

BALL ARITHMETIC*

Joris van der Hoeven

LIX, CNRS
École polytechnique
91128 Palaiseau Cedex
France

Email: vdhoeven@lix.polytechnique.fr

Web: <http://lix.polytechnique.fr/~vdhoeven>

March 24, 2010

The MATHEMAGIX project aims at the development of a “computer analysis” system, in which numerical computations can be done in a mathematically sound manner. A major challenge for such systems is to conceive algorithms which are both efficient, reliable and available at any working precision. In this paper, we survey several older and newer such algorithms. We mainly concentrate on the automatic and efficient computation of high quality error bounds, based on a variant of interval arithmetic which we like to call “ball arithmetic”.

KEYWORDS: Ball arithmetic, interval arithmetic, reliable computing, computable analysis

A.M.S. SUBJECT CLASSIFICATION: 65G20, 03F60, 65F99, 37-04

1. INTRODUCTION

Computer algebra systems are widely used today in order to perform mathematically correct computations with objects of algebraic or combinatorial nature. It is tempting to develop a similar system for analysis. On input, the user would specify the problem in a high level language, such as computable analysis [Wei00, BB85, Abe80, Tur36] or interval analysis [Moo66, AH83, Neu90, JKDW01, Kul08]. The system should then solve the problem in a certified manner. In particular, all necessary computations of error bounds must be carried out automatically.

There are several specialized systems and libraries which contain work in this direction. For instance, Taylor models have been used with success for the validated long term integration of dynamical systems [Ber98, MB96, MB04]. A fairly complete MATLAB interval arithmetic library for linear algebra and polynomial computations is INTLAB [Rum99b, Rum99a]. Another historical interval library, which continues to be developed is [ea67]. There exist several libraries for multiple precision arithmetic with correct or at least well specified rounding [HLRZ00, Hai95], as well as a library for multiple precision interval arithmetic [Rev01], and a library for computable real numbers [M0]. There are also libraries for very specific problems, such as the MPSOLVE library [BF00], which allows for the certified computation of the roots of multiple precision polynomials.

The MATHEMAGIX system [vdH+02b] intends to develop a new “computer analysis” system, distributed under a free licence and not depending on proprietary software such as MATLAB. Ultimately, we hope to cover most of the functionality present in the

*. This work has been supported by the ANR-09-JCJC-0098-01 MAGIX project, as well as a Digiteo 2009-36HD grant and Région Ile-de-France.

software mentioned above and add several new features, such as computable analytic functions [vdH07b] and transseries [vdH08b]. We also insist on performance and the possibility to easily switch between machine doubles and multiple precision arithmetic (without being penalized too much). In this paper, we will discuss several older and newer ideas for combining efficiency, reliability and multiple precision arithmetic (when necessary). Many algorithms are already present in the MATHEMAGIX system, or will be incorporated soon.

Most of the algorithms in this paper will be based on “ball arithmetic”, which provides a systematic tool for the automatic computation of error bounds. In section 3, we provide a short introduction to this topic and describe its relation with computable analysis. Although ball arithmetic is really a variant of interval arithmetic, we will give several arguments in section 4 why we actually prefer balls over intervals for most applications. Roughly speaking, balls should be used for the reliable approximation of numbers, whereas intervals are mainly useful for certified algorithms which rely on the subdivision of space. Although we will mostly focus on ball arithmetic in this paper, section 5 presents a rough overview of the other kinds of algorithms that are needed in a complete computer analysis system.

After basic ball arithmetic for real and complex numbers, the next challenge is to develop efficient algorithms for reliable linear algebra, polynomials, analytic functions, resolution of differential equations, etc. In sections 6 and 7 we will survey several basic techniques which can be used to implement efficient ball arithmetic for matrices and formal power series. Usually, there is a trade-off between efficiency and the quality of the obtained error bounds. One may often start with a very efficient algorithm which only computes rough bounds, or bounds which are only good in favourable well-conditioned cases. If the obtained bounds are not good enough, then we may switch to a more expensive and higher quality algorithm.

Of course, the main algorithmic challenge in the area of ball arithmetic is to reduce the overhead of the bound computations as much as possible with respect to the principal numeric computation. In favourable cases, this overhead indeed becomes negligible. For instance, for high working precisions p , the centers of real and complex balls are stored with the full precision, but we only need a small precision for the radii. Consequently, the cost of elementary operations ($+$, $-$, \cdot , \exp , etc.) on balls is dominated by the cost of the corresponding operations on their centers. Similarly, crude error bounds for large matrix products can be obtained quickly with respect to the actual product computation, by using norm bounds for the rows and columns of the multiplicands; see (21).

Another algorithmic challenge is to use fast algorithms for the actual numerical computations on the centers of the balls. In particular, it is important to use existing high performance libraries, such as BLAS, LINPACK, etc., whenever this is possible [Rum99a]. Similarly, we should systematically rely on asymptotically efficient algorithms for basic arithmetic, such as fast integer and polynomial multiplication [KO63, CT65, SS71]. There are several techniques to achieve this goal:

1. The representations of objects should be chosen with care. For instance, should we rather work with ball matrices or matricial balls (see sections 6.4 and 7.1)?
2. If the result of our computation satisfies an equation, then we may first solve the equation numerically and only perform the error analysis at the end. In the case of matrix inversion, this method is known as Hansen’s method; see section 6.2.
3. When considering a computation as a tree or dag, then the error tends to increase with the depth of the tree. If possible, algorithms should be designed so as to keep this depth small. Examples will be given in sections 6.4 and 7.5. Notice that this kind of algorithms are usually also more suitable for parallelization.

When combining the above approaches to the series of problems considered in this paper, we are usually able to achieve a constant overhead for the computation of sharp error bound computations. In favourable cases, the overhead becomes negligible. For particularly ill conditioned problems, we need a logarithmic overhead. It remains an open question whether we have been lucky or whether this is a general pattern.

In this paper, we will be easygoing on what is meant by “sharp error bound”. Regarding algorithms as functions $f: \mathbb{R}^k \rightarrow \mathbb{R}^l$; $x \mapsto y = f(x)$, an error $\delta_x \in \mathbb{R}^k$ in the input automatically gives rise to an error $\delta_y \approx J_f(x) \delta_x$ in the output. When performing our computations with bit precision p , we have to consider that the input error δ_x is at least of the order of $2^{-p} |x| \in (\mathbb{R}^{\geq})^k$ (where $|x|_i = |x_i|$ for all i). Now given an error bound $|\delta_x| \leq \varepsilon_x$ for the input, an error bound $|\delta_y| \leq \varepsilon_y$ is considered to be sharp if $\varepsilon_y \approx \max_{|\delta_x| \leq \varepsilon_x} |J_f(x) \delta_x|$. More generally, condition numbers provide a similar device for measuring the quality of error bounds. A detailed investigation of the quality of the algorithms presented in this paper remains to be carried out.

Of course, the development of asymptotically efficient ball arithmetic presupposes the existence of the corresponding numerical algorithms. This topic is another major challenge, which will not be addressed in this paper. Indeed, asymptotically fast algorithms usually suffer from lower numerical stability. Besides, switching from machine precision to multiple precision involves a huge overhead. Special techniques are required to reduce this overhead when computing with large compound objects, such as matrices or polynomials. It is also recommended to systematically implement single, double, quadruple and multiple precision versions of all algorithms (see [BHL00, BHL01] for a double-double and quadruple-double precision library). The MATHEMAGIX system has been designed so as to ease this task, since most current packages are C++ template libraries.

The aim of our paper is to provide a survey on how to implement an efficient library for ball arithmetic. Even though many of the ideas are classical in the interval analysis community, we do think that the paper sheds a new light on the topic. Indeed, we systematically investigate several issues which have received limited attention until now:

1. Using ball arithmetic instead of interval arithmetic, for a wide variety of center and radius types.
2. The trade-off between algorithmic complexity and quality of the error bounds.
3. The complexity of multiple precision ball arithmetic.

Throughout the text, we have attempted to provide pointers to existing literature, wherever appropriate. Unfortunately, we did have some difficulties to obtain copies of several seemingly relevant articles and intend to complete the bibliography in the future. It would be of great help if active researchers could systematically publish preprint versions of their papers on their websites.

2. NOTATIONS AND CLASSICAL FACTS

2.1. Floating point numbers and the IEEE 754 norm

We will denote by

$$\mathbb{D} = \mathbb{Z} 2^{\mathbb{Z}} = \{m 2^e : m, e \in \mathbb{Z}\}$$

the set of *dyadic numbers*. Given fixed bit precisions $p \in \mathbb{N}$ and $q \in \mathbb{N}$ for the unsigned mantissa (without the leading 1) and signed exponent, we also denote by

$$\mathbb{D}_{p,q} = \{m 2^e : m, e \in \mathbb{Z}, |m| < 2^{p+1}, 1 - 2^{q-1} < e + p < 2^q - 1\}$$

the corresponding set of floating point numbers. Numbers in $\mathbb{D}_{p,q}$ can be stored in $p + q + 1$ bit space and correspond to “ordinary numbers” in the IEEE 754 norm [ANS08]. For double precision numbers, we have $p = 51$ and $q = 12$. For multiple precision numbers, we usually take q to be the size of a machine word, say $q = 64$, and denote $\mathbb{D}_p = \mathbb{D}_{p,64}$.

The IEEE 754 norm also defines several special numbers. The most important ones are the infinities $\pm\infty$ and “not a number” NaN, which corresponds to the result of an invalid operation, such as $\sqrt{-2}$. We will denote

$$\mathbb{D}_{p,q}^* = \mathbb{D}_{p,q} \cup \{-\infty, +\infty, \text{NaN}\}.$$

Less important details concern the specification of signed zeros ± 0 and several possible types of NaNs. For more details, we refer to [ANS08].

The other main feature of the IEEE 754 is that it specifies how to round results of operations which cannot be performed exactly. There are four basic rounding modes \uparrow (upwards), \downarrow (downwards), \downarrow (nearest) and 0 (towards zero). Assume that we have an operation $f: U \rightarrow \mathbb{R}$ with $U \subseteq \mathbb{R}^k$. Given $x \in U \cap \mathbb{D}_{p,q}^k$ and $y = f(x)$, we define $f^\uparrow(x)$ to be the smallest number \tilde{y} in $\mathbb{D}_{p,q}^*$, such that $\tilde{y} \geq y$. More generally, we will use the notation $(f(x) + y)^\uparrow = f^\uparrow(x) +^\uparrow y$ to indicate that all operations are performed by rounding upwards. The other rounding modes are specified in a similar way. One major advantage of IEEE 754 arithmetic is that it completely specifies how basic operations are done, making numerical programs behave exactly in the same way on different architectures. Most current microprocessors implement the IEEE 754 norm for single and double precision numbers. A conforming multiple precision implementation also exists [HLRZ00].

2.2. Classical complexity results

Using Schönhage-Strassen multiplication [SS71], it is classical that the product of two p -bit integers can be computed in time $l(n) = O(n \log n \log \log n)$. A recent algorithm by Fürer [Für07] further improves this bound to $l(n) = O(n \log n 2^{O(\log^* n)})$, where \log^* denotes the iterator of the logarithm (we have $\log^* n = \log^* \log n + 1$). Other algorithms, such as division two $\leq n$ -bit integers, have a similar complexity $O(l(n))$. Given k prime numbers $p_1 < \dots < p_k$ and a number $0 \leq q < p_1 \dots p_k$, the computation of $q \bmod p_i$ for $i = 1, \dots, k$ can be done in time $O(l(n) \log n)$ using binary splitting [GG02, Theorem 10.25], where $n = \log(p_1 \dots p_k)$. It is also possible to reconstruct q from the remainders $q \bmod p_i$ in time $O(l(n) \log n)$.

Let \mathbb{K} be an effective ring, in the sense that we have algorithms for performing the ring operations in \mathbb{K} . If \mathbb{K} admits 2^p -th roots of unity for $2^p \geq n$, then it is classical [CT65] that the product of two polynomials $P, Q \in \mathbb{K}[z]$ with $\deg PQ < n$ can be computed using $O(n \log n)$ ring operations in \mathbb{K} , using the fast Fourier transform. For general rings, this product can be computed using $M(n) = O(n \log n \log \log n)$ operations [CK91], using a variant of Schönhage-Strassen multiplication.

A formal power series $f \in \mathbb{K}[[z]]$ is said to be *computable*, if there exists an algorithm which takes $n \in \mathbb{N}$ on input and which computes $f_n \in \mathbb{K}$. We denote by $\mathbb{K}[[z]]^{\text{com}}$ the set of computable power series. Many simple operations on formal power series, such as multiplication, division, exponentiation, etc., as well as the resolution of algebraic differential equations can be done up till order n using a similar time complexity $O(M(n))$ as polynomial multiplication [BK78, BCO+06, vdH06b]. Alternative *relaxed* multiplication algorithms of time complexity $R(n) = O(M(n))$ are given in [vdH02a, vdH07a]. In this case, the coefficients $(fg)_n$ are computed gradually and $(fg)_n$ is output as soon as f_0, \dots, f_n and g_0, \dots, g_n are known. This strategy is often most efficient for the resolution of implicit equations.

Let $\mathbb{K}^{m \times n}$ denote the set of $m \times n$ matrices with entries in \mathbb{K} . Given two $n \times n$ matrices $M, N \in \mathbb{K}^{n \times n}$, the naive algorithm for computing MN requires $O(n^\omega)$ ring operations with $\omega = 3$. The exponent ω has been reduced by Strassen [Str69] to $\omega = \log 7 / \log 2$. Using more sophisticated techniques, ω can be further reduced to $\omega < 2.376$ [Pan84, WC87]. However, this exponent has never been observed in practice.

2.3. Interval and ball notations

Given $x < y$ in a totally ordered set R , we will denote by $[x, y]$ the closed interval

$$[x, y] = \{z \in R : x \leq z \leq y\}. \quad (1)$$

Assume now that R is a totally ordered field and consider a normed vector space V over R . Given $c \in V$ and $r \in R^{\geq} = \{r \in R : r \geq 0\}$, we will denote by

$$\mathcal{B}(c, r) = \{x \in V : \|x - c\| \leq r\} \quad (2)$$

the closed ball with center c and radius r . Notice that we will never work with open balls in what follows. We denote the set of all balls of the form (2) by $\mathcal{B}(V, R)$. Given $x \in \mathcal{B}(V, R)$, we will denote by x_c and x_r the center resp. radius of x .

We will use the standard euclidean norm

$$\|x\| = \sqrt{\frac{|x_1|^2 + \dots + |x_n|^2}{n}}$$

as the default norm on \mathbb{R}^n and \mathbb{C}^n . Occasionally, we will also consider the max-norm

$$\|x\|_\infty = \max\{|x_1|, \dots, |x_n|\}.$$

For matrices $M \in \mathbb{R}^{m \times n}$ or $M \in \mathbb{C}^{m \times n}$, the default operator norm is defined by

$$\|M\| = \max_{\|x\|=1} \|Mx\|. \quad (3)$$

Occasionally, we will also consider the max-norm

$$\|M\|_\infty = \max_{i,j} |M_{i,j}|,$$

which satisfies

$$\|M\|_\infty \leq \|M\| \leq n \|M\|_\infty.$$

We also define the max-norm for polynomials $P \in \mathbb{R}[x]$ or $P \in \mathbb{C}[x]$ by

$$\|P\|_\infty = \max |P_i|.$$

For power series $f \in \mathbb{R}[[z]]$ or $f \in \mathbb{C}[[z]]$ which converge on the compact disk $\mathcal{B}(0, r)$, we define the norm

$$\|f\|_r = \max_{|z| \leq r} |f(z)|. \quad (4)$$

After rescaling $f(z) \mapsto f(rz)$, we will usually work with the norm $\|\cdot\| = \|\cdot\|_1$.

In some cases, it may be useful to generalize the concept of a ball and allow for radii in ordered rings. For instance, a ball $X \in \mathcal{B}(\mathbb{R}^n, \mathbb{R}^n)$ stands for the set

$$X = \{x \in \mathbb{R}^n \mid \forall 1 \leq i \leq n, |x_i - (X_c)_i| \leq (X_r)_i\}.$$

The sets $\mathcal{B}(\mathbb{R}^{m \times n}, \mathbb{R}^{m \times n})$, $\mathcal{B}(\mathbb{R}[x], \mathbb{R}[x])$, $\mathcal{B}(\mathbb{R}[[x]], \mathbb{R}[[x]])$, etc. are defined in a similar component wise way. Given $x \in \mathbb{R}^n$, it will also be convenient to denote by $|x|$ the vector with $|x|_i = |x_i|$. For matrices $M \in \mathbb{R}^{m \times n}$, polynomials $P \in \mathbb{R}[x]$ and power series $f \in \mathbb{R}[[x]]$, we define $|M|$, $|P|$ and $|f|$ similarly.

3. BALLS AND COMPUTABLE REAL NUMBERS

3.1. Ball arithmetic

Let \mathbb{V} be a normed vector space over \mathbb{R} and recall that $\mathcal{B}(\mathbb{V}, \mathbb{R})$ stands for the set of all balls with centers in \mathbb{V} and radii in \mathbb{R} . Given an operation $\varphi: \mathbb{V}^k \rightarrow \mathbb{V}$, the operation is said to lift to an operation

$$f^\bullet: \mathcal{B}(\mathbb{V}, \mathbb{R})^k \rightarrow \mathcal{B}(\mathbb{V}, \mathbb{R}),$$

if we have

$$f^\bullet(X_1, \dots, X_k) \supseteq \{f(x_1, \dots, x_k) : x_1 \in X_1, \dots, x_k \in X_k\},$$

for all $X_1, \dots, X_k \in \mathcal{B}(\mathbb{V}, \mathbb{R})$. For instance, both the addition $+: \mathbb{V}^2 \rightarrow \mathbb{V}$ and subtraction $-: \mathbb{V}^2 \rightarrow \mathbb{V}$ admits lifts

$$\mathcal{B}(x, r) +^\bullet \mathcal{B}(y, s) = \mathcal{B}(x + y, r + s) \quad (5)$$

$$\mathcal{B}(x, r) -^\bullet \mathcal{B}(y, s) = \mathcal{B}(x - y, r + s). \quad (6)$$

Similarly, if \mathbb{V} is a normed algebra, then the multiplication lifts to

$$\mathcal{B}(x, r) \cdot^\bullet \mathcal{B}(y, s) = \mathcal{B}(xy, (|x| + r)s + r|y|) \quad (7)$$

The lifts $+^\bullet$, $-^\bullet$, \cdot^\bullet , etc. are also said to be the ball arithmetic counterparts of addition, subtractions, multiplication, etc.

Ball arithmetic is a systematic way for the computation of error bounds when the input of a numerical operation is known only approximately. These bounds are usually not sharp. For instance, consider the mathematical function $f: x \in \mathbb{R} \mapsto x - x$ which evaluates to zero everywhere, even if x is only approximately known. However, taking $x = \mathcal{B}(2, 0.001)$, we have $x -^\bullet x = \mathcal{B}(0, 0.002) \neq \mathcal{B}(0, 0)$. This phenomenon is known as *overestimation*. In general, ball algorithms have to be designed carefully so as to limit overestimation.

In the above definitions, \mathbb{R} can be replaced by a subfield $R \subseteq \mathbb{R}$, \mathbb{V} by an R -vector space $V \subseteq \mathbb{V}$, and the domain of f by an open subset of \mathbb{V}^k . If V and R are *effective* in the sense that we have algorithms for the additions, subtractions and multiplications in V and R , then basic ball arithmetic in $\mathcal{B}(V, R)$ is again effective. If we are working with finite precision floating point numbers in $\mathbb{D}_{p,q}$ rather than a genuine effective subfield R , we will now show how to adapt the formulas (5), (6) and (7) in order to take into account rounding errors; it may also be necessary to allow for an infinite radius in this case.

3.2. Finite precision ball arithmetic

Let us detail what needs to be changed when using IEEE conform finite precision arithmetic, say $V = R = \mathbb{D}_{p,q}$. We will denote

$$\begin{aligned} \mathbb{B} &= \mathcal{B}(\mathbb{D}, \mathbb{D}) \\ \mathbb{B}_{p,q} &= \mathcal{B}(\mathbb{D}_{p,q}, \mathbb{D}_{p,q}) \\ \mathbb{B}[i]_{p,q} &= \mathcal{B}(\mathbb{D}_{p,q}[i], \mathbb{D}_{p,q}). \end{aligned}$$

When working with multiple precision numbers, it usually suffices to use low precision numbers for the radius type. Recalling that $\mathbb{D}_p = \mathbb{D}_{p,64}$, we will therefore denote

$$\begin{aligned}\mathbb{B}_p &= \mathcal{B}(\mathbb{D}_p, \mathbb{D}_{64}) \\ \mathbb{B}[i]_p &= \mathcal{B}(\mathbb{D}_p[i], \mathbb{D}_{64}).\end{aligned}$$

We will write $\epsilon = \epsilon_p = 2^{-p}$ for the machine accuracy and $\eta = \eta_{p,q} = 2^{2-2^q-p}$ for the smallest representable positive number in $\mathbb{D}_{p,q}$.

Given an operation $f: \mathbb{R}^k \rightarrow \mathbb{R}$ as in section 3.1, together with balls $X_i = \mathcal{B}(x_i, r_i)$, it is natural to compute the center y of

$$\mathcal{B}(y, s) = f^\bullet(X_1, \dots, X_s)$$

by rounding to the nearest:

$$y = f^\downarrow(x_1, \dots, x_k). \tag{8}$$

One interesting point is that the committed error

$$\delta = |y - f(x_1, \dots, x_k)|$$

does not really depend on the operation f itself: we have the universal upper bound

$$\begin{aligned}\delta &\leq \Delta(y) \\ \Delta(y) &:= (|y| + \uparrow \eta) \cdot \uparrow \epsilon.\end{aligned} \tag{9}$$

It would be useful if this *adjustment function* Δ were present in the hardware.

For the computation of the radius s , it now suffices to use the sum of $\Delta(y)$ and the theoretical bound formulas for the infinite precision case. For instance,

$$\mathcal{B}(x, r) + \bullet \mathcal{B}(y, s) = \mathcal{B}^\Delta(x + \uparrow y, r + \uparrow s) \tag{10}$$

$$\mathcal{B}(x, r) - \bullet \mathcal{B}(y, s) = \mathcal{B}^\Delta(x - \uparrow y, r + \uparrow s) \tag{11}$$

$$\mathcal{B}(x, r) \cdot \bullet \mathcal{B}(y, s) = \mathcal{B}^\Delta(x \cdot \uparrow y, [(|x| + r) s + r |y|]^\uparrow), \tag{12}$$

where \mathcal{B}^Δ stands for the ‘‘adjusted constructor’’

$$\mathcal{B}^\Delta(y, s) = \mathcal{B}(y, s + \uparrow \Delta(y)).$$

The approach readily generalizes to other ‘‘normed vector spaces’’ \mathbb{V} over $\mathbb{D}_{p,q}$, as soon as one has a suitable rounded arithmetic in \mathbb{V} and a suitable adjustment function Δ attached to it.

Notice that $\Delta(y) = \infty$, if $y = \infty$ or y is the largest representable real number in $\mathbb{D}_{p,q}$. Consequently, the finite precision ball computations naturally take place in domains of the form $\mathbb{B}_{p,q}^* = \mathcal{B}(\mathbb{D}_{p,q}^*, \mathbb{D}_{p,q}^*)$ rather than $\mathbb{B}_{p,q}$. Of course, balls with infinite radius carry no useful information about the result. In order to ease the reading, we will assume the absence of overflows in what follows, and concentrate on computations with ordinary numbers in $\mathbb{B}_{p,q}$. We will only consider infinities if they are used in an essential way during the computation.

Similarly, if we want ball arithmetic to be a natural extension of the IEEE 754 norm, then we need an equivalent of NaN. One approach consists of introducing a NaB (not a ball) object, which could be represented by $\mathcal{B}(\text{NaN}, \text{NaN})$ or $\mathcal{B}(0, -\infty)$. A ball function f^\bullet returns NaB if f returns NaN for one selection of members of the input balls. For instance, $\text{sqrt}^\bullet(\mathcal{B}(1, 3)) = \text{NaB}$. An alternative approach consists of the attachment of an additional flag to each ball object, which signals a possible invalid outcome. Following this convention, $\text{sqrt}^\bullet(\mathcal{B}(1, 3))$ yields $\mathcal{B}(1, 1)^{\text{NaN}}$.

3.3. Implementation details

Using the formulas from the previous section, it is relatively straightforward to implement ball arithmetic as a C++ template library, as we have done in the MATHEMAGIX system. However, in the case of multiple precision arithmetic this is far from optimal. Let us discuss several possible optimizations:

1. Multiple precision libraries such as MPFR [HLRZ00] suffer from a huge overhead when it comes to moderate (e.g. quadruple) precision computations. Since the radii are always stored in low precision, it is recommended to inline all computations on the radii. In the case of multiplication, this divides the number of function calls by four.
2. When computing with complex numbers $z \in \mathbb{D}_p[i]$, one may again save several function calls. Moreover, it is possible to regard z as an element of $\mathbb{Z}[i] 2^{\mathbb{Z}}$ rather than $(\mathbb{Z} 2^{\mathbb{Z}})[i]$, i.e. use a single exponent for both the real and imaginary parts of z . This optimization reduces the time spent on exponent computations and mantissa normalizations.
3. Consider a ball $\mathcal{B}(c, r) \in \mathbb{B}_p$ and recall that $c \in \mathbb{D}_p$, $r \in \mathbb{D}_{64}$. If $|c| < 2^{p-64} r$, then the $\lfloor \log_2 r + p - \log_2 |c| - 64 \rfloor$ least significant binary digits of c are of little interest. Hence, we may replace c by its closest approximation in $\mathbb{D}_{p'}$, with $p' = \lceil \log_2 |c| + 64 - \log_2 r \rceil$, and reduce the working precision to p' . Modulo slightly more work, it is also possible to share the exponents of the center and the radius.
4. If we don't need large exponents for our multiple precision numbers, then it is possible to use machine doubles $\mathbb{D}_{51,12}$ as our radius type and further reduce the overhead of bound computations.

When combining the above optimizations, it can be hoped that multiple precision ball arithmetic can be implemented almost as efficiently as standard multiple precision arithmetic. However, this requires a significantly higher implementation effort.

3.4. Computable numbers

Given $x \in \mathbb{R}$, $\tilde{x} \in \mathbb{D}$ and $\varepsilon \in \mathbb{D}^> = \{\varepsilon \in \mathbb{D} : \varepsilon > 0\}$, we say that \tilde{x} is an ε -*approximation* of x if $|\tilde{x} - x| \leq \varepsilon$. A real number $x \in \mathbb{R}$ is said to be *computable* [Tur36, Grz57, Wei00] if there exists an *approximation algorithm* which takes an absolute tolerance $\varepsilon \in \mathbb{D}^>$ on input and which returns an ε -approximation of x . We denote by \mathbb{R}^{com} the set of computable real numbers. We may regard \mathbb{R}^{com} as a data type, whose instances are represented by approximation algorithms.

In practice, it is more convenient to work with so called ball approximation algorithms: a real number is computable if and only if it admits a *ball approximation algorithm*, which takes a working precision $p \in \mathbb{N}$ on input and returns a *ball approximation* $\mathcal{B}(c_p, r_p) \in \mathbb{B}$ with $x \in \mathcal{B}(c_p, r_p)$ and $\lim_{p \rightarrow \infty} r_p = 0$. Indeed, assume that we have a ball approximation algorithm. In order to obtain an ε -approximation, it suffices to compute ball approximations at precisions $p = 32, 64, 128, \dots$ which double at every step, until $r_p \leq \varepsilon$. Conversely, given an approximation algorithm $\tilde{x}: \mathbb{D}^> \rightarrow \mathbb{D}$ of $x \in \mathbb{R}^{\text{com}}$, we obtain a ball approximation algorithm $p \mapsto \mathcal{B}(c_p, r_p)$ by taking $r_p = 2^{-p}$ and $c_p = \tilde{x}(r_p)$.

Given $x, y \in \mathbb{R}^{\text{com}}$ with ball approximation algorithms $\tilde{x}, \tilde{y}: \mathbb{N} \mapsto \mathbb{B}$, we may compute ball approximation algorithms for $x + y, x - y, xy$ simply by taking

$$\begin{aligned} (\widetilde{x + y})(p) &= \tilde{x}(p) + \bullet \tilde{y}(p) \\ (\widetilde{x - y})(p) &= \tilde{x}(p) - \bullet \tilde{y}(p) \\ (\widetilde{xy})(p) &= \tilde{x}(p) \cdot \bullet \tilde{y}(p). \end{aligned}$$

More generally, assuming a good library for ball arithmetic, it is usually easy to write a wrapper library with the corresponding operations on computable numbers.

From the efficiency point of view, it is also convenient to work with ball approximations. Usually, the radius r_p satisfies

$$\log_2 r_p = -p + o(1),$$

or at least

$$\log_2 r_p = -cp + o(1),$$

for some $c \in \mathbb{Q}^>$. In that case, doubling the working precision p until a sufficiently good approximation is found is quite efficient. An even better strategy is to double the “expected running time” at every step [vdH06a, KR06]. Yet another approach will be described in section 3.6 below.

The concept of computable real numbers readily generalizes to more general normed vector spaces. Let \mathbb{V} be a normed vector space and let \mathbb{V}^{dig} be an effective subset of *digital points* in \mathbb{V} , i.e. the elements of \mathbb{V}^{dig} admit a computer representation. For instance, if $\mathbb{V} = \mathbb{R}$, then we take $\mathbb{V}^{\text{dig}} = \mathbb{D}$. Similarly, if $\mathbb{V} = \mathbb{R}^{m \times n}$ is the set of real $m \times n$ matrices, with one of the matrix norms from section 2.3, then it is natural to take $\mathbb{V}^{\text{dig}} = (\mathbb{R}^{\text{dig}})^{m \times n} = \mathbb{D}^{m \times n}$. A point $x \in \mathbb{V}$ is said to be *computable*, if it admits an *approximation algorithm* which takes $\varepsilon \in \mathbb{D}^>$ on input, and returns an ε -*approximation* $\tilde{x} \in \mathbb{V}^{\text{dig}}$ of x (satisfying $\|\tilde{x} - x\| \leq \varepsilon$, as above).

3.5. Asymmetric computability

A real number x is said to be *left computable* if there exists an algorithm for computing an increasing sequence $\check{x}: \mathbb{N} \mapsto \mathbb{D}; n \mapsto \check{x}_n$ with $\lim_{n \rightarrow \infty} \check{x}_n = x$ (and \check{x} is called a left approximation algorithm). Similarly, x is said to be *right computable* if $-x$ is left computable. A real number is computable if and only if it is both left and right computable. Left computable but non computable numbers occur frequently in practice and correspond to “computable lower bounds” (see also [Wei00, vdH07b]).

We will denote by \mathbb{R}^{lcom} and \mathbb{R}^{rcom} the data types of left and right computable real numbers. It is convenient to specify and implement algorithms in computable analysis in terms of these data types, whenever appropriate [vdH07b]. For instance, we have computable functions

$$\begin{aligned} +: \mathbb{R}^{\text{lcom}} \times \mathbb{R}^{\text{lcom}} &\rightarrow \mathbb{R}^{\text{lcom}} \\ -: \mathbb{R}^{\text{lcom}} \times \mathbb{R}^{\text{rcom}} &\rightarrow \mathbb{R}^{\text{lcom}} \\ &\vdots \end{aligned}$$

More generally, given a subset $S \subseteq \mathbb{R}$, we say that $x \in S$ is left computable in S if there exists a left approximation algorithm $\check{x}: \mathbb{N} \rightarrow S$ for x . We will denote by S^{lcom} and S^{rcom} the data types of left and computable numbers in S , and define $S^{\text{com}} = S^{\text{lcom}} \cap S^{\text{rcom}}$.

Identifying the type of boolean numbers \mathbb{T} with $\{0, 1\}$, we have $\mathbb{T}^{\text{lcom}} = \mathbb{T}^{\text{rcom}} = \mathbb{T}$ as sets, but *not* as data types. For instance, it is well known that equality is non computable for computable real numbers [Tur36]. Nevertheless, equality *is* “ultimately computable” in the sense that there exists a computable function

$$=: \mathbb{R}^{\text{com}} \times \mathbb{R}^{\text{com}} \rightarrow \mathbb{T}^{\text{rcom}}.$$

Indeed, given $x, y \in \mathbb{R}^{\text{com}}$ with ball approximation algorithms \check{x} and \check{y} , we may take

$$(x \overset{\approx}{=} y)_n = \begin{cases} 1 & \text{if } (\check{x}_0 \cap \check{y}_0 \cap \dots \cap \check{x}_n \cap \check{y}_n) \neq \emptyset \\ 0 & \text{otherwise} \end{cases}$$

Similarly, the ordering relation \leq is ultimately computable.

This asymmetric point of view on equality testing also suggest a semantics for the relations $=$, \leq , etc. on balls. For instance, given balls $x, y \in \mathbb{B}$, it is natural to take

$$\begin{aligned} x = y &\iff x \cap y \neq \emptyset \\ x \neq y &\iff x \cap y = \emptyset \\ x \leq y &\iff \exists a \in x, b \in y, a \leq b \\ &\vdots \end{aligned}$$

These definitions are interesting if balls are really used as successive approximations of a real number. An alternative application of ball arithmetic is for modeling non-deterministic computations: the ball models a set of possible values and we are interested in the set of possible outcomes of an algorithm. In that case, the natural return type of a relation on balls becomes a “boolean ball”. In the area of interval analysis, this second interpretation is more common [ANS09].

3.6. Lipschitz ball arithmetic

Given a computable function $f: \mathbb{R}^{\text{com}} \rightarrow \mathbb{R}^{\text{com}}$, $x \in \mathbb{R}^{\text{com}}$ and $\varepsilon \in \mathbb{D}^>$, let us return to the problem of efficiently computing an approximation $\tilde{y} \in \mathbb{D}$ of $y = f(x)$ with $|\tilde{y} - y| < \varepsilon$. In section 3.4, we suggested to compute ball approximations of y at precisions which double at every step, until a sufficiently precise approximation is found. This computation involves an implementation $f^\bullet: \mathbb{B} \rightarrow \mathbb{B}$ of f on the level of balls, which satisfies

$$f^\bullet(X) \supseteq \{f(x) : x \in X\}, \quad (13)$$

for every ball $X \in \mathbb{B}$. In practice, f is often differentiable, with $f'(x) \neq 0$. In that case, given a ball approximation X of x , the computed ball approximation $Y = f^\bullet(X)$ of y typically has a radius

$$Y_r \sim |f'(x)| X_r, \quad (14)$$

for $X_r \rightarrow 0$. This should make it possible to directly predict a sufficient precision at which $Y_r \leq \varepsilon$. The problem is that (14) needs to be replaced by a more reliable relation. This can be done on the level of ball arithmetic itself, by replacing the usual condition (13) by

$$\begin{aligned} f^\bullet(X) &\supseteq \mathcal{B}(f(x_c), M X_r) \\ M &= \sup_{x \in X} |f'(x)|. \end{aligned} \quad (15)$$

Similarly, the multiplication of balls is carried out using

$$\mathcal{B}(x, r) \cdot \mathcal{B}(y, s) = \mathcal{B}(xy, (|x| + r)s + (|y| + s)r). \quad (16)$$

instead of (7). A variant of this kind of “Lipschitz ball arithmetic” has been implemented in [M0]. Although a constant factor is gained for high precision computations at regular points x , the efficiency deteriorates near singularities (i.e. the computation of $\sqrt{0}$).

4. BALLS VERSUS INTERVALS

In the area of reliable computation, interval arithmetic has for long been privileged with respect to ball arithmetic. Indeed, balls are often regarded as a more or less exotic variant of intervals, based on an alternative midpoint-radius representation. Historically, interval arithmetic is also preferred in computer science because it is easy to implement if floating point operations are performed with correct rounding. Since most modern microprocessors implement the IEEE 754 norm, this point of view is well supported by hardware.

Not less historically, the situation in mathematics is inverse: whereas intervals are the standard in computer science, balls are the standard in mathematics, since they correspond to the traditional ε - δ -calculus. Even in the area of interval analysis, one usually resorts (at least implicitly) to balls for more complex computations, such as the inversion of a matrix [HS67, Moo66]. Indeed, balls are more convenient when computing error bounds using perturbation techniques.

In this section, we will give a more detailed account on the respective advantages and disadvantages of interval and ball arithmetic.

4.1. Standardization

One advantage of interval arithmetic is that the IEEE 754 norm suggests a natural and standard implementation. Indeed, let f be a real function which is increasing on some interval I . Then the natural interval lift f^\sqcup of f is given by

$$f^\sqcup([l, h]) = [f^\downarrow(l), f^\uparrow(h)].$$

This implementation has the property that $f^\sqcup([l, h])$ is the smallest interval with end-points in $\mathbb{D}_{p,q} \cup \{\pm\infty\}$, which satisfies

$$f^\sqcup([l, h]) \supseteq \{f(x) : x \in [l, h]\}.$$

For not necessarily increasing functions f this property can still be used as a requirement for the “standard” implementation of f^\sqcup . For instance, this leads to the following implementation of the cosine function on intervals:

$$\cos^\sqcup([l, h]) = \begin{cases} [\cos^\downarrow l, \cos^\uparrow h] & \text{if } \lfloor l/2\pi \rfloor = \lfloor h/2\pi \rfloor \in 2\mathbb{Z} - 1 \\ [\cos^\downarrow h, \cos^\uparrow l] & \text{if } \lfloor l/2\pi \rfloor = \lfloor h/2\pi \rfloor \in 2\mathbb{Z} \\ [\min(\cos^\downarrow l, \cos^\downarrow h), 1] & \text{if } \lfloor l/2\pi \rfloor = \lfloor h/2\pi \rfloor - 1 \in 2\mathbb{Z} - 1 \\ [-1, \max(\cos^\uparrow l, \cos^\uparrow h)] & \text{if } \lfloor l/2\pi \rfloor = \lfloor h/2\pi \rfloor - 1 \in 2\mathbb{Z} \\ [-1, 1] & \text{if } \lfloor l/2\pi \rfloor < \lfloor h/2\pi \rfloor - 1 \end{cases}$$

Such a standard implementation of interval arithmetic has the convenient property that programs will execute in the same way on any platform which conforms to the IEEE 754 standard.

By analogy with the above approach for standardized interval arithmetic, we may standardize the ball implementation f^\bullet of f by taking

$$f^\bullet(\mathcal{B}(c, r)) = \mathcal{B}(f^\uparrow(c), s),$$

where the radius s is smallest in $\mathbb{D}_{p,q} \cup \{+\infty\}$ with

$$\mathcal{B}(f^\uparrow(c), s) \supseteq \{f(x) : x \in \mathcal{B}(c, r)\}.$$

Unfortunately, the computation of such an optimal s is not always straightforward. In particular, the formulas (10), (11) and (12) do not necessarily realize this tightest bound. In practice, it might therefore be better to achieve standardization by fixing once and for all the formulas by which ball operations are performed. Of course, more experience with ball arithmetic is required before this can happen.

4.2. Efficiency

The respective efficiencies of interval and ball arithmetic depend on the precision at which we are computing. For high precisions and most applications, ball arithmetic has the advantage that we can still perform computations on the radius at single precision. By contrast, interval arithmetic requires full precision for operations on both end-points. This makes ball arithmetic twice as efficient at high precisions.

When working at machine precision, the efficiencies of both approaches essentially depend on the hardware. *A priori*, interval arithmetic is better supported by current computers, since most of them respect the IEEE 754 norm, whereas the function Δ from (9) usually has to be implemented by hand. However, changing the rounding mode is often highly expensive (over hundred cycles). Therefore, additional gymnastics may be required in order to always work with respect to a fixed rounding mode. For instance, if \uparrow is our current rounding mode, then we may take

$$x + \downarrow y = -((-x) + \uparrow(-y)),$$

since the operation $x \mapsto -x$ is always exact (i.e. does not depend on the rounding mode). As a consequence, interval arithmetic becomes slightly more expensive. By contrast, when releasing the condition that centers of balls are computed using rounding to the nearest, we may replace (8) by

$$y = f^\uparrow(x_1, \dots, x_k) \tag{17}$$

and (9) by

$$\Delta(y) := (|y| + \uparrow \eta) \cdot \uparrow (2\epsilon).$$

Hence, ball arithmetic already allows us to work with respect to a fixed rounding mode. Of course, using (17) instead of (8) does require to rethink the way ball arithmetic should be standardized.

An alternative technique for avoiding changes in rounding mode exists when performing operations on compound types, such as vectors or matrices. For instance, when adding two vectors, we may first add all lower bounds while rounding downwards and next add the upper bounds while rounding upwards. Unfortunately, this strategy becomes more problematic in the case of multiplication, because different rounding modes may be needed depending on the signs of the multiplicands. As a consequence, matrix operations tend to require many conditional parts of code when using interval arithmetic, with increased probability of breaking the processor pipeline. On the contrary, ball arithmetic highly benefits from parallel architecture and it is easy to implement ball arithmetic for matrices on top of existing libraries: see [Rum99a] and section 6 below.

4.3. Quality

Besides the efficiency of ball and interval arithmetic for basic operations, it is also important to investigate the quality of the resulting bounds. Indeed, there are usually differences between the sets which are representable by balls and by intervals. For instance, when using the extended IEEE 754 arithmetic with infinities, then it is possible to represent $[1, \infty]$ as an interval, but not as a ball.

These differences get more important when dealing with complex numbers or compound types, such as matrices. For instance, when using interval arithmetic for reliable computations with complex numbers, it is natural to enclose complex numbers by rectangles $X + Yi$, where X and Y are intervals. For instance, the complex number $z = 1 + i$ may be enclosed by

$$z \in [1 - \epsilon, 1 + \epsilon] + [1 - \epsilon, 1 + \epsilon]i,$$

for some small number ϵ . When using ball arithmetic, we would rather enclose z by

$$z \in \mathcal{B}(1 + i, \epsilon).$$

Now consider the computation of $u = z^2$. The computed rectangular and ball enclosures are given by

$$\begin{aligned} u &\in [-2\varepsilon, 2\varepsilon] + [2 - 2\varepsilon, 2 + 2\varepsilon]i + o(\varepsilon) \\ u &\in \mathcal{B}(2, \sqrt{2}\varepsilon) + o(\varepsilon). \end{aligned}$$

Consequently, ball arithmetic yields a much better bound, which is due to the fact that multiplication by $1 + i$ turns the rectangular enclosure by 45 degrees, leading to an overestimation by a factor $\sqrt{2}$ when re-enclosing the result by a horizontal rectangle (see figure 1).

This is one of the simplest instances of the *wrapping effect* [Moo66]. For similar reasons, one may prefer to compute the square of the matrix

$$M = \begin{pmatrix} 0 & 1 \\ 1 & 0 \end{pmatrix} \tag{18}$$

in $\mathcal{B}(\mathbb{D}^{2 \times 2}, \mathbb{D})$ rather than $\mathcal{B}(\mathbb{D}, \mathbb{D})^{2 \times 2}$, while using the operator norm (3) for matrices. This technique highlights another advantage of ball arithmetic: we have a certain amount of flexibility regarding the choice of the radius type. By choosing a simple radius type, we do not only reduce the wrapping effect, but also improve the efficiency: when computing with complex balls in $\mathbb{B}[i]$, we only need to bound one radius instead of two for every basic operation. More precisely, we replace (8) and (9) by

$$\begin{aligned} y &= f^\uparrow(x_1, \dots, x_k) \\ &= (\operatorname{Re} f)^\uparrow(x_1, \dots, x_k) + (\operatorname{Im} f)^\uparrow(x_1, \dots, x_k) i \\ \Delta(y) &:= (\operatorname{abs}^\uparrow(y) + \uparrow \eta) \cdot \uparrow (\sqrt{2} \varepsilon). \end{aligned}$$

On the negative side, generalized norms may be harder to compute, even though a rough bound often suffices (e.g. replacing $\operatorname{abs}^\uparrow(y)$ by $|\operatorname{Re} y| + \uparrow |\operatorname{Im} y|$ in the above formula). In the case of matricial balls, a more serious problem concerns overestimation when the matrix contains entries of different orders of magnitude. Another more algorithmic technique for reducing the wrapping effect will be discussed in sections 6.4 and 7.5 below.

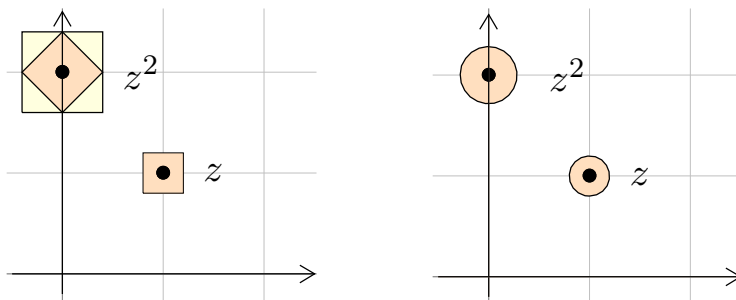


Figure 1. Illustration of the computation of z^2 using interval and ball arithmetic, for $z = 1 + i$.

4.4. Mathematical elegance

Even though personal taste in the choice between balls and intervals cannot be discussed, the elegance of the chosen approach for a particular application can partially be measured in terms of the human time needed which is needed to establish the necessary error bounds.

We have already seen that interval enclosures are particularly easy to obtain for monotonic real functions. Another typical algorithm where interval arithmetic is more convenient is the resolution of a system of equations using dichotomy. Indeed, it is easier to cut an n -dimensional block $[a_1, b_1] \times \dots \times [a_n, b_n]$ into 2^n smaller blocks than to perform a similar operation on balls.

For most other applications, ball representations are more convenient. Indeed, error bounds are usually obtained by perturbation methods. For any mathematical proof where error bounds are explicitly computed in this way, it is generally easy to derive a certified algorithm based on ball arithmetic. We will see several illustrations of this principle in the sections below.

Implementations of interval arithmetic often rely on floating point arithmetic with correct rounding. One may question how good correct rounding actually is in order to achieve reliability. One major benefit is that it provides a simple and elegant way to specify what a mathematical function precisely does at limited precision. In particular, it allows numerical programs to execute exactly in the same way on many different hardware architectures.

On the other hand, correct rounding does have a certain cost. Although the cost is limited for field operations and elementary functions [Mul06], the cost increases for more complex special functions, especially if one seeks for numerical methods with a constant operation count. For arbitrary computable functions on \mathbb{R}^{com} , correct rounding even becomes impossible. Another disadvantage is that correct rounding is lost as soon we perform more than one operation: in general, $g^\uparrow \circ f^\uparrow$ and $(g \circ f)^\uparrow$ do not coincide.

In the case of ball arithmetic, we only require an upper bound for the error, not necessarily the best possible representable one. In principle, this is just as reliable and usually more economic. Now in an ideal world, the development of numerical codes goes hand in hand with the systematic development of routines which compute the corresponding error bounds. In such a world, correct rounding becomes superfluous, since correctness is no longer ensured at the micro-level of hardware available functions, but rather at the top-level, *via* mathematical proof.

5. THE NUMERICAL HIERARCHY

Computable analysis provides a good high-level framework for the automatic and certified resolution of analytic problems. The user states the problem in a formal language and specifies a required absolute or relative precision. The program should return a numerical result which is certified to meet the requirement on the precision. A simple example is to compute an ε -approximation for π , for a given $\varepsilon \in \mathbb{D}^>$.

The MATHEMAGIX system [vdH+02b] aims the implementation of efficient algorithms for the certified resolution of numerical problems. Our ultimate goal is that these algorithms become as efficient as more classical numerical methods, which are usually non certified and only operate at limited precision. A naive approach is to systematically work with computable real numbers. Although this approach is convenient for theoretical purposes in the area of computable analysis, the computation with functions instead of ordinary floating point numbers is highly inefficient.

In order to address this efficiency problem, the MATHEMAGIX libraries for basic arithmetic on analytic objects (real numbers, matrices, polynomials, etc.) are subdivided into four layers of the so called *numerical hierarchy* (see figure 2). We will illustrate this decomposition on the problem of multiplying two $n \times n$ computable real matrices. The numerical hierarchy turns out to be a convenient framework for more complex problems as well, such as the analytic continuation of the solution to a dynamical system. As a matter of fact, the framework incites the developer to restate the original problem at the different levels, which is generally a good starting point for designing an efficient solution.

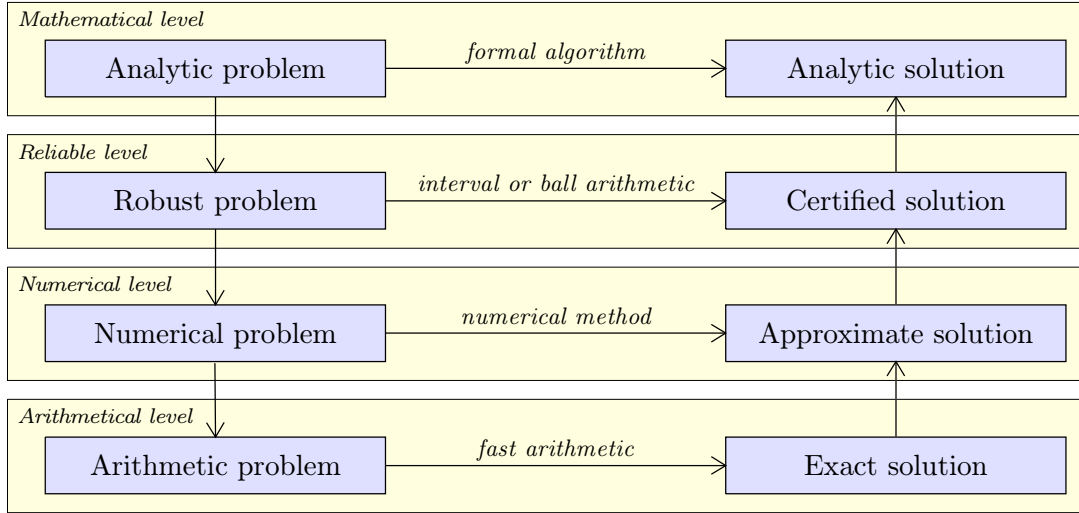


Figure 2. The numerical hierarchy.

Mathematical level. On the mathematical top level, we are given two computable real $n \times n$ matrices $A, B \in (\mathbb{R}^{\text{com}})^{n \times n}$ and an absolute error $\varepsilon \in \mathbb{D}^>$. The aim is to compute ε -approximations for all entries of the product $C = AB$.

The simplest approach to this problem is to use a generic formal matrix multiplication algorithm, using the fact that \mathbb{R}^{com} is an effective ring. However, as stressed above, instances of \mathbb{R}^{com} are really functions, so that ring operations in \mathbb{R}^{com} are quite expensive. Instead, when working at precision p , we may first compute ball approximations for all the entries of A and B , after which we form two approximation matrices $\tilde{A}, \tilde{B} \in \mathbb{B}^{n \times n}$. The multiplication problem then reduces to the problem of multiplying two matrices in $\mathbb{B}^{n \times n}$. This approach has the advantage that $O(n^3)$ operations on “functions” in \mathbb{R}^{com} are replaced by a single multiplication in $\mathbb{B}^{n \times n}$.

Reliable level. The aim of this layer is to implement efficient algorithms on balls. Whereas the actual numerical computation is delegated to the numerical level below, the reliable level should be able to perform the corresponding error analysis automatically.

When operating on non-scalar objects, such as matrices of balls, it is often efficient to rewrite the objects first. For instance, when working in fixed point arithmetic (i.e. all entries of A and B admit similar orders of magnitude), a matrix in $\mathbb{B}^{n \times n}$ may also be considered as a *matricial ball* in $\mathcal{B}(\mathbb{D}^{n \times n}, \mathbb{D})$ for the matrix norm $\|\cdot\|_\infty$. We multiply two such balls using the formula

$$\mathcal{B}(x, r) \cdot \mathcal{B}(y, s) = \mathcal{B}(xy, [n((\|x\|_\infty + r)s + r\|y\|_\infty)]^\uparrow),$$

which is a corrected version of (7), taking into account that the matrix norm $\|\cdot\|_\infty$ only satisfies $\|xy\|_\infty \leq n\|x\|_\infty\|y\|_\infty$. Whereas the naive multiplication in $\mathbb{B}^{n \times n}$ involves $4n^3$ multiplications in \mathbb{D} , the new method reduces this number to $n^3 + 2n^2$: one “expensive” multiplication in $\mathbb{D}^{n \times n}$ and two scalar multiplications of matrices by numbers. This type of tricks will be discussed in more detail in section 6.4 below.

Essentially, the new method is based on the isomorphism

$$\mathcal{B}(\mathbb{D}, \mathbb{D})^{n \times n} \cong \mathcal{B}(\mathbb{D}^{n \times n}, \mathbb{D}).$$

A similar isomorphism exists on the mathematical level:

$$(\mathbb{R}^{\text{com}})^{n \times n} \cong (\mathbb{R}^{n \times n})^{\text{com}}$$

As a variant, we directly may use the latter isomorphism at the top level, after which ball approximations of elements of $(\mathbb{R}^{n \times n})^{\text{com}}$ are already in $\mathcal{B}(\mathbb{D}^{n \times n}, \mathbb{D})$.

Numerical level. Being able to perform an automatic error analysis, the actual numerical computations are done at the numerical level. In our example, we should implement an efficient algorithm to multiply two matrices in $\mathbb{D}^{n \times n}$. In single or double precision, we may usually rely on highly efficient numerical libraries (BLAS, LAPACK, etc.). In higher precisions, new implementations are often necessary: even though there are efficient libraries for multiple precision floating point numbers, these libraries usually give rise to a huge overhead. For instance, when using MPFR [HLRZ00] at double precision, the overhead with respect to machine doubles is usually comprised between 10 and 100.

When working with matrices with multiple precision floating point entries, this overhead can be greatly reduced by putting the entries under a common exponent using the isomorphism

$$\mathbb{D}^{n \times n} = (\mathbb{Z} 2^{\mathbb{Z}})^{n \times n} \cong \mathbb{Z}^{n \times n} 2^{\mathbb{Z}}.$$

This reduces the matrix multiplication problem for floating point numbers to a purely arithmetic problem. Of course, this method becomes numerically unstable when the exponents differ wildly; in that case, row preconditioning of the first multiplicand and column preconditioning of the second multiplicand usually helps.

Arithmetic level. After the above succession of reductions, we generally end up with an arithmetic problem such as the multiplication of two $n \times n$ matrices in $\mathbb{Z}^{n \times n}$. The efficient resolution of this problem for all possible n and integer bit lengths p is again non-trivial.

Indeed, libraries such as GMP [Gra91] do implement Schönhage-Strassen's algorithm [SS71] algorithm for integer multiplication. However, the corresponding naive algorithm for the multiplication of matrices has a time complexity $O(l(p) n^3)$, which is far from optimal for large values of n .

Indeed, for large n , it is better to use multi-modular methods. For instance, choosing sufficiently many small primes $q_1, \dots, q_k < 2^{32}$ (or 2^{64}) with $q_1 \cdots q_k > 2 n 4^p$, the multiplication of the two integer matrices can be reduced to k multiplications of matrices in $(\mathbb{Z}/q_i \mathbb{Z})^{n \times n}$. Recall that a $2p$ -bit number can be reduced modulo all the q_i and reconstructed from these reductions in time $O(l(p) \log p)$. The improved matrix multiplication algorithm therefore admits a time complexity $O(p n^3 + n^2 l(p) \log p)$ and has been implemented in LINBOX [DGG+02b, DGG+02a], MATHEMAGIX and several other systems. FFT-based methods achieve similar practical time complexities $O(p n^3 + n^2 l(p))$ when n and p are of the same order of magnitude.

6. RELIABLE LINEAR ALGEBRA

In this section, we will start the study of ball arithmetic for non numeric types, such as matrices. We will examine the complexity of common operations, such as matrix multiplication, linear system solving, and the computation of eigenvectors. Ideally, the certified variants of these operations are only slightly more expensive than the non certified versions. As we will see, this objective can sometimes be met indeed. In general however, there is a trade-off between the efficiency of the certification and its quality, i.e. the sharpness of the obtained bound. As we will see, the overhead of bound computations also tends to diminish for increasing bit precisions p .

6.1. Matrix multiplication

Let us first consider the multiplication of two $n \times n$ double precision matrices

$$M, N \in \mathbb{B}_{51,12}^{n \times n}.$$

Naive strategy. The simplest multiplication strategy is to compute MN using the naive symbolic formula

$$(MN)_{i,k} = \sum_{j=1}^n M_{i,j} N_{j,k}. \quad (19)$$

Although this strategy is efficient for very small n , it has the disadvantage that we cannot profit from high performance BLAS libraries which might be available on the computer.

Revisited naive strategy. Reconsidering M and N as balls with matricial radii

$$M, N \in \mathcal{B}(\mathbb{D}_{51,12}^{n \times n}, \mathbb{D}_{51,12}^{n \times n}),$$

we may compute MN using

$$MN = \mathcal{B}([M_c N_c]^\uparrow, [|M_c| N_r + M_r (|N_c| + N_r) + n \epsilon |M_c| |N_c|]^\uparrow), \quad (20)$$

where $|M|$ is given by $|M|_{i,j} = |M_{i,j}|$. A similar approach was first proposed in [Rum99a]. Notice that the additional term $n \epsilon |M_c| \cdot^\uparrow |N_c|$ replaces $\Delta(MN)$. This extra product is really required: the computation of $(M_c N_c)_{i,j}$ may involve cancellations, which prevent a bound for the rounding errors to be read off from the end-result. The formula (20) does assume that the underlying BLAS library computes $M_c N_c$ using the naive formula (19) and correct rounding, with the possibility to compute the sums in any suitable order. Less naive schemes, such as Strassen multiplication [Str69], may give rise to additional rounding errors.

Fast strategy. The above naive strategies admit the disadvantage that they require four non certified $n \times n$ matrix products in $\mathbb{D}_{51,12}^{n \times n}$. If M and N are well-conditioned, then the following formula may be used instead:

$$MN = \mathcal{B}(M_c N_c, R) \quad (21)$$

$$R_{i,k} = [\| (M_{i,\cdot})_c \| \| (N_{\cdot,k})_r \| + \| (M_{i,\cdot})_r \| \| (N_{\cdot,k})_c \| + n \epsilon \| (M_{i,\cdot})_c \| \| (N_{i,\cdot})_c \|]^\uparrow, \quad (22)$$

where $M_{i,\cdot}$ and $N_{\cdot,k}$ stand for the i -th row of M and the k -th column of N . Since the $O(n)$ norms can be computed using only $O(n^2)$ operations, the cost of the bound computation is asymptotically negligible with respect to the $O(n^3)$ cost of the multiplication $M_c N_c$.

Hybrid strategy. For large $n \times n$ matrices, chances increase that M or N gets badly conditioned, in which case the quality of the error bound (22) decreases. Nevertheless, we may use a compromise between the naive and the fast strategies: fix a not too small constant K , such as $K \approx 16$, and rewrite M and N as $\lceil \frac{n}{K} \rceil \times \lceil \frac{n}{K} \rceil$ matrices whose entries are $K \times K$ matrices. Now multiply M and N using the revisited naive strategy, but use the fast strategy on each of the $K \times K$ block coefficients. Being able to choose K , the user has an explicit control over the trade-off between the efficiency of matrix multiplication and the quality of the computed bounds.

As an additional, but important observation, we notice that the user often has the means to perform an “*a posteriori* quality check”. Starting with a fast but low quality bound computation, we may then check whether the computed bound is suitable. If not, then we recompute a better bound using a more expensive algorithm.

High precision multiplication. Assume now that we are using multiple precision arithmetic, say $M, N \in \mathbb{B}_p^{n \times n}$. Computing MN using (20) requires one expensive multiplication in $\mathbb{D}_p^{n \times n}$ and three cheap multiplications in $\mathbb{D}_{64}^{n \times n}$. For large p , the bound computation therefore induces no noticeable overhead.

6.2. Matrix inversion

Assume that we want to invert an $n \times n$ ball matrix $M \in \mathbb{B}_p^{n \times n}$. This is a typical situation where the naive application of a symbolic algorithm (such as gaussian elimination or LR-decomposition) may lead to overestimation. An efficient and high quality method for the inversion of M is called Hansen's method [HS67, Moo66]. The main idea is to first compute the inverse of M_c using a standard numerical algorithm. Only at the end, we estimate the error using a perturbative analysis. The same technique can be used for many other problems.

More precisely, we start by computing an approximation $N_c \in \mathbb{D}_p^{n \times n}$ of $(M_c)^{-1}$. Putting $N = \mathcal{B}(N_c, 0)$, we next compute the product MN using ball arithmetic. This should yield a matrix of the form $1 - E$, where E is small. If E is indeed small, say $\|E\| \leq 2^{-p/2}$, then

$$\|(1 - E)^{-1} - 1 - E\| \leq \frac{\|E\|^2}{1 - \|E\|}. \quad (23)$$

Denoting by Ω_n the $n \times n$ matrix whose entries are all $\mathcal{B}(0, 1)$, we may thus take

$$(1 - E)^{-1} := 1 + E + \frac{\|E\|^2}{1 - \|E\|} \Omega_n.$$

Having inverted $1 - E$, we may finally take $M^{-1} = N (1 - E)^{-1}$. Notice that the computation of $\|E\|$ can be quite expensive. It is therefore recommended to replace the check $\|E\| \leq 2^{-p/2}$ by a cheaper check, such as $n \|E\|_\infty \leq 2^{-p/2}$.

Unfortunately, the matrix E is not always small, even if M is nicely invertible. For instance, starting with a matrix M of the form

$$M = J_{n,K} = \begin{pmatrix} 1 & K & & & \\ & 1 & \ddots & & \\ & & \ddots & K & \\ & & & & 1 \end{pmatrix},$$

with K large, we have

$$M^{-1} = \begin{pmatrix} 1 & -K & K^2 & \dots & (-K)^{n-1} \\ & 1 & -K & \ddots & \vdots \\ & & 1 & \ddots & K^2 \\ & & & \ddots & -K \\ & & & & 1 \end{pmatrix}.$$

Computing $(M_c)^{-1}$ using bit precision p , this typically leads to

$$\|E\| \approx K^{n-1} 2^{-p}.$$

In such cases, we rather reduce the problem of inverting $1 - E$ to the problem of inverting $1 - E^2$, using the formula

$$(1 - E)^{-1} = (1 + E)(1 - E^2)^{-1}. \quad (24)$$

More precisely, applying this trick recursively, we compute E^2, E^4, E^8, \dots until $\|E^{2^k}\|$ becomes small (say $\|E^{2^k}\| \leq 2^{-p/2}$) and use the formula

$$\begin{aligned} (1 - E)^{-1} &= (1 + E)(1 + E^2) \dots (1 + E^{2^{k-1}})(1 - E^{2^k})^{-1} \\ (1 - E^{2^k})^{-1} &:= 1 + E^{2^k} + \frac{\|E^{2^k}\|^2}{1 - \|E^{2^k}\|} \Omega_n. \end{aligned} \tag{25}$$

We may always stop the algorithm for $2^k > n$, since $(J_{n,K} - 1)^n = 0$. We may also stop the algorithm whenever $\|E^{2^k}\| \geq 1$ and $\|E^{2^k}\| \geq \|E^{2^{k-1}}\|$, since M usually fails to be invertible in that case.

In general, the above algorithm requires $O(\log m)$ ball $n \times n$ matrix multiplications, where m is the size of the largest block of the kind $J_{K,m}$ in the Jordan decomposition of M . The improved quality therefore requires an additional $O(\log n)$ overhead in the worst case. Nevertheless, for a fixed matrix M and $p \rightarrow \infty$, the norm $\|E\|$ will eventually become sufficiently small (i.e. $2^{-p/2}$) for (23) to apply. Again, the complexity thus tends to improve for high precisions. An interesting question is whether we can avoid the ball matrix multiplication MN altogether, if p gets really large. Theoretically, this can be achieved by using a symbolic algorithm such as Gaussian elimination or LR-decomposition using ball arithmetic. Indeed, even though the overestimation is important, it does not depend on the precision p . Therefore, we have $(M^{-1})_r = O(2^{-p})$ and the cost of the bound computations becomes negligible for large p .

6.3. Eigenproblems

Let us now consider the problem of computing the eigenvectors of a ball matrix $M \in \mathbb{B}_p^{n \times n}$, assuming for simplicity that the corresponding eigenvalues are non zero and pairwise distinct. We adopt a similar strategy as in the case of matrix inversion. Using a standard numerical method, we first compute a diagonal matrix $\Lambda \in \mathbb{D}_p^{n \times n}$ and an invertible transformation matrix $T \in \mathbb{D}_p^{n \times n}$, such that

$$T^{-1} M_c T \approx \Lambda. \tag{26}$$

The main challenge is to find reliable error bounds for this computation. Again, we will use the technique of small perturbations. The equation (26) being a bit more subtle than $M_c N_c \approx 1$, this requires more work than in the case of matrix inversion. In fact, we start by giving a numerical method for the iterative improvement of an approximate solution. A variant of the same method will then provide the required bounds. Again, this idea can often be used for other problems. The results of this section are work in common with B. MOURRAIN and PH. TREBUCHET; in [vdHMT], an even more general method is given, which also deals with the case of multiple eigenvalues.

Given \tilde{M} close to M_c , we have to find \tilde{T} close to T and $\tilde{\Lambda}$ close to Λ , such that

$$\tilde{T}^{-1} \tilde{M} \tilde{T} = \tilde{\Lambda}. \tag{27}$$

Putting

$$\begin{aligned} \tilde{T} &= T(1 + E) \\ \tilde{\Lambda} &= \Lambda(1 + \Delta) \\ N &= T^{-1} \tilde{M} T \\ H &= N - \Lambda, \end{aligned}$$

this yields the equation

$$(1 + E)^{-1} N (1 + E) = \Lambda(1 + \Delta).$$

Expansion with respect to E yields

$$\begin{aligned}
\frac{1}{1+E} N(1+E) &= N(1+E) - \frac{E}{1+E} N(1+E) \\
&= N + NE - EN + \frac{E^2}{1+E} N - \frac{E}{1+E} NE \\
&= N + [\Lambda, E] + [H, E] + \frac{E^2}{1+E} N - \frac{E}{1+E} NE \\
&= N + [\Lambda, E] + O([H, E]) + O(E^2).
\end{aligned} \tag{28}$$

Forgetting about the non-linear terms, the equation

$$[E, \Lambda] + \Lambda \Delta = H$$

admits a unique solution

$$\begin{aligned}
E_{i,j} &= \begin{cases} \frac{H_{i,j}}{\Lambda_{j,j} - \Lambda_{i,i}} & \text{if } j \neq i \\ 0 & \text{otherwise} \end{cases} \\
\Delta_{i,j} &= \begin{cases} \frac{H_{i,i}}{\Lambda_{i,i}} & \text{if } i = j \\ 0 & \text{otherwise} \end{cases}
\end{aligned}$$

Setting

$$\kappa = \kappa(\Lambda) = \max \left\{ \max \left\{ \frac{1}{|\Lambda_{j,j} - \Lambda_{i,i}|} : 1 \leq i < j \leq n \right\}, \max \left\{ \frac{1}{\Lambda_{i,i}} : 1 \leq i \leq n \right\} \right\},$$

it follows that

$$\max \{ \|E\|, \|\Delta\| \} \leq \kappa \sqrt{n} \|H\|. \tag{29}$$

Setting $T' = T(1+E)$, $\Lambda' = \Lambda(1+\Delta)$ and $H' = (T')^{-1} \tilde{M}T' - \Lambda'$, the relation (28) also implies

$$H' = [H, E] + \frac{E^2}{1+E} (\Lambda + H) - \frac{E}{1+E} (\Lambda + H) E.$$

Under the additional condition $\|E\| \leq \frac{1}{2}$, it follows that

$$\|H'\| \leq 3 \|H\| \|E\| + 4 \|E\|^2 \|\Lambda\|. \tag{30}$$

For sufficiently small H , we claim that iteration of the mapping $\Phi: (T, \Lambda) \mapsto (T', \Lambda')$ converges to a solution of (27).

Let us denote $(T^{(k)}, \Lambda^{(k)}) = \Phi^k(T, \Lambda)$, $H^{(k)} = (T^{(k)})^{-1} \tilde{M}T^{(k)} - \Lambda^{(k)}$ and let $(E^{(k)}, \Delta^{(k)})$ be such that $T^{(k+1)} = T^{(k)}(1+E^{(k)})$ and $\Lambda^{(k+1)} = \Lambda^{(k)}(1+\Delta^{(k)})$. Assume that

$$\|H\| \leq \frac{1}{80n\kappa^2\|\Lambda\|} \tag{31}$$

and let us prove by induction over k that

$$\|H^{(k)}\| \leq 2^{-k} \|H\| \tag{32}$$

$$\|\Lambda^{(k)}\| \leq 2 \|\Lambda\| \tag{33}$$

$$\begin{aligned}
\max \{ \|E^{(k)}\|, \|\Delta^{(k)}\| \} &\leq 2\sqrt{n}\kappa \|H^{(k)}\| \\
&\leq \frac{1}{40\sqrt{n}\kappa\|\Lambda\|2^k}
\end{aligned} \tag{34}$$

This is clear for $k=0$, so assume that $k > 0$. In a similar way as (30), we have

$$\|H^{(k)}\| \leq 3 \|H^{(k-1)}\| \|E^{(k-1)}\| + 4 \|E^{(k-1)}\|^2 \|\Lambda^{(k-1)}\|. \tag{35}$$

Using the induction hypotheses and $\kappa \|\Lambda\| \geq 1$, it follows that

$$\begin{aligned} \|H^{(k)}\| &\leq (3 + 16\sqrt{n}\kappa\|\Lambda\|) \|E^{(k-1)}\| \|H^{(k-1)}\| \\ &\leq \frac{1}{2} \|H^{(k-1)}\|, \end{aligned}$$

which proves (32). Now let $\Sigma^{(k)}$ be such that

$$\begin{aligned} \Lambda^{(k)} &= \Lambda (1 + \Delta^{(0)}) \dots (1 + \Delta^{(k-1)}) \\ &= \Lambda (1 + \Sigma^{(k)}). \end{aligned}$$

From (34), it follows that

$$\|\Sigma^{(k)}\| \leq \frac{1}{4\kappa\|\Lambda\|}.$$

On the one hand, this implies (33). On the other hand, it follows that

$$\kappa(\Lambda^{(k)}) \leq 2\kappa,$$

whence (29) generalizes to (34). This completes the induction and the linear convergence of $\|H^{(k)}\|$ to zero. In fact, the combination of (34) and (35) show that we even have quadratic convergence.

Let us now return to the original bound computation problem. We start with the computation of $H = T^{-1}MT - \Lambda$ using ball arithmetic. If the condition (31) is met (using the most pessimistic rounding), the preceding discussion shows that for every $\tilde{M} \in M$ (in the sense that $\tilde{M}_{i,j} \in M_{i,j}$ for all i, j), the equation (27) admits a solution of the form

$$\begin{aligned} \tilde{T} &= T(1 + \tilde{\Upsilon}) = T(1 + E^{(0)})(1 + E^{(1)}) \dots \\ \tilde{\Lambda} &= \Lambda(1 + \tilde{\Sigma}) = \Lambda(1 + \Delta^{(0)})(1 + \Delta^{(1)}) \dots, \end{aligned}$$

with

$$\max\{\|E^{(k)}\|, \|\Delta^{(k)}\|\} \leq 2^{1-k}\sqrt{n}\kappa\|H\|,$$

for all k . It follows that

$$\max\{\|\tilde{\Upsilon}\|, \|\tilde{\Sigma}\|\} \leq \eta := 6\sqrt{n}\kappa\|H\|.$$

We conclude that

$$\begin{aligned} \tilde{T} &\in T(1 + \eta\Omega_n) \\ \tilde{\Lambda} &\in \mathcal{B}(1, \eta)\Lambda. \end{aligned}$$

We may thus return $(T(1 + \eta\Omega_n), \mathcal{B}(1, \eta)\Lambda)$ as the solution to the original eigenproblem associated to the ball matrix M .

The reliable bound computation essentially reduces to the computation of three matrix products and one matrix inversion. At low precisions, the numerical computation of the eigenvectors is far more expensive in practice, so the overhead of the bound computation is essentially negligible. At higher precisions p , the iteration $(T, \Lambda) \mapsto \Phi(T, \Lambda)$ actually provides an efficient way to double the precision of a numerical solution to the eigenproblem at precision $p/2$. In particular, even if the condition (31) is not met initially, then it usually can be enforced after a few iterations and modulo a slight increase of the precision. For fixed M and $p \rightarrow \infty$, it also follows that the numerical eigenproblem essentially reduces to a few matrix products. The certification of the end-result requires a few more products, which induces a constant overhead. By performing a more refined error analysis, it is probably possible to make the cost certification negligible, although we did not investigate this issue in detail.

6.4. Matricial balls versus ball matrices

In section 4.3, we have already seen that matricial balls in $\mathcal{B}(\mathbb{D}^{n \times n}, \mathbb{D})$ often provide higher quality error bounds than ball matrices in $\mathcal{B}(\mathbb{D}, \mathbb{D})^{n \times n}$ or essentially equivalent variants in $\mathcal{B}(\mathbb{D}^{n \times n}, \mathbb{D}^{n \times n})$. However, ball arithmetic in $\mathcal{B}(\mathbb{D}^{n \times n}, \mathbb{D})$ relies on the possibility to quickly compute a sharp upper bound for the operator norm $\|M\|$ of a matrix $M \in \mathbb{D}^{n \times n}$. Unfortunately, we do not know of any really efficient algorithm for doing this.

One expensive approach is to compute a reliable singular value decomposition of M , since $\|M\|$ coincides with the largest singular value. Unfortunately, this usually boils down to the resolution of the eigenproblem associated to $M^* M$, with a few possible improvements (for instance, the dependency of the singular values on the coefficients of M is less violent than in the case of a general eigenproblem).

Since we only need the largest singular value, a faster approach is to reduce the computation of $\|M\|$ to the computation of $\|M^* M\|$, using the formula

$$\|M\| = \sqrt{\|M^* M\|}.$$

Applying this formula k times and using a naive bound at the end, we obtain

$$\|M\| \leq \sqrt[k]{n \| (M^* M)^{2^{k-1}} \|_\infty}.$$

This bound has an accuracy of $\approx k - O(\log_2 n)$ bits. Since $M^* M$ is symmetric, the $k - 1$ repeated squarings of $M^* M$ only correspond to about $\frac{k}{2}$ matrix multiplications. Notice also that it is wise to renormalize matrices before squaring them, so as to avoid overflows and underflows.

The approach can be speeded up further by alternating steps of tridiagonalization and squaring. Indeed, for a symmetric tridiagonal matrix D , the computation of D^2 and its tridiagonalization only take $O(n)$ steps instead of $O(n^3)$ for a full matrix product. After a few $k = O(n^2)$ steps of this kind, one obtains a good approximation μ of $\|D\|$. One may complete the algorithm by applying an algorithm with quadratic convergence for finding the smallest eigenvalues of $D - \mu$. In the lucky case when D has an isolated maximal eigenvalue, a certification of this method will provide sharp upper bounds for $\|D\|$ in reasonable time.

Even after the above improvements, the computation of sharp upper bounds for $\|M\|$ remains quite more expensive than ordinary matrix multiplication. For this reason, it is probably wise to avoid ball arithmetic in $\mathcal{B}(\mathbb{D}^{n \times n}, \mathbb{D})$ except if there are good reasons to expect that the improved quality is really useful for the application in mind.

Moreover, when using ball arithmetic in $\mathcal{B}(\mathbb{D}^{n \times n}, \mathbb{D}^{n \times n})$, it is often possible to improve algorithms in ways to reduce overestimation. When interpreting a complete computation as a dag, this can be achieved by minimizing the depth of the dag, i.e. by using an algorithm which is better suited for parallelization. Let us illustrate this idea for the computation of the k -th power M^k of a matrix M . When using Horner's method (multiply the identity matrix k times by M), we typically observe an overestimation of $O(k)$ bits (as for the example (18)). If we use binary powering, based on the rule

$$M^k = M^{\lfloor k/2 \rfloor} M^{\lceil k/2 \rceil}, \tag{36}$$

then the precision loss drops down to $O(\log k)$ bits. We will encounter a less obvious application of the same idea in section 7.5 below.

7. RELIABLE POWER SERIES ARITHMETIC

7.1. Ball series versus serial balls

There are two typical applications of power series $f \in \mathbb{R}[[z]]$ or $f \in \mathbb{C}[[z]]$ with certified error bounds. When f occurs as a generating function in a counting problem or random object generator, then we are interested in the computation of the coefficients f_n for large n , together with reliable error bounds. A natural solution is to systematically work with computable power series with ball coefficients in $\mathbb{B}_p[[z]]^{\text{com}}$. For many applications, we notice that p is fixed, whereas $n \gg p$ may become very large.

The second typical application is when $f \in \mathbb{C}[[z]]$ is the local expansion of an analytic function on a disk $\mathcal{B}(0, \rho)$ and we wish to evaluate f at a point z with $|z| < \rho$. The geometric decrease of $|f_n z^n|$ implies that we will need only $n = O(p)$ coefficients of the series. In order to bound the remaining error using Cauchy's formula, we do not only need bounds for the individual coefficients f_n , but also for the norm $\|f\|_\rho$ defined in (4). Hence, it is more natural to work with serial balls in $\mathcal{B}(\mathbb{D}_p[\mathbb{i}][[z]]^{\text{com}}, \mathbb{D}_{64})$, while using the $\|\cdot\|_\rho$ norm. Modulo a rescaling $f(z) \mapsto f(\rho z)$, it will be convenient to enforce $\rho = 1$. In order to compute sharp upper bounds $\|f\|$ for $\|f\| = \|f\|_1$, it will also be convenient to have an algorithm which computes bounds $\|f_n\|$ for the tails

$$f_{;n} = f_n z^n + f_{n+1} z^{n+1} + \dots$$

Compared to the computation of the corresponding head

$$f_{;n} = f_0 + \dots + f_{n-1} z^{n-1},$$

we will show in section 7.4 that the computation of such a tail bound is quite cheap.

Again the question arises how to represent $f_{;n}$ in a reliable way. We may either store a global upper bound for the error, so that $f_{;n} \in \mathcal{B}(\mathbb{D}_p[\mathbb{i}][z], \mathbb{D}_{64})$, or compute individual bounds for the errors, so that $f_{;n} \in \mathbb{B}[\mathbb{i}]_p[z]$. If our aim is to evaluate f at a point z with $|z| \approx 1$, then both representations $P \in \mathbb{B}[\mathbb{i}]_p[z]$ and $\hat{P} = \mathcal{B}(P_c, \|P_r\|_\infty) \in \mathcal{B}(\mathbb{D}_p[\mathbb{i}][z], \mathbb{D}_{64})$ give rise to evaluations $P(z), \hat{P}(z) \in \mathbb{B}[\mathbb{i}]_p$ with equally precise error bounds. Since the manipulation of global error bounds is more efficient, the corresponding representation should therefore be preferred in this case. In the multivariate case, one has the additional benefit that ‘‘small’’ coefficients $f_i z^i$ (e.g. $|f_i| \leq \epsilon_p \|f\|$) can simply be replaced by a global error $\mathcal{B}(0, |f_i|)$, thereby increasing the sparsity of f . On the other hand, individual error bounds admit the advantage that rescaling $f(z) \mapsto f(\lambda z)$ is cheap. If we suddenly find out that f is actually convergent on a larger disk and want to evaluate f at a point z with $|z| > 1$, then we will not have to recompute the necessary error bounds for $f_{;n}$ from scratch.

Serial ball representations similar to what has been described above are frequently used in the area of Taylor models [MB96, MB04] for the validated long term integration of dynamical systems. In the case of Taylor models, there is an additional twist: given a dynamical system of dimension d , we not only compute a series expansion with respect to the time t , but also with respect to small perturbations $\varepsilon_1, \dots, \varepsilon_d$ of the initial conditions. In particular, we systematically work with power series in several variables. Although such computations are more expensive, the extra information may be used in order to increase the sharpness of the computed bounds. A possible alternative is to compute the expansions in $\varepsilon_1, \dots, \varepsilon_d$ only up to the first order and to use binary splitting for the multiplication of the resulting Jacobian matrices on the integration path. This approach will be detailed in section 7.5.

7.2. Reliable multiplication of series and polynomials

In order to study the reliable multiplication of series f and g , let us start with the case when $f, g \in \mathcal{B}(\mathbb{D}_p[[z]]^{\text{com}}, \mathbb{D}_{64})$, using the sup-norm on the unit disk. We may take

$$h = fg = \mathcal{B}(\widetilde{f_c g_c}, [(\|f_c\| \|g_r\| + \|f_r\| (\|g_c\| + \|g_r\|) + \delta)^\uparrow]),$$

where $h_c = \widetilde{f_c g_c}$ stands for a δ -approximation of $f_c g_c$. Since h_c is really a numeric algorithm for the computation of its coefficients, the difficulty resides in the fact that δ has to be chosen once and for all, in such a way that the bound $\|h_c - f_c g_c\|$ will be respected at the limit. A reasonable choice is $\delta = \epsilon_p \|f_c\| \|g_c\|$. We next distribute this error over the infinity of coefficients: picking some $\alpha < 1$, each coefficient $(h_c)_n$ is taken to be an $[(1 - \alpha) \alpha^n \delta]$ -approximation of $f_c g_c$. Of course, these computation may require a larger working precision than p . Nevertheless, f and g are actually convergent on a slightly larger disk $\mathcal{B}(0, \rho)$. Picking $\alpha = 1/\rho$, the required increase of the working precision remains modest.

Let us now turn our attention to the multiplication of two computable series $f, g \in \mathbb{B}_p[[z]]^{\text{com}}$ with ball coefficients. Except for naive power series multiplication, based on the formula $(fg)_n = \sum_{i+j=n} f_i g_j$, most other multiplication algorithms (whether relaxed or not) use polynomial multiplication as a subalgorithm. We are thus left with the problem of finding an efficient and high quality algorithm for multiplying two polynomials $P, Q \in \mathbb{B}_p[z]$ of degrees $< n$. In order to simplify the reading, we will assume that $P_0, P_{n-1}, Q_0, Q_{n-1}$ are all non-zero.

As in the case of matrix multiplication, there are various approaches with different qualities, efficiencies and aptitudes to profit from already available fast polynomial arithmetic in $\mathbb{D}_p[z]$. Again, the naive $O(n^2)$ approach provides almost optimal numerical stability and qualities for the error bounds. However, this approach is both slow from an asymptotic point of view and unable to rely on existant multiplication algorithms in $\mathbb{D}_p[z]$.

If the coefficients of P and Q are all of the same order of magnitude, then we may simply convert P and Q into polynomial balls in $\mathcal{B}(\mathbb{D}_p[z], \mathbb{D}_{64})$ for the norm $\|\cdot\|_\infty$ and use the following crude formula for their multiplication:

$$PQ = \mathcal{B}(\widetilde{P_c Q_c}, [n \|P_c\|_\infty \|Q_r\|_\infty + n \|P_r\|_\infty (\|Q_c\|_\infty + \|Q_r\|_\infty) + \delta]^\uparrow), \quad (37)$$

where $\widetilde{P_c Q_c}$ stands for a δ -approximation of $P_c Q_c$. In other words, we may use any efficient multiplication algorithm in $\mathbb{D}_p[z]$ for the approximation of $P_c Q_c$, provided that we have a means to compute a certified bound δ for the error.

In our application where P and Q correspond to ranges of coefficients in the series f and g , we usually have $P_i \approx P_0 \rho_f^{-i}$ and $Q_i \approx Q_0 \rho_g^{-i}$ for the convergence radii ρ_f and ρ_g of P and Q . In order to use (37), it is therefore important to scale $P(z) \mapsto P(z/\rho)$ and $Q(z) \mapsto Q(z/\rho)$ for a suitable ρ . If we are really interested in the evaluation of fg at points z in a disk $\mathcal{B}(0, r)$, then we may directly take $\rho = r$. In we are rather interested in the coefficients of fg , then $\rho = \min(\rho_f, \rho_g)$ is the natural choice. However, since ρ_f and ρ_g are usually unknown, we first have to compute suitable approximations for them, based on the available coefficients of P and Q . A good heuristic approach is to determine indices $i < n/2 \leq j$ such that

$$\left| \frac{P_k}{P_i} \right|^{j-i} \leq \left| \frac{P_j}{P_i} \right|^{k-i},$$

for all k , and to take

$$\rho_f \approx \left| \frac{P_i}{P_j} \right|^{\frac{1}{j-i}}$$

as the approximation for ρ_f . Recall that the numerical Newton polygon of N_P is the convex hull of all points $(i, \log |P_i| - \lambda) \in \mathbb{R}^2$ with $\lambda \geq 0$. Consider the edge of N_P whose projection on the x -axis contains $n/2$. Then i and j are precisely the extremities of this projection, so they can be computed in linear time.

For large precisions $n = O(p)$, the scaling algorithm is both very efficient and almost of an optimal quality. For small p and large n , there may be some precision loss which depends on the nature of the smallest singularities of f and g . Nevertheless, for many singularities, such as algebraic singularities, the precision loss is limited to $O(\log n)$ bits. For a more detailed discussion, we refer to [vdH02a, Section 6.2].

In other applications, where P and Q are not obtained from formal power series, it is usually insufficient to scale using a single factor ρ . This is already the case when multiplying $P = 1 + z$ and $Q = 1 + \delta z$ for small $\delta \ll \epsilon_p$, since the error bound for $(PQ)_2 = \delta$ exceeds ϵ_p . One possible remedy is to “precondition” P and Q according to their numerical Newton polygons and use the fact that N_{PQ} is close to the Minkowski sum $N_P + N_Q$.

More precisely, for each i , let $(N_P)_i > 0$ denote the number such that $(i, \log (N_P)_i)$ lies on one of the edges of N_P . Then $(N_P)_i$ is of the same order of magnitude as P_i , except for indices for which P_i is small “by accident”. Now consider the preconditioned relative error

$$\nu_P = \max_i \frac{(P_i)_r}{(N_P)_i}$$

and similarly for Q . Then

$$[(PQ)_i]_r \leq n(\nu_P + \nu_Q + \nu_P \nu_Q)(N_P + N_Q)_i, \tag{38}$$

if PQ is computed using infinite precision ball arithmetic. Assuming a numerically stable multiplication algorithm for the centers, as proposed in [vdH08a], and incorporating the corresponding bound for the additional errors into the right hand side of (38), we thus obtain an efficient way to compute an upper bound for $(PQ)_r$.

Notice that the numerical Newton polygon N_P has a close relation to the orders of magnitudes of the roots of P . Even though the error bounds for some “accidentally small” coefficients $(PQ)_i$ may be bad for the above method, the error bounds have a good quality if we require them to remain valid for small perturbations of the roots of PQ .

7.3. Reliable series division and exponentiation

The algorithms from the previous section can in particular be used for the relaxed multiplication of two ball power series $f, g \in \mathbb{B}_p[[z]]^{\text{com}}$. In particular, using the implicit equation

$$g = 1 + zfg, \tag{39}$$

this yields a way to compute $g = (1 - zf)^{-1}$. Unfortunately, the direct application of this method leads to massive overestimation. For instance, the computed error bounds for the inverses of

$$g = \frac{1}{1 - 3z - 2z^2}$$

$$h = \frac{1}{1 - 3z + 2z^2}$$

coincide. Indeed, even when using a naive relaxed multiplication, the coefficients of g and h are computed using the recurrences

$$g_n = 3g_{n-1} + 2g_{n-2}$$

$$h_n = 3h_{n-1} - 2h_{n-2},$$

but the error bound ϵ_n for g_n and h_n is computed using the same recurrence

$$\epsilon_n = 3\epsilon_{n-1} + 2\epsilon_{n-2},$$

starting with $\varepsilon_0 \approx \epsilon_p$. For $n \rightarrow \infty$, it follows that $g_n \sim \lambda \alpha^{-n}$ and $\varepsilon_n \sim \lambda \epsilon_p \alpha^{-n}$ for some λ , where $\alpha \approx 0.281$ is the smallest root of $1 - 3\alpha - 2\alpha^2$. Hence the error bound for g_n is sharp. On the other hand, $h_n \sim \mu \beta^n$ for some μ , where $\beta = \frac{1}{2}$ is the smallest root of $1 - 3\beta + 2\beta^2$. Hence, the error bound ε_n is $n \log_2(\beta/\alpha)$ bits too pessimistic in the case of h . The remedy is similar to what we did in the case of matrix inversion. We first introduce the series

$$\begin{aligned}\varphi &= \frac{1}{1 - z f_c} \in \mathbb{D}_p[[z]]^{\text{com}} \\ \psi &= \frac{[\varphi(zf - 1)]_1}{z} \in \mathbb{B}_p[[z]]^{\text{com}}\end{aligned}$$

and next compute g using

$$g = \frac{\varphi}{1 - z\psi},$$

where $1 - z\psi$ is inverted using the formula (39). This approach has the advantage of being compatible with relaxed power series expansion and it yields high quality error bounds.

Another typical operation on power series is exponentiation. Using relaxed multiplication, we may compute the exponential g of an infinitesimal power series $f \in z \mathbb{B}_p[[z]]^{\text{com}}$ using the implicit equation

$$g = 1 + \int f' g, \quad (40)$$

where \int stands for ‘‘distinguished integration’’ in the sense that $(\int h)_0 = 0$ for all h . Again, one might fear that this method leads to massive overestimation. As a matter of fact, it usually does not. Indeed, assume for simplicity that $f_r = 0$, so that $f \in \mathbb{D}_p[[z]]$. Recall that $|f|$ denotes the power series with $|f|_n = |f_n|$. Roughly speaking, the error bound for g_n , when computed using the formula (40) will be the coefficient ε_n of the power series $\varepsilon = \exp(\epsilon_p |f|)$. Since $|f|$ has the same radius of convergence as f , it directly follows that the bit precision loss is sublinear $o(n)$. Actually, the dominant singularity of $|f|$ often has the same nature as the dominant singularity of f . In that case, the computed error bounds usually become very sharp. The observation generalizes to the resolution of linear differential equations, by taking a square matrix of power series for f .

7.4. Automatic tail bounds

In the previous sections, we have seen various reliable algorithms for the computation of the expansion $f_{;n}$ of a power series $f \in \mathbb{C}[[z]]$ up till a given order n . Such expansions are either regarded as ball polynomials in $\mathbb{B}[i]_p[z]$ or as polynomial balls in $\mathcal{B}(\mathbb{D}_p[i][z], \mathbb{D}_{64})$. Assuming that f is convergent on the closed unit disk $\mathcal{B}(0, 1)$, it remains to be shown how to compute tail bounds $\llbracket f_{;n} \rrbracket$. We will follow [vdH07b]. Given a ball polynomial P , then we notice that reasonably sharp upper and lower bounds $\llbracket P \rrbracket$ and $\lll P \lll$ for P can be obtained efficiently by evaluating P at $O(\deg P)$ primitive roots of unity using fast Fourier transforms [vdH07b, Section 6.2].

We will assume that the series f is either an explicit series, the result of an operation on other power series or the solution of an implicit equation. Polynomials are the most important examples of explicit series. Assuming that f is a polynomial of degree $< d$, we may simply take

$$\llbracket f_{;n} \rrbracket = \begin{cases} \llbracket z^{-n} f_{;d} \rrbracket & \text{if } n < d \\ 0 & \text{otherwise} \end{cases}, \quad (41)$$

where

$$f_{;d} = f_n z^n + \dots + f_{d-1} z^{d-1}.$$

For simple operations on power series, we may use the following bounds:

$$\llbracket (f + g);_n \rrbracket = \llbracket f;_n \rrbracket + \llbracket g;_n \rrbracket \quad (42)$$

$$\llbracket (f - g);_n \rrbracket = \llbracket f;_n \rrbracket + \llbracket g;_n \rrbracket \quad (43)$$

$$\llbracket (fg);_n \rrbracket = \llbracket f;_n; \rrbracket (\llbracket g;_n \rrbracket + \llbracket g_n; \rrbracket) + \llbracket f;_n \rrbracket \llbracket g_n; \rrbracket + \llbracket (f;_n g;_n)_n; \rrbracket \quad (44)$$

$$\llbracket (ff)_n; \rrbracket = \frac{1}{n+1} \llbracket f_n; \rrbracket, \quad (45)$$

where

$$\llbracket (f;_n g;_n)_n; \rrbracket \leq \sum_{k=0}^{n-1} |f_k| \left(\sum_{l=n-k}^{n-1} |g_l| \right)$$

can be computed in time $O(n)$. Due to possible overestimation, division has to be treated with care. Given an infinitesimal power series ε and

$$f = \frac{1}{1 - \varepsilon},$$

so that

$$f_n; = \frac{1 + \varepsilon f;_n - f;_n}{1 - \varepsilon},$$

we take

$$\llbracket f_n; \rrbracket = \frac{\llbracket (\varepsilon f;_n)_n; \rrbracket}{\llbracket 1 - \varepsilon;_n + \mathcal{B}(0, \llbracket \varepsilon_n; \rrbracket) \rrbracket}, \quad (46)$$

where $\llbracket (\varepsilon f;_n)_n; \rrbracket$ is computed using (44).

Let $\Phi(f)$ be an expression which is constructed from f and polynomials, using the ring operations and distinguished integration (we exclude division in order to keep the discussion simple). Assume that each coefficient $\Phi(f)_n$ only depends on the previous coefficients f_0, \dots, f_{n-1} of f and not on f_n, f_{n-1}, \dots . Then the sequence $0, \Phi(0), \Phi^2(0), \dots$ tends to the unique solution f of the implicit equation

$$f = \Phi(f), \quad (47)$$

Moreover, $\Phi^k(0)_n = f_n$ for any $k > n$. In (39) and (40), we have already seen two examples of implicit equations of this kind.

The following technique may be used for the computation of tail bounds $\llbracket f_n; \rrbracket$. Given $c \in \mathbb{D}^{\geq}$ and assuming that $\llbracket f_n; \rrbracket \leq c$, we may use the rules (41–46) in order to compute a “conditional tail bound” $\llbracket \Phi(f)_n; |c \rrbracket$ for $\llbracket \Phi(f)_n; \rrbracket$. Mathematically speaking, this bound has the property that for any power series g with $g;_n = f;_n$ and $\llbracket g_n; \rrbracket \leq c$, we have $\llbracket \Phi(g)_n; \rrbracket \leq \llbracket \Phi(g)_n; |c \rrbracket$. If $\llbracket \Phi(g)_n; |c \rrbracket \leq c$, then it follows that $\llbracket \Phi^{(k)}(f;_n);_n \rrbracket \leq c$ for all k . Given any disk $\mathcal{B}(0, \rho)$ with $\rho < 1$, it follows that $\llbracket \Phi^{(k)}(f;_n) - f \rrbracket_{\rho} \leq 2c\rho^k / (1 - \rho)$ for any $k \geq n$, since $\Phi^{(k)}(f;_n) - f = O(z^k)$. In other words, $\Phi^{(k)}(f;_n)$ uniformly converges to f on $\mathcal{B}(0, \rho)$. Therefore, $\llbracket f_n; \rrbracket_{\rho} \leq c$ and $\llbracket f_n; \rrbracket \leq c$, by letting $\rho \rightarrow 1$.

Let us now consider the mapping $\varphi: c \mapsto \llbracket \Phi(f)_n; |c \rrbracket$ and assume that Φ involves no divisions. When computing $\llbracket \Phi(f)_n; |c \rrbracket$ using infinite precision and the rules (41–44), we notice that φ is a real analytic function, whose power series expansion contains only positive coefficients. Finding the smallest c with $\varphi(c) \leq c$ thus reduces to finding the smallest fixed point c_{fix} of φ , if such a fixed point exists. We may use the secant method:

$$\begin{aligned} c_0 &:= 0 \\ c_1 &:= \varphi(c_0) \\ c_{k+2} &:= c_k + \frac{\varphi(c_k) - c_k}{c_{k+1} - \varphi(c_{k+1}) + \varphi(c_k) - c_k} (c_{k+1} - c_k) \end{aligned}$$

If $c_{k+1} < c_k$ for some k or if k exceeds a given threshold K , then the method fails and we set $\|f_n\| = +\infty$. Otherwise, c_k converges quadratically to c_{fix} . As soon as $|c_{k+1}/c_k - 1| < \varepsilon$, for some given $\varepsilon > 0$, we check whether $\varphi(\tilde{c}_{\text{fix}}) \leq \tilde{c}_{\text{fix}}$ for $\tilde{c}_{\text{fix}} = 2c_{k+1} - c_k$, in which case we stop. The resulting \tilde{c}_{fix} is an approximation of c_{fix} with relative accuracy $\varepsilon > 0$.

Assuming that $\Psi(f)_n$ has been computed for every subexpression of Φ , we notice that the computation of $\|\Phi(f)_n; c\|$ only requires $O(ns)$ floating point operations, where s is the size of Φ as an expression. More generally, the evaluation of $\|\Phi(f)_n; c_i\|$ for k different hypotheses c_1, \dots, c_k only requires $O(nsk)$ operations, since the heads $\Psi(f)_n$ do not need to be recomputed. Our algorithm for the computation of \tilde{c}_{fix} therefore requires at most $O(nsk)$ floating point operations. Taking $K = o(R(n)/n)$, it follows that the cost of the tail bound computation remains negligible with respect to the series expansion itself.

The approach generalizes to the case when f is a vector or matrix of power series, modulo a more involved method for the fixed point computation [vdH07b, Section 6.3]. If f is indeed convergent on $\mathcal{B}(0, 1)$, then it can also be shown that φ indeed admits a fixed point if n is sufficiently large.

7.5. Long term integration of dynamical systems

Let us now consider a dynamical system

$$f = f(0) + \int \Phi(f), \quad (48)$$

where f is a vector of d unknown complex power series, $f(0) \in \mathbb{B}[i]_p^d$, and Φ an expression built up from the entries of f and polynomials in $\mathbb{D}_p[i][z]$ using the ring operations. Given $t \in \mathbb{D}[i] \setminus \{0\}$, denote by $f_{z \times t}$ the scaled power series $f(tz)$ and by f_{z+t} the shifted power series $f(t+z)$, assuming that f converges at t . Also denote by $\Phi_{z \times t}$ and Φ_{z+t} the expressions which are obtained when replacing each polynomial P in Φ by $P_{z \times t}$ resp. P_{z+t} . Then $f_{z \times t}$ and f_{z+t} satisfy

$$f_{z \times t} = f(0) + t \int \Phi_{z \times t}(f_{z \times t}) \quad (49)$$

$$f_{z+t} = f(t) + \int \Phi_{z+t}(f_{z+t}) \quad (50)$$

Since (48) is a particular case of an implicit equation of the form (47), we have an algorithm for the computation of tail bounds $\|f_n\| \in (\mathbb{D}_{64}^{\geq} \cup \{\infty\})^d$ for f on $\mathcal{B}(0, 1)$. Modulo rescaling (49), we may also compute tail bounds $\|f_n\|_\rho$ on more general disks $\mathcal{B}(0, \rho)$.

Assume that we want to integrate (48) along the real axis, as far as possible, and performing all computations with an approximate precision of p bits. Our first task is to find a suitable initial step size t and evaluate f at t . Since we require a relative precision of approximately p bits, we roughly want to take $t \approx \rho_f/2$, where ρ_f is the radius of convergence of f , and evaluate the power series f up to order $n \approx p$. We thus start by the numerical computation of f_n ; and the estimation $\tilde{\rho}_f$ of ρ_f , based on the numerical Newton polygon of f_n . Setting $t := \tilde{\rho}_f/2$, we next compute a tail bound $\|f_n\|_t$. This bound is considered to be of an acceptable quality if

$$\|f_n\|_t \leq \epsilon_p \|f_n\|_t.$$

In the contrary case, we keep setting $t := t/2$ and recomputing $\|f_n\|_t$, until we find a bound of acceptable quality. It can be shown [vdH07b] that this process ultimately terminates, when p is sufficiently large. We thus obtain an appropriate initial step size t which allows us to compute a δ -approximation of $f(t)$ with $\delta = 2\epsilon_p \|f_n\|_t + \uparrow(f_n(t))_r$.

In principle, we may now perform the translation $z \mapsto t + z$ and repeat the analytic continuation process using the equation (50). Unfortunately, this approach leads to massive overestimation. Indeed, if the initial condition $f(0)$ is given with relative precision η , then the relative precisions of the computed coefficients f_k are typically a non trivial factor times larger than η , as well as the next “initial condition” $f(t)$ at t . Usually, we therefore lose $O(1)$ bits of precision at every step.

One remedy is the following. Let $\Delta = \Delta_{0,t}$ be the analytic continuation operator which takes an initial condition $c \in \mathbb{C}^d$ on input and returns the evaluation $f(t) \in \mathbb{C}^d$ of the unique solution to (48) with $f(0) = c$. Now instead of computing $f_{;n}$ using a ball initial condition $f(0) \in \mathbb{B}[i]_p^d$, we rather use its center $f(0)_c$ as our initial condition and compute $\Delta(f(0)_c)$ using ball arithmetic. In order to obtain a reliable error bound for $f(t)$, we also compute the Jacobian J_Δ of Δ at $f(0)_c$ using ball arithmetic, and take

$$f(t) = \Delta(f(0)_c) + \mathcal{B}(0, [J_\Delta(f(0)_c) f(0)_r]^\uparrow).$$

The Jacobian can either be computed using the technique of automatic differentiation [BS83] or using series with coefficients in a jet space of order one.

With this approach $\Delta(f(0)_c)$ is computed with an almost optimal p -bit accuracy, and $J_\Delta(f(0)_c)$ with an accuracy which is slightly worse than the accuracy of the initial condition. In the lucky case when $J_\Delta(f(0)_c)$ is almost diagonal, the accuracy of $f(t)$ will therefore be approximately the same as the accuracy of $f(0)$. However, if $J_\Delta(f(0)_c)$ is not diagonal, such as in (18), then the multiplication $J_\Delta(f(0)_c) f(0)_r$ may lead to overestimation. This case may already occur for simple linear differential equations such as

$$\begin{pmatrix} f \\ g \end{pmatrix} = \begin{pmatrix} 1 \\ 0 \end{pmatrix} + \int \begin{pmatrix} 0 & -1 \\ 1 & 0 \end{pmatrix} \begin{pmatrix} f \\ g \end{pmatrix} \quad (51)$$

and the risk is again that we lose $O(1)$ bits of precision at every step.

Let us describe a method for limiting the harm of this manifestation of the wrapping effect. Consider the analytic continuations of f at successive points $0 = t_0 < t_1 < \dots < t_k$ and denote

$$J_i = J_{\Delta_{t_{i-1}, t_i}}(f(t_{i-1})_c), \quad i = 1, \dots, k.$$

Instead of computing the error at t_i due to perturbation of the initial condition $f(0)$ using the formula

$$[f(t)_r]^{\text{per}(0)} = J_k(J_{k-1}(\dots J_1 f(0)_r \dots)),$$

we may rather use the formula

$$[f(t)_r]^{\text{per}(0)} = (J_k J_{k-1} \dots J_1) f(0)_r,$$

where matrix chain products are computed using a variant of binary powering (36):

$$J_j \dots J_i = (J_j \dots J_{\lfloor (i+j)/2 \rfloor + 1}) (J_{\lfloor (i+j)/2 \rfloor} \dots J_i).$$

In order to be complete, we also have to take into account the additional error δ_i , which occurs during the computation of $\Delta_{t_{i-1}, t_i}(f(t_{i-1})_c)$. Setting $\delta_0 = f(0)_r$, we thus take

$$f(t_i)_r = \delta_i + J_i \delta_{i-1} + \dots + (J_k \dots J_1) \delta_0. \quad (52)$$

When using this algorithm, at least in the case of simple systems such as (51), the precision loss at step k will be limited to $O(\log k)$ bits. Notice that we benefit from the fact that the Jacobian matrices remain accurate as long as the initial conditions remain accurate.

Although we have reduced the wrapping effect, the asymptotic complexity of the above algorithm is cubic or at least quadratic in k when evaluating (52) using a binary splitting version of Horner's method. Let us now describe the final method which requires only $O(k \log k)$ matrix multiplications up till step k . For $0 \leq i < j$, we define

$$\begin{aligned}\delta_{i,j} &= \delta_{j-1} + J_{j-1} \delta_{j-2} + \cdots + (J_{j-1} \cdots J_{i+1}) \delta_i, \\ J_{i,j} &= J_{j-1} \cdots J_i,\end{aligned}$$

where $J_0 = 1$. At stage $i = 2^{e_1} + \cdots + 2^{e_l}$ with $e_1 > \cdots > e_l$, we assume that

$$\begin{aligned}\delta_{[j]} &:= \delta_{[i;j]} := \delta_{2^{e_1+\dots+2^{e_{j-1}}}, 2^{e_1+\dots+2^{e_j}}} \\ J_{[j]} &:= J_{[i;j]} := J_{2^{e_1+\dots+2^{e_{j-1}}}, 2^{e_1+\dots+2^{e_j}}}\end{aligned}$$

are stored in memory for all $1 \leq j \leq l$. From these values, we may compute

$$f(t_{i-1})_r = \delta_{[l]} + J_{[l]} (\delta_{[l-1]} + \cdots J_{[3]} (\delta_{[2]} + J_{[2]} \delta_{[1]}) \cdots)$$

using only $O(l) = O(\log i)$ matrix multiplications. Now assume that we go to the next stage $i+1$. If m is maximal such that 2^m divides $i+1$, then $i+1 = 2^{e_1} + \cdots + 2^{e_{l-m}} + 2^m$. Consequently, $\delta_{[i+1;j]} = \delta_{[i;j]}$ and $J_{[i+1;j]} = J_{[i;j]}$ for $j < l' := l - m + 1$ and

$$\begin{aligned}\delta_{[i+1;l']} &= \delta_i + J_i (\delta_{[i;l]} + \cdots J_{[i;l'+2]} (\delta_{[i;l'+1]} + J_{[i;l'+1]} \delta_{[i;l']}) \cdots) \\ J_{[i+1;l']} &= J_i J_{[i;l]} \cdots J_{[i;l']}.\end{aligned}$$

Hence, the updated lists $\delta_{[i+1;j]}$ and $J_{[i+1;j]}$ can be computed using $O(m) = O(\log i)$ matrix multiplications. Furthermore, we only need to store $O(\log i)$ auxiliary matrices at step i .

7.6. Discussion

There is a vast literature on validated integration of dynamical systems and reduction of the wrapping effect [Moo65, Moo66, Nic85, Loh88, GS88, Neu93, K8, Loh01, Neu02, MB04]. We refer to [Loh01] for a nice review. Let us briefly discuss the different existing approaches.

The wrapping effect was noticed in the early days of interval arithmetic [Moo65] and local coordinate transformations were proposed as a remedy. The idea is to work as much as possible with respect to coordinates in which all errors are parallel to axes. Hence, instead of considering blocks $x \in \mathbb{B}_p^d$, we rather work with parallelepipeds $x = c + T\mathcal{B}(0, r)$, with $c \in \mathbb{D}_p^d$, $T \in \mathbb{D}_p^{d \times d}$ and $r \in \mathbb{D}_p^d$. A natural choice for T after k steps is $T = J_k \cdots J_1$, but more elaborate choices may be preferred [Loh88]. Other geometric shapes for the enclosures have been advanced in the literature, such as ellipsoids [Neu93], which are also invariant under linear transformations, and zonotopes [K8]. However, as long as we integrate (48) using a straightforward iterative method, and even if we achieve a small average loss $\varepsilon \ll 1$ of the bit precision at a single step, the precision loss after k steps will be of the form $k\varepsilon$.

The idea of using dichotomic algorithms in order to reduce the wrapping effect was first described in the case of linear differential equations [GS88]. The previous section shows how to adapt that technique to non linear equations. We notice that the method may very well be combined with other geometric shapes for the enclosures: this will help to reduce the precision loss to $(\log k)\varepsilon$ instead of $k\varepsilon$ in the above discussion.

Notice that there are two common misunderstandings about the dichotomic method. Contrary to what we have shown above (see also [GS88] in the linear case), it is sometimes believed that we need to keep all matrices J_1, \dots, J_k in memory and that we have to reevaluate products $J_j \cdots J_i$ over and over again. Secondly, one should not confuse *elimination* and *reduction* of the wrapping effect: if M is the 2×2 rotation matrix by a 30° angle, then none of its repeated squarings M^{2^l} will be the identity matrix, so every squaring $(M^{2^l})^2$ will involve a wrapping. Even though we have *not* eliminated the wrapping effect, we *did* achieve to reduce the number of wrappings to l instead of 2^l .

Taylor models are another approach for the validated long term integration of dynamical systems [EKW84, EMOK91, MB96, MB04]. The idea is to rely on higher k -th order expansions of the continuation operator Δ . This allows for real algebraic enclosures which are determined by polynomials of degree k . Such enclosures are *a priori* more precise for non linear problems. However, the method requires us to work in order k jet spaces in d variables; the mere storage of such a jet involves $\binom{d+k}{d}$ coefficients. From a theoretical point of view it is also not established that Taylor methods eliminate the wrapping effect by themselves. Nevertheless, Taylor models can be combined with any of the above methods and the non linear enclosures seem to be more adapted in certain cases. For a detailed critical analysis, we refer to [Neu02].

Let us finally investigate how sharp good error bounds actually may be. If $\rho_{f(0)} = \|f(0)_r\|/\|f(0)\|_c$ denotes the relative precision of the initial condition at the start and $\rho_{f(t)} = \|f(t)_r\|/\|f(t)\|_c$ the relative precision of the final result, then it is classical that

$$\begin{aligned} \rho_{f(t)} &\geq \kappa(J_\Delta(f(0))) \rho_{f(0)}, \\ \kappa(J_\Delta(f(0))) &= \|J_\Delta(f(0))\| \|J_\Delta(f(0))^{-1}\|, \end{aligned}$$

where $\kappa(J_\Delta(f(0)))$ is the condition number of the matrix $J_\Delta(f(0))$. We propose to define the condition number $\kappa(\Phi, f(0), 0, t)$ for the integration problem (48) on $[0, t]$ by

$$K = \kappa(\Phi, f(0), 0, t) = \max_{0 \leq t_1 \leq t_2 \leq t} \kappa(J_{\Delta_{t_1, t_2}}(f(t_1))).$$

Indeed, without using any particular mathematical properties of Φ or f , we somehow have to go through the whole interval $[0, t]$. Of course, it may happen that Φ is indeed special. For instance, if $\Phi = Mf$ for a matrix M with a special triangular or diagonal form, then special arguments may be used to improve error bounds and more dedicated condition numbers will have to be introduced.

However, in general, we suspect that a precision loss of $\log K$ cannot be avoided. If K gets large, the only real way to achieve long term stability is to take $p > 2 \log K$ (say), whence the interest of efficient multiple precision and high order ball algorithms. It seems to us that the parallelepiped method should manage to keep the precision loss bounded by $\log K$, for $p > 2 \log K$ and $\varepsilon \approx \epsilon_p$. The algorithm from section 7.5 (without coordinate transformations) only achieves a $\log k \log K$ in the worst case, although a $\log K$ bound is probably obtained in many cases of interest. We plan to carry out a more detailed analysis once we will have finished a first efficient multiple precision implementation.

BIBLIOGRAPHY

- [Abe80] O. Aberth. *Computable analysis*. McGraw-Hill, New York, 1980.
- [AH83] G. Alefeld and J. Herzberger. *Introduction to interval analysis*. Academic Press, New York, 1983.
- [ANS08] ANSI/IEEE. IEEE standard for binary floating-point arithmetic. Technical report, ANSI/IEEE, New York, 2008. ANSI-IEEE Standard 754-2008. Revision of IEEE 754-1985, approved on June 12, 2008 by IEEE Standards Board.
- [ANS09] ANSI/IEEE. IEEE interval standard working group - p1788. <http://grouper.ieee.org/groups/1788/>, 2009.
- [BB85] E. Bishop and D. Bridges. *Foundations of constructive analysis*. Die Grundlehren der mathematischen Wissenschaften. Springer, Berlin, 1985.
- [BCO+06] A. Bostan, F. Chyzak, F. Ollivier, B. Salvy, É. Schost, and A. Sedoglavic. Fast computation of power series solutions of systems of differential equation. preprint, april 2006. submitted, 13 pages.

- [Ber98] M. Berz. Cosy infinity version 8 reference manual. Technical Report MSUCL-1088, Michigan State University, East Lansing, 1998.
- [BF00] D. A. Bini and G. Fiorentino. Design, analysis, and implementation of a multiprecision polynomial rootfinder. *Numerical Algorithms*, 23:127–173, 2000.
- [BHL00] D.H. Bailey, Y. Hida, and X.S. Li. QD, a C/C++/Fortran library for double-double and quad-double arithmetic. <http://crd.lbl.gov/~dhbailey/mpdist/>, 2000.
- [BHL01] D.H. Bailey, Y. Hida, and X.S. Li. Algorithms for quad-double precision floating point arithmetic. In *15th IEEE Symposium on Computer Arithmetic*, pages 155–162, June 2001.
- [BK78] R.P. Brent and H.T. Kung. Fast algorithms for manipulating formal power series. *Journal of the ACM*, 25:581–595, 1978.
- [BS83] Walter Baur and Volker Strassen. The complexity of partial derivatives. *Theor. Comput. Sci.*, 22:317–330, 1983.
- [CK91] D.G. Cantor and E. Kaltofen. On fast multiplication of polynomials over arbitrary algebras. *Acta Informatica*, 28:693–701, 1991.
- [CT65] J.W. Cooley and J.W. Tukey. An algorithm for the machine calculation of complex Fourier series. *Math. Computat.*, 19:297–301, 1965.
- [DGG+02a] J.-G. Dumas, T. Gautier, M. Giesbrecht, P. Giorgi, B. Hovinen, E. Kaltofen, B. D. Saunders, W. J. Turner, and G. Villard. Linbox: A generic library for exact linear algebra. In *First Internat. Congress Math. Software ICMS*, pages 40–50, Beijing, China, 2002.
- [DGG+02b] J.-G. Dumas, T. Gautier, M. Giesbrecht, P. Giorgi, B. Hovinen, E. Kaltofen, B. D. Saunders, W. J. Turner, and G. Villard. Project linbox: Exact computational linear algebra. <http://www.linalg.org/>, 2002.
- [ea67] U.W. Kulisch et al. Scientific computing with validation, arithmetic requirements, hardware solution and language support. <http://www.math.uni-wuppertal.de/~xsc/>, 1967.
- [EKW84] J.-P. Eckmann, H. Koch, and P. Wittwer. A computer-assisted proof of universality in area-preserving maps. *Memoirs of the AMS*, 47(289), 1984.
- [EMOK91] J.-P. Eckmann, A. Malaspinas, and S. Oliffson-Kamphorst. A software tool for analysis in function spaces. In K. Meyer and D. Schmidt, editors, *Computer aided proofs in analysis*, pages 147–166. Springer, New York, 1991.
- [Für07] M. Fürer. Faster integer multiplication. In *Proceedings of the Thirty-Ninth ACM Symposium on Theory of Computing (STOC 2007)*, pages 57–66, San Diego, California, 2007.
- [GG02] J. von zur Gathen and J. Gerhard. *Modern Computer Algebra*. Cambridge University Press, 2-nd edition, 2002.
- [Gra91] T. Granlund et al. GMP, the GNU multiple precision arithmetic library. <http://www.swox.com/gmp>, 1991.
- [Grz57] A. Grzegorzcyk. On the definitions of computable real continuous functions. *Fund. Math.*, 44:61–71, 1957.
- [GS88] T.N. Gambill and R.D. Skeel. Logarithmic reduction of the wrapping effect with application to ordinary differential equations. *SIAM J. Numer. Anal.*, 25(1):153–162, 1988.
- [Hai95] B. Haible. CLN, a Class Library for Numbers. <http://www.ginac.de/CLN/cln.html>, 1995.
- [HLRZ00] G. Hanrot, V. Lefèvre, K. Ryde, and P. Zimmermann. MPFR, a C library for multiple-precision floating-point computations with exact rounding. <http://www.mpfr.org>, 2000.
- [HS67] E. Hansen and R. Smith. Interval arithmetic in matrix computations, part ii. *Siam J. Numer. Anal.*, 4:1–9, 1967.
- [JKDW01] L. Jaulin, M. Kieffer, O. Didrit, and E. Walter. *Applied interval analysis*. Springer, London, 2001.
- [K8] W. Kühn. Rigorously computed orbits of dynamical systems without the wrapping effect. *Computing*, 61:47–67, 1998.
- [KO63] A. Karatsuba and J. Ofman. Multiplication of multidigit numbers on automata. *Soviet Physics Doklady*, 7:595–596, 1963.

- [**KR06**] V. Kreinovich and S. Rump. Towards optimal use of multi-precision arithmetic: a remark. Technical Report UTEP-CS-06-01, UTEP-CS, 2006.
- [**Kul08**] U.W. Kulisch. *Computer Arithmetic and Validity. Theory, Implementation, and Applications*. Number 33 in Studies in Mathematics. de Gruyter, 2008.
- [**Loh88**] R. Lohner. *Einschließung der Lösung gewöhnlicher Anfangs- und Randwertaufgaben und Anwendungen*. PhD thesis, Universität Karlsruhe, 1988.
- [**Loh01**] R. Lohner. On the ubiquity of the wrapping effect in the computation of error bounds. In U. Kulisch, R. Lohner, and A. Facius, editors, *Perspectives of enclosure methods*, pages 201–217, Wien, New York, 2001. Springer.
- [**M0**] N. Müller. iRRAM, exact arithmetic in C++. <http://www.informatik.uni-trier.de/iRRAM/>, 2000.
- [**MB96**] K. Makino and M. Berz. Remainder differential algebras and their applications. In M. Berz, C. Bischof, G. Corliss, and A. Griewank, editors, *Computational differentiation: techniques, applications and tools*, pages 63–74, SIAM, Philadelphia, 1996.
- [**MB04**] K. Makino and M. Berz. Suppression of the wrapping effect by Taylor model-based validated integrators. Technical Report MSU Report MSUHEP 40910, Michigan State University, 2004.
- [**Moo65**] R.E. Moore. Automatic local coordinate transformations to reduce the growth of error bounds in interval computation of solutions to ordinary differential equation. In L.B. Rall, editor, *Error in Digital Computation*, volume 2, pages 103–140. John Wiley, 1965.
- [**Moo66**] R.E. Moore. *Interval analysis*. Prentice Hall, Englewood Cliffs, N.J., 1966.
- [**Mul06**] Jean-Michel Muller. *Elementary Functions, Algorithms and Implementation*. Birkhauser Boston, Inc., 2006.
- [**Neu90**] A. Neumaier. *Interval methods for systems of equations*. Cambridge university press, Cambridge, 1990.
- [**Neu93**] A. Neumaier. The wrapping effect, ellipsoid arithmetic, stability and confidence regions. *Computing Supplementum*, 9:175–190, 1993.
- [**Neu02**] A. Neumaier. Taylor forms - use and limits. *Reliable Computing*, 9:43–79, 2002.
- [**Nic85**] K. Nickel. How to fight the wrapping effect. In Springer-Verlag, editor, *Proc. of the Intern. Symp. on interval mathematics*, pages 121–132, 1985.
- [**Pan84**] V. Pan. *How to multiply matrices faster*, volume 179 of *Lect. Notes in Math*. Springer, 1984.
- [**Rev01**] N. Revol. MPFI, a multiple precision interval arithmetic library. <http://perso.ens-lyon.fr/nathalie.revol/software.html>, 2001.
- [**Rum99a**] S.M. Rump. Fast and parallel interval arithmetic. *BIT*, 39(3):534–554, 1999.
- [**Rum99b**] S.M. Rump. INTLAB - INTerval LABoratory. In Tibor Csendes, editor, *Developments in Reliable Computing*, pages 77–104. Kluwer Academic Publishers, Dordrecht, 1999. <http://www.ti3.tu-harburg.de/rump/>.
- [**SS71**] A. Schönhage and V. Strassen. Schnelle Multiplikation grosser Zahlen. *Computing*, 7:281–292, 1971.
- [**Str69**] V. Strassen. Gaussian elimination is not optimal. *Numer. Math.*, 13:352–356, 1969.
- [**Tur36**] A. Turing. On computable numbers, with an application to the Entscheidungsproblem. *Proc. London Maths. Soc.*, 2(42):230–265, 1936.
- [**vdH02a**] J. van der Hoeven. Relax, but don't be too lazy. *JSC*, 34:479–542, 2002.
- [**vdH+02b**] J. van der Hoeven et al. Mathemagix, 2002. <http://www.mathemagix.org>.
- [**vdH06a**] J. van der Hoeven. Computations with effective real numbers. *TCS*, 351:52–60, 2006.
- [**vdH06b**] J. van der Hoeven. Newton's method and FFT trading. Technical Report 2006-17, Univ. Paris-Sud, 2006. Submitted to JSC.
- [**vdH07a**] J. van der Hoeven. New algorithms for relaxed multiplication. *JSC*, 42(8):792–802, 2007.
- [**vdH07b**] J. van der Hoeven. On effective analytic continuation. *MCS*, 1(1):111–175, 2007.
- [**vdH08a**] J. van der Hoeven. Making fast multiplication of polynomials numerically stable. Technical Report 2008-02, Université Paris-Sud, Orsay, France, 2008.

-
- [vdH08b] J. van der Hoeven. Meta-expansion of transseries. In Y. Boudabbous and N. Zaguia, editors, *ROGICS '08: Relations Orders and Graphs, Interaction with Computer Science*, pages 390–398, Mahdia, Tunisia, May 2008.
- [vdHMT] J. van der Hoeven, B. Mourrain, and Ph. Trébuchet. Efficient and reliable resolution of eigenproblems. In preparation.
- [WC87] Winograd and Coppersmith. Matrix multiplication via arithmetic progressions. In *Proc. of the 19th Annual Symposium on Theory of Computing*, pages 1–6, New York City, may 25–27 1987.
- [Wei00] K. Weihrauch. *Computable analysis*. Springer-Verlag, Berlin/Heidelberg, 2000.