



HAL
open science

Sharp error bounds for complex floating-point inversion

Claude-Pierre Jeannerod, Nicolas Louvet, Jean-Michel Muller, Antoine Plet

► **To cite this version:**

Claude-Pierre Jeannerod, Nicolas Louvet, Jean-Michel Muller, Antoine Plet. Sharp error bounds for complex floating-point inversion. *Numerical Algorithms*, 2016, 73 (3), pp.735-760. 10.1007/s11075-016-0115-x . ensl-01195625v2

HAL Id: ensl-01195625

<https://ens-lyon.hal.science/ensl-01195625v2>

Submitted on 19 Feb 2016

HAL is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.

Sharp error bounds for complex floating-point inversion

Claude-Pierre Jeannerod ·
Nicolas Louvet · Jean-Michel Muller ·
Antoine Plet

the date of receipt and acceptance should be inserted later

Abstract We study the accuracy of the classic algorithm for inverting a complex number given by its real and imaginary parts as floating-point numbers. Our analyses are done in binary floating-point arithmetic, with an unbounded exponent range and in precision p ; we also assume that the basic arithmetic operations $(+, -, \times, /)$ are rounded to nearest, so that the roundoff unit is $u = 2^{-p}$. We bound the largest relative error in the computed inverse either in the componentwise or in the normwise sense. We prove the componentwise relative error bound $3u$ for the complex inversion algorithm (assuming $p \geq 4$), and we show that this bound is asymptotically optimal (as $p \rightarrow \infty$) when p is even, and sharp when using one of the basic IEEE 754 binary formats with an odd precision ($p = 53, 113$). This componentwise bound obviously leads to the same bound $3u$ for the normwise relative error. However, we prove that the smaller bound $2.707131u$ holds (assuming $p \geq 24$) for the normwise relative error, and we illustrate the sharpness of this bound for the basic IEEE 754 binary formats ($p = 24, 53, 113$) using numerical examples.

Keywords floating-point arithmetic · rounding error analysis · complex inversion

Claude-Pierre Jeannerod
Inria, Laboratoire LIP (CNRS, ENSL, Inria, UCBL), Université de Lyon, France.
E-mail: claude-pierre.jeannerod@inria.fr

Nicolas Louvet
UCBL, Laboratoire LIP (CNRS, ENSL, Inria, UCBL), Université de Lyon, France.
E-mail: nicolas.louvet@ens-lyon.fr

Jean-Michel Muller
CNRS, Laboratoire LIP (CNRS, ENSL, Inria, UCBL), Université de Lyon, France.
E-mail: jean-michel.muller@ens-lyon.fr

Antoine Plet
ENS de Lyon, Laboratoire LIP (CNRS, ENSL, Inria, UCBL), Université de Lyon, France.
E-mail: antoine.plet@ens-lyon.fr

1 Introduction

This paper deals with the accuracy of the inversion of a nonzero complex number given by its real and imaginary parts as floating-point numbers. We assume that the underlying floating-point arithmetic has radix 2 and precision $p \geq 2$, and we also assume an unbounded exponent range, which means that our results apply to practical floating-point calculations according to the IEEE 754 standard [6] as long as underflow and overflow do not occur.

Given a nonzero complex number $a + ib$, its inverse satisfies

$$z = R + iI, \quad R = \frac{a}{a^2 + b^2}, \quad I = -\frac{b}{a^2 + b^2}. \quad (1)$$

Assuming a and b are floating-point numbers and denoting by RN a round-to-nearest function, we focus in this paper on the approximation $\hat{z} = \hat{R} + i\hat{I}$ that can be computed classically in floating-point arithmetic according to

$$\hat{R} = \text{RN}\left(\frac{a}{\text{RN}(\text{RN}(a^2) + \text{RN}(b^2))}\right) \quad (2)$$

for the real part, and with a similar expression for the imaginary part \hat{I} . This scheme corresponds to Algorithm 1 below.

Algorithm 1 Inversion of a nonzero complex floating-point number $a + ib$.

```

 $s_a \leftarrow \text{RN}(a^2)$ 
 $s_b \leftarrow \text{RN}(b^2)$ 
 $s \leftarrow \text{RN}(s_a + s_b)$ 
 $\hat{R} \leftarrow \text{RN}(a/s)$ 
 $\hat{I} \leftarrow \text{RN}(-b/s)$ 
return  $\hat{R} + i\hat{I}$ 

```

We provide an accuracy analysis of Algorithm 1, for both the componentwise relative error $E_C = \max(|R - \hat{R}|/|R|, |I - \hat{I}|/|I|)$ and the normwise relative error $E_N = |z - \hat{z}|/|z|$. In each case, we bound the *largest* error value by a function $B(p)$ depending only on the precision p , and study the tightness of that bound. In this context, we typically distinguish between three levels of quality:

- If we can show that there exist some inputs $a + ib$ parametrized by p and for which the bound is attained for every $p \geq p_0$ (for a given $p_0 \geq 2$), then we say that the bound is *optimal*.
- If we can show that there exist some inputs parametrized by p and for which the relative error $E(p)$ satisfies $E(p)/B(p) \rightarrow 1$ as $p \rightarrow \infty$, then we say that the bound is *asymptotically optimal*.
- In some cases, we did not manage to establish (asymptotic) optimality, but have found input examples—either parametrized by p , or just for some values of p of practical interest (like those corresponding to the basic IEEE

754 formats)—for which $E(p)$ is very close to $B(p)$. In this case, we say that the bound is *sharp*. (See [13] for a similar use of the word “sharp”.)

The componentwise relative error generated by Algorithm 1 can easily be bounded as $E_C \leq 3u + O(u^2)$, where $u = 2^{-p}$ is the unit roundoff. Our first contribution is to show that the term $O(u^2)$ can in fact be removed, which leads to the simpler bound $E_C \leq 3u$ (assuming $p \geq 4$). Furthermore, when p is even, we show that this bound is asymptotically optimal by providing floating-point numbers a and b parametrized by p and for which $E_C \geq 3u - \frac{31}{2}u^{\frac{3}{2}} + O(u^2)$; when p is odd, we show that the bound $3u$ is sharp, especially for the corresponding basic IEEE 754 binary formats ($p = 53, 113$).

The normwise relative error bound $E_N \leq 3u + O(u^2)$ can be found in [4, p. 30], and a direct application of our componentwise error analysis leads further to $E_N \leq 3u$. The second main contribution of our paper is to show that if $p \geq 10$ then the following smaller bound holds: $E_N < \gamma u + 9u^2$, where γ is an explicit constant in (2.70712, 2.70713). When using for example the IEEE 754 binary32 format ($p = 24$), this implies $E_N < 2.707131u$. The techniques and the case distinction we use to prove this bound are inspired from [13], but we also extensively use real analysis and differentiation for the treatment of each case. We provide numerical examples to show that the bound we obtain is sharp for the basic IEEE 754 formats ($p = 24, 53, 113$).

Several authors [10, 11, 8, 2] have suggested ways of avoiding spurious overflows and underflows in complex division, and some of them may be used in Algorithm 1. Of course, if the computation introduces further rounding errors, which is the case for example in Smith’s method [10], then our error bounds may not hold anymore. However, following the technique developed by Priest in [8], it is possible to scale a and b by a power of two in order to avoid spurious overflows and underflows without introducing new rounding errors: in that case, our analyses are valid if neither overflow nor underflow occurs during the computation. Nonetheless, we do not deal with scaling techniques here, and focus only on the largest error assuming an unbounded exponent range.

Outline. Section 2 is devoted to the componentwise relative error analysis of Algorithm 1, and Section 3 to its normwise relative error analysis. We conclude in Section 4 with some remarks on the implications of these error analyses for complex floating-point division. The technical parts of the proofs that can be skipped at first reading are gathered in Appendix A.

Assumptions and notation. For any real number t , we denote by $\text{RN}(t)$ the binary floating-point number that is nearest to t , with a tie-breaking strategy preserving the following properties:

- $\text{RN}(2^k t) = 2^k \text{RN}(t)$ for any integer k ;
- $\text{RN}(-t) = -\text{RN}(t)$.

In particular, either the *roundTiesToEven* or the *roundTiesToAway* rounding direction attribute defined in the IEEE 754 standard [6] can be used.

Throughout the paper, we also rely on the following relative error bound [7, p. 232]: for any real number t ,

$$\text{RN}(t) = t(1 + \epsilon) \quad \text{with} \quad |\epsilon| \leq \frac{u}{1+u}. \quad (3)$$

Note that (3) implies the well-known inequality $|\text{RN}(t) - t| \leq u|t|$; see [5, p. 38].

Finally, we use the notation $\text{ufp}(t)$ (*unit in the first place*, introduced in [9]) to denote the weight of the most significant bit of t : if $t \neq 0$ then $\text{ufp}(t)$ is the unique integer power of two such that $1 \leq |t|/\text{ufp}(t) < 2$, and $\text{ufp}(0) = 0$. The usual ulp function (*unit in the last place*) is related to the ufp function through the relation $\text{ulp}(t) = 2u \cdot \text{ufp}(t)$, so that

$$|t - \text{RN}(t)| \leq \frac{1}{2}\text{ulp}(t) = \text{ufp}(t)u. \quad (4)$$

2 Componentwise error bound

In this section, we focus on the componentwise relative error of Algorithm 1. We note first that since $a+ib$ is nonzero, $R = a/(a^2+b^2)$ and $I = -b/(a^2+b^2)$ cannot both be zero, and that if one of them is zero then the returned result is very accurate. Assume for example that $R = 0$ (the case $I = 0$ is similar). In that case, $a = 0$ and $I = -1/b$. Using the bound in (3), it is then easily checked that the values \widehat{R} and \widehat{I} returned by the algorithm are as follows:

- $\widehat{R} = 0$, which means that the real part is computed exactly;
- $\widehat{I} = -\text{RN}(b/\text{RN}(b^2))$ and the relative error on the imaginary part is bounded by $2u$ (and thus smaller than the bound we are going to give in the general case).

Therefore, the rest of this section is devoted to analyzing $E_C = \max(|R - \widehat{R}|/|R|, |I - \widehat{I}|/|I|)$ for R and I nonzero. Repeated applications of the bound in (3) give immediately $E_C \leq 3u + \mathcal{O}(u^2)$. We show below that if $p \geq 4$ then the $\mathcal{O}(u^2)$ term can in fact be removed, leading to the simpler bound $3u$.

To do this, we prove that if $p \neq 3$ then the relative error bound $u/(1+u)$ in (3) can be replaced by $u/(1+3u)$ when evaluating a square $\text{RN}(a^2)$ instead of a general product. (When $p = 3$, it is easily checked that the bound $u/(1+u)$ is attained when squaring the floating-point numbers $3/2 \cdot 2^e$, $e \in \mathbb{Z}$.) This slight refinement will turn out to be enough to show that Algorithm 1 satisfies $E_C \leq 3u$.

Lemma 1 *Let a be a floating-point number. If $p \neq 3$ then $|a^2 - (2+2u)| \geq 4u^2$.*

Proof If $|a| < 1$ then $|a^2 - (2+2u)| > 1+2u$, and the result follows from the fact that $1+2u > 4u^2$ when $p > 0$. Assume now that $|a| \geq 1$. To handle this case, we show first that

$$a^2 = 2 + 2u \quad \Rightarrow \quad p = 3. \quad (5)$$

Since $|a|$ is a floating-point number not smaller than 1, there exists a positive integer A such that $|a| = A \cdot 2^{1-p} = A \cdot 2u$. The equality $a^2 = 2 + 2u$ is thus equivalent to $A^2 = (2^p + 1) \cdot 2^{p-1}$ and, using the (unique) decomposition $A = (2B + 1) \cdot 2^C$ with $B, C \in \mathbb{N}$, it can also be rewritten $(2B + 1)^2 \cdot 2^{2C} = (2^p + 1) \cdot 2^{p-1}$. Now, $p > 0$ implies that $2^p + 1$ is odd and at least 3, so $B \neq 0$ and $(2B + 1)^2 = 2^p + 1$. The latter equality can be rewritten as $4B(B + 1) = 2^p$ and its unique solution over $\mathbb{N}_{>0}^2$ is $(B, p) = (1, 3)$, so (5) follows.

If $p \neq 3$ then, by (5) we have $a^2 \neq 2 + 2u$, that is, $A^2 \neq (2^p + 1) \cdot 2^{p-1}$. Since the latter inequality involves only integers, it is equivalent to $|A^2 - (2^p + 1) \cdot 2^{p-1}| \geq 1$ and thus to $|a^2 - (2 + 2u)| \geq 4u^2$.

Lemma 2 *Let a be a floating-point number. If $p \neq 3$ then $\text{RN}(a^2) = a^2(1 + \epsilon)$ with $|\epsilon| \leq u/(1 + 3u)$.*

Proof We can assume that $1 \leq a < 2$. If $a = 1$ then $\text{RN}(a^2) = a^2$ and the result is clear. If $1 < a < \sqrt{2}$ then it follows from a being a floating-point number that $p \geq 4$ and that a belongs to the non-empty interval $[1 + 2u, \sqrt{2})$. Consequently, $1 + 4u < a^2 < 2$ and thus $|\epsilon| \leq u \text{ufp}(a^2)/a^2 = u/a^2 < u/(1 + 4u)$. Finally, if $\sqrt{2} < a < 2$ then $2 < a^2 < 4$ and, by Lemma 1, it suffices to consider the following four subcases:

- If $2 < a^2 \leq 2 + 2u - 4u^2$ then $\text{RN}(a^2) = 2$ and, therefore,

$$|\epsilon| = 1 - \frac{2}{a^2} \leq 1 - \frac{2}{2 + 2u - 4u^2} \leq \frac{u}{1 + 3u}.$$

- If $2 + 2u + 4u^2 \leq a^2 < 2 + 4u$ then $\text{RN}(a^2) = 2 + 4u$ and, therefore,

$$|\epsilon| = \frac{2 + 4u}{a^2} - 1 \leq \frac{2 + 4u}{2 + 2u + 4u^2} - 1 \leq \frac{u}{1 + 3u}.$$

- If $2 + 4u \leq a^2 < 2 + 6u$ then $\text{RN}(a^2) = 2 + 4u$ and, therefore,

$$|\epsilon| = 1 - \frac{2 + 4u}{a^2} \leq 1 - \frac{2 + 4u}{2 + 6u} = \frac{u}{1 + 3u}.$$

- If $2 + 6u \leq a^2 < 4$ then $\text{ufp}(a^2) = 2$ and $|\epsilon| \leq 2u/a^2 \leq 2u/(2 + 6u) = u/(1 + 3u)$.

Theorem 1 *If $p \geq 4$ then the componentwise relative error for Algorithm 1 satisfies $E_C \leq 3u$.*

Proof Due to the symmetry of Algorithm 1, it suffices to show that $|R - \widehat{R}| \leq 3u|R|$. From (3) and Lemma 2 we have

$$s_a = a^2(1 + \epsilon_a), \quad s_b = b^2(1 + \epsilon_b), \quad s = (s_a + s_b)(1 + \epsilon_s), \quad \widehat{R} = \frac{a}{s}(1 + \epsilon_R)$$

with $|\epsilon_a|, |\epsilon_b| \leq u/(1 + 3u)$ and $|\epsilon_s|, |\epsilon_R| \leq u/(1 + u)$. Hence

$$\widehat{R} = \frac{a}{a^2(1 + \epsilon_a) + b^2(1 + \epsilon_b)} \cdot \frac{1 + \epsilon_R}{1 + \epsilon_s}$$

and, using $R = a/(a^2 + b^2)$, we deduce that $\varphi R \leq \widehat{R} \leq \varphi' R$ with

$$\varphi := \frac{1 - \frac{u}{1+u}}{\left(1 + \frac{u}{1+3u}\right)\left(1 + \frac{u}{1+u}\right)} \quad \text{and} \quad \varphi' := \frac{1 + \frac{u}{1+u}}{\left(1 - \frac{u}{1+3u}\right)\left(1 - \frac{u}{1+u}\right)}.$$

It is easily checked that $\varphi > 1 - 3u$ and $\varphi' = 1 + 3u$, which completes the proof.

We conclude this section by showing that the componentwise bound $E_C \leq 3u$ is sharp. More precisely, when the precision p is even, the following example shows that the componentwise error bound $3u$ is asymptotically optimal as $p \rightarrow \infty$. Assuming an even $p \geq 12$, let us consider the following binary floating-point numbers in precision p :

$$\begin{aligned} a &= 2^{\frac{p}{2}-1} + 5 \cdot 2^{-2} + 2^{-\frac{p}{2}+2}, \\ b &= 2^{p-1} + 2^{\frac{p}{2}-1} + 1. \end{aligned}$$

With these values as inputs of Algorithm 1, we have (the details are provided in Appendix A.1):

$$\begin{aligned} s_a &= 2^{p-2} + 5 \cdot 2^{\frac{p}{2}-2} + 11 \cdot 2^{-1}, \\ s_b &= 2^{2p-2} + 2^{\frac{3p}{2}-1} + 3 \cdot 2^{p-1}, \\ s &= 2^{2p-2} + 2^{\frac{3p}{2}-1} + 2^{p+1}. \end{aligned}$$

From this we deduce

$$\frac{a}{s} = 2^{-\frac{3p}{2}+1} + 2^{-2p} - 2^{-\frac{5p}{2}+1} - 2^{-3p+2} + \mathcal{O}\left(2^{-\frac{7p}{2}}\right),$$

and $\text{ulp}\left(\frac{a}{s}\right) = 2^{-\frac{5p}{2}+2}$. Then, defining the floating-point number τ by

$$\tau = 2^{-\frac{3p}{2}+1} + 2^{-2p} - 2^{-\frac{5p}{2}+2},$$

it can be checked that

$$\left| \frac{a}{s} - \tau \right| = \frac{2^{-\frac{5p}{2}+1} + 2^{-\frac{7p}{2}+5}}{1 + 2^{-\frac{p}{2}+1} + 2^{-p+3}} < \frac{1}{2} \text{ulp}\left(\frac{a}{s}\right).$$

Hence $\widehat{R} = \text{RN}\left(\frac{a}{s}\right) = \tau$, which together with $R = a/(a^2 + b^2)$ leads to

$$\frac{R - \widehat{R}}{R} = 3u - \frac{31}{2}u^{\frac{3}{2}} + \mathcal{O}(u^2).$$

As a consequence, in this example the componentwise relative error in the computed \widehat{z} is at least $3u - \frac{31}{2}u^{\frac{3}{2}} + \mathcal{O}(u^2)$, which shows the asymptotic optimality (as $p \rightarrow \infty$) of the bound when p is even.

When p is odd, we have not found an input set parametrized by the precision to prove the asymptotic optimality of the error bound $3u$. However, we illustrate the sharpness of the bound by numerical examples in Table 1.

p	Inputs a and b	E_C/u
15	$a = 16732$ $b = 23252 \cdot 2^3$	2.93047...
17	$a = 66078$ $b = 93014 \cdot 2^8$	2.96359...
19	$a = 131435$ $b = 370969 \cdot 2^8$	2.98509...
53	$a = 4508053433127332$ $b = 6369149602646415 \cdot 2^{16}$	2.97894...
113	$a = 5192393427440123027423416459819356$ $b = 7343016638055329519853569740503421 \cdot 2^{16}$	2.97647...

Table 1 Examples with p odd and a componentwise relative error close to $3u$.

3 Normwise error bound

In this section, we are interested in the normwise relative error of Algorithm 1, that is,

$$E_N = \sqrt{a^2 + b^2} \sqrt{(R - \widehat{R})^2 + (I - \widehat{I})^2}.$$

The analysis is done in radix 2 and precision p , and we assume that overflows and underflows never occur. If we apply directly the componentwise bound obtained in Section 2, we end up with the normwise error bound $E_N \leq 3u$. In this section, we establish the following result, which achieves a smaller bound by keeping track of the correlations between the various rounding errors committed by the algorithm.

Theorem 2 *If $p \geq 10$ then the normwise relative error for Algorithm 1 satisfies $E_N \leq \gamma u + 9u^2$, where γ is defined by*

$$\gamma = \frac{\sqrt{8778980525057 + 16793600(8\sqrt{2} - \sqrt{127}) - 550842155008\sqrt{254}}}{8192(16 - \sqrt{254})}, \quad (6)$$

and is such that $\gamma \in (2.70712, 2.70713)$.

If $p \geq 10$, $E_N < 2.70713u + 9u^2$ is therefore a rigorous bound for the normwise error of Algorithm 1. It should also be noticed that the second order term in the error bound can be absorbed by the first order term, at the cost of a slight overestimation: for example, for $p \geq 24$, we have $9u = 9 \cdot 2^{-24} < 10^{-6}$ so that $E_N < 2.707131u$. The numerical examples listed in Table 2 show that the error bound of Theorem 2 is sharp for the basic IEEE 754 formats ($p = 24, 53, 113$).

3.1 Preliminaries

The first step in the error analysis of Algorithm 1 is to reduce the input domain. Since the function RN is symmetric with respect to zero, the signs of

p	Inputs a and b	E_N/u
24	$a = 11863283$ $b = 11865457 \cdot 2^{12}$	2.69090...
53	$a = 4503599709991314$ $b = 6369051770002436 \cdot 2^{26}$	2.70679...
113	$a = 2^{112}$ $b = 7343016637207171132572330391109909 \cdot 2^{56}$	2.70559...

Table 2 Examples with a normwise relative error close to γu .

a and b are not relevant and we can assume that both a and b are nonnegative. Swapping the inputs a and b does not affect the relative error; moreover, if $a = 0$, then a simple analysis, based on (3), leads to the upper bound $2u$ for E_N , so we can assume that $0 < a \leq b$. Finally, multiplying or dividing by two both a and b does not affect either the relative error, and we can restrict the analysis to the case $1 \leq b < 2$.

From the definition of the ufp function and this input range reduction, we know that $\text{ufp}(b^2) \in \{1, 2\}$. Moreover, b is a floating-point number, so that $1 \leq b \leq 2 - 2u$ and thus $1 \leq b^2 \leq 4 - 4u$. Since $4 - 4u$ is a floating-point number, and using the monotonicity of the rounding function RN, we deduce that $1 \leq s_b < 4$. Using again the monotonicity of RN, we also deduce that $0 < s_a < 4$. Hence $1 < s_a + s_b < 8$, which implies $\text{ufp}(s_a + s_b) \in \{1, 2, 4\}$.

We now define δ_a , δ_b , δ_s , δ_R , and δ_I as follows:

$$\begin{aligned} s_a &= a^2 + \delta_a u, & |\delta_a| &\leq \text{ufp}(a^2), \\ s_b &= b^2 + \delta_b u, & |\delta_b| &\leq \text{ufp}(b^2), \\ s &= s_a + s_b + \delta_s u, & |\delta_s| &\leq \text{ufp}(s_a + s_b), \\ \widehat{R} &= \frac{a}{s} + \delta_R u, & |\delta_R| &\leq \text{ufp}\left(\frac{a}{s}\right), \\ \widehat{I} &= -\left(\frac{b}{s} + \delta_I u\right), & |\delta_I| &\leq \text{ufp}\left(\frac{b}{s}\right). \end{aligned}$$

Let us also define $\delta = \delta_a + \delta_b + \delta_s$ and $\epsilon = \frac{|\delta|}{a^2 + b^2}$, so that $|\delta|u$ and ϵu are the absolute and relative errors, respectively, in the evaluation of $a^2 + b^2$. Since $0 < a \leq b$, $\text{ufp}(b^2) \leq 2$ and $\text{ufp}(s_a + s_b) \leq 4$, we deduce that $|\delta| \leq 8$. Moreover, it can be deduced from (3) that $\epsilon \leq 2$. (This bound on ϵ already appeared in [3, p. 1471].)

With these notations, we have

$$R - \widehat{R} = \frac{a}{s(a^2 + b^2)} \delta u - \delta_R u,$$

and since a similar expression holds for $I - \widehat{I}$, we arrive at

$$\frac{E_N^2}{u^2} = (a^2 + b^2) (\delta_R^2 + \delta_I^2) - 2 \frac{\delta(a\delta_R + b\delta_I)}{a^2 + b^2 + \delta u} + \left(\frac{\delta}{a^2 + b^2 + \delta u} \right)^2.$$

Then, using the triangular inequality, we obtain

$$\begin{aligned} \frac{E_N^2}{u^2} &\leq (a^2 + b^2) \left(\text{ufp}\left(\frac{a}{s}\right)^2 + \text{ufp}\left(\frac{b}{s}\right)^2 \right) \\ &\quad + 2 \frac{|\delta| \left(\text{ufp}\left(\frac{a}{s}\right)a + \text{ufp}\left(\frac{b}{s}\right)b \right)}{a^2 + b^2 - |\delta|u} + \left(\frac{\delta}{a^2 + b^2 - |\delta|u} \right)^2. \end{aligned}$$

For $p \geq 2$, $\epsilon u < 1$ and we use the equality $\frac{1}{a^2 + b^2 - |\delta|u} = \frac{1}{a^2 + b^2} \left(1 + \frac{\epsilon}{1 - \epsilon u} u \right)$ and the inequality $\left(1 + \frac{\epsilon}{1 - \epsilon u} u \right)^2 \leq 1 + \frac{2\epsilon}{(1 - \epsilon u)^2} u$ to get

$$E_N^2 \leq f_2(a, b)u^2 + f_3(a, b)u^3, \quad (7)$$

with

$$\begin{aligned} f_2(a, b) &= (a^2 + b^2) \left(\text{ufp}\left(\frac{a}{s}\right)^2 + \text{ufp}\left(\frac{b}{s}\right)^2 \right) \\ &\quad + 2 \frac{|\delta| \left(\text{ufp}\left(\frac{a}{s}\right)a + \text{ufp}\left(\frac{b}{s}\right)b \right)}{a^2 + b^2} + \left(\frac{\delta}{a^2 + b^2} \right)^2 \end{aligned} \quad (8)$$

and

$$f_3(a, b) = 2 \left(\text{ufp}\left(\frac{a}{s}\right)a + \text{ufp}\left(\frac{b}{s}\right)b \right) \frac{\epsilon^2}{1 - \epsilon u} + \frac{2\epsilon^3}{(1 - \epsilon u)^2}.$$

From (4), we have

$$\text{ufp}\left(\frac{a}{s}\right)a + \text{ufp}\left(\frac{b}{s}\right)b \leq \frac{a^2 + b^2}{s} \leq \frac{a^2 + b^2}{a^2 + b^2 - |\delta|u} = \frac{1}{1 - \epsilon u},$$

and since $0 \leq \epsilon \leq 2$, it follows that $f_3(a, b) \leq \frac{2\epsilon^2(1+\epsilon)}{(1-\epsilon u)^2} < 25$ for $p \geq 10$. Moreover, if f_2 is upper bounded by κ , we can conclude from (7) that

$$E_N \leq \sqrt{\kappa}u + \frac{25}{2\sqrt{\kappa}}u^2. \quad (9)$$

3.2 Taking care of some corner cases

We can first roughly bound f_2 using the inequality $\text{ufp}(t) \leq |t|$, valid for any real t , which will allow us to conclude in some particular cases and to further reduce the input domain. From (8) we have

$$\begin{aligned} f_2(a, b) &\leq \left(\frac{a^2 + b^2}{a^2 + b^2 - |\delta|u} \right)^2 + 2 \frac{|\delta| (a^2 + b^2)}{(a^2 + b^2)(a^2 + b^2 - |\delta|u)} + \left(\frac{\delta}{a^2 + b^2} \right)^2 \\ &= \left(1 + \epsilon + \frac{\epsilon}{1 - \epsilon u} u \right)^2. \end{aligned}$$

This last bound is increasing with respect to ϵ and u (*i.e.*, decreasing with respect to the precision p). Therefore, if $\epsilon \leq 1 + \frac{\sqrt{2}}{2} + u$, and as soon as $p \geq 5$, we have $f_2(a, b) \leq (2 + \frac{\sqrt{2}}{2} + 3u)^2$ and, from (9),

$$E_N \leq \left(2 + \frac{\sqrt{2}}{2}\right) u + 8u^2. \quad (10)$$

Below are five cases that lead to the inequality $\epsilon \leq 1 + \frac{\sqrt{2}}{2} + u$, so they can be ignored in the following analysis.

- If $a = b$, then $s_a = s_b$ and $s = s_a + s_b$ so that $\delta_s = 0$ and one can check that $\epsilon \leq 1$. In this case, the previous bound (10) holds and we can continue the analysis assuming that

$$a < b. \quad (11)$$

- If $b = 1$, then $s_b = b^2 = 1$ and $\delta_b = 0$. Moreover, from (11) we have $a < 1$, so that $s_a < 1$, which implies $\text{ufp}(1 + s_a) = 1$ and $\epsilon \leq 1$. Again, the bound (10) holds and we can continue the analysis assuming that $1 < b$. In fact, since b is a floating-point number, we can assume that

$$1 + 2u \leq b. \quad (12)$$

- If $a = 1$, then $\delta_a = 0$ and we can distinguish three cases. If $\text{ufp}(b^2) = 1$ then $\text{ufp}(1 + s_b) = 2$ and $\epsilon \leq \frac{3}{2}$. If $\text{ufp}(b^2) = 2$ then either $\text{ufp}(1 + s_b) = 2$ which implies $\epsilon \leq \frac{4}{3}$, or $\text{ufp}(1 + s_b) = 4$ and then $\epsilon \leq \frac{3}{2} + u$. In all these cases, (10) holds, hence we can assume now that

$$a \neq 1. \quad (13)$$

- If $a^2 + b^2 < \text{ufp}(s_a + s_b)$, then we have $(s_a + s_b) - \text{ufp}(s_a + s_b) < (\delta_a + \delta_b)u \leq (a^2 + b^2)u < \text{ufp}(s_a + s_b)u = \frac{1}{2}\text{ulp}(s_a + s_b)$. Since $\text{ufp}(s_a + s_b)$ is a floating-point number, we can deduce that $s = \text{RN}(s_a + s_b) = \text{ufp}(s_a + s_b)$ hence $\epsilon \leq 1$ and (10) holds. In the following, we can then assume that

$$\text{ufp}(s_a + s_b) \leq a^2 + b^2. \quad (14)$$

- One last case is when $s_a + s_b \geq \sqrt{2} \text{ufp}(s_a + s_b)$. In this case, $\epsilon \leq 1 + \frac{\sqrt{2}}{2} + u$ and the previous bound (10) holds. Therefore, we now assume that

$$s_a + s_b < \sqrt{2} \text{ufp}(s_a + s_b). \quad (15)$$

3.3 Overview of the case analysis

The analysis goes through the possible values of $\text{ufp}(s_a + s_b)$, which are 1, 2, and 4. In each case, we first deduce upper bounds for $\text{ufp}(b^2)$, $\text{ufp}(\frac{a}{s})$, and $\text{ufp}(\frac{b}{s})$. This leads to a new function g , which is greater than or equal to f_2 , and which depends on a and b as well as on a third parameter, e , defined as the unique integer such that

$$\text{ufp}(a^2) = 2^{-e}.$$

The function g does not involve floating-point operations anymore and can be seen as a continuous and differentiable function over real inputs. We then look for an upper bound on this function over a restricted domain D containing all the floating-point numbers we are interested in. For this latter step, we mainly use real analysis, especially partial derivatives. In some cases, we can maximize with respect to a and b at the same time. The last step is always to maximize with respect to e , using the change of variable $x = 2^{-e}$ and considering x as a continuous variable.

The analysis is split into seven cases depending on the values of some ufp functions involved in the definition (8) of f_2 . Note that, since $a^2 < 4$, we have $e \geq -1$. In each case but the last one, we end up with a bound smaller than or equal to $(2 + \frac{\sqrt{2}}{2})^2$ for f_2 , from which we conclude using (9) that $E_N \leq (2 + \frac{\sqrt{2}}{2})u + 5u^2$. The last case is similar although we have a slightly larger bound $\gamma^2 + 20u$ for f_2 (we have $2 + \frac{\sqrt{2}}{2} = 2.70710\dots$, while $\gamma = 2.70712\dots$), which leads to $E_N \leq \gamma u + 9u^2$. The table below summarizes the bounds in each case, under the assumptions (11) to (15).

$\text{ufp}(s_a + s_b)$	$\text{ufp}(b^2)$	e	$\text{ufp}(\frac{a}{s})$	f_2	E_N
1	1	≥ 2	$\leq 2^{-\frac{e}{2}}$	6.565	$2.6u$
4	2	$= -1$	$\leq \frac{1}{4}$	$(2 + \frac{\sqrt{2}}{2})^2$	$(2 + \frac{\sqrt{2}}{2})u + 5u^2$
		≥ 0	$\leq 2^{-2-\frac{e}{2}}$	$(\frac{7}{4} + \frac{\sqrt{2}}{2})^2$	$2.5u$
2	1	≥ 1	$\leq 2^{-1-\frac{e}{2}}$	$(\frac{7}{4} + \frac{\sqrt{3}}{2})^2$	$2.65u$
		$= 0$	$\leq \frac{1}{4}$	$(\frac{5}{2})^2$	$\frac{5}{2}u + 5u^2$
	2	≥ 1	$\leq 2^{-\frac{3+e}{2}}$	$(2 + \frac{\sqrt{2}}{2})^2$	$(2 + \frac{\sqrt{2}}{2})u + 5u^2$
		$\geq 2, \text{ even}$	$= 2^{-1-\frac{e}{2}}$	$\gamma^2 + 20u$	$\gamma u + 9u^2$

We give all the details of the analysis of the first case. For the other cases, we only give a sketch of the analysis, while deferring the details to Appendices A.2 to A.7.

3.4 Case $\text{ufp}(s_a + s_b) = 1$

In this case, we can deduce from (15) that $1 \leq s_a + s_b < \sqrt{2}$. As a consequence, we must have $b < \sqrt{2}$ (otherwise we would have $s_a + s_b > 2$), hence

$$\text{ufp}(b^2) = 1.$$

Since $s_a < \sqrt{2} - 1 < \frac{1}{2}$ and $s_a = \text{RN}(a^2)$, we have $a^2 < \frac{1}{2}$, and

$$e \geq 2.$$

Moreover, we know from (12) that $b \geq 1+2u$ so we have $b^2 \geq b(1+2u) \geq b+2u$, which is a floating-point number because $\text{ufp}(b) = 1$. Consequently $s_b \geq b+2u$ and $s \geq s_a + s_b - u \geq s_a + b + u > b$, hence $\frac{b}{s} < 1$, which implies

$$\text{ufp}\left(\frac{b}{s}\right) \leq \frac{1}{2}.$$

Finally, $s = \text{RN}(s_a + s_b) \geq 1$ so $\frac{a}{s} \leq a < 2^{\frac{1-e}{2}}$ and

$$\text{ufp}\left(\frac{a}{s}\right) \leq 2^{-\frac{e}{2}}.$$

Therefore, using (8) we deduce in this case that $f_2(a, b) \leq g_1(a, b, e)$, with

$$g_1(a, b, e) := (a^2 + b^2) \left(2^{-e} + \frac{1}{4}\right) + 2 \frac{(2 + 2^{-e}) \left(2^{-\frac{e}{2}} a + \frac{b}{2}\right)}{a^2 + b^2} + \left(\frac{2 + 2^{-e}}{a^2 + b^2}\right)^2.$$

Let us now characterize explicitly the domain over which we will bound $g_1(a, b, e)$. First, we know that $2^{-\frac{e}{2}} \leq a < 2^{\frac{1-e}{2}}$. Next, since $s_a + s_b < \sqrt{2}$ and $s_a > 0$, we have $s_b < \sqrt{2}$, so that $b^2 < \sqrt{2} + u$ and $1 < b < \sqrt{\sqrt{2} + u}$. Finally, we have $a^2 + b^2 \leq s_a + \text{ufp}(a^2)u + s_b + \text{ufp}(b^2)u < \sqrt{2} + \frac{5}{4}u$, which concludes the domain analysis: it suffices to look for an upper bound for g_1 over the domain

$$D_1 := \left\{ (a, b, e) \mid 2^{-\frac{e}{2}} \leq a < 2^{\frac{1-e}{2}}, 1 \leq b < \sqrt{\sqrt{2} + u}, a^2 + b^2 < \sqrt{2} + \frac{5}{4}u, e \geq 2 \right\}.$$

We now compute the partial derivatives of g_1 with respect to a and b ,

$$\frac{\partial g_1}{\partial a} = 2a \left(2^{-e} + \frac{1}{4}\right) + \frac{2 + 2^{-e}}{a^2 + b^2} 2^{1-\frac{e}{2}} - 4a \frac{(2 + 2^{-e}) \left(2^{-\frac{e}{2}} a + \frac{b}{2}\right)}{(a^2 + b^2)^2} - 4a \frac{(2 + 2^{-e})^2}{(a^2 + b^2)^3},$$

$$\frac{\partial g_1}{\partial b} = 2b \left(2^{-e} + \frac{1}{4}\right) + \frac{2 + 2^{-e}}{a^2 + b^2} - 4b \frac{(2 + 2^{-e}) \left(2^{-\frac{e}{2}} a + \frac{b}{2}\right)}{(a^2 + b^2)^2} - 4b \frac{(2 + 2^{-e})^2}{(a^2 + b^2)^3},$$

and the next step is to prove that they are both negative over the domain D_1 . Since $\frac{1}{b} \frac{\partial}{\partial b} g_1(a, b, e) - \frac{1}{a} \frac{\partial}{\partial a} g_1(a, b, e) = \frac{2+2^{-e}}{a^2+b^2} \left(\frac{1}{b} - \frac{1}{a} 2^{1-\frac{e}{2}}\right) < 0$ over D_1 , it is sufficient to prove that $\frac{\partial}{\partial a} g_1(a, b, e) < 0$. Since $2a \frac{2+2^{-e}}{a^2+b^2} > 0$, we can rewrite this inequality as

$$\frac{(2^{-e} + \frac{1}{4})(a^2 + b^2)}{2 + 2^{-e}} + \frac{2^{-\frac{e}{2}}}{a} < 2 \frac{2^{-\frac{e}{2}} a + \frac{b}{2}}{a^2 + b^2} + 2 \frac{2 + 2^{-e}}{(a^2 + b^2)^2}.$$

This inequality follows from the following three relations:

$$\begin{aligned} \frac{(2^{-e} + \frac{1}{4})(a^2 + b^2)}{2 + 2^{-e}} + \frac{2^{-\frac{e}{2}}}{a} &< \frac{\sqrt{2} + \frac{5}{4}u}{4} + 1 \quad \text{for } (a, b, e) \in D_1, \\ \frac{\sqrt{2} + \frac{5}{4}u}{4} + 1 &< \frac{1}{\sqrt{2} + \frac{5}{4}u} + \frac{4}{(\sqrt{2} + \frac{5}{4}u)^2} \quad \text{for } p \geq 3, \\ \frac{1}{\sqrt{2} + \frac{5}{4}u} + \frac{4}{(\sqrt{2} + \frac{5}{4}u)^2} &< 2 \frac{2^{-\frac{e}{2}} a + \frac{b}{2}}{a^2 + b^2} + 2 \frac{2 + 2^{-e}}{(a^2 + b^2)^2} \quad \text{for } (a, b, e) \in D_1. \end{aligned}$$

Since both $\frac{\partial g_1}{\partial a}$ and $\frac{\partial g_1}{\partial b}$ are negative over D_1 , since $(a, b, e) \in D_1$ implies $a \geq 2^{-\frac{e}{2}}$ and $b \geq 1$, and since $(2^{-\frac{e}{2}}, 1, e) \in D_1$, we deduce that $g_1(a, b, e) \leq g_1(2^{-\frac{e}{2}}, 1, e) =: h_1(x)$, with $x = 2^{-e}$ and

$$h_1(x) = (x+1)\left(x + \frac{1}{4}\right) + \frac{(x+2)(2x+1)}{x+1} + \left(\frac{x+2}{x+1}\right)^2.$$

Since $e \geq 2$, we have $0 < x \leq \frac{1}{4}$ and

$$h_1'(x) = \frac{8x^4 + 37x^3 + 63x^2 + 43x + 1}{4(x+1)^3} > 0.$$

Overall, we thus have $f_2(a, b) \leq g_1(a, b, e) \leq h_1(x) \leq h_1(\frac{1}{4}) = 6.565$.

3.5 Case $\text{ufp}(s_a + s_b) = 4$

From (15) and (11), we know that $4 \leq s_a + s_b < 4\sqrt{2}$ and $s_a < s_b$. As a consequence, we have $2 < s_b$ which implies $2 < b^2$, so that

$$\text{ufp}(b^2) = 2 \quad \text{and} \quad \sqrt{2} < b \leq 2 - 2u.$$

Since 4 is a floating-point number, we have $s = \text{RN}(s_a + s_b) \geq 4$ and $\frac{b}{s} \leq \frac{b}{4} < \frac{1}{2}$ hence

$$\text{ufp}\left(\frac{b}{s}\right) \leq \frac{1}{4}.$$

In the same way, $\frac{a}{s} \leq \frac{a}{4} < 2^{-\frac{3+e}{2}}$ so that

$$\text{ufp}\left(\frac{a}{s}\right) \leq 2^{-2-\frac{e}{2}}.$$

We now distinguish two subcases, namely $e = -1$ and $e \geq 0$.

3.5.1 Subcase $e = -1$

We have $\text{ufp}\left(\frac{a}{s}\right) \leq 2^{-\frac{3}{2}}$, hence $\text{ufp}\left(\frac{a}{s}\right) \leq \frac{1}{4}$, thus we deduce from (8) that $f_2(a, b) \leq g_2(a, b)$ with

$$g_2(a, b) := \frac{a^2 + b^2}{8} + \frac{4(a+b)}{a^2 + b^2} + \left(\frac{8}{a^2 + b^2}\right)^2.$$

From (15), we know that $s_a + s_b < 4\sqrt{2}$, which implies $a^2 + b^2 < 4\sqrt{2} + 4u$. The domain of interest is then given by

$$D_2 := \{(a, b) \mid \sqrt{2} \leq a \leq b < 2, a^2 + b^2 < 4\sqrt{2} + 4u\}.$$

Computing the partial derivatives of g_2 with respect to a and b , and proving that they are both negative over the domain D_2 (detailed computations are in §A.2), we end up with $f_2(a, b) \leq g_2(\sqrt{2}, \sqrt{2}) = \left(2 + \frac{\sqrt{2}}{2}\right)^2$.

3.5.2 Subcase $e \geq 0$

Since $|\delta| \leq \text{ufp}(a^2) + \text{ufp}(b^2) + \text{ufp}(s_a + s_b) = 6 + 2^{-e}$ and $\text{ufp}(\frac{a}{s}) \leq 2^{-2-\frac{e}{2}}$, from (8) we get $f_2(a, b) \leq g_3(a, b, e)$ with

$$g_3(a, b, e) := \frac{(a^2 + b^2)(2^{-e} + 1)}{16} + \frac{(6 + 2^{-e})(2^{-\frac{e}{2}}a + b)}{2(a^2 + b^2)} + \left(\frac{6 + 2^{-e}}{a^2 + b^2}\right)^2.$$

From (14), $a^2 + b^2$ is lower bounded by 4, and we restrict the analysis of g_3 to the domain

$$D_3 := \{(a, b, e) \mid 2^{-\frac{e}{2}} \leq a \leq 2^{\frac{1-e}{2}}, \sqrt{2} \leq b < 2, 4 \leq a^2 + b^2 < 4\sqrt{2} + 4u, e \geq 0\}.$$

First, it can be checked that the partial derivative of g_3 with respect to b is negative over D_3 (details are in §A.3). Since $b \geq \sqrt{4 - a^2}$, and $(a, b, e) \in D_3$ implies $(a, \sqrt{4 - a^2}, e) \in D_3$, we deduce that $g_3(a, b, e) \leq g_3(a, \sqrt{4 - a^2}, e)$, where

$$g_3(a, \sqrt{4 - a^2}, e) = \frac{2^{-e} + 1}{4} + \frac{(6 + 2^{-e})(2^{-\frac{e}{2}}a + \sqrt{4 - a^2})}{8} + \frac{(6 + 2^{-e})^2}{16}.$$

We then compute $\frac{\partial}{\partial a} g_3(a, \sqrt{4 - a^2}, e) = \frac{6 + 2^{-e}}{8} \left(2^{-\frac{e}{2}} - \frac{a}{\sqrt{4 - a^2}}\right)$, which is non-negative because $a^2 \leq \frac{2a^2}{1 + 2^{-e}} \leq \frac{4 \cdot 2^{-e}}{1 + 2^{-e}}$. Since $(2^{\frac{1-e}{2}}, \sqrt{4 - 2^{1-e}}, e) \in D_3$, we have $g_3(a, b, e) \leq g_3(2^{\frac{1-e}{2}}, \sqrt{4 - 2^{1-e}}, e) =: h_3(x)$, with $x = 2^{-e}$ and

$$h_3(x) = \frac{x + 1}{4} + \frac{(6 + x)(\sqrt{2x} + \sqrt{4 - 2x})}{8} + \frac{(6 + x)^2}{16}.$$

Since

$$h_3'(x) = 1 + \frac{x}{8} \left(1 + \sqrt{2}\right) + \frac{\sqrt{4 - 2x}}{8} + \frac{x + 6}{8} \left(\sqrt{2} - \frac{1}{\sqrt{4 - 2x}}\right)$$

is positive for $0 < x \leq 1$, we deduce $f_2(a, b) \leq h_3(1) = \left(\frac{7}{4} + \frac{\sqrt{2}}{2}\right)^2 = 6.037\dots$

3.6 Case $\text{ufp}(s_a + s_b) = 2$

From (14) we have $2 \leq a^2 + b^2$, and from (15) we have $2 \leq s_a + s_b < 2\sqrt{2}$ hence

$$e \geq 0.$$

Since 2 is a floating-point number, we know that $s \geq 2$. Therefore $\frac{a}{s} < 2^{-\frac{1+e}{2}}$, hence

$$\text{ufp}\left(\frac{a}{s}\right) \leq 2^{-1-\frac{e}{2}}, \quad (16)$$

and $\frac{b}{s} < 1$ so that

$$\text{ufp}\left(\frac{b}{s}\right) \leq \frac{1}{2}.$$

We handle separately the two possible values, 1 and 2, for $\text{ufp}(b^2)$.

3.6.1 Subcase $\text{ufp}(b^2) = 1$

We distinguish the cases $e \geq 1$ and $e = 0$.

- *Subsubcase $e \geq 1$:* From (8) we have $f_2(a, b) \leq g_4(a, b, e)$ with

$$g_4(a, b, e) := \frac{(a^2 + b^2)(2^{-e} + 1)}{4} + \frac{(3 + 2^{-e})(2^{-\frac{e}{2}}a + b)}{a^2 + b^2} + \left(\frac{3 + 2^{-e}}{a^2 + b^2}\right)^2.$$

From (14), we know that $a^2 + b^2$ is lower bounded by 2. On the other hand, we have $a^2 + b^2 \leq s_a + s_b + (\text{ufp}(a^2) + \text{ufp}(b^2))u < 2\sqrt{2} + 2u$ and $1 < b < \sqrt{2}$, hence we can restrict the analysis to the domain

$$D_4 := \{(a, b, e) \mid 2^{-\frac{e}{2}} \leq a < 2^{\frac{1-e}{2}}, 1 < b < \sqrt{2}, 2 \leq a^2 + b^2 < 2\sqrt{2} + 2u, e \geq 1\}.$$

We can first compute the partial derivative of g_4 with respect to b and prove it is negative over D_4 for $p \geq 4$ (see the details in §A.4). Since $(a, \sqrt{2 - a^2}, e)$ is in D_4 , we deduce that $g_4(a, b, e) \leq g_4(a, \sqrt{2 - a^2}, e)$, and we have

$$g_4(a, \sqrt{2 - a^2}, e) = \frac{2^{-e} + 1}{2} + \frac{(3 + 2^{-e})(2^{-\frac{e}{2}}a + \sqrt{2 - a^2})}{2} + \frac{(3 + 2^{-e})^2}{4}.$$

Next, we can compute the derivative of $g_4(a, \sqrt{2 - a^2}, e)$ with respect to a (see §A.4) and check that the maximum is attained at $a_0 = 2^{-\frac{e}{2}} \sqrt{\frac{2}{1 + 2^{-e}}}$, so that $g_4(a, b, e) \leq g_4(a_0, \sqrt{2 - a_0^2}, e) =: h_4(x)$ with

$$h_4(x) = \frac{x + 1}{2} + \frac{3 + x}{2} \left(x \sqrt{\frac{2}{1 + x}} + \sqrt{2 - \frac{2x}{1 + x}} \right) + \frac{(3 + x)^2}{4}.$$

Since $h_4'(x) > 0$ for $0 < x \leq \frac{1}{2}$, we conclude that $f_2(a, b) \leq g_4(a, b, e) \leq h_4(\frac{1}{2}) = (\frac{7}{4} + \frac{\sqrt{3}}{2})^2$.

- *Subsubcase $e = 0$:* According to (13), we assume that $1 < a$, so that $\text{ufp}(b^2) = \text{ufp}(a^2) = 1$. It follows that $s \geq s_a + s_b - 2u \geq a^2 + b^2 - 4u$, hence $\frac{a}{s} \leq \frac{a}{a^2 + b^2 - 4u}$. Since a and b are both floating-point numbers, and from (11), we know that $b \geq a + 2u$ so that $b^2 - 4u > a^2$. By computing its partial derivative, it can then be checked that $\frac{a}{a^2 + b^2 - 4u}$ is increasing with respect to a , which implies $\frac{a}{s} \leq \frac{b - 2u}{(b - 2u)^2 + b^2 - 4u}$. This last expression is decreasing with respect to b , and since $b \geq 1 + 2u$ we deduce $\frac{a}{s} \leq \frac{1}{2(1 + 2u^2)} < \frac{1}{2}$. Thus,

$$\text{ufp}\left(\frac{a}{s}\right) \leq \frac{1}{4}.$$

In the same way, it can be derived from $\frac{b}{s} \leq \frac{b}{a^2 + b^2 - 4u}$ that

$$\text{ufp}\left(\frac{b}{s}\right) \leq \frac{1}{4}.$$

Combining these bounds on $\text{ufp}(\frac{a}{s})$ and $\text{ufp}(\frac{b}{s})$ with (8) gives $f_2(a, b) \leq g_5(a, b)$, where

$$g_5(a, b) := \frac{a^2 + b^2}{8} + \frac{2(a+b)}{a^2 + b^2} + \frac{16}{(a^2 + b^2)^2}.$$

Hence it remains to bound $g_5(a, b)$ over the domain D_5 defined by

$$D_5 := \{(a, b) \mid 1 \leq a \leq b < \sqrt{2} \text{ and } a^2 + b^2 < 2\sqrt{2} + 2u\}.$$

In this domain, we have $\frac{\partial}{\partial b} g_5(a, b) < 0$ (details are in §A.5), so that $g_5(a, b) \leq g_5(a, a) = \frac{a^2}{4} + \frac{4}{a^4} + \frac{2}{a}$, which is maximal for $a = 1$. Therefore, we deduce that $f_2(a, b) \leq g_5(a, b) \leq g_5(1, 1) = (\frac{5}{2})^2$.

3.6.2 Subcase $\text{ufp}(b^2) = 2$

In this paragraph, $a^2 < 1$ (otherwise we would have $s_a + s_b \geq 2 + 1$ while from (15) we have $s_a + s_b < 2\sqrt{2}$), hence $e \geq 1$. We split the inequality (16) into two possible cases. Either $\text{ufp}(\frac{a}{s}) < 2^{-1-\frac{e}{2}}$ which implies $\text{ufp}(\frac{a}{s}) \leq 2^{-\frac{3+e}{2}}$, or $\text{ufp}(\frac{a}{s}) = 2^{-1-\frac{e}{2}}$ in which case e is even.

• *Subsubcase* $\text{ufp}(\frac{a}{s}) < 2^{-1-\frac{e}{2}}$: We deduce from (8) and $|\delta| \leq 4 + 2^{-e}$ that $f_2(a, b) \leq g_6(a, b, e)$ with

$$g_6(a, b, e) := \frac{(a^2 + b^2)(2^{-1-e} + 1)}{4} + \frac{(4 + 2^{-e})(2^{-\frac{1+e}{2}}a + b)}{a^2 + b^2} + \left(\frac{4 + 2^{-e}}{a^2 + b^2}\right)^2.$$

We can compute the derivatives of g_6 (details are provided in §A.6) with respect to a and b and prove that they are negative over the domain

$$D_6 := \{(a, b, e) \mid 2^{-\frac{e}{2}} \leq a < 2^{\frac{1-e}{2}}, \sqrt{2} \leq b < 2, \\ 2 \leq a^2 + b^2 < 2\sqrt{2} + (2 + 2^{-e})u, e \geq 1\}.$$

For $(a, b, e) \in D_6$, we deduce that $g_6(a, b, e) \leq g_6(2^{-\frac{e}{2}}, \sqrt{2}, e) =: h_6(x)$ with

$$h_6(x) = \frac{(x+2)(\frac{x}{2}+1)}{4} + \frac{\sqrt{2}(4+x)(\frac{x}{2}+1)}{x+2} + \left(\frac{4+x}{x+2}\right)^2, \quad x = 2^{-e}.$$

We can maximize $h_6(x)$ for $0 < x \leq \frac{1}{2}$, which leads to $f_2(a, b) \leq h_6(0) = (2 + \frac{\sqrt{2}}{2})^2$.

• *Subsubcase* $\text{ufp}(\frac{a}{s}) = 2^{-1-\frac{e}{2}}$: In this case, e is even, hence $e \geq 2$. We have $f_2(a, b) \leq g_7(a, b, e)$ with

$$g_7(a, b, e) := \frac{(a^2 + b^2)(2^{-e} + 1)}{4} + \frac{(4 + 2^{-e})(2^{-\frac{e}{2}}a + b)}{a^2 + b^2} + \left(\frac{4 + 2^{-e}}{a^2 + b^2}\right)^2.$$

The lower bound $2^{-\frac{e}{2}}$ for a does not lead to a sufficiently tight bound for f_2 in this case: to get a better bound, we exploit further the hypothesis $\text{ufp}(\frac{a}{s}) = 2^{-1-\frac{e}{2}}$. This gives $s2^{-1-\frac{e}{2}} \leq a$, which implies $a^2 - 2^{1+\frac{e}{2}}a + b^2 + \delta u \leq 0$, hence

$$a \geq 2^{\frac{e}{2}} - \sqrt{2^e - 2 + (4 + 2^{-e})u} = a_0 + \eta(u)$$

with

$$a_0 = 2^{\frac{e}{2}} - \sqrt{2^e - 2}, \quad \eta(u) < 0, \quad |\eta(u)| \in \mathcal{O}(u).$$

Therefore, we analyze g_7 over the domain

$$D_7 := \{(a, b, e) \mid a_0 + \eta(u) \leq a < 2^{\frac{1-e}{2}}, \sqrt{2} \leq b < 2, \\ 2 \leq a^2 + b^2 < 2\sqrt{2} + (2 + 2^{-e})u, e \geq 2, e \text{ even}\}.$$

First, we can compute the partial derivative of g_7 with respect to b and prove (see Appendix A.7) that it is negative over the domain D_7 , hence we know that $g_7(a, b, e) \leq g_7(a, \sqrt{2}, e)$.

It can be checked that $a_0 + \eta(u)$ belongs to $[2^{-1-\frac{e}{2}}, 2^{\frac{1-e}{2}}]$ and that $g_7(a, \sqrt{2}, e)$ is decreasing with respect to a over $[2^{-1-\frac{e}{2}}, 2^{\frac{1-e}{2}}]$; see Appendix A.7. We then deduce that $g_7(a, \sqrt{2}, e) \leq g_7(a_0 + \eta(u), \sqrt{2}, e)$, for any $(a, b, e) \in D_7$.

Next, it can be proved that $g_7(a_0 + \eta(u), \sqrt{2}, e) \leq g_7(a_0, \sqrt{2}, e) + 20u$ (again, the details are provided in §A.7). As a consequence, for any $(a, b, e) \in D_7$ we have $g_7(a, b, e) \leq g_7(a_0, \sqrt{2}, e) + 20u$.

The last step is to bound $g_7(a_0, \sqrt{2}, e)$ for e an even positive integer. With $y = \sqrt{1 - 2^{1-e}}$, we have $g_7(a_0, \sqrt{2}, e) =: h_7(y)$ with $h_7(y)$ a rational function in y . The variable y belongs to $[\sqrt{2}/2, 1]$, and $h_7'(y) = \frac{P(y)}{32(y+1)^2}$ with

$$P(y) = 3y^7 + 11y^6 - 5y^5 - (12\sqrt{2} + 85)y^4 - (32\sqrt{2} + 143)y^3 \\ - (23 - 8\sqrt{2})y^2 + (64\sqrt{2} + 113)y + 36\sqrt{2} + 33.$$

Using Descartes' rule of signs, one can check that P has exactly one root in the interval $[\sqrt{2}/2, 1]$, and since the evaluation of P is positive at $\sqrt{1 - 2^{-5}}$ and negative at $\sqrt{1 - 2^{-7}}$, we deduce that h_7 is increasing over $[\sqrt{2}/2, \sqrt{1 - 2^{-5}}]$ and decreasing over $[\sqrt{1 - 2^{-7}}, 1]$. Comparing the values of h_7 at the points $\sqrt{1 - 2^{-5}}$ and $\sqrt{1 - 2^{-7}}$, we conclude that $h_7(\sqrt{1 - 2^{-7}})$ is an upper bound for h_7 .

Finally, it can be checked that $h_7(\sqrt{1 - 2^{-7}}) = \gamma^2$ hence we get $f_2(a, b) \leq \gamma^2 + 20u$. From (9), we derive the final upper bound $\gamma u + 9u^2$ for E_N , which concludes the proof of Theorem 2.

4 Implications for complex floating-point division

Let us conclude with some remarks about complex division. The conventional complex division algorithm for computing an approximation $\hat{z} = \hat{R} + i\hat{T}$ of $(a + ib)/(c + id)$ in floating-point arithmetic consists in evaluating the real part as

$$\hat{R} = \text{RN}\left(\frac{\text{RN}(\text{RN}(ac) + \text{RN}(bd))}{\text{RN}(\text{RN}(c^2) + \text{RN}(d^2))}\right) \quad (17)$$

and using a similar scheme for the imaginary part. An approximate quotient \hat{z} can also be obtained by first computing an inverse of $c + id$ using Algorithm 1, and then multiplying it by $a + ib$ by means of the classic complex multiplication algorithm. Note that both algorithms require 3 additions/subtractions, 6 multiplications, and 2 divisions.

Normwise relative accuracy analyses of the method based on (17) can be found in [4, 12, 5]. To our knowledge, the best known upper bound for the normwise relative error generated by this method is $(3 + \sqrt{5})u + \mathcal{O}(u^2) \approx 5.2u$: as noted in [1], this bound can be derived from the bound $\sqrt{5}u$ from [3] on the normwise relative error for the classic complex multiplication algorithm. On the other hand, it can be checked using Theorem 2 and, again, the bound $\sqrt{5}u$ from [3] that the algorithm combining inversion and multiplication admits the smaller normwise error bound $(\gamma + \sqrt{5})u + \mathcal{O}(u^2) \approx 4.9u$. The following examples of complex quotients in precision $p = 11$ show that in both cases the largest normwise relative error cannot be bounded by $\gamma u + \mathcal{O}(u^2) \approx 2.7u$ as for inversion:

- with $a + ib = 1575 + i 1419$ and $c + id = 1457 + i 1480$, using (17) gives $|\hat{z} - z|/(u|z|) = 4.67973\dots$;
- dividing $1506 + i 1512$ by $1491 + i 1504$ using the inversion-multiplication approach leads to $|\hat{z} - z|/(u|z|) = 4.34446\dots$

However, these examples are not sufficient to conclude about the sharpness of the bounds $(3 + \sqrt{5})u + \mathcal{O}(u^2)$ and $(\gamma + \sqrt{5})u + \mathcal{O}(u^2)$, and further investigation is needed to understand the accuracy of complex floating-point division.

Acknowledgements

We thank the associate editor and the anonymous reviewers for their helpful comments and suggestions.

References

1. Baudin, M.: Error bounds of complex arithmetic (2011). Available at http://forge.scilab.org/upload/compdiv/files/complexerrorbounds_v0.2.pdf
2. Baudin, M., Smith, R.L.: A robust complex division in Scilab (2012). Available at <http://arxiv.org/abs/1210.4539>

3. Brent, R., Percival, C., Zimmermann, P.: Error bounds on complex floating-point multiplication. *Mathematics of Computation* **76**, 1469–1481 (2007)
4. Champagne, W.P.: On finding roots of polynomials by hook or by crook. Master's thesis, University of Texas (1964)
5. Higham, N.J.: *Accuracy and Stability of Numerical Algorithms*, second edn. SIAM, Philadelphia, PA, USA (2002)
6. IEEE Computer Society: IEEE Standard for Floating-Point Arithmetic. IEEE Standard 754-2008 (2008). Available at <http://ieeexplore.ieee.org/servlet/opac?punumber=4610933>
7. Knuth, D.E.: *The Art of Computer Programming, Volume 2, Seminumerical Algorithms*, third edn. Addison-Wesley, Reading, MA, USA (1998)
8. Priest, D.M.: Efficient scaling for complex division. *ACM Transactions on Mathematical Software* **30**(4) (2004)
9. Rump, S.M., Ogita, T., Oishi, S.: Accurate floating-point summation, Part I: Faithful rounding. *SIAM Journal on Scientific Computing* **31**(1), 189–224 (2008)
10. Smith, R.L.: Algorithm 116: Complex division. *Communications of the ACM* **5**(8), 435 (1962)
11. Stewart, G.W.: A note on complex division. *ACM Transactions on Mathematical Software* **11**(3), 238–241 (1985)
12. Wilkinson, J.H.: *The Algebraic Eigenvalue Problem*. Oxford University Press (1965)
13. Ziv, A.: Sharp ULP rounding error bound for the hypotenuse function. *Mathematics of Computation* **68**(227), 1143–1148 (1999)

A Details omitted in the proofs

A.1 Asymptotic optimality of the componentwise error bound

We briefly detail the computations of s_a , s_b and s in the example parametrized by p given in Section 2. We assume that $p \geq 12$ is even, and we recall that

$$\begin{aligned} a &= 2^{\frac{p}{2}-1} + 5 \cdot 2^{-2} + 2^{-\frac{p}{2}+2}, \\ b &= 2^{p-1} + 2^{\frac{p}{2}-1} + 1. \end{aligned}$$

- Computation of $s_a = \text{RN}(a^2)$:

$$\begin{aligned} a^2 &= 2^{p-2} + 5 \cdot 2^{\frac{p}{2}-2} + 11 \cdot 2^{-1} + 2^{-4} + 10 \cdot 2^{-\frac{p}{2}} + 2^{-p+4} \\ \text{ulp}(a^2) &= 2^{-1} \\ \tilde{s}_a &:= 2^{p-2} + 5 \cdot 2^{\frac{p}{2}-2} + 11 \cdot 2^{-1} \\ |a^2 - \tilde{s}_a| &= 2^{-4} + 10 \cdot 2^{-\frac{p}{2}} + 2^{4-p} \\ &\leq 2^{-4} + 10 \cdot 2^{-6} + 2^{-8} \\ &< 2^{-2} = \frac{1}{2} \text{ulp}(a^2) \end{aligned}$$

Hence $s_a = \tilde{s}_a$.

- Computation of $s_b = \text{RN}(b^2)$:

$$\begin{aligned} b^2 &= 2^{2p-2} + 2^{\frac{3p}{2}-1} + 2^p + 2^{p-2} + 2^{\frac{p}{2}} + 1 \\ \tilde{s}_b &:= 2^{2p-2} + 2^{\frac{3p}{2}-1} + 3 \cdot 2^{p-1} \\ \text{ulp}(b^2) &= 2^{p-1} \\ |b^2 - \tilde{s}_b| &= 2^{p-2} - 2^{\frac{p}{2}} - 1 \\ &< 2^{p-2} = \frac{1}{2} \text{ulp}(b^2) \end{aligned}$$

Hence $s_b = \tilde{s}_b$.

- Computation of $s = \text{RN}(s_a + s_b)$:

$$\begin{aligned} s_a + s_b &= 2^{2p-2} + 2^{\frac{3p}{2}-1} + 3 \cdot 2^{p-1} + 2^{p-2} + 5 \cdot 2^{\frac{p}{2}-2} + 11 \cdot 2^{-1} \\ \tilde{s} &= 2^{2p-2} + 2^{\frac{3p}{2}-1} + 2^{p+1} \\ \text{ulp}(s_a + s_b) &= 2^{p-1} \\ |s_a + s_b - \tilde{s}| &= 2^{p-2} - 5 \cdot 2^{\frac{p}{2}-2} - 11 \cdot 2^{-1} \\ &< 2^{p-2} = \frac{1}{2} \text{ulp}(s_a + s_b) \end{aligned}$$

Hence $s = \tilde{s}$.

A.2 Partial derivatives of g_2

Computing the partial derivatives of g_2 with respect to a and b gives

$$\begin{aligned}\frac{\partial g_2}{\partial a} &= \frac{a}{4} + \frac{4}{a^2 + b^2} - \frac{8a(a+b)}{(a^2 + b^2)^2} - \frac{256a}{(a^2 + b^2)^3}, \\ \frac{\partial g_2}{\partial b} &= \frac{b}{4} + \frac{4}{a^2 + b^2} - \frac{8b(a+b)}{(a^2 + b^2)^2} - \frac{256b}{(a^2 + b^2)^3}.\end{aligned}$$

First, we know that $b > a$ so $\frac{1}{b} \frac{\partial}{\partial b} g_2(a, b) < \frac{1}{a} \frac{\partial}{\partial a} g_2(a, b)$. We just have to prove that $\frac{\partial}{\partial a} g_2(a, b) < 0$, that is,

$$\frac{(a^2 + b^2)^2}{4} + \frac{4(a^2 + b^2)}{a} < 8(a + b) + \frac{256}{a^2 + b^2}.$$

Since for $(a, b) \in D_2$ we have $\sqrt{2} < a, b$, and $a^2 + b^2 < 4\sqrt{2} + 4u$, it is enough to check that

$$\frac{(4\sqrt{2} + 4u)^2}{4} + \frac{4(4\sqrt{2} + 4u)}{\sqrt{2}} < 16\sqrt{2} + \frac{256}{4\sqrt{2} + 4u},$$

which holds for $p \geq 2$.

A.3 Partial derivative of g_3

We compute the partial derivative of g_3 with respect to b , and check that this derivative is negative over the domain D_3 . We have

$$\frac{\partial g_3}{\partial b} = \frac{b(2^{-e} + 1)}{8} + \frac{6 + 2^{-e}}{2(a^2 + b^2)} - b \frac{(6 + 2^{-e})(2^{-\frac{e}{2}}a + b)}{(a^2 + b^2)^2} - 4b \frac{(6 + 2^{-e})^2}{(a^2 + b^2)^3},$$

and we check that

$$\frac{b(2^{-e} + 1)}{8} + \frac{6 + 2^{-e}}{2(a^2 + b^2)} < b \frac{(6 + 2^{-e})(2^{-\frac{e}{2}}a + b)}{(a^2 + b^2)^2} + 4b \frac{(6 + 2^{-e})^2}{(a^2 + b^2)^3}.$$

Multiplying both sides by $\frac{(a^2 + b^2)^2}{b(6 + 2^{-e})}$ and since $1 \leq b$, it is enough to prove that

$$\frac{(2^{-e} + 1)(a^2 + b^2)^2}{8(6 + 2^{-e})} + \frac{a^2 + b^2}{2} < 2^{-\frac{e}{2}}a + b + 4 \frac{6 + 2^{-e}}{a^2 + b^2}.$$

This follows from the following sequence of three inequalities

$$\begin{aligned}\frac{(2^{-e} + 1)(a^2 + b^2)^2}{8(6 + 2^{-e})} + \frac{a^2 + b^2}{2} &< \frac{2(4\sqrt{2} + 4u)^2}{48} + \frac{4\sqrt{2} + 4u}{2}, \\ \frac{2(4\sqrt{2} + 4u)^2}{48} + \frac{4\sqrt{2} + 4u}{2} &< 4 \frac{6}{4\sqrt{2} + 4u} + 1 \quad \text{for } p \geq 3, \\ 4 \frac{6}{4\sqrt{2} + 4u} + 1 &< 4 \frac{6 + 2^{-e}}{a^2 + b^2} + 2^{-\frac{e}{2}}a + b.\end{aligned}$$

A.4 Partial derivatives of g_4

The partial derivative of g_4 with respect to b is given by

$$\frac{\partial g_4}{\partial b} = \frac{b(2^{-e} + 1)}{2} + \frac{3 + 2^{-e}}{a^2 + b^2} - 2b \frac{(2^{-\frac{\varepsilon}{2}}a + b)(3 + 2^{-e})}{(a^2 + b^2)^2} - 4b \frac{(3 + 2^{-e})^2}{(a^2 + b^2)^3}.$$

We want to prove that $\frac{\partial}{\partial b}g_4(a, b, e) < 0$ or, equivalently, that

$$\frac{(a^2 + b^2)^2(2^{-e} + 1)}{2(3 + 2^{-e})} + \frac{a^2 + b^2}{b} < 2(2^{-\frac{\varepsilon}{2}}a + b) + 4\frac{3 + 2^{-e}}{a^2 + b^2}.$$

This inequality can be derived from the following ones:

$$\begin{aligned} \frac{(a^2 + b^2)^2(2^{-e} + 1)}{2(3 + 2^{-e})} + \frac{a^2 + b^2}{b} &< \frac{2(2\sqrt{2} + 2u)^2}{6} + 2\sqrt{2} + 2u, \\ \frac{(2\sqrt{2} + 2u)^2}{3} + 2\sqrt{2} + 2u &< 2 + \frac{12}{2\sqrt{2} + 2u} \quad \text{for } p \geq 4, \\ 2 + \frac{12}{2\sqrt{2} + 2u} &< 2(2^{-\frac{\varepsilon}{2}}a + b) + 4\frac{3 + 2^{-e}}{a^2 + b^2}. \end{aligned}$$

The partial derivative of $g_4(a, \sqrt{2 - a^2}, e)$ with respect to a is

$$\frac{\partial}{\partial a}g_4(a, \sqrt{2 - a^2}, e) = \frac{3 + 2^{-e}}{2} \left(2^{-\frac{\varepsilon}{2}} - \frac{a}{\sqrt{2 - a^2}} \right),$$

which is zero if $a = a_0$ with $a_0 = 2^{-\frac{\varepsilon}{2}}\sqrt{\frac{2}{1 + 2^{-e}}}$, positive if $a < a_0$, and negative if $a > a_0$.

A.5 Partial derivative of g_5

We have

$$\frac{\partial g_5}{\partial b} = \frac{b}{4} + \frac{2}{a^2 + b^2} - \frac{4(a + b)}{(a^2 + b^2)^2}b - \frac{64}{(a^2 + b^2)^3}b,$$

and it can be checked that this partial derivative is negative using the following inequalities:

$$\begin{aligned} \frac{(a^2 + b^2)^2}{4} + \frac{2}{b}(a^2 + b^2) &< \frac{(2\sqrt{2} + 2u)^2}{4} + 2(2\sqrt{2} + 2u), \\ \frac{(2\sqrt{2} + 2u)^2}{4} + 2(2\sqrt{2} + 2u) &< 8 + \frac{64}{2\sqrt{2} + 2u} \quad \text{for } p \geq 2, \\ 8 + \frac{64}{2\sqrt{2} + 2u} &< 4(a + b) + \frac{64}{a^2 + b^2}. \end{aligned}$$

A.6 Partial derivatives of g_6

The partial derivatives of g_6 with respect to a and b are

$$\frac{\partial g_6}{\partial a} = \frac{a}{4}(2^{-e} + 2) + \frac{4 + 2^{-e}}{a^2 + b^2} 2^{-\frac{1+\varepsilon}{2}} - 2a \frac{(2^{-\frac{1+\varepsilon}{2}}a + b)(4 + 2^{-e})}{(a^2 + b^2)^2} - 4a \frac{(4 + 2^{-e})^2}{(a^2 + b^2)^3}$$

and

$$\frac{\partial g_6}{\partial b} = \frac{b}{4} (2^{-e} + 2) + \frac{4 + 2^{-e}}{a^2 + b^2} - 2b \frac{\left(2^{-\frac{1+e}{2}} a + b\right) (4 + 2^{-e})}{(a^2 + b^2)^2} - 4b \frac{(4 + 2^{-e})^2}{(a^2 + b^2)^3}.$$

For $(a, b, e) \in D_6$, it can be checked that $\frac{\partial g_6}{\partial a}(a, b, e) < 0$ and $\frac{\partial g_6}{\partial b}(a, b, e) < 0$ as follows. Note first that $2^{-\frac{e}{2}} \leq a$ implies

$$\frac{4 + 2^{-e}}{a^2 + b^2} 2^{-\frac{1+e}{2}} \leq \frac{4 + 2^{-e}}{\sqrt{2}(a^2 + b^2)} a,$$

and that $\sqrt{2} \leq b$ implies

$$\frac{4 + 2^{-e}}{a^2 + b^2} \leq \frac{4 + 2^{-e}}{\sqrt{2}(a^2 + b^2)} b.$$

Thus, the same expression can be used as an upper bound for both $\frac{1}{a} \frac{\partial g_6}{\partial a}$ and $\frac{1}{b} \frac{\partial g_6}{\partial b}$. Then, multiplying it by $\frac{(a^2 + b^2)^2}{4 + 2^{-e}}$, it is enough to prove that

$$\frac{(a^2 + b^2)^2 (2^{-1-e} + 1)}{2(4 + 2^{-e})} + \frac{a^2 + b^2}{\sqrt{2}} < 2 \left(2^{-\frac{1+e}{2}} a + b\right) + 4 \frac{4 + 2^{-e}}{a^2 + b^2}.$$

This last inequality follows from the following three ones:

$$\begin{aligned} \frac{(a^2 + b^2)^2 (2^{-1-e} + 1)}{2(4 + 2^{-e})} + \frac{a^2 + b^2}{\sqrt{2}} &< \frac{\left(2\sqrt{2} + \left(2 + \frac{1}{2}\right)u\right)^2 \left(\frac{1}{4} + 1\right)}{8} + 2 + \frac{2 + \frac{1}{2}}{\sqrt{2}}u, \\ \frac{\left(2\sqrt{2} + \left(2 + \frac{1}{2}\right)u\right)^2 \left(\frac{1}{4} + 1\right)}{8} + 2 + \frac{2 + \frac{1}{2}}{\sqrt{2}}u &< 2\sqrt{2} + \frac{16}{2\sqrt{2} + \left(2 + \frac{1}{2}\right)u} \quad \text{for } p \geq 2, \end{aligned}$$

and

$$2\sqrt{2} + \frac{16}{2\sqrt{2} + \left(2 + \frac{1}{2}\right)u} < 2 \left(2^{-\frac{1+e}{2}} a + b\right) + 4 \frac{4 + 2^{-e}}{a^2 + b^2}.$$

A.7 Analysis of g_7

In this section, we provide some details about the analysis of g_7 that were omitted in §3.6.2.

- Let us first maximize g_7 with respect to b . We have

$$\frac{\partial g_7}{\partial b} = \frac{b}{2} (2^{-e} + 1) + \frac{4 + 2^{-e}}{a^2 + b^2} - 2b \frac{\left(2^{-\frac{e}{2}} a + b\right) (4 + 2^{-e})}{(a^2 + b^2)^2} - 4b \frac{(4 + 2^{-e})^2}{(a^2 + b^2)^3}.$$

We want to prove that $\frac{\partial}{\partial b} g_7(a, b, e) < 0$ over D_7 . Multiplying by $\frac{(a^2 + b^2)^2}{(4 + 2^{-e})b}$ and using the inequality $\frac{1}{b} < 1$, we only need to prove that

$$\frac{(a^2 + b^2)^2 (2^{-e} + 1)}{2(4 + 2^{-e})} + a^2 + b^2 < 2 \left(2^{-\frac{e}{2}} a + b\right) + 4 \frac{4 + 2^{-e}}{a^2 + b^2}.$$

Since $e \geq 2$, we can derive this inequality for $p \geq 2$ from the three following ones using the definition of D_7 :

$$\frac{(a^2 + b^2)^2 (2^{-e} + 1)}{2(4 + 2^{-e})} + a^2 + b^2 < \frac{\left(2\sqrt{2} + \left(2 + \frac{1}{4}\right)u\right)^2 \left(\frac{1}{4} + 1\right)}{8} + 2\sqrt{2} + \left(2 + \frac{1}{4}\right)u,$$

$$\frac{\left(2\sqrt{2} + \left(2 + \frac{1}{4}\right)u\right)^2 \left(\frac{1}{4} + 1\right)}{8} + 2\sqrt{2} + \left(2 + \frac{1}{4}\right)u < 2\sqrt{2} + \frac{16}{2\sqrt{2} + \left(2 + \frac{1}{4}\right)u},$$

and

$$2\sqrt{2} + \frac{16}{2\sqrt{2} + \left(2 + \frac{1}{4}\right)u} < 2\left(2^{-\frac{e}{2}}a + b\right) + 4\frac{4 + 2^{-e}}{a^2 + b^2}.$$

Therefore, g_7 is decreasing with respect to b , and for all (a, b, e) in D_7 , $g_7(a, b, e) \leq g_7(a, \sqrt{2}, e)$.

- We now maximize $g_7(a, \sqrt{2}, e)$ with respect to a . Let us recall that in D_7 ,

$$a \geq a_0 + \eta(u) = 2^{\frac{e}{2}} - \sqrt{2^e - 2 + (4 + 2^{-e})u};$$

and prove that

$$a_0 + \eta(u) \geq 2^{-1-\frac{e}{2}}.$$

Using the notation $x = 2^{-e}$, the inequality $a_0 + \eta(u) \geq 2^{-1-\frac{e}{2}}$ is equivalent to $\left(\frac{1}{4} - u\right)x \geq -1 + 4u$ which holds for $p \geq 2$ since $\frac{1}{4} - u \geq 0 \geq -1 + 4u$.

Moreover, we have

$$\begin{aligned} \frac{(a^2 + 2)^2}{a(4 + 2^{-e})} \frac{\partial}{\partial a} g_7(a, \sqrt{2}, e) &= \frac{(2^{-e} + 1)(a^2 + 2)^2}{2(4 + 2^{-e})} + \frac{2^{-\frac{e}{2}}}{a} (a^2 + 2) \\ &\quad - 2\left(2^{-\frac{e}{2}}a + \sqrt{2}\right) - 4\frac{4 + 2^{-e}}{a^2 + 2}, \end{aligned}$$

with $\frac{(a^2 + 2)^2}{a(4 + 2^{-e})} > 0$ for $a \in I := \left[2^{-1-\frac{e}{2}}, 2^{\frac{1-e}{2}}\right]$. For $e \geq 2$ and $a \in I$, we have

$$\frac{(a^2 + 2)^2}{a(4 + 2^{-e})} \frac{\partial}{\partial a} g_7(a, \sqrt{2}, e) < \frac{125}{128} + 5 - 2\sqrt{2} - \frac{32}{5} < 0.$$

As a consequence, $g_7(a, \sqrt{2}, e)$ is decreasing with respect to a over I , and since $a_0 + \eta(u) \in I$, the maximum of $g_7(a, \sqrt{2}, e)$ for $a \in [a_0 + \eta(u), 2^{\frac{1-e}{2}}]$ is $g_7(a_0 + \eta(u), \sqrt{2}, e)$.

Thus, for (a, b, e) in D_7 , we have $g_7(a, b, e) \leq g_7(a, \sqrt{2}, e) \leq g_7(a_0 + \eta(u), \sqrt{2}, e)$.

- Let us prove that $g_7(a_0 + \eta(u), \sqrt{2}, e) \leq g_7(a_0, \sqrt{2}, e) + 20u$. For this purpose, we first show that $|\eta(u)| < 2u$. We have

$$|\eta(u)| = \sqrt{2^e - 2 + (4 + 2^{-e})u} - \sqrt{2^e - 2}$$

and a short calculation shows that $|\eta(u)| < 2u$. It can also be checked that $a_0 < 2^{\frac{1-e}{2}}$ using again $x = 2^{-e}$ and a short calculation. Since $e \geq 2$, this implies $a_0 \leq \frac{\sqrt{2}}{2} < 1$. Let us now consider

$$\lambda_0(u) = \frac{1}{(a_0 + \eta(u))^2 + 2}.$$

We have

$$\lambda_0(u) = \frac{1}{a_0^2 + 2} - \frac{2a_0 + \eta(u)}{(a_0^2 + 2)((a_0 + \eta(u))^2 + 2)} \eta(u),$$

and using $-2u < \eta(u) \leq 0$, we deduce

$$\lambda_0(u) < \frac{1}{a_0^2 + 2} + a_0 u.$$

Moreover, we have

$$\lambda_0(u)^2 = \left(\frac{1}{a_0^2 + 2}\right)^2 - \frac{4a_0}{(a_0^2 + 2)^2((a_0 + \eta(u))^2 + 2)}\eta(u) + \frac{(2a_0 + \eta(u))^2 - 2((a_0 + \eta(u))^2 + 2)}{(a_0^2 + 2)^2((a_0 + \eta(u))^2 + 2)^2}\eta(u)^2,$$

and using both $-2u < \eta(u) \leq 0$ and $a_0 < 1$, we also deduce

$$\lambda_0(u)^2 < \left(\frac{1}{a_0^2 + 2}\right)^2 + a_0 u.$$

From the definition of g_7 , using $0 < a_0 + \eta(u) < a_0$ and the previous upper bounds on $\lambda_0(u)$ and $\lambda_0(u)^2$, we obtain

$$g_7(a_0 + \eta(u), \sqrt{2}, e) < \frac{(a_0^2 + 2)(2^{-e} + 1)}{4} + (4 + 2^{-e}) \left(2^{-\frac{e}{2}} a_0 + \sqrt{2}\right) \left(\frac{1}{a_0^2 + 2} + a_0 u\right) + (4 + 2^{-e})^2 \left(\frac{1}{(a_0^2 + 2)^2} + a_0 u\right),$$

and $g_7(a_0 + \eta(u), \sqrt{2}, e) < g_7(a_0, \sqrt{2}, e) + (4 + 2^{-e}) \left(2^{-\frac{e}{2}} a_0 + \sqrt{2} + 4 + 2^{-e}\right) a_0 u$. The inequality $g_7(a_0 + \eta(u), \sqrt{2}, e) < g_7(a_0, \sqrt{2}, e) + 20u$ then follows from $e \geq 2$ and $a_0 \leq \frac{\sqrt{2}}{2}$.

• Finally, we check that the function h_7 is increasing over $\left[\frac{\sqrt{2}}{2}, \sqrt{1 - 2^{-5}}\right]$ and decreasing over $\left[\sqrt{1 - 2^{-7}}, 1\right]$. We have $h_7(y) = \frac{H(y)}{64(y+1)}$ with

$$H(y) = y^7 + 3y^6 - 7y^5 - (8\sqrt{2} + 45)y^4 - (16\sqrt{2} + 53)y^3 + (64\sqrt{2} + 113)y^2 + (144\sqrt{2} + 315)y + 72\sqrt{2} + 249.$$

Hence $h_7'(y) = \frac{P(y)}{32(y+1)^2}$ where P is the polynomial

$$P(y) = 3y^7 + 11y^6 - 5y^5 - (12\sqrt{2} + 85)y^4 - (32\sqrt{2} + 143)y^3 - (23 - 8\sqrt{2})y^2 + (64\sqrt{2} + 113)y + 36\sqrt{2} + 33.$$

This polynomial has 0 or 2 positive roots according to Descartes' rule of signs (there are two sign changes in the sequence of coefficients). Moreover,

$$P(y + 1) = 3y^7 + 32y^6 + 124y^5 + (160 - 12\sqrt{2})y^4 - (208 + 80\sqrt{2})y^3 - (784 + 160\sqrt{2})y^2 - (640 + 64\sqrt{2})y - 96 + 64\sqrt{2},$$

with only one sign change, so there is exactly one root of P greater than 1 and at most one root of P in $\left[\frac{\sqrt{2}}{2}, 1\right]$. Since $P(\sqrt{1 - 2^{-5}}) > 0$ and $P(\sqrt{1 - 2^{-7}}) < 0$, we deduce that $P(y)$ is positive for $y \in \left[\frac{\sqrt{2}}{2}, \sqrt{1 - 2^{-5}}\right]$, and negative for $y \in \left[\sqrt{1 - 2^{-7}}, 1\right]$, which implies that h_7 is increasing over the former interval, and decreasing over the latter.