proof.

$$jT_{N}(h; 1=2)j \qquad \frac{hx}{p} \overset{¥}{\underset{m=N+1}{a}} p \frac{e^{-t\frac{x}{m}}}{x^{4} + t\frac{4}{m}} \qquad \frac{1}{2p} \frac{e^{-t\frac{x}{m}}}{x^{4} + t\frac{4}{N+1}} \qquad 2he^{-t\frac{2}{N+1}} + 2h \overset{¥}{\underset{m=N+2}{a}} e^{-t\frac{2}{m}} \frac{1}{e^{-t\frac{2}{m}}} \frac{1}{e^{-t\frac{2}{m}$$

To arrive at the last line we have used that, for x > 0,

$$2 \sum_{x}^{Z} e^{t^{2}} dt = \frac{e^{x^{2}}}{x} \sum_{x}^{Z} \frac{e^{t^{2}}}{t^{2}} dt < \frac{e^{x^{2}}}{x}:$$
 (2.40)

At this point we make a choice of h to approximately equalise  $D_h(x)$  in Theorem 2.3.1 and the bound on  $T_N(h; 1=2)$  in Proposition 2.3.1, choosing h so that  $p=h = t_{N+1} = (N + 1=2)h$ , giving that

$$h = \frac{p}{p = (N + 1 = 2)};$$
 (2.41)

in which case  $t_{N+1} = A_N = {p \over (N+1=2)p}$ , and  $t_k = t_k$ , where  $t_k$  is defined by (2.13). Making this choice of h we see that
**Theorem 2.3.2.** For  $h = {p \over p=(N+1=2)}$  so that  $H = p = h = A_N$  we have that

$$jE_N(x)j$$
  $h_N(x) := D_h(jxj) + \frac{(2p+1)jxj}{2pA_N}e^{-A_N^2}$  (2.42)

where

**Theorem 2.3.3.** For x > 0,

$$jF(x) = jE_N(x)j = h_N(x) + h_N(x) = h_N(x) + h_N(x) = \frac{e^{-pN}}{N+1=2};$$
 for x 2 R; (2.44)

where

$$c_{N} = \frac{20^{P} \overline{2}e^{-p=2}}{9p - 1 - e^{-2A_{N}^{2}}} - 1 + 2^{P} \overline{p} e^{-bA_{N}^{2}} + \frac{(2p + 1)e^{-p=2}}{2^{P} \overline{2}p^{3=2}A_{N}};$$

which decreases as N increases, with

$$c_1 = 0.825 \text{ and } \lim_{N! \neq} c_N = \frac{20^p \overline{2}e^{-p=2}}{9p} = 0.208:$$
 (2.45)

**Proof.** It is easy to see that  $D_h(x)$  is increasing on  $[0; \frac{5}{4}^p \overline{2}A_N)$  and decreasing on  $[\frac{5}{4}^p \overline{2}A_N; \mathbf{Y})$ . Further, where  $D_h(\frac{5}{4}^p \overline{2}A_N)$  denotes the limiting value of  $D_h(x)$  as  $x ! = \frac{5}{4}^p \overline{2}A_N$  from below, since  $2A_N^{-1} > e^{-A_N^2}$ ,

$$\begin{array}{rcl} D_{h} & \frac{5}{4}^{p} \, \overline{2} A_{N} & = & \frac{20^{p} \, \overline{2} \, e^{-A_{N}^{2}}}{9^{p} \, \overline{p} \, A_{N} & 1 & e^{-2A_{N}^{2}}} & 1 + 2^{p} \, \overline{p} \, e^{-b \, A_{N}^{2}} \\ & > & \frac{20^{p} \, \overline{2} \, e^{-A_{N}^{2}}}{9^{p} \, \overline{p} \, A_{N} & 1 & e^{-2A_{N}^{2}}} + \frac{e^{-5A_{N}^{2}=2}}{1 & e^{-5A_{N}^{2}=2}} = D_{h} & \frac{5}{4}^{p} \, \overline{2} \, A_{N} & \vdots \end{array}$$

Similarly,  $xD_h(x)$  is increasing on  $[0; \frac{5}{4}^p \overline{2}A_N)$  and decreasing on  $[\frac{5}{4}^p \overline{2}A_N; ¥)$ . Thus, for x > 0,

$$D_{h}(x) \quad D_{h} \quad \frac{5}{4}^{p} \overline{2} A_{N} \quad \text{and} \quad x D_{h}(x) \quad \frac{5}{4}^{p} \overline{2} A_{N} D_{h} \quad \frac{5}{4}^{p} \overline{2} A_{N} :$$
 (2.46)

Moreover,

$$q \frac{x}{\overline{x^4 + A_N^4}} = \frac{1}{2A_N} \text{ and } q \frac{x^2}{\overline{x^4 + A_N^4}} < 1; \text{ for } x > 0:$$
 (2.47)

Combining (2.42), (2.46) and (2.47) we reach the result.

**Remark 2.3.1.** We have shown the boun(2s42) and (2.44) for x > 0, but the symmetries (2.17) and (2.18) imply that  $E_N(x) = E_N(x)$ , so that (2.42) and (2.44) hold also for x < 0, and, by continuity, also for  $\neq 0$  (and in fact  $E_N(0) = h_N(0) = 0$ ).

The following result from [3, Theorem 4] will be used to bound the relative error of  $F_N(x)$ .

Lemma 2.3.4. For the Fresnel integral F(x) we have that

$$\begin{cases} 8 \\ \frac{1}{2+2^{p}} \overline{px}; & \text{for } x = 0 \\ jF(x)j \\ \vdots \\ \frac{1}{2}; & \text{for } x = 0; \end{cases}$$
(2.48)

**Theorem 2.3.5.** For the Fresnel integral  $\mathbf{F}(x)$  and its approximation  $\mathbf{F}(x)$  we have that

$$\frac{\mathbf{j}\mathbf{F}(\mathbf{x}) \quad \mathbf{F}_{N}(\mathbf{x})\mathbf{j}}{\mathbf{j}\mathbf{F}(\mathbf{x})\mathbf{j}} \quad \frac{\mathbf{h}_{N}(\mathbf{x})}{\mathbf{j}\mathbf{F}(\mathbf{x})\mathbf{j}} \quad \stackrel{\mathbf{0}}{\stackrel{\mathbf{k}}{\stackrel{\mathbf{k}}{\stackrel{\mathbf{k}}{\mathbf{j}}}} \stackrel{\mathbf{k}}{\stackrel{\mathbf{k}}{\stackrel{\mathbf{k}}{\mathbf{j}}}} \stackrel{\mathbf{k}}{\stackrel{\mathbf{k}}{\mathbf{j}}} \stackrel{\mathbf{k}}{\mathbf{j}} \stackrel{\mathbf{k}}{\mathbf{j}$$

where

$$c_{N} = \frac{10^{p} \overline{2} \ 4 + 5^{p} \overline{2p} A_{N} \ 1 + 2^{p} \overline{p} e^{bA_{N}^{2}}}{9^{p} \overline{p} e^{p=2} A_{N} \ 1 \ e^{2A_{N}^{2}}} + \frac{(2p+1)}{p e^{p=2} A_{N}} \ \frac{1}{p \overline{2} A_{N}} + \frac{p}{\overline{p}} \overline{p} :$$

which decreases as N increases, with c10:4 and  $\lim_{N!} c_N = 100e^{-p=2}=9$  2:3.

**Proof.** Combining (2.42), (2.46), (2.47) and (2.48) we see, for x > 0, that

$$\frac{h_N(x)}{jF(x)j} = 2 + \frac{5}{2}^p$$

This implies (2.49) for x > 0. The bound for x = 0 follows immediately from (2.48), (2.44) and Remark 2.3.1.

The above estimates use (2.42) and (2.43)

0 arg(z) p=2; moreover, it is clear from (2.12) that the same holds for  $F_N(z)$  and hence for  $E_N(z)$ . Thus (2.52) implies that (2.44) holds for 0 arg(z) p=2, and (2.17) and (2.18) then imply that (2.44) holds also for p arg(z) 3p=4.

It is clear from the derivations above that, if h is given by (2.41), then I (h; 1=2) also satisfies the bound (2.44), i.e.,

$$jF(z) = I (h; 1=2)j \quad c_N p = \frac{e^{-pN}}{N+1=2};$$
 (2.53)

this holding in the first instance for real z, then for imaginary z, and finally for all z in the first and third quadrants. The bound (2.44) cannot hold in the second or fourth quadrant because  $E_N(z) = F(z) - F_N(z)$  has poles there. This issue does not hold for F(z) - I - (h; 1=2), which is an entire function, but (2.53) cannot hold in the whole compleW 0 0 12Tf 7.531 0 Td [(tcomple).8n7

Thus, for z = x + iy in the second and fourth quadrants with  $A_N = (2^p \overline{2})$ ,

$$jF(z) = F_N(z)j \quad \hat{c}_N e^{-xy} \frac{e^{-pN}}{N+1=2};$$
 (2.55)

where

$$\hat{c}_{N} := c_{N} + \frac{p \bar{2}(2p+1)}{p^{3=2} \exp(p=2)} \frac{p \bar{2}(2p+1)}{N+1=2}$$
(2.56)

The sequence  $\hat{c}_N$  is decreasing with  $\hat{c}_1$  1:14 and  $\lim_{N!} \pm \hat{c}_N = \lim_{N!} \pm c_N$  0:208.

We observe above that the bou(2d44) on  $E_N(z) = F(z) - F_N(z)$  holds for all complexz in the rst and third quadrants of the complex plane, and on the boundaries of those quadrants, the real and imaginary axes, while the bou(2d55) holds in the second and fourth quadrants for  $jIm(z)j - A_N = (2^p \overline{2})$ . A signi cant implication of these bounds is that they imply that the coef cients in the Maclaurin series  $\overline{bf}_N(z)$  are close to those  $d\overline{f}(z)$ . Precisely, at least for  $jzj < A_N = \overline{2}$ ,

$$F(z) = \overset{*}{\underset{n=0}{a}} a_n z^n \text{ and } F_N(z) = \overset{*}{\underset{n=0}{a}} b_n z^n;$$

with a<sub>h</sub>

and are given explicitly in (2.14) and (2.15). We note the similarity betwee (2.14) and (2.15) and the formulae [46, (7.5.3)-(7.5.4)]

$$C(x) = \frac{1}{2} + f(x) \sin \frac{1}{2}px^2 \quad g(x) \cos \frac{1}{2}px^2 ; \qquad (2.60)$$

$$S(x) = \frac{1}{2} f(x) \cos \frac{1}{2}px^2 \quad g(x) \sin \frac{1}{2}px^2 ;$$
 (2.61)

which expres  $\mathfrak{L}(x)$  and  $\mathfrak{L}(x)$  in terms of the auxiliary functions  $\mathfrak{L}(x)$  and  $\mathfrak{L}(x)$ , for the Fresnel integrals [46, §7.2(iv)]. Indeed, it follows from [46, (7.7.10)-(7.7.11)] that, for x > 0, f(x) and  $\mathfrak{L}(x)$  have the integral representations

$$f(x) = \frac{p_{\overline{p}}x^3}{2} \int_{0}^{Z} \frac{e^{t^2}}{\frac{p_{\overline{p}}x^{2}}{2} + t^4} dt \text{ and } g(x) = p_{\overline{p}} \frac{x}{\overline{p}} \int_{0}^{Z} \frac{t^2 e^{t^2}}{\frac{p_{\overline{p}}x^{2}}{2} + t^4} dt;$$

and, recalling that  $A_N$  is linked to the quadrature step-size through 1), it is clear that, for x > 0,  $p p x_{a_N} p x^2 = A_N$  and  $p p x_{b_N} p x^2 = A_N$  can be viewed as quadrature approximations to these integrals.

The approximation (2.14) and (2.15) inherit the accuracy  $dF_N(x)$  on the real line: from (2.58) and (2.59) we see, for R, that

$$jC(x) \quad C_N(x)j \quad \stackrel{p}{\overline{2}}jE_N(\stackrel{p}{\overline{p=2}}x)j \text{ and } jS(x) \quad S_N(x)j \quad \stackrel{p}{\overline{2}}jE_N(\stackrel{p}{\overline{p=2}}x)j: \quad (2.62)$$

where  $E_N(x) = F(x)$   $F_N(x)$ . Thus the error bounds of the previous section can be applied. In particular, from(2.44) and (2.50) it follows that both  $C(x) = C_N(x)j$  and  $S(x) = S_N(x)j$  are

$$2c_{N} p \frac{e^{pN}}{2N+1}$$
; for x 2 R; (2.63)

and

$$p \overline{p} \tilde{c}_{N} j x j \frac{e^{pN}}{2N+1};$$
 for jxj  $p \overline{N+1=2}$ : (2.64)

Here  $c_N < 0.83$  and  $\tilde{c}_N < 0.18$  are the decreasing sequences of positive numbers de ned by (2.14) and (2.51), respectively.

These bounds show that (x) and  $S_N(x)$  are exponentially convergent by 4, uni-.955iv

the power series [46, §7.6(i)]

$$C(x) = \overset{\text{``}}{\underset{n=0}{\overset{\text{(}}{a}}} \frac{(-1)^{n} \frac{1}{2} p^{-2n} x^{4n+1}}{(2n)!(4n+1)}; \quad S(x) = \overset{\text{``}}{\underset{n=0}{\overset{\text{(}}{a}}} \frac{(-1)^{n} \frac{1}{2} p^{-2n+1} x^{4n+3}}{(2n+1)!(4n+3)}; \quad (2.65)$$

It follows from the analyticity of  $F_N(x)$  in F6603J/F6r[(å)]TJ/F69 8.9664 t 0d  $B_R$  cussed TJ/F65

equally spaced numbers between 0 and 1,000. The average elapsed times were 11.1 and 15.6 seconds, respectively, so that F(x,12) is almost 50% faster.

In Figure 2.2 we see that the theoretical error bounds are upper bounds as claimed, and that these bounds appear to capture the x-dependence of the errors fairly well, for example that  $E_N(x) = O(x)$  as  $x ! 0 = O(x^{-1})$  as x ! 4, and that  $E_N(x)$  reaches a maximum at about  $x = \frac{p}{2}A_N = \frac{p}{p(2N+1)}(-7.7 \text{ when } N = 9)$ .

Turning to C(x) and S(x), in Figure 2.3 we have plotted the maximum values of the absolute and relative errors in  $S_N(x)$  and  $C_N(x)$ , computed using fresnelCS in Table A.2. As accurate values for C(x) and S(x) we use  $C_{20}(x)$  and  $S_{20}(x)$  for x > 1:5 while, for 0 < x < 1:5 (following [52]) we approximate by the series (2.65) truncated after 15 terms, evaluated by the Horner algorithm. Exponential convergence is seen in Figure 2.3: the absolute errors are 4:5 10 <sup>16</sup> for N 11, the maximum relative error in  $C_N(x)$  is 3:6 10 <sup>15</sup> for N = 11 but that in  $S_N(x)$  as large as 2:7 10 <sup>13</sup>. These errors may be entirely acceptable, but the truncated power series (2.65) must achieve smaller errors for small x and is cheaper to evaluate. (Evaluating at 10<sup>7</sup> equally spaced points between 0 and 1:5 takes 2.9 times longer in Matlab with fresnelCS than evaluating 15 terms of both the series (2.65) via Horner's algorithm.)





**Fig. 2.2** Left hand side: Absolute error,  $jF(x) = F_N(x)j( )$ , and its upper bound  $h_N(x)$  given by (2.42) ( ), plotted against x. Right hand side: Relative error,  $jF(x) = F_N(x)j=jF(x)j( )$ , and its upper bound  $2(1 + \frac{1}{p}x)h_N(x)$  ( ), plotted against x. In both figures N = 9 and the exact value for F(x) is approximated by  $F_{20}$ 

### Chapter 3

## The Faddeeva function

#### 3.1 Introduction

The complex error function is defined by [46, (7.2.1)]

$$\operatorname{erf}(z) = \frac{p^2}{\overline{p}} \int_{0}^{z} e^{t^2} dt;$$

quadrants can be obtained using the symmetries [50, (3.1) and (3.2)]

w(z) = e<sup>z<sup>2</sup></sup> w(z) and w(
$$\overline{z}$$
) =  $\overline{w(z)}$ : (3.5)

Chiarella and Reichel [15] and Matta and Reichel [43] first proposed to compute erfc(z) for complex z by I (h; 0) given by (1.18) with H = p=h starting from the integral representation, which follows from (3.4), that

erfc(z) = 
$$\frac{z e^{z^2}}{p} \sum_{i=1}^{z^2} \frac{e^{t^2}}{z^2 + t^2} dt; \quad \text{Re}(z) > 0:$$
 (3.6)

Hunter and Regan [31] discussed the stability of these approximations when z is near one of the quadrature points, and proposed to use the formula I (h; 0), if jf (y=h) 0:5j 0:25, otherwise to use formula I (h; 1=2) given by (1.18) with H = p = h, where y = Im(z) and

$$f(t) = t [t] 2 [0; 1)$$
 (3.7)

is the function that gives the fractional part of t. This criterion and proposal is our main starting point for the methods developed in this chapter to approximate w(z).

There are a number of other effective schemes for computation of w(z), and we briefly summarise here the best of these. Gautschi [22] proposed an approximation for complex z based on continued fractions and this approximation is the basis of ACM TOM Algorithm 680 in Poppe and Wijers [50] which achieves a relative error of 10<sup>-14</sup> over nearly all the complex plane by Taylor expansions of degree up to 20 in an ellipse around the origin, convergents of up to order 20 of continued fractions outside a larger ellipse, and a more expensive mix of Taylor expansion and continued fraction calculations in between.

Weideman [62] proposed a rational approximation (the derivation starts from the integral representation (3.4)) to compute w(z), for Im(z) > 0. The approximation proposed is

w(z) 
$$P \frac{1}{\overline{p}(L iz)} + \frac{2}{(L iz)^2} \mathop{a}\limits^{N}_{n=0}^{1} a_{n+1} \frac{L + iz}{L iz}^{n};$$
 (3.8)

where the size of N controls the accuracy of the approximation,  $L = 2^{-1=4}N^{1=2}$  and the coefficients are computed as

$$\mathbf{a}_{n} := \frac{1}{2M} \sum_{j=M+1}^{M} (L^{2} + t_{j}^{2}) e^{-t_{j}^{2}} e^{-inq_{j}}; \quad n = 1; ...; N;$$
(3.9)

with M = 2N,  $t_j = L \tan(q_j=2)$  and  $q_j = p j=M$  for  $j = M + 1; \dots; M = 1$ . Weideman [62] argued that, for intermediate values of jzj, and as measured by operation counts, the work required to compute w(z) to 10<sup>-14</sup> relative accuracy is much smaller for the approximation (3.8) than for ACM TOMS Algorithm 680n [50].

**Remark 3.1.1.** Weideman [62] also compared his method to the modi ed trapezium rule approximation developed in [3, 31] and commented that the trapezium rule very accurate, provided for given z and N the optimal step-size h is selected. It is not easy, however, to determine this optimal h a priori." As we will see shortly, we address this comment

erfcx(y) =  $e^{y^2}$ erf(y) and  $S_1 := \stackrel{*}{a}_{k=1}^{*} \frac{1}{a^2k^2 + y^2} e^{(a^2k^2 + x^2)};$   $S_2 := \stackrel{*}{a}_{k=1}^{*} \frac{1}{a^2k^2 + y^2} e^{(ak + x)^2};$   $S_3 := \stackrel{*}{a}_{k=1}^{*} \frac{1}{a^2k^2 + y^2} e^{(ak - x)^2};$   $S_4 := \stackrel{*}{a}_{k=1}^{*} \frac{ak}{a^2k^2 + y^2} e^{(ak - x)^2};$   $TJ/TJ/F66 8277(2.965.977 3.455103)]TJ/F66 \stackrel{*}{a}_{k=1}^{*}1.958388/1287f939f 11.4p5 -116.8620.755(v)552 T056 2= Tf 11.4p5$ 

The authors have supplied us with the fatlab implementation of this methods [4] in the form of a Matlab function Faddeyeva\_v2(z,M), where the parameter is the number of accurate signi cant gures required, and the code enforces a choil veio f the range 4 M 13. In this Matlab implementation the choic = 1=2 is made and the sums in (3.14) are truncated, the number of terms retained depending in a complicated way on Zagloul and Ali [63] argued, using numerical calculations, that the approxima (Boh1), with appropriate choices for and truncation of (3.14)

$$A_{m} := \frac{p \overline{p}(2m \ 1)}{2^{2M}h} \sum_{n=N}^{N} e^{a^{2} = 4 \ n^{2}h^{2}} \sin \frac{p(2m \ 1)(nh + a = 2)}{2^{M}h} ; \qquad (3.19)$$

and

$$B_{m} := p \frac{i}{\overline{p} 2^{M-1}} \mathop{a}\limits^{N}_{n=} e^{a^{2}=4} e^{n^{2}h^{2}} \cos \frac{p(2m-1)(nh+a=2)}{2^{M}h} : \qquad (3.20)$$

Abrarov and Quine [1] argued, based on numerical calculations, that the approximation (3.16) is more accurate and faster (using the same number of summation terms in (3.16) as in (3.8)) than the approximation (3.8). We will be investigating these claims in Section §3.4 and we will be comparing the efficiency (accuracy and speed) of  $w_N(z)$  given in (3.21) with the approximations (3.8), (3.11) and (3.16).

We end this introduction by outlining the remainder of this chapter. Section 3.2 gives summary of the main results; §3.3 is concerned with the proposed approximation and its error bounds and §3.4 explores, using the theoretical and numerical calculations, the accuracy of our approximation in comparison with the approximations (3.8), (3.11) and (3.17).

#### 3.2 Summary of the main results

The main contributions of this chapter are: (i) to propose a family of approximations to w(z), based on the truncated modified trapezium rules defined in (1.22) adopting (at least for 0 arg(z) < p=4) the proposals of Hunter and Regan [31], but making explicit the choice of the step-size h as a function of N, the number of quadrature points addressing the criticism in Remark 3.1.1 by Weideman [62]; (ii) to prove completely explicit and rigorous bounds on both the absolute and relative errors as a function of N, uniform in z = x + iy, with x; y 0; and (iii) to demonstrate through the bounds and numerical experiments the high accuracy and efficiency of our approximation in comparison with the approximations (3.8), (3.12), (3.13) and (3.17).

The proposed approximation to w(z) for z = x + iy, with x; y = 0, is

$$\begin{cases} 8 \\ \gtrless \\ I_N(h; 1=2); & y \mod (x; p=h); \\ W_N(z) := & I_N(h; 0); & y < x \text{ and } jf (x=h) \quad 1=2j \quad 1=4; \\ \end{cases}$$
 x 0 Td [(;)]T90 G [-2]

where f is defined by (3.7),

$$I_{N}(h; 1=2) := \frac{2ihz}{p} \mathop{a}\limits^{N}_{k=0} \frac{e^{t_{k}^{2}}}{z^{2} - t_{k}^{2}}; \qquad (3.22)$$

$$I_N(h; 1=2) := \frac{2e^{-z^2}}{1+e^{-2ipz=h}} + I_N(h; 1=2);$$
 (3.23)

$$I_{N}(h;0) := \frac{2e^{-z^{2}}}{1 - e^{-2ipz=h}} + \frac{ih}{pz} + \frac{2ihz}{p} \sum_{k=1}^{N} \frac{e^{-t_{k}^{2}}}{z^{2} - t_{k}^{2}}; \qquad (3.24)$$

$$h = \frac{r}{\frac{p}{N+1}}; \quad t_k := (k+1=2)h \text{ and } t_k := kh:$$
 (3.25)

The main error estimate that we prove is

**Theorem 3.2.1.** Suppose  $v_N(z)$  is given by (3.21). Then, for z = x + iy with

The approximation  $w_N$  is proven in Theorem 3.2.1 (where we give completely explicit error bounds) to converge exponentially, uniformly in the first quadrant with respect to both absolute and relative errors, and this predicted rate of exponential convergence is observed in numerical experiments in Section §3.4 below (we know of no other rigorous error bounds for approximations for w(z) in the whole quadrant Re(z); Im(z) = 0).

This approximation is straightforward to code. Listing A.3 shows the Matlab code used to evaluate  $w_N$  for all the computations in this paper.

The approximation  $w_N$  is very competitive in accuracy and operation counts with other methods, as discussed in Section §3.4.

#### 3.3 The proposed approximation and its error bounds

In this section we derive the approximation  $w_N(z)$  given by (3.21) and its error bounds which demonstrate that the absolute and relative errors are both converging exponentially as N (the number of quadrature points) increases.

We can rewrite (3.4) as

$$w(z) = \sum_{i=1}^{Z} f(t) dt;$$
 (3.32)

where

$$f(t) = e^{t^2} F(t)$$
 and  $F(t) = \frac{iz}{p(z^2 - t^2)}$ : (3.33)

Note that the function  $e^{t^2}F(t)$  is even and meromorphic with simple poles at t = z. The residues at these two simple poles are

$$R_1 = \text{Res}(f; z) = \frac{ie^{-z^2}}{2p}$$
 and  $R_2 = \text{Res}(f; z) = -R_1$ : (3.34)

Using (1.16) and Remark 1.2.2, we have

$$C(h;a) = \frac{2e^{z^2}}{1 e^{2ip(a+z=h)}} \text{ so that } jC(h;a)j \quad \frac{2e^{2py=h}}{1 e^{2py=h}}e^{y^2 x^2}:$$
(3.35)

Applying the trapezium rule (1.6) to the integral in (3.32) leads to

$$I(h;a) = h \mathop{a}\limits_{k2Z} \frac{ize^{(k-a)^2h^2}}{p(z^2 (k-a)^2h^2)}:$$
(3.36)

Let

$$I(h;a) := I(h;a) + C(h;a);$$
 for  $a = 0; 1=2;$  (3.37)

where C(h; a) and I(h; a) are given by (3.35) and (3.36)

we have

$$jw(z) = I(h;a)j = \frac{2^{p} \overline{p} M_{H}(F) e^{H^{2} 2pH=h}}{1 e^{2pH=h}};$$
 (3.44)

where F is given by (3.33) and

$$M_{H}(F) := \sup_{t \ge R} jF(t + iH)j:$$
(3.45)

For H > 0 and z = t + iH, we have

$$jz^2$$
  $z^2j = jz$   $zjjz + zj$   $j$   $y$   $Hjjy + Hj = H^2$   $y^2;$ 

and hence we have, for H = p = h, that

$$jw(z) \quad I \quad (h;a)j \quad d_1(y) := \frac{2^p \overline{2}y e^{p^2 = h^2}}{\overline{p}(p^2 = h^2 \quad y^2) \quad 1 \quad e^{-2p^2 = h^2}}:$$
 (3.46)

Similarly and using the bound in (3.35) for C(h; a), we have for y  $\frac{5}{4}$ H, that

$$jw(z) = I(h;a)j = d_1(y) + jC(h;a)j = d_3(y)$$
: (3.47)

Select e in the range (0; H) and consider the case that jy Hj < e. We can easily show that

w(z) I (h; a) = 
$$\int_{C_H}^{Z} f(z)(1 g(z)) dz;$$
 (3.48)

where f is given by (3.33),  $g(z) = i \cot(pz=h+ap)$  and the contour  $C_H$ , passing above the pole of f at z = z, is the union of  $C_H$  and g, where  $C_H = ft + iH : t 2 R$  and j(t + iH) = zj > eg and  $g = fz + ee^{iq} : q_0 = q = p = q_0g$ , where  $q_0 = sin^{-1}((H = y)=e) = 2(p=2;p=2)$ .

For z 2  $C_H$ , it holds that

$$jz^{2} z^{2}j = jz zjjz + zj ejy + Hj:$$
 (3.49)

Thus, using (1.8), similarly to (3.46) we deduce that

<sup>Z</sup>  
<sub>C<sub>H</sub></sub> f(z)(1 g(z)) dz 
$$p \frac{2^{p} \overline{2} y e^{p^{2} = h^{2}}}{\overline{p} e(p = h + y) 1 e^{2p^{2} = h^{2}}}$$
: (3.50)

To bound the integral over g we note, for z = X + iY 2 g, that (3.49) is true and Y H. Further,

$$je^{z^2}j = e^{P_z}$$

#### where

$$P = Y^{2} X^{2}$$
  
= y^{2} x^{2} e^{2} \cos(2q) + 2e^{p} \overline{y^{2} + x^{2}} \sin(q \tan^{1}(y=

where

$$D_{h} \frac{p}{h} = \frac{4^{p} \overline{2}h + 2^{p} \overline{p}e^{-bp^{2} \pm h^{2}}}{p^{3=2} + 1 - e^{-2p^{2} \pm h^{2}}}; \qquad (3.56)$$

and b is given by (3.43).

**Proof.** It is easy to show, using (3.39), that  $D_h(y)$  and  $yD_h(y)$  are increasing functions of y for 0 y < p=h, in particular

$$D_{h} \quad \frac{3p}{4h} = \frac{3^{p} \overline{2}h}{14p(1 e^{2p^{2}=h^{2}})} < D_{h} \quad \frac{p}{h} = \frac{4^{p} \overline{2}h \ 1 + 2^{p} \overline{p}e^{bp^{2}=h^{2}}}{p^{3=2} \ 1 e^{2p^{2}=h^{2}}}: \quad (3.57)$$

Also we have, using (3.53), that

$$\frac{jw(z) \ I \ (h;a)j}{jw(z)j} \qquad (1 + {}^{p} \overline{p}jzj)jw(z) \ I \ (h;a)j (1 + {}^{p} \overline{2p}y)jw(z) \ I \ (h;a)j; \qquad (3.58)$$

and the two results follow.

In the following proposition we bound jw(z) = I(h; a)j and jw(z) = I(h; a)j=jw(z)j.

**Proposition 3.3.3.** Suppose that(h; a) is given by(3.36). Then, form > 0 and z = x + iy with x 0 and y max(x; p=h), we have

jw(z) I(h;a)j  $D_h \frac{5p}{4h} + \frac{2e^{1=4}}{1 e^{2p^2=h^2}}e$ 

where

where  

$$M_{H}(F) := \sup_{t2R} jF(t+iH)j \qquad \frac{p}{2y} \frac{1}{p(y^{2} - H^{2})}: \qquad (3.63)$$
Since  $\frac{y}{y^{2} - H^{2}}$  and  $\frac{y^{2}}{y^{2} - H^{2}}$  are both decreasing functions of y on (H;¥), we have  
 $\frac{y}{y^{2} - H^{2}} = \frac{H + e}{e^{2} + 2eH} = \frac{H + e}{2eH} = \frac{5}{8e} \text{ and } \frac{y^{2}}{y^{2} - H^{2}} = \frac{25}{32e}H: \qquad (3.64)$ 
Thus, we have  

$$-P = -\frac{1}{2e}$$

jw(z) I(h;a)j 
$$\frac{5^{4} \overline{2}}{4^{4}}$$

Similarly, using (3.53) and since  $yD_h(y)$  and  $\frac{2ye^{y^2} + 2py=h}{1 + e^{-2py=h}}$  are both increasing functions of y for H y < e, we have that

$$\frac{jw(z) \quad I(h;a)j}{jw(z)j} \qquad (1 + p \frac{1}{2p}y)(jw(z) \quad I(h;a)j + jC(h;a)j)$$
$$1 + \frac{5^{p} \frac{1}{2p} y^{3-2}}{2p^{3-2}}$$

where  $D_h = \frac{5p}{4h}$  is given by (3.61).

Proof. De ne

 $E_h(z) = w(z)$  I (h; a) and  $e_h(z) = E_h(z)=w(z)$ ;

on G := f z 2 C :  $0 < \arg(z) < p=4g$ . Sincew(z) and (h; a) are both entire functions of and, using(3.53), w(z)  $\in$  0 for all z 2 G,  $E_h(z)$  and  $e_h(z)$  are analytic or G and continuous on its closure. From the asymptotic expansion  $\psi(z)$  in the complex plane (se@2, (2.6)]) it follows that w(z) ! 0 as jzj ! ¥, uniformly for  $0 < \arg(z) < p=4$ . Moreover it follows from (3.37) and (3.35) that the same holds for (h; a) and hence fo  $E_h(z)$ . Thus we have, using Lemma 3.3.1, that

$$supjE_{h}(z)j = supjE_{h}(z)j:$$

$$z_{2}G \qquad z_{2}MG$$

Let  $z = re^{ip=4}$  with r 0. Then, using Proposition 3.3.1, we have that

Now, for z2 G, using (3.53) and (3.71),

$$je_h(z)j$$
 (1+<sup>p</sup> $\overline{p}jzj)jE_h(z)j$  Pe<sup>jzj</sup>;

where  $P := M D_h \frac{5p}{4h}$  e  $p^2 = h^2$  and  $M := max(1 + p \overline{p}jzj)e^{jzj}$ , for z 2 G. Thus we have, using Lemma 3.3.1, that

$$\sup_{z \in G} je_h(z)j = \sup_{z \in \P} je_h(z)j:$$
(3.75)

Let  $z = re^{ip=4}$  with r = 0. Then, we have, using Proposition 3.3.1, that  $yD_h(y)$  is increasing on  $0; \frac{5}{4}\frac{p}{h}$  and decreasing on  $\frac{5}{4}\frac{p}{h}; \neq$  with  $D_h \frac{5}{4}\frac{p}{h} > D_h \frac{5}{4}\frac{p}{h}$ ; thus we have

$$je_{h}(z)j = 1 + \frac{5^{p} \bar{2}p^{3=2}}{4h}^{!} D_{h} = \frac{5p}{4h} e^{-p^{2}=h^{2}}$$
 (3.76)

Let z = x + ie with 0 < e < p = h. Then we have, using (3.53) and Proposition 1.2.4, that

$$je_{h}(z)j \qquad (1+\stackrel{p}{\overline{p}jzj})jE_{h}(z)j \\ \frac{2jzj(1+\stackrel{p}{\overline{p}jzj})e^{p^{2}=h^{2}}}{p(1 e^{2p^{2}=h^{2}})} \stackrel{Z}{\xrightarrow{}} \frac{e^{t^{2}}}{iz^{2}} \frac{e^{t^{2}}}{(t+ip=h)^{2}j}dt;$$

Taking the limit  $e! 0^+$ , since both sides in the above bound are continuous for 0 < e < p=h, we obtain

$$je_{h}(x)j = \frac{2x(1 + p \overline{p}x)e^{-p^{2} \pm h^{2}}}{p(1 - e^{-2p^{2} \pm h^{2}})} \sum_{i=1}^{i} G(t) dt; \quad x = 0; \quad (3.77)$$

where

$$G(t) = \frac{e^{t^2}}{jx^2 (t + ip = h)^2 j}$$
:

Note

$$Z_{*} G(t) dt = Z_{*} G(t) dt + Z_{*} G(t) d$$

Since, for x 0 and t 2 R,

$$jx^{2}$$
  $(t + ip=h)^{2}j = jx$  t  $i(p=h)jx + t + i(p=h)j$   $\frac{p}{h}^{q} \frac{q}{x^{2} + (p=h)^{2}};$ 

we have

$$\sum_{x=2}^{Z} G(t) dt = \frac{h}{p} \frac{Z}{\overline{x^2 + (p=h)^2}} \sum_{x=2}^{Z} e^{t^2} dt = \frac{hxe^{x^2=4}}{p} \frac{hxe^{x^2=4}}{\overline{x^2 + (p=h)^2}};$$
 (3.79)

$$\sum_{3x=2}^{Z} G(t) dt \quad \frac{p}{p} \frac{h}{\overline{x^2 + p^2 = h^2}} \sum_{3x=2}^{Z} e^{t^2} dt \quad \frac{h e^{9x^2 = 4}}{3px};$$
(3.80)

with f given by (3.33) and t<sub>k</sub> and t<sub>k</sub> are given by (3.25).

2<sup>p</sup> –

We will call the error in approximating I(h; a) by  $I_N(h; a)$  the truncation error, given by

$$T_{N}(h;a) := 2h \mathop{a}\limits_{k=N+1}^{4} f((k+a)h):$$
(3.88)

**Proposition 3.3.5.** Suppose k is given by (3.25) and jz  $t_k j$  h=4 for k = N + 1; N + 2; ... and z= x + iy with 0 y < x. Then, for h> 0,

For the second summation we have that

$$2h \overset{\texttt{¥}}{\overset{\texttt{a}}{\texttt{h}}} \frac{e^{t_{k}^{2}}}{jz \ t_{k}j} \qquad \frac{4}{h} \quad 2h \overset{\texttt{¥}}{\overset{\texttt{a}}{\texttt{a}}} e^{t_{k}^{2}}$$

$$\frac{4}{h} \quad 2he^{t_{M}^{2}} + 2h \overset{\texttt{¥}}{\overset{\texttt{a}}{\texttt{a}}} e^{t_{k}^{2}}$$

$$\frac{4}{h} \quad 2he^{t_{M}^{2}} + 2h \overset{\texttt{¥}}{\overset{\texttt{a}}{\texttt{a}}} e^{t_{k}^{2}}$$

$$\frac{4}{h} \quad 2he^{t_{M}^{2}} + 2 \overset{\texttt{Z}}{\overset{\texttt{X}=\mathsf{M}+1}{t_{M}}} e^{t^{2}} dt$$

$$\frac{4}{h} \quad \frac{1 + 2ht_{M}}{t_{M}} e^{t_{M}^{2}}$$

$$\frac{4}{h} \quad \frac{1 + 2ht_{M}}{q_{X}} e^{t_{M}^{2}}$$

Note that  $(1 + 2ht)e^{t^2}$  is a decreasing function of t for t  $t_0$ , where  $t_0 := 2h=(1 + p^2)^2$  and  $t_0 < h < t_{N+1}$ . Thus we have that

$$2h \overset{*}{\overset{a}{a}}_{k=M} \frac{e^{t_{k}^{2}}}{jz t_{k}j} \qquad \frac{4}{h} \frac{1+2ht_{N+1}}{qx} e^{t_{N+1}^{2}}$$
(3.93)

We have, using 
$$q \frac{1}{\overline{x^2 + t_{N+1}^2}} = \frac{1}{t_{N+1}}$$
 and (3.91), (3.92) and (3.93), that

$$jT_N(h;0)j = \frac{P \overline{2}(1+2ht_{N+1})}{pt_{N+1}} = \frac{1}{(1-q)t_{N+1}} + \frac{4}{hq} = t_{N+1}^2$$

Choose q such that

$$\frac{1}{(1 q) t_{N+1}} = \frac{4}{hq};$$

i.e.

$$q = \frac{4t_{N+1}}{h+4t_{N+1}}$$
:

Then we have that

$$jT_{N}(h;0)j = \frac{2^{P} \overline{2}(1+2ht_{N+1})(h+4t_{N+1})}{pht_{N+1}^{2}} e^{t_{N+1}^{2}}$$
(3.94)

Similarly, we have, using  $q \frac{x}{x^2 + t_{N+1}^2}$  1 and (3.91), (3.92) and (3.93), that

xjT<sub>N</sub>(h;0)j 
$$\frac{2^{p} \overline{2}(1+2ht_{N+1})(h+4t_{N+1})}{pht_{N+1}} e^{t_{N+1}^{2}}$$
: (3.95)

In a similar way we can prove the following result for T(h; 1=2).

**Proposition 3.3.6.** Suppose is given by (3.25) and jz  $t_k j$  h=4 for k = N + 1; N + 2; ... and z= x + iy with 0 y < x. Then, for h> 0,

$$jT_N(h; 1=2)j$$
  $\frac{2^{p} \overline{2}(1+2ht_{N+1})(h+4t_{N+1})}{pht_{N+1}^2} e^{t_{N+1}^2}; and$  (3.96)

$$\frac{jT_{N}(h; 1=2)j}{jw(z)j} \qquad (1 + \frac{p}{2p}T_{N+1})jT_{N}(h; 1=2)j:$$
(3.97)

Proof. Suppose that 0 < q < 1, then we have, using (3.88) with a = 1=2, that  $p_{-}$ 

jT<sub>N</sub>(h;1=2)j

Then we have that

$$jT_{N}(h; 1=2)j = \frac{2^{p} \overline{2}(1+2ht_{N+1})(h+4t_{N+1})}{pht_{N+1}^{2}}e^{t_{N+1}^{2}}$$
(3.101)

Similarly, we have, using1

**Proposition 3.3.7.** Suppose a = 0 or a = 1=2 and z = x + iy with y = x = 0. Then, for h > 0,

$$jT_N(h;a)j = \frac{(1+2ht_{N+1})}{pt_{N+1}^2}e^{-t_{N+1}^2}; and$$
 (3.108)

$$\frac{jT_{N}(h;a)j}{jw(z)j} \qquad \frac{(1+2ht_{N+1})(1+2^{p}\overline{p}t_{N+1})}{pt_{N+1}^{2}}e^{-t_{N+1}^{2}}: \qquad (3.109)$$

**Proof.** Suppose  $t_k$  and  $t_k$  be given by (3.25) and F(t) is given by (3.33). Then, for z = x + iy with y = x = 0,

$$jz^2$$
  $t_k^2 j^2 = y^4 + t_k^4 + x^4 + 2x^2y^2 + 2t_k^2(y^2 - x^2) j z^2 - t_k^2 j^2$ 

Thus, we have

$$jT_N(h;a)j$$
  $2h \overset{\texttt{Y}}{\overset{\texttt{a}}{a}} e^{t_k^2} jF(t_k)j;$ 

and, using (3.53),

$$\frac{jT_{N}(h;a)j}{jw(z)j} = (1 + {}^{p}\overline{p}jzj) 2h \overset{¥}{\overset{a}{a}} e^{t_{k}^{2}}jF(t_{k})j \\ (1 + {}^{p}\overline{2p}y) 2h \overset{¥}{\overset{a}{a}} e^{t_{k}^{2}}jF(t_{k})j ; y 0:$$

Since

$$jz^2$$
  $t_k^2 j^2 = y^4 + t_k^4 + x^4 + 2x^2 y^2 + 2t_k^2 (y^2 - x^2) - y^4 + t_k^4;$ 

$$jT_{N}(h;a)j \qquad \frac{2^{p} \overline{2}hy}{p} \overset{\texttt{¥}}{\overset{\texttt{a}}_{k=N+1}} q \frac{e^{-t_{k}^{2}}}{\overline{y^{4} + t_{k}^{4}}} \\ -\frac{p \overline{2}y}{\overline{y^{4} + t_{N+1}^{4}}} 2he^{-t_{N+1}^{2} + 2} \overset{\texttt{Z}}{\overset{\texttt{¥}}_{t_{N+1}}} e^{-t^{2}}dt \\ -\frac{p \overline{2}y(1 + 2ht_{N+1})}{p \overline{2}y(1 + 2ht_{N+1})} e^{-t_{N+1}^{2}} e^{-t_{N+1}^{2}}:$$

Moreover

$$q \frac{y}{y^4 + t_{N+1}^4} = p \frac{1}{2t_{N+1}}$$
 and  $q \frac{y^2}{y^4 + t_{N+1}^4} = 1$ :

The Faddeeva function

# 3.59

# 3.Proof.
Since

$$jw(z) = I_N(h;a)j j w(z) = I(h;a)j + jT_N(h;a)j;$$

and

$$\frac{jw(z) \quad I_N(h;a)j}{jw(z)j} \quad \frac{jw(z) \quad I(h;a)j}{jw(z)j} + \frac{jT_N(h;a)j}{jw(z)j};$$

the first result follows by combining (3.125) and (3.127) and the second result follows by combining (3.126) and (3.128).  $\hfill \Box$ 

**Remark 3.3.3.** Using Proposition 3.3.3 and Theorem 3.3.3, we can easily show; four iy with  $0 < x \le p=h$ , that

$$jw(z) I_N(h;a)j b_N e^{pN}$$
 (3.129)

and

$$\frac{jw(z) \quad I_N(h;a)j}{jw(z)j} \quad b_N^{\ p} \overline{N+1} e^{-pN}; \qquad (3.130)$$

where  $b_N$  and  $b_N$  are given by(3.122) and (3.123), respectively.

### 3.4 Numerical results

In this section we show numerical calculations that illustrate and confirm the theoretical (this section wilvo g 0620 mbining (p100 (i) and the second r48 TJ/ [([(i) 91) 25(w)-ining (rejults (Theorems 3.3.2)]w Th02 11.9552 Tf 5.977 0 Td1[(ET[(j) 91) 25(w)-wt 10.005 3.391 Td]]

(iii) the approximationw<sub>N</sub>, with N 14, is signi cantly more accurate than the approximation (3.8) from Weideman [62];

Figure 3.2 below shows that  $x_N(z)$  is very accurate  $a_j z_j = 0$ , and with N as small a the computed relative error is 10<sup>-12</sup>, which con rms the calculations in Figure 3.1.

We will comment now on the accuracy and the efficiency of computiv(z) using the approximation  $w_N(z)$  given by (3.21) and its codew(z,N) in Listing A.3 in comparison with the approximation (\$3.8), (3.11) and (3.16) and their codes. We do not have access to exact values forw(z) and so we use four different accurate approximation (\$2):

- (i) Our own approximation  $w_N(z)$  with N = 20 computed by the call (z,20) to the code in Listing A.3;
- (ii) Weideman's approximatio((8.8)) with N = 40 (this choice of N gives maximum accuracy for this approximation), implemented by the **case** f(z,40) in Table 1 [62];
- (iii) The approximation(3.11) of Zagloul and Ali [63], implemented in the Matlab code [64], supplied to us by the author, computed by the Eakldeyeva\_v2(z,M) with M = 13 (the maximum value permitted by the code), where the number of accurate signi cant gures required, which must be in the range 4M 13;
- (iv) The approximation(3.16) of Abrarov and Quine1[], implemented as the the atlab function comperf(z) of Abrarov and Quine1[, Appendix], which uses the method (3.16) with a = 2:75 and M = 5.

The maximum absolute errors and computation times are shown in Table 3.1 (using Matlab (R2015a) on a laptop with Intel core i7-4510U 2.00 GHz processor/( $p_{i}$ ) in Table A.3, cef(z,40) in Table 1 of Weideman6[2], comperf(z) of Abrarov and Quine 1, Appendix] and the method of M. Zaghloul and A. Al63] as implemented inFaddeyeva\_v2(z,13) of [64]. The calculations are implemented f $p_{i} = 10^{p}e^{iq}$ , with p = -6(0:0006)6 and q = 0(p=400)p=2 giving in total 4020201 values. It can be seen from Table 3.1 that the approximationw<sub>N</sub> given by(3.21), with N as small as 1, is as accurate as most accurate version of the approximatio(8.11) in Zagloul and Ali [63] as implemented inf[4] with a = 1=2 and M = 13, and signi cantly more accurate than tMetatlab code of Abrarov and Quine [1] based or(3.16) with a = 2:75 and M = 5, and at least as accurate as Weideman's approximation(3.8) with N = 40

Algorithm	Maximum absolute error	Computation time in seconds			
w(z,11)	1:11 10 <sup>15</sup>	0.64			
cef(z,40)	1:30 10 <sup>15</sup>	1.46			
comperf(z)	5:53 10 <sup>10</sup>	0.90			
Faddeyeva_v2(z,13)	3:92 10 <sup>15</sup>	0.51			

**Table 3.1** Accuracy and computation times of the Matlab codes of the approximations (3.21), (3.8), (3.11) and (3.16).



Fig. 3.2 The surfaces of the absolute (top) and relative (bottom) errors of the approximation  $w_N(z)$  given by (3.21) with N = 9, where the exact value of w(z) is computed by  $w_{20}(z)$ .

# Chapter 4

# The 2D impedance half-space Green's function for the Helmholtz equation

## 4.1 Introduction

This chapter is concerned with the problem of calculating sound propagation from a monofrequency coherent line source above an impedance plane. The interest in this problem has been motivated by the development of boundary element methods (BEMs) for the calculation of outdoor sound propagation for many applications (e.g. [26], [11], [12] and [13]). These  $d^0 = j\mathbf{r} + \mathbf{r}_0^0 \mathbf{j}$  be the distance from the image source to the receiver and  $\mathbf{r} = \mathbf{k} d^0$ , where k is the wave number that satisfies  $\mathbf{k} = 2\mathbf{p} = \mathbf{l}$  where l is the wavelength.

The problem is to calculate the acoustic pressure at **r**, denoted by  $G_b(r; r_0)$ , due to the source at  $r_0$ , where b is the normalised admittance of the impedance plane with Re(b) > 0.  $G_b(r; r_0)$  (the Green's function) satisfies the following conditions:

(i) the Helmholtz equation, that is

 $\tilde{N}^{2}G_{b}(\mathbf{r};\mathbf{r}_{0}) + k^{2}G_{b}(\mathbf{r};\mathbf{r}_{0}) =$ 

#### 4.1 Introduction

deform the path L to the steepest descent path (see [35], [8] and [14]), and obtain [14]

$$P_{b}(\mathbf{r};\mathbf{r}_{0}) = P_{b}^{(G)} + P_{b}^{(s)};$$
(4.11)

where

$$P_{b}^{(G)} = \frac{b e^{ir}}{p} \sum_{i=1}^{Z} e^{-r t^{2}} F(t) dt; \qquad (4.12)$$

with

$$\mathsf{F}(\mathsf{t}) := \quad \mathsf{p}\frac{\mathsf{b} + \mathsf{g}(1 + \mathsf{i}\mathsf{t}^2)}{\overline{\mathsf{t}^2 \ 2\mathsf{i}}(\mathsf{t}^2 \ \mathsf{z}_1^2)(\mathsf{t}^2 \ \mathsf{z}_2^2)}; \qquad \frac{\mathsf{p}}{2} < \arg^{\mathsf{p}} \ \overline{\mathsf{t}^2 \ 2\mathsf{i}} < \frac{\mathsf{p}}{2}; \tag{4.13}$$

$$z_1 := \stackrel{p}{\underset{q}{\longrightarrow}} \frac{\overline{a_+}}{a_+}; \qquad \frac{p}{4} < \arg^p \frac{\overline{a_+}}{a_+} < \frac{3p}{4}; \qquad (4.14)$$

$$a := 1 + bg = 1 b^2 - 1 g^2; Re = 1 b^2 - 0;$$
 (4.16)

and

$$\mathsf{P}_{\mathsf{b}}^{(\mathsf{s})} := \frac{\mathsf{b}\,\mathsf{e}^{\mathsf{i}\mathsf{r}}}{\mathsf{p}} \frac{\mathsf{p}\,\mathsf{e}^{-\mathsf{i}\mathsf{r}\,\mathsf{a}_{+}}}{2} \,\mathsf{d}_{\mathsf{s}}; \tag{4.17}$$

where

$$\begin{array}{cccc} 8 \\ \gtrless & 2; & \text{Im} b < 0; \text{Re} a_{+} < 0 \\ d_{s} := & 1; & \text{Im} b < 0; \text{Re} a_{+} = 0; \\ \Re & 0; & \text{otherwise} \end{array}$$
 (4.18)

The integral representation (4.12) from [14] will be the starting point for our proposed approximation of  $P_b$ .

Numerical computation of the solution of the problem (

has been widely cited and applied (e.g. [41], [25], [7], [47] and [39]) as a well-established method for solving this problem. In particular, it is used in many papers as an efficient method for the solution of outdoor sound propagation problems via the BEM (e.g. [32], [34], [49], [51]). The following representations for  $P_b(r; r_0)$  is derived and used in [14]:

$$P_{b}(\mathbf{r};\mathbf{r}_{0}) = \frac{b e^{ir}}{p} \int_{0}^{Z_{*}} t^{-1=2} e^{-rt} f(t) dt; \quad Im(b) > 0 \quad or \quad Re(a_{+}) > 0; \quad (4.19)$$

and

$$P_{b}(\mathbf{r};\mathbf{r}_{0}) = \frac{b e^{jr}}{p} \frac{Z_{4}}{ct > TJ/F95 \ 8.96642 \ Tf \ 15525508 \ 4.938 \ Td \ []TJ/F102 \ 11.955b}$$

and H is the Heaviside step function defined by

$$H(t) := \begin{cases} 0 \\ \gtrless \\ 1; \\ 1=2; \\ 0; \\ t < 0; \end{cases}$$
(4.32)

and

$$d_{+}^{(1)} := \begin{cases} 8 \\ \gtrless 2e^{2ipz_{1}=h}; \\ 1+e^{2ipz_{1}=h}; \\ \end{Bmatrix} y_{1} < 0; \\ y_{1} = 0; \\ 2; \\ y_{1} > 0; \end{cases}$$
(4.33)

La Porte [38] proved a bound on  $jP_b = P_b^{h;N;H}j$  derived largely from Proposition 1.2.4 and using, for F given by (4.13), that

$$M_{H}(F) := \sup_{x \ge R; j \neq j = H} jF(x + iy)j \qquad p \frac{jbj + g}{1 - H^{2}} M_{I};$$
(4.34)

where

$$\mathbf{M} := \max 3; \frac{2\max(\mathbf{x}_1^2; \mathbf{x}_2^2) + 2}{\mathbf{j}\mathbf{H}^2 + \mathbf{y}_1^2 \mathbf{j}\mathbf{j}\mathbf{H}^2 + \mathbf{y}_2^2 \mathbf{j}}; \qquad (4.35)$$

and  $x_j = \operatorname{Re}(z_j)$  and  $y_j = \operatorname{Im}(z_j)$  for j = 1; 2.

La Porte [38] showed, using numerical calculations, that the approximation in (4.28) achieves with N = 11 higher accuracy than the approximation (4.25), with n = 40 and m = 22, in Chandler-Wilde and Hothersall [14] for 0:5 r 8:54, 0 g 1 and 0:1 j bj 1.

This chapter of the thesis builds on the work of La Porte [38] but extends this work significantly. The main issues with the approximation  $P_{a_{\pm}}^{h;N;H}$  in (4.28) are that: (i) the approximation formula blows up if the simple pole at  $z_1 = P_{a_{\pm}}^{h;N;H}$  coincides with a quadrature point at kh and is inaccurate in floating point arithmetic when  $z_1$  is close to kh; (ii) the expression (4.31) blows up when a = 2 and is inaccurate in floating point arithmetic when a is close to 2; and (iii) the bound (4.34) blows up when  $H = Im(z_1)$  or  $H = Im(z_1)$ . In this chapter of the thesis we address all these issues: we propose an approximation which is stable for numerical calculations for r > 0, 0 = 1 and b with Re(b) > 0; we prove a rigorous and uniform error bound for this approximation; and finally we show through systematic numerical experiments that this approximation is at least as accurate as the approximation (4.28) in La Porte [38] and is more accurate and more efficient than the approximation of Chandler-Wilde and Hothersall [14].

Recently, O'Neil et al. [47] propose a method of computing  $P_b(r;r_0)$ , for 0 b 1, based on the following representation for  $P_b(r;r_0)$ :

$$P_{b}(r; r_{0}) = I_{1} + I_{2};$$

where

$$I_{1} := \frac{ikb}{2p} \sum_{0}^{Z_{1}} H_{0}^{(1)}(kjr \ \mathbf{e}_{0}j) e^{ikbh} dh; \ \mathbf{e}_{0} = (x_{0}; (y_{0} + h));$$

and

$$I_{2} := \frac{ikb}{2p} \frac{Z}{4} \frac{e}{p} \frac{p^{p} \frac{1}{12} \frac{1}{k^{2}} (y+y_{0})}{p^{p} \frac{1}{12} \frac{1}{k^{2}} \frac{1}{k^{2$$

experiments to demonstrate the accuracy of the proposed approxining tion comparison with the approximations of Chandler-Wilde and Hothersall [14] and La Porte [38].

q  $\frac{\text{Let P}_{b}(r;r_{0}) \text{ be given by equation}(4.11)-(4.18) \text{ and } H := \min(0:9; \mathbf{A}_{N}) \text{ with } \mathbf{A}_{N} := \frac{p_{0}(N+1)}{2p(N+1)}$ , and recall that

with

$$\begin{split} & C_{N} := (jbj+1)^{4} \frac{384^{p} \overline{10}(4jbj+7)(1+4^{p} \overline{pr})r^{3=2}}{p^{3=2}(N+1)^{2} 1 e^{2p(N+1)=p^{n} \overline{3}}} + 20 1 + \frac{1}{R_{N}} 5; \quad (4.46) \\ & \mathfrak{E}_{N} := (jbj+1)^{u} \frac{781^{p} \overline{10p}(4jbj+7)(1+4^{p} \overline{pr})}{p \overline{r} 1 e^{0:9p=h_{N}} (N+1)^{1=3}} + \mathfrak{K}_{N} \quad \text{and} \\ & \mathfrak{K}_{N} := 8K_{N} 1 + \frac{r^{1=3}}{ap^{1=3}H^{1=3}(N+1)^{1=3}}; \quad (4.47) \end{split}$$

where  $K_N$  is given by (4.120)

Remark 4.3.1. The advantage of choosing the branch cut  $foir(a_{+} 2)$  as in (4.51) is that a cut from 2 to + ¥ on the positive real axis in the -plane is implied. This is convenient, sincea<sub>+</sub> 2 is impossible unless 1. Thus,  $iin(a_{+} 2)$ , considered as a function **b**f, is analytic in the cut half-plane.

The formulas(4.51)and(4.52)for  $R_1$  and  $R_2$  are not numerically stable in oating point arithmetic wherb and g are close to zero, close **f**oor whenb = g. We simplify them in the following lemma to make them more stable in numerical calculations.

Lemma 4.3.1. For 0 g 1 and b in the cut half-plane, we have that

$$R_{1} = \frac{ie^{ira_{+}}}{4} \frac{1}{1} \frac{b^{2}}{b^{2}};$$
(4.54)

$$R_2 = \frac{10^{10}}{4} \frac{1}{1} \frac{b^2}{b^2} W;$$
(4.55)

where

Q

where Wis given by (4.56),

$$d_{+} := \begin{cases} 3 \\ 2e^{2ip(a+z_{1}=h)}; \\ 1+e^{2ip(a+z_{1}=h)}; \\ 2; \\ 2; \end{cases} y_{1} < 0; \qquad (4.73)$$

and **H** is the Heaviside step function given by (4.32).

#### 4.3.1 Bounding the discretisation error

This section is concerned with bounding, for h > 0 and a = 0 or a = 1=2,

E(h;a) := I I(h;a);

where I and I (h; a) are given by (4.50) and (4.70), respectively.

Since F given by (4.13) is meromorphic for jIm(t)j < 1, we will be defining H throughout this chapter as

$$H := \min \ 0.9; \frac{p}{rh} :$$
 (4.74)

Then, we have the following result.

**Proposition 4.3.1.** Let h > 0 and  $H := \min 0:9; \frac{p}{r h}$ . Then

jE (h;a)j 
$$\frac{D(H)e^{rH^2=4 pH=h}}{1 e^{pH=h}};$$
 (4.75)

where

$$D(H) := \frac{512^{p} \overline{10p}(jbj+1)(4jbj+7)(1+4^{p} \overline{pr})}{p \overline{r} H^{4}} + \frac{2p}{j1 b^{2} j^{1=2}}:$$
 (4.76)

**Proof.** Let  $z_1 = x_1 + iy_1$  and  $z_2 = x_2 + iy_2$  be given by (4.14) and (4.15), respectively. Select e 2 (0; H=4) and consider the case jH j y<sub>1</sub>jj e and jH y<sub>2</sub>j e. Then, using Proposition 1.2.4, we have that

jE (h;a)j 
$$p \frac{2^{P} \overline{p} M_{H}(F)}{\overline{r}(1 e^{2pH=h})} e^{rH^{2} 2pH=h};$$
 (4.77)

and using equation (4.34) and noting  $x_j^2$  j  $z_j j^2$  2+ 2jbj with j = 1;2, it holds that

$$M_{H}(F) \qquad \frac{P}{10} \frac{10(jbj+1)}{P(1+H)} \max \quad 3; \frac{2\max(x_{1}^{2}; x_{2}^{2}) + 3}{e^{2}(jy_{1}j+H)(y_{2}+H)}$$

$$P_{10}(jbj+1) \max \quad 3; \frac{7+4jbj}{e^{2}(jy_{1}j+H)(y_{2}+H)} + 4e)$$

$$P_{10}(jbj+1) \max \quad 3; \frac{7+4jbj}{e^{2}(H-4e)^{2}} : \qquad (4.78)$$

We consider now the case jH j  $y_1$ jj < e or jH  $y_2$ j < e. Let D be the region in the complex plane defined by

$$D := f z: 0 < Im(z) < Hgn \lim_{j=1;2}^{l} B_e(z_j);$$
(4.79)

where, for j = 1; 2,

$$B_{e}(z_{j}) := \begin{cases} f z : jz \quad z_{j}j < eg; & \text{if } j \text{Im}(z_{j}) & \text{H}j < 2e; \\ \vdots & f; & \text{otherwise} \end{cases}$$
(4.80)

and let, for j = 1; 2,

$$g_j = f z 2 \ \P D : jz \quad z_j j = eg \quad \text{and} \quad G_H = f z 2 \ \P D : z = t + iH; \ t \ 2 \ Rg;$$

where  $\P D$  is the boundary of D. Then we can show, recalling that g, F and C(h; a) are given by (1.7), (4.13) and (4.72), respectively, that

$$\begin{array}{c} z \\ jE(h;a)j \\ G_{H} \end{array} \stackrel{r}{=} e^{rz^{2}}F(z)(1 \ g(z))dz + a^{2}a \\ j=1 \ g_{j} \end{aligned} e^{rz^{2}}F(z)(1 \ g(z))dz + jC(h;a)j:$$

$$\begin{array}{c} (4.81) \end{array}$$

If  $H = (jy_1j) + (jy_1j) + (jy_2) + ($ 

$$jC(h;a)j \qquad \frac{p}{2j1 \ b^2 j^{1=2}} \quad \frac{2e^{2pjy_1j=h}}{1 \ e^{2pjy_1j=h}} + \frac{2e^{2py_2=h}}{1 \ e^{2py_2=h}}$$
$$\frac{p}{2j1 \ b^2 j^{1=2}} \quad \frac{4e^{2p(H \ 4e)=h}}{1 \ e^{2p(H \ 4e)=h}}$$

and

$$\frac{2^{p} \overline{p} M_{H}(F)}{\overline{r}(1 e^{2pH=h})} e^{rH^{2} 2pH=h} = \frac{512^{p} \overline{10p}(jbj+1)(7jbj+4)}{p} e^{rH^{2}=4} e^{H=h}$$
 (4.94)

Thus, the result follows by combining, with e = H=8, (4.82), (4.92) and (4.94).

#### 4.3.2 Bounding the truncation error

This section will give bounds on the truncation error  $T_N(h; a)$  as defined in (1.24) for a = 0 or a = 1=2. We will present two results on the truncation errors  $T_N(h; 0)$  and  $T_N(h; 1=2)$  and then we propose a scheme for choosing the step-size h. This scheme will be used to simplify further the bounds on  $T_N(h; 0)$  and  $T_N(h; 1=2)$ .

Recall that  $z_1 = x^{+7} d$  [(=)]TJ/F69 11.J/F66633.886 Td [(1)]TJ/F102 19 055211.J/F66633+7 0 T

Combining the above inequalities, we find that

$$jF(t_k)j \qquad \frac{8x_2(jbj+1)^{-1} + t_k^4}{h^4 + t_k^4(t_k+jx_1j)(t_k+x_2)} \qquad (4.101)$$

$$\frac{\delta(JDJ+T) - T + \tau_k}{h(t_k + jx_1j)}$$
(4.102)

$$\frac{8(jbj+1)(1+t_k)}{h(t_k+jx_1j)};$$
(4.103)

where the last line comes from

$$\frac{1+t}{\overline{2}} \quad (1+t^4)^{1=4} \quad (1+t^2)^{2} \quad ^{1=4} \quad (1+t):$$

Also, note that

$$\frac{d}{dt} \quad \frac{1+t}{jx_1j+t} = \frac{jx_1j}{(t+jx_1j)^2};$$

thus we have that

$$jF(t_{k})j = \frac{8}{h}(jbj+1) \begin{cases} \frac{1+t_{N+1}}{jx_{1}j+t_{N+1}}; & \text{if } jx_{1}j = 1; \\ \vdots & \vdots \\ 1; & \text{otherwise }; \end{cases}$$
(4.104)

but

$$\frac{1+t_{N+1}}{jx_1j+t_{N+1}} = 1+\frac{1}{t_{N+1}} ;$$

and hence the result follows.

**Proposition 4.3.2.** Leth > 0, N 2 N, F(t) be given by (4.13) and  $t_k = kh \text{ with } jt_k z_1 j$  h=4 for k = N + 1; N + 2; :::. Then, for

$$T_{N}(h;0) := 2h \overset{\texttt{Y}}{\underset{k=N+1}{\overset{\texttt{P}}{a}}} e^{-rt \overset{2}{\underset{k}{\overset{\texttt{P}}{c}}}} F(t_{k});$$

we have

$$jT_N(h;0)j = \frac{8(jbj+1)(1+2hrt_{N+1})}{hrt_{N+1}} + \frac{1}{t_{N+1}} e^{-rt_{N+1}^2}$$
(4.105)

Proof. Using Lemma 4.3.2 we find that

$$jT_{N}(h;0)j = \frac{8M_{N}(jbj+1)}{h} + 2h \overset{¥}{\overset{a}_{k=N+1}} e^{rt \frac{2}{k}}$$

$$= \frac{8M_{N}(jbj+1)}{h} + 2h e^{rt \frac{2}{N+1}} + 2h \overset{¥}{\overset{a}_{k=N+2}} e^{rt \frac{2}{k}}$$

$$= \frac{8M_{N}(jbj+1)}{h} + 2h e^{rt \frac{2}{N+1}} + 2 e^{rt \frac{2}{k}} e^{rt^{2}} dt$$

$$= \frac{8M_{N}(jbj+1)}{h} + 2h e^{rt \frac{2}{N+1}} + \frac{e^{rt \frac{2}{N+1}}}{rt \frac{2}{N+1}}$$

$$= \frac{8M_{N}(jbj+1)(1+2hrt \frac{2}{N+1})}{hrt \frac{2}{N+1}} e^{rt \frac{2}{N+1}}$$

To arrive at the last line we have used that, for x > 0 and r > 0,

$$2\sum_{x}^{Z} e^{rt^{2}} dt = 2 \quad \frac{e^{rx^{2}}}{2rx} \quad \frac{Z}{x} \frac{e^{rt^{2}}}{2rt^{2}} dt \quad < \frac{e^{rx^{2}}}{rx}:$$
(4.106)

**Remark 4.3.2.** We can show in a similar way, for x = (k + 1 = 2)h with  $jt_k z_1 j$  h=4 and k = N + 1; N + 2; ..., that

$$jF(t_k)j = \frac{8}{h}(jbj+1) + \frac{1}{t_{N+1}}$$
: (4.107)

Also, since  $h_{+1} = t_{N+1} + h_{=2}$ , it holds that

$$1 + \frac{1}{t_{N+1}} \qquad 1 + \frac{1}{t_{N+1}}$$
;

and hence we have that

$$jT_N(h; 1=2)j = \frac{8(jbj+1)(1+2hrt_{N+1})}{hrt_{N+1}} + \frac{1}{t_{N+1}} e^{-rt_{N+1}^2} (4.108)$$

#### 4.3.3 Choices of the step-size h

This section is concerned with proposing explicit recommendations on how to choose the step-size h, following the recommendations in La Porte [38].

For r > 0, H := min(0:9; p=(r h)) and  $t_{N+1} = (N + 1)h$  with N 2 N, we define two possible choices,  $h_N$  and  $h_N$ , for the step-size h. For both we choose the step-size to satisfy the right hand equations in (4.109) and (4.111) below, i.e. to equalise the exponents in our

Lemma 4.3.3. If b > 0 and a given by (4.114), then

$$\frac{1}{1+3b}$$
 a  $\frac{1}{1+b}$ : (4.115)

Given r > 0 and N 2 N we choose h > 0 as follows.

**Remark 4.3.3.** Let  $H := \min(0:9; \mathbf{A}_N)$  with  $\mathbf{A}_N := \begin{pmatrix} q & \frac{1}{2p(N+1)=(p-\overline{3}r)} \\ p & \frac{1}{3} \end{pmatrix}$ , and set

$$\begin{array}{c} 8 & s & -p \\ 8 & h_{N} := & \frac{p}{3p} \\ \frac{1}{2r(N+1)}; & f_{N} & 0:9 \\ h := & & \\ h_{N} = a & \frac{pH}{r(N+1)^{2}} \\ \end{array}$$
 (4.116)

where a2 [1=(1+3b); 1=(1+b)] is given by

$$\mathbf{a} = {}^{\mathbf{s}} \frac{\mathbf{r}}{\frac{1}{2} + \frac{1}{4} + \mathbf{b}^{3}} + {}^{\mathbf{s}} \frac{\mathbf{r}}{\frac{1}{2} + \frac{1}{4} + \mathbf{b}^{3}}; \qquad (4.117)$$

and

$$b = \frac{r^{2=3}H^{4=3}}{12p^{2=3}(N+1)^{2=3}}$$
(4.118)

The following result bounds the expression

$$\frac{1+2hr t_{N+1}}{hr t_{N+1}}$$

for the choice of h given in Remark 4.3.3 which will be used to simplify further the bound (4.105) in Proposition 4.3.2.

Lemma 4.3.4. Let r > 0, N 2 N and  $\mathbf{A}_{N} := \begin{pmatrix} q \\ 2p(N+1)=(p \\ \overline{3}r) \end{pmatrix}$ , and h be given as in Remark 4.3.3. Then,  $f \phi_{N+1} = (N+1)h$ ,

$$\frac{1+2hrt_{N+1}}{rht_{N+1}} \stackrel{\substack{8\\ \gtrless}}{\underset{(N+1)}{5}} \frac{5}{2}; \quad \text{if } \mathcal{R}_{N} \quad 0:9; \quad (4.119)$$

where

$$K_{N} := \frac{2}{(N+1)^{1=3}} + \frac{2}{r^{1=3}p^{2=3}H^{2=3}} + \frac{rH^{2}}{8p^{2}(N+1)^{4=3}}:$$
(4.120)  
Proof. For  $h = h_{N} = \frac{r}{2r(N+1)}$ , we have  

$$\frac{1+2rht_{N+1}}{rht_{N+1}} = 2 + \frac{1}{r(N+1)(h_{N})^{2}} = 2 + p\frac{2}{3p}$$
(2.5):  
For  $h = h_{N} = a \frac{pH}{r(N+1)^{2}} + \frac{1}{1} = (2 + 1)^{-1}$ 

and hence the rst bound follows.

Now we consider the case = 0:9 and h =  $h_N = a \frac{pH}{r (N+1)^2}$ . Using(

 $P_{b;N}$  given by (4.36) in comparison with the approximations (4.25) and (4.28). Systematic numerical calculations are implemented for  $q_0 = 0^{\circ}(10^{\circ})90^{\circ}$ , jbj = 0:1(0:1)0:999 and  $arg(b) = 89^{\circ}(8:9^{\circ})89^{\circ}$ , and the Faddeeva function in  $P_{n;m}^{(2)}$  given by (4.24) is computed by Wiedeman's approximation (3.8), implemented by the call cef(z,40) in Table 1 [62].

For convenience, we denote in this section the approximation (4.28) in La Porte [38] by  $P_N^{(1)}$  and our approximation  $P_{b;N}$  given by (4.36) by  $P_N^{(2)}$ . We do not have access to exact values for  $P_b$  and so using different accurate approximations to  $P_b$ :

- (i) Our approximation P<sub>b;N</sub> given by (4.36) with N = 100, computed by the Matlab code in Listing A.4;
- (ii) Chandler-Wilde and Hothersall's approximation P<sub>100;100</sub> given by (4.25) computed by a Matlab

(ii) With N = 11, our approximation  $P_b$ 

 $r = kd^{\circ}$ 

$r = kd^{0}$	d <sup>o</sup>	Ea	ipprox	$E_{9}^{(4)}$		E <sup>(4)</sup>		E <sup>(4)</sup>	
0:5	0:0796	5:8	10 4	1:7	10 <sup>5</sup>	3:0	10 <sup>6</sup>	9:8	10 <sup>9</sup>
0:75	0:119	8:1	10 <sup>5</sup>	2:7	10 <sup>6</sup>	7:1	10 <sup>7</sup>	2:5	10 <sup>9</sup>
1:125	0:179	7:1	10 6	1:3	10 <sup>6</sup>	2:8	10 7	4:9	10 <sup>10</sup>
1:688	0:269	3:5	10 7	5:3	10 <sup>7</sup>	9:4	10 <sup>8</sup>	7:6	10 <sup>11</sup>
2:531	0:403	8:3	10 <sup>9</sup>	1:8	10 <sup>7</sup>	2:7	10 <sup>8</sup>	8:6	10 <sup>12</sup>
3:793	0:604	8:4	10 <sup>11</sup>	5:2	10 <sup>8</sup>	6:1	10 <sup>9</sup>	7:1	10 <sup>13</sup>
5:70	0:906	7:0	10 <sup>13</sup>	1:3	10 <sup>8</sup>	1:1	10 <sup>9</sup>	4:0	10 <sup>14</sup>
8:54	1:36	4:0	10 <sup>13</sup>	2:5	10 <sup>9</sup>	1:7	10 <sup>10</sup>	6:7	10 <sup>15</sup>
12:814	2:039	4:0	10 <sup>13</sup>	5:2	10 <sup>10</sup>	2:2	10 <sup>11</sup>	1:1	10 <sup>14</sup>
19:222	3:059	3:9	10 <sup>13</sup>	9:7	10 <sup>11</sup>	2:7	10 <sup>12</sup>	3:2	10 <sup>15</sup>
28:833	4:589	3:9	10 <sup>13</sup>	1:9	10 <sup>11</sup>	3:1	10 <sup>13</sup>	5:0	10 <sup>15</sup>
43:249	6:883	3:9	10 <sup>13</sup>	4:3	10 <sup>12</sup>	3:7	10 <sup>14</sup>	3:7	10 <sup>15</sup>
64:873	10:325	3:9	10 <sup>13</sup>	1:7	10 <sup>11</sup>	4:1	10 <sup>14</sup>	6:0	10 <sup>15</sup>
97:31	15:487	3:9	10 <sup>13</sup>	1:5	10 <sup>11</sup>	6:8	10 <sup>14</sup>	7:8	10 <sup>15</sup>
145:96	23:230	3:9	10 <sup>13</sup>	7:4	10 <sup>12</sup>	5:8	10 <sup>14</sup>	6:1	10 <sup>15</sup>
218:95	34:847	3:9	10 <sup>13</sup>	1:0	10 <sup>11</sup>	5:0	10 <sup>14</sup>	6:2	10 <sup>15</sup>
328:42	51:633	3:9	10 <sup>13</sup>	4:1	10 <sup>12</sup>	2:2	10 <sup>14</sup>	1:3	10 <sup>14</sup>
492:63	78:404	3:9	10 <sup>13</sup>	3:0	10 <sup>12</sup>	1:8	10 <sup>14</sup>	2:9	10 <sup>15</sup>
738:95	117:608	3:9	10 <sup>13</sup>	2:9	10 <sup>12</sup>	1:2	10 <sup>14</sup>	7:6	10 <sup>15</sup>
1108:4	176:407	3:9	10 <sup>13</sup>	2:0	10 <sup>12</sup>	9:9	10 <sup>15</sup>	4:9	10 <sup>15</sup>

 Table 4.3 Maximum values of  $E_{approx}$  and  $E_N^{(4)}$  given by (4.129) and (4.133), respectively, with N = 9;11;21, for  $q_0 = 0^{\circ}(10^{\circ})90^{\circ}$ , jbj = 0:1(0:1)0:999 and arg(b) =  $89^{\circ}(8:9^{\circ})89^{\circ}$ .



**Fig. 4.2** Accuracy of our approximation (4.36) and its upper bound (4.43), as a function of N, in comparison with La Porte's approximation (4.28).



Fig. 4.3 Accuracy of our approximation (4.36), as a function of r, in comparison with La Porte's approximation (4.28).
# Chapter 5

In Chapter 4, building on the works of Chandler-Wilde and Hothersall [14] and La Porte [38], we extended and improved the approximation of La Porte [38] by proposing a more stable (in floating point arithmetic) approximation of the 2D impedance half-space Green's function of the Helmholtz equation. We proved a uniform bound on the absolute error of this approximation and we showed, using systematic numerical calculations, that our approximation is more accurate and more efficient than the approximation of Chandler-Wilde and Hothersall [14].

We have achieved our objectives in this thesis and we hope that the presented approximations will be of great benefit for the wide range of applications of these three special functions.

#### 5.2 Further work

It was shown in this thesis that the truncated modified trapezium rule given by (1.23) is an accurate and efficient method to approximate three special functions which can be written as integrals of the form  $z_{\rm rule}$ 

$$I := \sum_{i=1}^{2} e^{rt^{2}} F(t); dt; \text{ for } r > 0;$$
 (5.1)

where **F** is an even meromorphic function with simple poles in a strip surrounding the real line. It is of interest to investigate further to what extent the methods of this thesis are applicable to other special functions. In particular, we summarize below suggested extensions to the work of this thesis, motivated by our theoretical and numerical results, as follows:

(i) The Voigt function, denoted by V(x; y), is defined as V(x; y) = Re(w(z)), and its derivatives satisfy that

$$\frac{\P V}{\P x} = 2 \operatorname{Re}$$

(iii) Additionally, it is interesting to investigate to what extent the methods of Chapter 4 are applicable to the 3D impedance half-space Green's function for the Helmholtz equation [14], to the 2D case of an in nite periodic array of point sources above an impedance plane [28], and the related important 2D case of an in nite periodic array of point sources in free space(1). In all three cases integral representations of the f(5). In are relevant with F meromorphic.

### References

- Abrarov, S. and Quine, B. M. (2015). Sampling by incomplete cosine expansion of the sinc function: Application to the Voigt/complex error function. Applied Mathematics and Computation 258:425–435.
- [2] Abromowitz, M. and Stegun, I. A. (1968). Handbook of Mathematical FunctionBover.
- [3] Alazah, M., Chandler-Wilde, S. N., and La Porte, S. (2014). Computing Fresnel integrals via modified trapezium rules. Numerische Mathematik 28(4):635–661.
- [4] Allasia, G. and Besenghi, R. (1986). Numerical calculation of incomplete gamma functions by the trapezoidal rule. Numerische Mathemațiso(4):419–428.
- [5] Bialecki, B. (1989). A modified sinc quadrature rule for functions with poles near the arc of integration. BIT Numerical Mathematics 9(3):464–476.
- [6] Bowman, J., Senior, T., and Uslenghi, P. (1969). Electromagnetic and acoustic scattering by simple shapesAmsterdam: North Holland.
- [7] Brambley, E. and Gabard, G. (2014). Reflection of an acoustic line source by an impedance surface with uniform flow. Journal of Sound and Vibration 333(21):5548– 5565.
- [8] Chandler-Wilde, S. N. (1988). Ground Effects in Environmental Sound PropagatiehD thesis, University of Bradford.
- [9] Chandler-Wilde, S. N. (2016). Private communication.
- [10] Chandler-Wilde, S. N., Hewett, D., Langdon, S., and Twigger, A. (2015). A high frequency boundary element method for scattering by a class of nonconvex obstacles. Numerische Mathematik29(4):647–689.
- [11] Chandler-Wilde, S. N. and Hothersall, D. (1985). Sound propagation above an inhomogeneous impedance plane. Journal of Sound and Vibration 8(4):475–491.
- [12] Chandler-Wilde, S. N. and Hothersall, D. (1988a). Integral equations in traffic noise simulation. In

- [14] Chandler-Wilde, S. N. and Hothersall, D. (1995). Efficient calculation of the Green function for acoustic propagation above a homogeneous impedance plane. Journal of Sound and Vibration180(5):705–724.
- [15] Chiarella, C. and Reichel, A. (1968). On the evaluation of integrals related to the error function. Mathematics of Computatio@2(101):137–143.
- [16] Cody, W. (1968). Chebyshev approximations for the Fresnel integrals. Mathematics of Computation, 22(102):450–453.
- [17] Conway, J. B. (1978). Functions of one complex variableSpringer.
- [18] Davis, P. J. and Rabinowitz, P. (2007). Methods of numerical integrationDover.
- [19] Durán, M., Hein, R., and Nédélec, J.-C. (2007). Computing numerically the Green's function of the half-plane Helmholtz operator with impedance boundary conditions. Numerische Mathematik07(2):295–314.
- [20] Fettis, H. E. (1955). Numerical calculation of certain definite integrals by Poisson's summation formula. Mathematical Tables and Other Aids to Computation 85–92.
- [21] Filippi, P. (1983). Extended sources radiation and Laplace type integral representation: Application to wave propagation above and within layered media. Journal of Sound and Vibration, 91(1):65–84.
- [22] Gautschi, W. (1970). Efficient computation of the complex error function. SIAM Journal on Numerical Analysis (1):187–198.
- [23] Gil, A., Segura, J., and Temme, N. M. (2002). Computing complex Airy functions by numerical quadrature. Numerical Algorithms30(1):11–23.
- [24] Goodwin, E. (1949). The evaluation of integrals of the form  $\stackrel{R_{+}}{\underset{k}{\longrightarrow}} f(x)e^{-x^2}dx$ . Mathematical Proceedings of the Cambridge Philosophical Society(02):241–245.
- [25] Grubeša, S., Jambrošić, K., and Domitrović, H. (2012). Noise barriers with varying cross-section optimized by genetic algorithms. Applied Acoustics73(11):1129–1137.
- [26] Habault, D. (1985). Sound propagation above an inhomogeneous plane: boundary integral equation methods. Journal of Sound and Vibratior100(1):55–67.
- [27] Heald, M. A. (1985). Rational approximations for the Fresnel integrals. Mathematics of Computation,44(170):459–461.
- [28] Horoshenkov, K. V. and Chandler-Wilde, S. N. (2002). Efficient calculation of twodimensional periodic and waveguide acoustic Green's functions. The Journal of the Acoustical Society of America11(4):1610–1622.
- [29] Hunter, D. (1964). The calculation of certain Bessel functions. Mathematics of Computation, 18(85):123–128.
- [30] Hunter, D. (1968). The evaluation of a class of functions defined by an integral. Mathematics of Computatio@2(102):440–444.

- [31] Hunter, D. and Regan, T. (1972). A note on the evaluation of the complementary error function. Mathematics of Computatio 26(118):539–541.
- [32] Jean, P. and Gabillet, Y. (2000). Using a boundary element approach to study small screens close to rails. Journal of sound and vibration 231(3):673–679.
- [33] Jiménez-Mier, J. (2001). An approximation to the plasma dispersion function. Journal of Quantitative Spectroscopy and Radiative Trans**76(**3):273–284.
- [34] Kang, J. (2002). Numerical modelling of the sound fields in urban streets with diffusely reflecting boundaries. Journal of Sound and Vibration 258(5):793–813.
- [35] Kawai, T., Hidaka, T., and Nakajima, T. (1982). Sound propagation above an impedance boundary. Journal of Sound and Vibratio®3(1):125–138.
- [36] Kress, R. (1998). Numerical analysisSpringer-Verlag, New York.
- [37] Krommer, A. R. and Ueberhuber, C. W. (1998). Computational integrationSIAM.
- [38] La Porte, S. (2007). Modi ed Trapezium Rule Methods for the Ef cient Evaluation of Green's Functions in Acoustic PhD thesis, Brunel University.
- [39] Li, J., Sun, G., and Zhang, R. (2016). The numerical solution of scattering by infinite rough interfaces based on the integral equation method. Computers & Mathematics with Applications 71(7):1491–1502.
- [40] Linton, C. (1998). The Green's function for the two-dimensional Helmholtz equation in periodic domains. Journal of Engineering Mathematics 3(4):377–401.
- [41] Liu, S. and Li, K. M. (2012). Efficient computation of the sound fields above a layered porous ground. The Journal of the Acoustical Society of America1(6):4389–4398.
- [42] Luke, Y. L. (1969). Special functions and their approximations lume 2. Academic press.
- [43] Matta, F. and Reichel, A. (1971). Uniform computation of the error function and other related functions. Mathematics of Computatiopages 339–344.
- [44] McNamee, J. (1964). Error-bounds for the evaluation of integrals by the Euler-Maclaurin formula and by Gauss-type formulae. Mathematics of Computation 8(87):368– 381.
- [45] Mori, M. (1983). A method for evaluation of the error function of real and complex variable with high relative accuracy. Publications of the Research Institute for Mathematical Sciences19(3):1081–1094.
- [46] Olver, F. W. (2010). NIST Handbook of Mathematical Functions Hardback and CD-ROM Cambridge University Press.
- [47] O'Neil, M., Greengard, L., and Pataki, A. (2014). On the efficient representation of the half-space impedance Green's function for the Helmholtz equation. Wave Motion, 51(1):1–13.

# Appendix A

# Matlab codes

#### A.1 Matlab codes to compute Fresnel integrals

```
Listing A.1 Matlab code to evaluate F_N(x) given by (2.12)
1 function f = fresnel(x, N)
2 select = x > = 0;
3 f = zeros(size(x));
4 if any(select), f(select) = F(x(select),N); end
5 if any(~select), f(~select) = 1 F( x(~select),N); end
6 function f = F(x, N)
7 h = sqrt(pi/(N+0.5));
8 t = h * ((N: 1:1) 0.5); AN = pi/h;
9 t_2 = t \cdot t_1; t_4 = t_2 \cdot t_2; et_2 = exp(t_2);
10 rooti = exp(i * pi/4);
11 z = rooti *x; x2 = x. *x; x4 = x2. *x2; z2 = i *x2;
12 S = (et2(1)./(x4+t4(1))).*(z2+t2(1));
13 for n = 2:N
       S = S + (et2(n)) / (x4+t4(n))) . * (z2+t2(n));
14
15 end
  ez = exp((2 * AN * i * rooti) * x);
16
17
```

```
106
```

**Listing A.2** Matlab code to evaluate  $C_N(x)$  and  $S_N(x)$  given by (2.14) and (2.15)

```
1 function [C,S] = fresnelCS(x,N)
2 h = sqrt(pi/(N+0.5));
3 t = h*((N: 1:1) 0.5); AN = pi/h; rootpi = sqrt(pi);
4 t2 = t.
```

### A.2 Matlab code to compute Faddeeva function

```
34 end
35 function f = w3(z,N)
36 z2 = z.*z; az = (2i/AN)*z;
37 a = h*((N: 1:1)+0.5); a2 = a.^2; et2 = exp( a2);
38 S1 = et2(1)./(z2 a2(1));
39 for n = 2 : N
40 S1 = S1 + et2(n)./(z2 a2(n));
41 end
42 h0 = 0.5*h;
43 S0 = exp( h0.^2)./(z2 h0.^2);
44 f = az.*(S0 + S1);
45 end
46 end
```

### A.3 Matlab code to compute $P_{\!\! b}$

```
65 if V1 == V2
   V = 1;
66
67 else
      V = 1;
68
69 end
70 Cm = 2*V*exp(1i*rho.*am).*heaviside(Himag(z2))./(1+exp
     ( 2*1i*pi*z2./h));
71 TC = pi * (Cp + Cm) . / (2 * sqrt(1 beta .^2));
72 t = h \cdot ((N: 1:1) + 0.5); t2 = t \cdot 2; h0 = 0.5 \cdot h; et2 =
                                                                exp
      ( t2.*rho);
73 s1 = beta + gamma.*(1+1i*t2); s2 = sqrt(t2 2*1i);
74 \ s3 = t2 \ 1i * ap; \ s4 = t2 \ 1i * am;
75 S1 = et2(1).*s1(1)./(s2(1).*s3(1).*s4(1));
76 for n = 2 : N
       S1 = S1 + et2(n) \cdot s1(n) \cdot (s2(n) \cdot s3(n) \cdot s4(n));
77
78 end
79 A =
        (beta + gamma.*(1+1i*h0.^2)).*exp(rho.*h0.^2);
80 B = sqrt (h0.^2 2*1i).*(h0^2 z1.^2).*(h0^2
                                                      z2.^2);
81 I = (beta.*exp(1i*rho)/pi).*2*h.*(S1 + A./B);
82 f2 = I + (beta.*exp(1i*rho)/pi).*TC;
83 end
84 end
```