



HAL
open science

Efficient implementation of interval matrix multiplication

Hong Diep Nguyen

► **To cite this version:**

Hong Diep Nguyen. Efficient implementation of interval matrix multiplication. Para 2010: State of the Art in Scientific and Parallel Computing, Jun 2010, Reykjavik, Iceland. inria-00469472

HAL Id: inria-00469472

<https://inria.hal.science/inria-00469472v1>

Submitted on 1 Apr 2010

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

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

Efficient implementation of interval matrix multiplication

NGUYEN Hong Diep
INRIA
LIP(UMR 5668 CNRS - ENS de Lyon - INRIA - UCBL)
Université de Lyon
hong.diep.nguyen@ens-lyon.fr

April 1, 2010

Abstract

The straightforward implementation of interval matrix product suffers from poor efficiency, far from the performances of highly optimized floating-point implementations. In this paper, we show how to reduce the interval matrix multiplication to 9 floating-point matrix products - for performance issues - without sacrificing the quality of the result. We show that, compared to the straightforward implementation, the overestimation factor is at most 1.18.

1 Introduction

Interval arithmetic is a means to obtain guaranteed results: enclosures of the results are computed. Nonetheless, it suffers from lower performance than non-guaranteed floating-point computations. Indeed, the factor of performance between interval arithmetic and floating-point arithmetic in theory is four. It is even much worse in practice, especially when a large number of operations is involved.

In this paper we will study the case of matrix multiplication. Our implementation, based on the natural algorithm, provides good results at a cost of very low performance. The idea is to exploit existing libraries which are well optimized for floating-point matrix operations such as BLAS, ATLAS, etc. Rump proposed a fast algorithm which uses only four floating-point matrix products [2]. It returns a result wider than the result obtained by replacing each floating-point operation between two numbers by its interval counterpart: the factor of overestimation in the worst case of this algorithm is 1.5. This paper proposes a new algorithm which costs **nine** floating-point matrix products with the factor of overestimation in the worst case of 1.18.

This paper is organized as follows. Section 2 briefly presents interval arithmetic with some basic operations. These are extended to matrix operations which are studied in section 3. Section 4 reminds the idea of Rump algorithm. Section 5 is devoted to our proposed algorithm.

Notations

In this paper, bold-face lower-case letters represent scalar intervals and bold-face upper-case letters represent interval matrices. Below are some additional notations.

- $[o]$ an expression computed by interval arithmetic,
- $(o)_\downarrow, (o)_\uparrow$ expressions computed by floating-point arithmetic with downward and upward rounding mode respectively,
- $[i, s]$ an interval bounded by i and s ,
- $\{m, r\}$ an interval whose mid-point is m and radius is r ,
- $\underline{\mathbf{a}}, \overline{\mathbf{a}}$ lower and upper bound of \mathbf{a} ,
- $\text{mag}(\mathbf{a})$ maximal magnitude of \mathbf{a} : $\text{mag}(\mathbf{a}) \stackrel{def}{=} \max\{|a|, a \in \mathbf{a}\} = \max(|\underline{\mathbf{a}}|, |\overline{\mathbf{a}}|)$

2 Interval Arithmetic

Interval are used to represent connected closed sets of real values. Interval arithmetic defines operations between intervals. The result of an operation between intervals is also an interval containing all the possible results between all possible pairs of real values taken from input intervals: $\mathbf{r} = \mathbf{a}[o]\mathbf{b} = \{a \circ b, a \in \mathbf{a}, b \in \mathbf{b}\}$.

Due to the monotonicity property, the sum of two intervals $\mathbf{r} = \mathbf{a} + \mathbf{b}$ can be computed by the sum of their two respective lower and upper bounds:

$$\underline{\mathbf{r}} = \underline{\mathbf{a}} + \underline{\mathbf{b}} \quad \overline{\mathbf{r}} = \overline{\mathbf{a}} + \overline{\mathbf{b}} \quad (1)$$

Interval multiplication $\mathbf{r} = \mathbf{a} * \mathbf{b}$ is formed by taking the minimum and maximum value of the four products between two pairs of bounds of the two input intervals.

$$\begin{cases} \underline{\mathbf{r}} &= \min(\underline{\mathbf{a}} * \underline{\mathbf{b}}, \underline{\mathbf{a}} * \overline{\mathbf{b}}, \overline{\mathbf{a}} * \underline{\mathbf{b}}, \overline{\mathbf{a}} * \overline{\mathbf{b}}) \\ \overline{\mathbf{r}} &= \max(\underline{\mathbf{a}} * \underline{\mathbf{b}}, \underline{\mathbf{a}} * \overline{\mathbf{b}}, \overline{\mathbf{a}} * \underline{\mathbf{b}}, \overline{\mathbf{a}} * \overline{\mathbf{b}}) \end{cases} \quad (2)$$

Hence, in theory the factor, in terms of performance, between interval arithmetic and floating-point arithmetic is 4.

Implementation

Intervals are implemented on computers using floating-point numbers. To ensure the inclusion property of the results, rounding errors must be taken into account.

For interval addition $\mathbf{r} = \mathbf{a} + \mathbf{b}$, the lower bound must be computed with downward rounding mode, and the upper bound with upward rounding mode.

$$\underline{\mathbf{r}} = (\underline{\mathbf{a}} + \underline{\mathbf{b}})_\downarrow \quad \bar{\mathbf{r}} = (\bar{\mathbf{a}} + \bar{\mathbf{b}})_\uparrow \quad (3)$$

Interval product $\mathbf{r} = \mathbf{a} * \mathbf{b}$ is computed following (2) by four floating-point products. Nevertheless, to tackle rounding errors, each floating-point product must be computed twice, once with upward and once with downward rounding mode. Thus in total, it requires eight floating-point products.

$$\begin{cases} \underline{\mathbf{r}} &= \min((\underline{\mathbf{a}} * \underline{\mathbf{b}})_\downarrow, (\underline{\mathbf{a}} * \bar{\mathbf{b}})_\downarrow, (\bar{\mathbf{a}} * \underline{\mathbf{b}})_\downarrow, (\bar{\mathbf{a}} * \bar{\mathbf{b}})_\downarrow) \\ \bar{\mathbf{r}} &= \max((\underline{\mathbf{a}} * \underline{\mathbf{b}})_\uparrow, (\underline{\mathbf{a}} * \bar{\mathbf{b}})_\uparrow, (\bar{\mathbf{a}} * \underline{\mathbf{b}})_\uparrow, (\bar{\mathbf{a}} * \bar{\mathbf{b}})_\uparrow) \end{cases} \quad (4)$$

We can reduce the number of floating-point products by inspecting the sign of each component. But testing the sign is costly also.

Particular cases

If \mathbf{a} and \mathbf{b} are centered in zero

ie $\underline{\mathbf{a}} = -\bar{\mathbf{a}}$ and $\underline{\mathbf{b}} = -\bar{\mathbf{b}}$ then $\underline{\mathbf{a}} + \underline{\mathbf{b}} = -(\bar{\mathbf{a}} + \bar{\mathbf{b}})$. Hence, the sum $\mathbf{r} = \mathbf{a} + \mathbf{b}$ can be computed by:

$$\bar{\mathbf{r}} = (\bar{\mathbf{a}} + \bar{\mathbf{b}})_\uparrow \quad \underline{\mathbf{r}} = -\bar{\mathbf{r}} \quad (5)$$

If \mathbf{b} is centered in zero

then $\max(x * \underline{\mathbf{b}}, x * \bar{\mathbf{b}}) = |x| * \bar{\mathbf{b}}$ and $\min(x * \underline{\mathbf{b}}, x * \bar{\mathbf{b}}) = -|x| * \bar{\mathbf{b}}$ for all $x \in \mathbb{R}$. Hence according to (2), $\mathbf{r} = \mathbf{a} * \mathbf{b}$ can be computed by:

$$\begin{aligned} \bar{\mathbf{r}} &= \max(|\underline{\mathbf{a}}|, |\bar{\mathbf{a}}|) * \bar{\mathbf{b}} \\ &= \text{mag}(\mathbf{a}) * \bar{\mathbf{b}} \\ \underline{\mathbf{r}} &= -\max(|\underline{\mathbf{a}}|, |\bar{\mathbf{a}}|) * \bar{\mathbf{b}} \\ &= -\bar{\mathbf{r}} \end{aligned}$$

Using floating-point arithmetic, \mathbf{r} can be computed by only one floating-point product in upward rounding mode:

$$\bar{\mathbf{r}} = (\text{mag}(\mathbf{a}) * \bar{\mathbf{b}})_\uparrow \quad \underline{\mathbf{r}} = -\bar{\mathbf{r}} \quad (6)$$

3 Interval matrix operations

Let's now study the case of matrix operations. Suppose that each interval matrix here is of dimension $n \times n$ and is represented by two floating-point matrices, one for its lower bound and the other for its upper bound.

The addition of two interval matrices $\mathbf{C} = \mathbf{A} + \mathbf{B}$ can be computed by performing scalar additions between corresponding elements of the two matrices. Thus

$$\underline{\mathbf{C}} = (\underline{\mathbf{A}} + \underline{\mathbf{B}})_\downarrow \quad \bar{\mathbf{C}} = (\bar{\mathbf{A}} + \bar{\mathbf{B}})_\uparrow$$

For the case of interval matrix multiplication, each element of the result matrix is computed by:

$$\mathbf{C}_{i,j} = \sum_k \mathbf{A}_{i,k} * \mathbf{B}_{k,j} \quad (7)$$

Using this natural formula, the computation of each product element requires n interval multiplications and n interval additions, or, following (3) and (4), it requires $8n$ floating-point products and $2n$ floating-point additions. Hence, the overall cost in theory is $8n_*^3 + 2n_+^3$ floating-point operations. Nonetheless, this cost does not take into account neither \min, \max functions nor the cost of rounding mode changes.

Particular cases

A is centered in zero

ie $\mathbf{A}_{i,k}$ is centered in zero for all i, k . From (6) we get

$$\begin{aligned} \mathbf{A}_{i,k} * \mathbf{B}_{k,j} &= \left[-\overline{\mathbf{A}}_{i,k} * \text{mag}(\mathbf{B})_{k,j}, \overline{\mathbf{A}}_{i,k} * \text{mag}(\mathbf{B})_{k,j} \right] \\ \Rightarrow \mathbf{C}_{i,j} &= \left[-\sum_k \overline{\mathbf{A}}_{i,k} * \text{mag}(\mathbf{B})_{k,j}, \sum_k \overline{\mathbf{A}}_{i,k} * \text{mag}(\mathbf{B})_{k,j} \right] \\ \Rightarrow \overline{\mathbf{C}} &= (\overline{\mathbf{A}} * \text{mag}(\mathbf{B}))_{\uparrow} \\ \underline{\mathbf{C}} &= -\overline{\mathbf{C}} \end{aligned}$$

Hence $\mathbf{A} * \mathbf{B}$ can be computed by one floating-point matrix product with upward rounding mode.

A is non-negative

It means that $0 \leq \underline{\mathbf{A}}_{i,k}$ and $0 \leq \overline{\mathbf{A}}_{i,k}$ for all i, k . Hence

$$\begin{cases} \max(\underline{\mathbf{A}}_{i,k} * \underline{\mathbf{B}}_{k,j}, \underline{\mathbf{A}}_{i,k} * \overline{\mathbf{B}}_{k,j}) &= \underline{\mathbf{A}}_{i,k} * \overline{\mathbf{B}}_{k,j} \\ \min(\underline{\mathbf{A}}_{i,k} * \underline{\mathbf{B}}_{k,j}, \underline{\mathbf{A}}_{i,k} * \overline{\mathbf{B}}_{k,j}) &= \underline{\mathbf{A}}_{i,k} * \underline{\mathbf{B}}_{k,j} \\ \max(\overline{\mathbf{A}}_{i,k} * \underline{\mathbf{B}}_{k,j}, \overline{\mathbf{A}}_{i,k} * \overline{\mathbf{B}}_{k,j}) &= \overline{\mathbf{A}}_{i,k} * \overline{\mathbf{B}}_{k,j} \\ \min(\overline{\mathbf{A}}_{i,k} * \underline{\mathbf{B}}_{k,j}, \overline{\mathbf{A}}_{i,k} * \overline{\mathbf{B}}_{k,j}) &= \overline{\mathbf{A}}_{i,k} * \underline{\mathbf{B}}_{k,j} \end{cases}$$

Denote by $x^{pos} = \max(x, 0)$ and $x^{neg} = \min(x, 0)$ then

$$\begin{aligned} \begin{cases} \max(\underline{\mathbf{A}}_{i,k} * \overline{\mathbf{B}}_{k,j}, \overline{\mathbf{A}}_{i,k} * \overline{\mathbf{B}}_{k,j}) &= \overline{\mathbf{A}}_{i,k} * \overline{\mathbf{B}}_{k,j}^{pos} + \underline{\mathbf{A}}_{i,k} * \overline{\mathbf{B}}_{k,j}^{neg} \\ \min(\underline{\mathbf{A}}_{i,k} * \underline{\mathbf{B}}_{k,j}, \overline{\mathbf{A}}_{i,k} * \underline{\mathbf{B}}_{k,j}) &= \underline{\mathbf{A}}_{i,k} * \underline{\mathbf{B}}_{k,j}^{pos} + \overline{\mathbf{A}}_{i,k} * \underline{\mathbf{B}}_{k,j}^{neg} \end{cases} \\ \Rightarrow \begin{cases} \overline{\mathbf{A}}_{i,k} * \overline{\mathbf{B}}_{k,j} &= \overline{\mathbf{A}}_{i,k} * \overline{\mathbf{B}}_{k,j}^{pos} + \underline{\mathbf{A}}_{i,k} * \overline{\mathbf{B}}_{k,j}^{neg} \\ \underline{\mathbf{A}}_{i,k} * \underline{\mathbf{B}}_{k,j} &= \underline{\mathbf{A}}_{i,k} * \underline{\mathbf{B}}_{k,j}^{pos} + \overline{\mathbf{A}}_{i,k} * \underline{\mathbf{B}}_{k,j}^{neg} \end{cases} \\ \Rightarrow \begin{cases} \overline{\mathbf{C}} &= (\overline{\mathbf{A}} * \overline{\mathbf{B}}^{pos} + \underline{\mathbf{A}} * \overline{\mathbf{B}}^{neg})_{\uparrow} \\ \underline{\mathbf{C}} &= (\underline{\mathbf{A}} * \underline{\mathbf{B}}^{pos} + \overline{\mathbf{A}} * \underline{\mathbf{B}}^{neg})_{\downarrow} \end{cases} \end{aligned}$$

In this case, an interval matrix product can be computed by 4 floating-point matrix products.

Similarly, if \mathbf{A} is non positive, 4 floating-point matrix products suffice to compute an interval matrix product.

A does not contain zero in its interior

We can then split \mathbf{A} into positive and negative part by

$$\begin{aligned} \mathbf{a}^{pos} &= \begin{cases} \{a \in \mathbf{a} \mid a \geq 0\} & \text{if } \bar{\mathbf{a}} \geq 0 \\ 0 & \text{if } \bar{\mathbf{a}} < 0 \end{cases} \\ \mathbf{a}^{neg} &= \begin{cases} \{a \in \mathbf{a} \mid a \leq 0\} & \text{if } \underline{\mathbf{a}} \leq 0 \\ 0 & \text{if } \underline{\mathbf{a}} > 0 \end{cases} \end{aligned}$$

\mathbf{A} does not contain zero in its interior, thus either $\mathbf{A}_{i,j}^{pos} = \mathbf{A}_{i,j}$, $\mathbf{A}_{i,j}^{neg} = 0$ or $\mathbf{A}_{i,j}^{pos} = 0$, $\mathbf{A}_{i,j}^{neg} = \mathbf{A}_{i,j}$. Hence

$$\begin{aligned} \mathbf{A}_{i,j} * \mathbf{B}_{j,k} &= \mathbf{A}_{i,j} * \mathbf{B}_{j,k}^{pos} + \mathbf{A}_{i,j} * \mathbf{B}_{j,k}^{neg} \\ \Rightarrow \mathbf{A} * \mathbf{B} &= \mathbf{A}^{pos} * \mathbf{B} + \mathbf{A}^{neg} * \mathbf{B} \end{aligned}$$

In total, it costs eight floating-point matrix products.

For these three particular cases, if all operations are performed in infinite precision or if there is no rounding error, then the computed result is exact.

In general, when there is no assumption about the input intervals then it would not be efficient to use the natural algorithm. One solution is to exploit floating-point matrix multiplication because there are available libraries which are well optimised for floating-point matrix operations.

4 Rump's algorithm

Rump proposes an algorithm which makes use of floating-point operations for speed [2, 3]. This algorithm is based on the midpoint-radius representation of intervals.

Let $\mathbf{A} = \{m_A, r_A\}$ and $\mathbf{B} = \{m_B, r_B\}$ be two interval matrices, with m_A, m_B their midpoints and r_A, r_B their radius respectively. The product $\mathbf{A} * \mathbf{B}$ is enclosed by an interval matrix \mathbf{C} whose midpoint and radius are computed by:

$$\begin{aligned} m_C &= m_A * m_B \\ r_C &= (|m_A| + r_A) * r_B + r_A * |m_B| \end{aligned}$$

In fact, because of rounding errors, $m_A * m_B$ cannot be computed exactly. Thus it must be computed with both upward and downward rounding mode to obtain an enclosure of the midpoint. The radius r_C must also be computed with upward rounding to ensure the inclusion property.

$$\begin{aligned} \mathbf{mC} &= \left[(m_A * m_B)_\downarrow, (m_B * m_B)_\uparrow \right] \\ r_C &= ((|m_A| + r_A) * r_B + r_A * |m_B|)_\uparrow \end{aligned}$$

Finally, the result can be easily converted back to endpoints representation by two floating-point additions.

$$\overline{\mathbf{C}} = (\overline{\mathbf{mC}} + r_C)_\uparrow \quad \underline{\mathbf{C}} = (\underline{\mathbf{mC}} - r_C)_\downarrow$$

In total, this algorithm uses **four** floating-point products.

Over-estimation

Without taking into account of rounding errors, Rump's algorithm always provides over-estimating results. The factor of over-estimation in the worst case is 1.5.

For example, by definition $[0, 2] * [0, 4] \stackrel{def}{=} [0, 8]$.

Meanwhile, Rump's algorithm gives

$$\begin{aligned} [0, 2] * [0, 4] &= \{1, 1\} * \{2, 2\} \\ &= \{1 * 2, (1 + 1) * 2 + 1 * 2\} \\ &= \{2, 6\} \\ &= [-4, 8] \end{aligned}$$

The over-estimation factor in this case is $12/8 = 1.5$.

5 Proposition

Rump algorithm can be considered as decomposing \mathbf{A} and \mathbf{B} into a sum of two components representing its midpoint and radius respectively. $\mathbf{A} * \mathbf{B}$ is then replaced by its development, which is a sum of four sub-products. Due to the sub-distributive property of interval product, the result yielded by this sum is an enclosure of the original product.

Our idea is to find another decomposition such that sub-products can be efficiently computed, and the over-estimation is small.

As we can see in Section 3, an interval product can be efficiently computed when one of the two multipliers is centered in zero or does not contain zero.

Proposition 1. *Let \mathbf{A} be an interval matrix. If \mathbf{A} is decomposed into two interval matrices \mathbf{A}^0 and \mathbf{A}^* which satisfy:*

- $\mathbf{A}_{i,j}^0 = 0, \mathbf{A}_{i,j}^* = \mathbf{A}_{i,j}$ if $\underline{\mathbf{A}}_{i,j} * \overline{\mathbf{A}}_{i,j} \geq 0$,
- $\mathbf{A}_{i,j}^0 = [\underline{\mathbf{A}}_{i,j}, -\underline{\mathbf{A}}_{i,j}], \mathbf{A}_{i,j}^* = [0, \underline{\mathbf{A}}_{i,j} + \mathbf{A}_{i,j}]$ if $\underline{\mathbf{A}}_{i,j} < 0 < |\underline{\mathbf{A}}_{i,j}| \leq \overline{\mathbf{A}}_{i,j}$,
- $\mathbf{A}_{i,j}^0 = [-\overline{\mathbf{A}}_{i,j}, \overline{\mathbf{A}}_{i,j}], \mathbf{A}_{i,j}^* = [\underline{\mathbf{A}}_{i,j} + \mathbf{A}_{i,j}, 0]$ if $\underline{\mathbf{A}}_{i,j} < 0 < \overline{\mathbf{A}}_{i,j} < |\underline{\mathbf{A}}_{i,j}|$

then

- \mathbf{A}^0 is centered in zero,
- \mathbf{A}^* does not contain zero in its interior, and

- $\mathbf{A}^0 + \mathbf{A}^* = \mathbf{A}$.

Proof. Easily deduced from the formula. \square

Proposition 2. Let \mathbf{A} and \mathbf{B} be two interval matrices and $(\mathbf{A}^0, \mathbf{A}^*)$ a decomposition of \mathbf{A} by Proposition 1. Let \mathbf{C} be an interval matrix computed by

$$\mathbf{C} = \mathbf{A}^0 * \mathbf{B} + \mathbf{A}^* * \mathbf{B} \quad (8)$$

Then $\mathbf{A} * \mathbf{B}$ is contained in \mathbf{C} . Denote $\mathbf{C} \stackrel{def}{=} \mathbf{A} \otimes \mathbf{B}$.

Proof. Following Proposition 1:

$$\mathbf{A} = \mathbf{A}^0 + \mathbf{A}^*$$

Interval multiplication is sub-distributive, hence

$$\begin{aligned} \mathbf{A} * \mathbf{B} &\subseteq \mathbf{A}^0 * \mathbf{B} + \mathbf{A}^* * \mathbf{B} \\ &\subseteq \mathbf{C} \end{aligned}$$

\square

\mathbf{A}^0 is centered in zero and \mathbf{A}^* does not contain zero, so according to Section 3, $\mathbf{A}^0 * \mathbf{B}$ and $\mathbf{A}^* * \mathbf{B}$ can be computed using 1 and 8 floating-point matrix products respectively. Hence, the overall cost is 9 floating-point matrix products. The following section will study the over-estimation factor of this operation in the worst case.

Over-estimation factor

Let us first consider the product of two scalar intervals \mathbf{a} and \mathbf{b} with \mathbf{a} being decomposed, following Proposition 1, into two parts: $\mathbf{a} = \mathbf{a}^0 + \mathbf{a}^*$. Suppose that all calculations here are performed in infinite precision. It means that rounding errors will not be taken into account.

If \mathbf{a} is centered in zero

then $\mathbf{a}^* = 0 \rightarrow \mathbf{a} * \mathbf{b} = \mathbf{a}^0 * \mathbf{b}$. Thus the result is exact.

If \mathbf{a} does not contain zero in its interior

then $\mathbf{a}^0 = 0 \rightarrow \mathbf{a} * \mathbf{b} = \mathbf{a}^* * \mathbf{b}$. Thus the result is exact too.

If \mathbf{a} contains zero

It means that $\underline{\mathbf{a}} < 0 < \bar{\mathbf{a}}$. The case $\bar{\mathbf{a}} < |\underline{\mathbf{a}}|$ is contrary to the case $\bar{\mathbf{a}} > |\underline{\mathbf{a}}|$, hence without loss of generality we suppose that $\bar{\mathbf{a}} > |\underline{\mathbf{a}}|$. In this case $\mathbf{a}^0 = [\underline{\mathbf{a}}, -\underline{\mathbf{a}}]$ and $\mathbf{a}^* = [0, \bar{\mathbf{a}} + \underline{\mathbf{a}}]$. Following Proposition 2

$$\begin{aligned} \mathbf{a} \otimes \mathbf{b} &= \mathbf{a}^0 * \mathbf{b} + \mathbf{a}^* * \mathbf{b} \\ &= [\underline{\mathbf{a}}, -\underline{\mathbf{a}}] * \mathbf{b} + [0, \bar{\mathbf{a}} + \underline{\mathbf{a}}] * \mathbf{b} \\ &= [\underline{\mathbf{a}} * \text{mag}(\mathbf{b}), -\underline{\mathbf{a}} * \text{mag}(\mathbf{b})] + [0, \bar{\mathbf{a}} + \underline{\mathbf{a}}] * \mathbf{b} \end{aligned}$$

- If $\mathbf{b} \geq 0$, then $\text{mag}(\mathbf{b}) = \bar{\mathbf{b}}$

$$\begin{aligned} \mathbf{a} \otimes \mathbf{b} &= [\underline{\mathbf{a}} * \text{mag}(\mathbf{b}), -\underline{\mathbf{a}} * \text{mag}(\mathbf{b})] + [0, \bar{\mathbf{a}} + \underline{\mathbf{a}}] * \mathbf{b} \\ &= [\underline{\mathbf{a}} * \bar{\mathbf{b}}, -\underline{\mathbf{a}} * \bar{\mathbf{b}}] + [0, (\bar{\mathbf{a}} + \underline{\mathbf{a}}) * \bar{\mathbf{b}}] \\ &= [\underline{\mathbf{a}} * \bar{\mathbf{b}}, \bar{\mathbf{a}} * \bar{\mathbf{b}}] \end{aligned}$$

Meanwhile, $\mathbf{b} \geq 0 \rightarrow \mathbf{a} * \mathbf{b} = [\underline{\mathbf{a}} * \bar{\mathbf{b}}, \bar{\mathbf{a}} * \bar{\mathbf{b}}]$. Hence $\mathbf{a} \otimes \mathbf{b} = \mathbf{a} * \mathbf{b}$. It means that there is no overestimation.

- Idem for the case $\mathbf{b} \leq 0$.
- If $\bar{\mathbf{b}} > 0 > \underline{\mathbf{b}}$ and $\bar{\mathbf{b}} \geq |\underline{\mathbf{b}}|$ then $\text{mag}(\mathbf{b}) = \bar{\mathbf{b}}$. Hence

$$\begin{aligned} \mathbf{a} \otimes \mathbf{b} &= [\underline{\mathbf{a}} * \text{mag}(\mathbf{b}), -\underline{\mathbf{a}} * \text{mag}(\mathbf{b})] + [0, \bar{\mathbf{a}} + \underline{\mathbf{a}}] * \mathbf{b} \\ &= [\underline{\mathbf{a}} * \bar{\mathbf{b}}, -\underline{\mathbf{a}} * \bar{\mathbf{b}}] + [(\bar{\mathbf{a}} + \underline{\mathbf{a}}) * \underline{\mathbf{b}}, (\bar{\mathbf{a}} + \underline{\mathbf{a}}) * \bar{\mathbf{b}}] \\ &= [(\bar{\mathbf{a}} + \underline{\mathbf{a}}) * \underline{\mathbf{b}} + \underline{\mathbf{a}} * \bar{\mathbf{b}}, \bar{\mathbf{a}} * \bar{\mathbf{b}}] \end{aligned}$$

$$\begin{aligned} \Rightarrow \text{diam}(\mathbf{a} \otimes \mathbf{b}) &= \bar{\mathbf{a}} * \bar{\mathbf{b}} - ((\bar{\mathbf{a}} + \underline{\mathbf{a}}) * \underline{\mathbf{b}} + \underline{\mathbf{a}} * \bar{\mathbf{b}}) \\ &= \bar{\mathbf{a}} * \bar{\mathbf{b}} * (1 - \underline{\mathbf{a}}/\bar{\mathbf{a}} - \underline{\mathbf{b}}/\bar{\mathbf{b}} - \underline{\mathbf{a}}/\bar{\mathbf{a}} * \underline{\mathbf{b}}/\bar{\mathbf{b}}) \end{aligned}$$

As $\bar{\mathbf{a}} > 0 > \underline{\mathbf{a}} > -\bar{\mathbf{a}}$ and $\bar{\mathbf{b}} > 0 > \underline{\mathbf{b}} > -\bar{\mathbf{b}}$

$$\begin{aligned} \mathbf{a} * \mathbf{b} &= [\min(\bar{\mathbf{a}} * \underline{\mathbf{b}}, \underline{\mathbf{a}} * \bar{\mathbf{b}}), \bar{\mathbf{a}} * \bar{\mathbf{b}}] \\ \Rightarrow \text{diam}(\mathbf{a} * \mathbf{b}) &= \bar{\mathbf{a}} * \bar{\mathbf{b}} - \min(\bar{\mathbf{a}} * \underline{\mathbf{b}}, \underline{\mathbf{a}} * \bar{\mathbf{b}}) \\ &= \bar{\mathbf{a}} * \bar{\mathbf{b}} * (1 - \min(\underline{\mathbf{a}}/\bar{\mathbf{a}}, \underline{\mathbf{b}}/\bar{\mathbf{b}})) \end{aligned}$$

Denote by $\begin{cases} M &= \max(|\underline{\mathbf{a}}|/\bar{\mathbf{a}}, |\underline{\mathbf{b}}|/\bar{\mathbf{b}}) \\ m &= \min(|\underline{\mathbf{a}}|/\bar{\mathbf{a}}, |\underline{\mathbf{b}}|/\bar{\mathbf{b}}) \end{cases}$ then

$$\begin{cases} 0 < m \leq M \leq 1 \\ \min(\underline{\mathbf{a}}/\bar{\mathbf{a}}, \underline{\mathbf{b}}/\bar{\mathbf{b}}) &= -M \\ \underline{\mathbf{a}}/\bar{\mathbf{a}} + \underline{\mathbf{b}}/\bar{\mathbf{b}} &= -m - M \\ \underline{\mathbf{a}}/\bar{\mathbf{a}} * \underline{\mathbf{b}}/\bar{\mathbf{b}} &= m * M \end{cases}$$

The factor of overestimation is computed by:

$$\begin{aligned}
\frac{\text{diam}(\mathbf{a} \otimes \mathbf{b})}{\text{diam}(\mathbf{a} * \mathbf{b})} &= \frac{1 - \underline{\mathbf{a}}/\bar{\mathbf{a}} - \underline{\mathbf{b}}/\bar{\mathbf{b}} - \underline{\mathbf{a}}/\bar{\mathbf{a}} * \underline{\mathbf{b}}/\bar{\mathbf{b}}}{1 - \min(\underline{\mathbf{a}}/\bar{\mathbf{a}}, \underline{\mathbf{b}}/\bar{\mathbf{b}})} \\
&= \frac{1 + M + m - Mm}{1 + M} \\
&= \frac{1 + M + m(1 - M)}{1 + M} \\
&\leq \frac{1 + M + M(1 - M)}{1 + M}
\end{aligned}$$

Inspecting the last function of unknown M between 0 and 1, we have that its maximum is $4 - 2\sqrt{2} \approx 1.18$ and this maximum is reached for $m = M = \sqrt{2} - 1$, or $\underline{\mathbf{a}}/\bar{\mathbf{a}} = \underline{\mathbf{b}}/\bar{\mathbf{b}} = 1 - \sqrt{2}$.

- Idem for the case $\bar{\mathbf{b}} > 0 > \underline{\mathbf{b}}$ and $\bar{\mathbf{b}} < |\underline{\mathbf{b}}|$

Let's now extend to the case of matrix multiplication. Each element of the result matrix is computed by

$$\begin{aligned}
\mathbf{C}_{i,j} &= \sum_k \mathbf{A}_{i,k}^0 * \mathbf{B}_{k,j} + \sum_k \mathbf{A}_{i,k}^* * \mathbf{B}_{k,j} \\
&= \sum_k (\mathbf{A}_{i,k}^0 * \mathbf{B}_{k,j} + \mathbf{A}_{i,k}^* * \mathbf{B}_{k,j})
\end{aligned}$$

Moreover

$$\begin{aligned}
\text{diam}(\mathbf{a} + \mathbf{b}) &= (\bar{\mathbf{a}} + \bar{\mathbf{b}}) - (\underline{\mathbf{a}} + \underline{\mathbf{b}}) \\
&= (\bar{\mathbf{a}} - \underline{\mathbf{a}}) + (\bar{\mathbf{b}} - \underline{\mathbf{b}}) \\
&= \text{diam}(\mathbf{a}) + \text{diam}(\mathbf{b}) \\
\Rightarrow \text{diam}(\mathbf{C}_{i,j}) &= \text{diam} \left(\sum_k (\mathbf{A}_{i,k}^0 * \mathbf{B}_{k,j} + \mathbf{A}_{i,k}^* * \mathbf{B}_{k,j}) \right) \\
&= \sum_k \text{diam}(\mathbf{A}_{i,k}^0 * \mathbf{B}_{k,j} + \mathbf{A}_{i,k}^* * \mathbf{B}_{k,j})
\end{aligned}$$

With the assumption of no rounding error, the over-estimation factor in the worst case for this scalar product is $4 - 2\sqrt{2}$. Hence

$$\begin{aligned}
\text{diam}(\mathbf{C}_{i,j}) &\leq (4 - 2\sqrt{2}) * \sum_k \text{diam}(\mathbf{A}_{i,k} * \mathbf{B}_{k,j}) \\
&\leq (4 - 2\sqrt{2}) * \text{diam} \left(\sum_k \mathbf{A}_{i,k} * \mathbf{B}_{k,j} \right) \\
\Rightarrow \text{diam}(\mathbf{C}) &\leq (4 - 2\sqrt{2}) * \text{diam}(\mathbf{A} * \mathbf{B})
\end{aligned}$$

Hence, the over-estimation factor in the worst case of interval matrix product is also $4 - 2\sqrt{2} \approx 1.18$.

6 Conclusion

The algorithm presented in this paper implements the product of interval matrices using floating-point operations. It constitutes a trade-off between the performances of optimized floating-point libraries, and a slight overestimation of the result. We have proven that this overestimation factor is at most 1.18, which means that the width of computed result is always less than 1.18 times the width of exact result. In particular, if one of the two multipliers does not contain zero in its interior then the computed result is exact using exact arithmetic.

The performance of the algorithm given above relies entirely on the performance of the employed floating-point library. Such a library must support directed rounding modes to be of use. This requirement is not really an issue, since libraries such as LAPACK or ATLAS, which are renowned for their performances, support directed roundings.

However, the use of directed rounding modes restricts the algorithm that can be used as a basic algorithm to the natural one. Indeed, fast algorithms for matrix products [1], which have a complexity below $\mathcal{O}(n^3)$, do not preserve the monotonicity of operations (they use subtraction as well as addition and multiplication) and they cannot be used. Furthermore, such algorithms rely on algebraic properties such as $(x + y) - y = x$, which do not hold in interval arithmetic.

References

- [1] M. Ceberio and V. Kreinovich. Fast multiplication of interval matrices (interval version of strassen's algorithm). *Reliable Computing*, 10(3):241–243, 2004.
- [2] S. M. Rump. Fast and parallel interval arithmetic. *BIT*, 39(3):534–554, 1999.
- [3] S. M. Rump. Computer-assisted proofs and self-validating methods. In B. Einarsson, editor, *Handbook on Accuracy and Reliability in Scientific Computation*, chapter 10, pages 195–240. SIAM, 2005.