

# Formal Proofs of Rounding Error Bounds

With application to an automatic positive definiteness check.

Pierre Roux

Received: date / Accepted: date

**Abstract** Floating-point arithmetic is a very efficient solution to perform computations in the real field. However, it induces rounding errors making results computed in floating-point differ from what would be computed with reals. Although numerical analysis gives tools to bound such differences, the proofs involved can be painful, hence error prone. We thus investigate the ability of a proof assistant like Coq to mechanically check such proofs. We demonstrate two different results involving matrices, which are pervasive among numerical algorithms, and show that a large part of the development effort can be shared between them.

**Keywords** floating-point arithmetic · rounding error · numerical analysis · proof assistant · Coq · matrices · Cholesky decomposition

## 1 Introduction

Floating-point arithmetic is a very efficient solution to perform computations in the real field  $\mathbb{R}$ . Unfortunately, intermediate results of computations need to be rounded to fit in the floating-point format used. Due to this rounding errors, final results of computations differ from what would have been obtained by computing in the real field  $\mathbb{R}$ , although both results usually remain pretty close.

Fortunately, each rounding can only introduce a bounded error. By combining these atomic errors, one can get a bound on the error affecting the final result. Numerical analysis [10] thus aims at bounding these differences between results of numerical algorithms using floating-point or real arithmetic. Using such mathematical properties, rigorous results can be obtained despite the use of floating-point arithmetic [14], ensuring that no disastrous rounding error can happen during the computation.

---

This work was done while the author was a visiting researcher at LRI, Inria Saclay – Île-de-France.

---

Pierre Roux  
ISAE, ONERA  
E-mail: pierre.roux@onera.fr

More precisely, denoting  $f$  a function in the real field  $\mathbb{R}$  and  $\tilde{f}$  its actually computed floating-point counterpart, we have  $\tilde{f}(x) = f(x) + e$ . The value  $e$  is called the *forward error* and is expected to be negligible in front of  $f(x)$ . When  $\tilde{f}(x) = f(x + d)$ ,  $d$  is called a *backward error*. The goal is then to prove some bound  $b$ , preferably as small as possible, such that  $|e| \leq b(x)$  (respectively  $|d| \leq b(x)$ ).

Proofs of this kind of mathematical results are hard to automate when they involve an arbitrary number of operations and are therefore mostly done by hand. However, they can be particularly painful and repetitive which make them specially error prone. That is why we want to investigate the ability of a proof assistant, namely Coq [1, 6], to check them. Matrices being pervasive in numerical algorithms, we will particularly focus on them.

Formal proofs of error bounds have already been performed with proof assistants such as HOL [9] or Coq [3] or with automatic tools such as Gappa [7]. Yet, to the extent of author's knowledge, those work only address results with a fixed number of basic arithmetic operations whereas algorithms with an arbitrary, parameterized, number of operations are targeted in this paper.

The paper is organized as follows. The remainder of this section first introduces our motivating example (Section 1.1) and a secondary example (Section 1.2), then gives basic properties of floating-point arithmetic (Section 1.3) and eventually details a simple proof about summations (Section 1.4). Section 2 then gives the detailed specification of floating-point arithmetic used while Section 3 shows how error terms can be combined. Section 4 eventually details proofs of numerical analysis results involving matrices, Section 5 deals with overflows, that were ignored until there, and Section 6 concludes.

### 1.1 Motivating Example: Cholesky Decomposition

We will use as motivating example throughout this paper a Cholesky decomposition which is a typical example of numerical algorithm involving matrices. Checking positive definiteness of matrices is one common use of Cholesky decomposition. A symmetric matrix  $A \in \mathbb{R}^{n \times n}$  is said to be positive semi-definite, written  $A \succeq 0$ , when for all  $x \in \mathbb{R}^n$ ,  $x^T A x \geq 0$ .

To prove that a scalar  $a \in \mathbb{R}$  is non negative, one can exhibit some  $r \in \mathbb{R}$  such that  $a = r^2$  (typically  $r = \sqrt{a}$ ). Similarly, one can prove that a symmetric matrix  $A \in \mathbb{R}^{n \times n}$  is positive semi-definite by exposing a matrix  $R$  such that  $A = R^T R$  (for  $x^T (R^T R) x = (R x)^T (R x) = \|R x\|_2^2 \geq 0$  for all  $x \in \mathbb{R}^n$ ). The Cholesky decomposition algorithm, shown on Figure 1, can compute such a matrix  $R$ , thus playing on matrices a similar role to the square root on reals. Indeed, taking as input a symmetric matrix  $A \in \mathbb{R}^{n \times n}$ , if the algorithm runs to completion, without ever attempting to take the square root of a negative value or perform a division by zero, it returns a matrix  $R$  such that  $A = R^T R$  (this can be proven by induction on  $n$ ). Thus, the mere termination of the Cholesky decomposition without error is enough to prove  $A \succeq 0$ . Conversely, the decomposition will terminate without error for all symmetric positive definite matrix  $A$  (i.e., such that for all  $x \in \mathbb{R}^n$ , if  $x \neq 0$  then  $x^T A x > 0$ ).

```

R := 0;
for j from 0 to  $n - 1$  do
  for i from 0 to  $j - 1$  do
     $R_{i,j} := (A_{i,j} - \sum_{k=0}^{i-1} R_{k,i}R_{k,j}) / R_{i,i};$ 
  od
   $R_{j,j} := \sqrt{A_{j,j} - \sum_{k=0}^{j-1} R_{k,j}^2};$ 
od

```

**Fig. 1** Cholesky decomposition: from a symmetric matrix  $A \in \mathbb{R}^{n \times n}$ , attempts to compute  $R$  such that  $A = R^T R$ . The resulting matrix  $R$  is upper-triangular (its elements  $R_{i,j}$  with  $i \leq j$  are the only one assigned after the initialization to 0). The elements  $R_{i,j}$  are computed from left to right (outer loop) and from top to bottom (inner loop), each element being assigned exactly once a value obtained from the corresponding element  $A_{i,j}$  of the input  $A$  and previously computed elements of  $R$  (i.e., elements on its top left). Due to the three nested loops (one being hidden in the  $\Sigma$  notation), the algorithm performs  $\Theta(n^3)$  arithmetic operations (that is  $n\sqrt{n}$  operations if  $n$  is the size of the input).

This would work perfectly if the algorithm were run using real arithmetic. However, performing it with floating-point arithmetic, it could run to completion while  $A \not\geq 0$ , due to rounding errors. That is, there could be  $x \in \mathbb{R}^n$  such that  $x^T A x < 0$ . But rounding errors remain bounded, so that there exists a bound  $c \in \mathbb{R}$  such that, if the floating-point Cholesky decomposition of  $A$  succeeds, then for all  $x \in \mathbb{R}^n$ ,  $x^T A x \geq -c \|x\|_2^2$ . That is  $A + cI \geq 0$ . The successful floating-point Cholesky decomposition of  $A - cI$  eventually proves<sup>1</sup> that  $A \geq 0$ . Moreover, such a constant  $c$  can be easily computed from simple characteristics of  $A$  and the floating-point arithmetic format used.

It can be noticed that the method is not complete. In particular, if  $A$  is positive semi-definite but not positive-definite (i.e., there is  $x \neq 0$  such that  $x^T A x = 0$ ), then  $A - cI$  is not positive semi-definite and the Cholesky decomposition will certainly fail. Thus, rather than just proving that a matrix  $A$  is positive semi-definite  $A \geq 0$ , the method could as well prove that  $A$  is positive definite  $A \succ 0$ . This is done by choosing  $c$  such that if the decomposition succeeds, then for all  $x$ ,  $x^T A x > -c \|x\|_2^2$  (strict inequality instead of the wide one above). This ability to prove only strict, and not wide, inequalities is a common characteristic of this kind of numerical methods, due to the rounding errors [14].

## 1.2 Secondary Example: Ellipsoidal Invariants for Linear Controllers

Figure 2 shows a code modeling a typical linear controller. One may want to prove that, as long as the value of  $u \in \mathbb{R}^p$  remains bounded ( $\|u\|_\infty \leq 1$ ), the value of  $x$  in  $\mathbb{R}^n$  will remain bounded along any execution of this code.

This is the case if and only if there exists a positive definite matrix  $P \in \mathbb{R}^{n \times n}$  and a scalar  $\lambda \in \mathbb{R}$  such that “ $x^T P x \leq \lambda$ ” is an inductive loop invariant. That is “ $x^T P x \leq \lambda$ ” holds before entering the loop (i.e., for  $x = 0$ , which amount to  $\lambda \geq 0$ ) and if it holds before an iteration of the loop (i.e., if  $x^T P x \leq \lambda$ ) then it still holds

<sup>1</sup> Indeed,  $A - cI$  leads to another bound  $c'$  but  $c' \leq c$  (c.f., Corollary 4.4, page 15 for details).

```

x := 0;
while true do
  // u models input read from sensors
  u := ?(-1, 1); // random value in  $\mathbb{R}^p$  with  $\|u\|_\infty \leq 1$ 
  x := Ax + Bu;
  // send orders to actuators and wait for next period
od

```

**Fig. 2** Model of a typical linear controller: such a controller can be used to maintain a physical system in a wanted state, when one can get some physical quantities of the system thanks to some sensors and has to act accordingly on the system through some actuators. The body of the infinite loop is repeatedly executed at periodic instants of time (for instance, every 10 milliseconds). At each iteration, the variable  $u$  first receives values from the sensors (this value being bounded, for instance between  $-1$  and  $1$ ). Then, the variable  $x$  is updated according to its value from the previous iteration and the value of  $u$ . Finally, the value of  $x$  could be used to send some orders to the actuators.

after the iteration (i.e.,  $(Ax + Bu)^T P (Ax + Bu) \leq \lambda$  for any  $u$  such that  $\|u\|_\infty \leq 1$ ). Denoting  $x = (x_0, \dots, x_{n-1})$ , the expression  $x^T P x$  is a polynomial of degree 2 in the  $n$  variables  $x_0, \dots, x_{n-1}$  and when  $P$  is positive definite, the set  $\{x \in \mathbb{R}^n \mid x^T P x \leq \lambda\}$  is, geometrically speaking, an ellipsoid. Since the variable  $x$  remains in the bounded set  $\mathcal{E}$  along any execution of the program, it remains bounded. All this is illustrated on Figure 3.

The goal will be to show that an invariant “ $x^T P x \leq \lambda$ ” for an ideal program executed with real arithmetic (i.e.,  $Ax + Bu$  computed in the real field) remains valid, under some conditions, for the actual program using floating-point arithmetic (i.e.,  $Ax + Bu$  computed with floating-point arithmetic).

### 1.3 Definitions and Basic Properties

**Definition 1.1**  $\mathbb{F} \subset \mathbb{R}$  denotes the set of floating point values,  $\circ : \mathbb{R} \rightarrow \mathbb{F}$  a rounding function (toward  $+\infty$  or to nearest for instance) and  $\text{fl}(e) \in \mathbb{F}$  the floating point evaluation of expression  $e$  from left to right<sup>2</sup>.

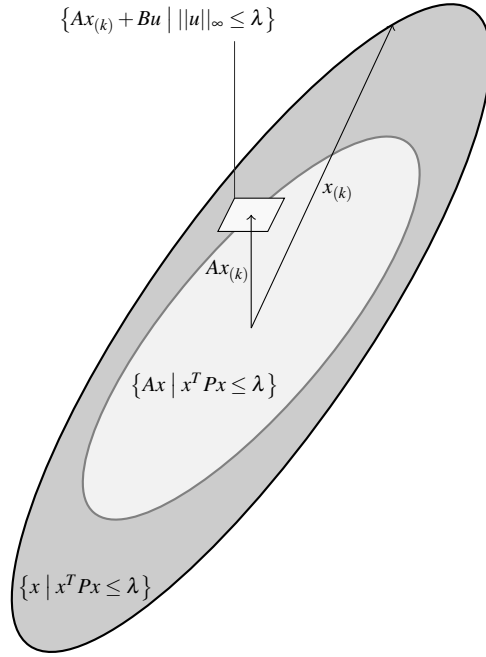
*Example 1.2* Assuming 1, 2 and 3 are floating-point values,  $\text{fl}(1 + 2 + 3)$  denotes the value  $\circ(\circ(1 + 2) + 3)$ .

In practice, a floating-point value  $f \in \mathbb{F}$  is encoded as a mantissa  $m \in \mathbb{Z}$  and an exponent  $e \in \mathbb{Z}$  such that  $f = m\beta^e$  where  $\beta$  is the radix (commonly 2). Since  $f$  is encoded with a constant number of bits,  $m$  lies in a bounded range. More precisely, if the mantissa  $m$  is encoded with a precision of  $prec$  bits:  $|m| < \beta^{prec}$ . To fully exploit the available precision,  $m$  and  $e$  can be chosen such that  $|m| \geq \beta^{prec-1}$ . Such numbers are called *normalized*.

The exponent  $e$  also lies in a bounded range. We will initially treat  $e$  as unbounded above which means that we ignore potential overflows<sup>3</sup>. In contrary,  $e$  is bounded below:  $e \geq e_{\min}$ . This leads to an underflow phenomena: numbers close to 0 (i.e., with

<sup>2</sup> Order of evaluation matters since floating point operations are not associative.

<sup>3</sup> They will be taken into account later (c.f., Section 5).



**Fig. 3** “ $x^T P x \leq \lambda$ ” is an inductive invariant: for any  $x_{(k)}$  is in the dark gray ellipse (i.e.,  $x_{(k)}^T P x_{(k)} \leq \lambda$ ),  $Ax_{(k)}$  is in the light gray one. The white box represents the potential values of  $Bu$ . Eventually,  $Ax_{(k)} + Bu$  is in the original dark gray ellipse (i.e.,  $(Ax_{(k)} + Bu)^T P (Ax_{(k)} + Bu) \leq \lambda$ ).

absolute value less than  $\beta^{e_{\min} + prec - 1}$ ) cannot be represented with the usual precision of  $prec$  bits. However, these numbers can be represented with a decreased precision when they are not too close from 0 (i.e., with absolute value larger than  $\beta^{e_{\min}}$ ). Such numbers are called *denormalized*.

Thus, a real number  $x \in \mathbb{R}$  not too close from 0 can be represented by a normalized number  $\circ(x) \in \mathbb{F}$  such that  $|x - \circ(x)| \leq \beta^{1-prec} |x|$  (or  $(\beta^{1-prec}/2) |x|$  if  $\circ$  is a rounding to nearest). Otherwise, a denormalized number has to be used, and only  $|x - \circ(x)| \leq \beta^{e_{\min}}$  can be guaranteed (or  $\beta^{e_{\min}}/2$  if  $\circ$  is a rounding to nearest).

**Definition 1.3**  $\text{eps} \in \mathbb{R}$  and  $\text{eta} \in \mathbb{R}$  are constants, depending on the floating-point format used, such that for all  $x \in \mathbb{R}$ ,  $\circ(x) \in \mathbb{F}$  satisfies either  $|x - \circ(x)| \leq \text{eps} |x|$  or  $|x - \circ(x)| \leq \text{eta}$ .

*Example 1.4* For the IEEE754 [11] binary64<sup>4</sup> format with  $\circ$  a rounding to nearest, we have  $\text{eps} = 2^{-53}$  ( $\simeq 10^{-16}$ ) and  $\text{eta} = 2^{-1075}$  ( $\simeq 10^{-323}$ ).

These constants allow to bound the rounding errors of the basic arithmetic operations.

<sup>4</sup> Usual implementation of the type `double` in C.

*Property 1.5* For all  $x, y \in \mathbb{F}$

$$\begin{aligned} \exists \delta \in \mathbb{R}, |\delta| \leq \text{eps} \wedge \text{fl}(x \diamond y) &= (1 + \delta)(x \diamond y), & \text{for } \diamond \in \{+, -\} \\ \exists \delta, \eta \in \mathbb{R}, |\delta| \leq \text{eps} \wedge |\eta| \leq \text{eta} \wedge \text{fl}(x \diamond y) &= (1 + \delta)(x \diamond y) + \eta, & \text{for } \diamond \in \{\times, /\} \\ \exists \delta \in \mathbb{R}, |\delta| \leq \text{eps} \wedge \text{fl}(\sqrt{x}) &= (1 + \delta)\sqrt{x}. \end{aligned}$$

For multiplication and division, the above property is a direct consequence of Definition 1.3. For addition, subtraction and square root, it also takes into account that, under mild conditions on the floating-point format, these operations are exact in case of underflow.

#### 1.4 Simple Example: the Sum

The previous bounds on rounding errors of basic operations can be combined to get bounds on the error of larger expressions, as for instance a summation in the following classic result [10, 14].

We first define  $\gamma_n$  that will be extensively used in the following.

**Definition 1.6** For all  $n \in \mathbb{N}$ ,  $\gamma_n := \frac{n \text{eps}}{1 - n \text{eps}}$ .

These terms are particularly useful to accumulate relative errors  $(1 + \delta)$  thanks, among others, to the following property.

**Lemma 1.7** (*phi\_gamma*) For all  $i, j \in \mathbb{N}$ ,  $\delta \in \mathbb{R}^{j-i}$ , if  $2(j-i)\text{eps} < 1$  and for all  $k$ ,  $|\delta_k| \leq \text{eps}$ , then

$$\exists \theta, |\theta| \leq \gamma_{j-i} \wedge \prod_{k=i}^{j-1} (1 + \delta_k) = 1 + \theta.$$

**Theorem 1.8** (*fsum\_l2r\_err\_abs*) For all  $n \in \mathbb{N}$ ,  $a \in \mathbb{F}^n$ , if  $2(n-1)\text{eps} < 1$ , then

$$\left| \text{fl} \left( \sum_{i=0}^{n-1} a_i \right) - \sum_{i=0}^{n-1} a_i \right| \leq \gamma_{n-1} \sum_{i=0}^{n-1} |a_i|.$$

The following proof is a good example of how equalities with error terms are used to derive such results. It gives a first overview of the properties and lemmas required about these error terms.

*Proof* We have by direct application of Property 1.5

$$\text{fl} \left( \sum_{i=0}^{n-1} a_i \right) = \text{fl} \left( \text{fl} \left( \sum_{i=0}^{n-2} a_i \right) + x_{n-1} \right) = (1 + \delta_n) \left( \text{fl} \left( \sum_{i=0}^{n-2} a_i \right) + x_{n-1} \right)$$

for some  $\delta_{n-1} \in \mathbb{R}$ ,  $|\delta_{n-1}| \leq \text{eps}$ . Then

$$\text{fl} \left( \sum_{i=0}^{n-1} a_i \right) = (1 + \delta_{n-1}) \left( (1 + \delta_{n-2}) \left( \text{fl} \left( \sum_{i=0}^{n-3} a_i \right) + x_{n-2} \right) + x_{n-1} \right)$$

for some  $\delta_{n-2} \in \mathbb{R}$ ,  $|\delta_{n-2}| \leq \text{eps}$ . By an immediate induction

$$\text{fl} \left( \sum_{i=0}^{n-1} a_i \right) = \left( \prod_{j=1}^{n-1} (1 + \delta_j) \right) a_0 + \sum_{i=1}^{n-1} \left( \prod_{j=i}^{n-1} (1 + \delta_j) \right) a_i \quad (1)$$

for some  $\delta_j \in \mathbb{R}$ ,  $|\delta_j| \leq \text{eps}$ . Then by Lemma 1.7, for all  $i \in \llbracket 1, n-1 \rrbracket$ , there exist  $\theta_{n-i} \in \mathbb{R}$  such that  $|\theta_{n-i}| \leq \gamma_{n-i}$  and  $\prod_{j=i}^{n-1} (1 + \delta_j) = 1 + \theta_{n-i}$ , hence

$$\text{fl} \left( \sum_{i=0}^{n-1} a_i \right) = (1 + \theta_{n-1}) a_0 + \sum_{i=1}^{n-1} (1 + \theta_{n-i}) a_i = \sum_{i=0}^{n-1} a_i + \left( \theta_{n-1} a_0 + \sum_{i=1}^{n-1} \theta_{n-i} a_i \right). \quad (2)$$

Finally, since for all  $i \in \llbracket 1, n-1 \rrbracket$ ,  $|\theta_{n-i}| \leq \gamma_{n-i} \leq \gamma_{n-1}$ , there exists  $\theta \in \mathbb{R}$  such that  $|\theta| \leq \gamma_{n-1}$  and

$$\text{fl} \left( \sum_{i=0}^{n-1} a_i \right) = \sum_{i=0}^{n-1} a_i + \theta \sum_{i=0}^{n-1} |a_i| \quad (3)$$

which enables to conclude.

One may notice that this result admits somewhat simpler proofs by direct induction. This more complicated proof was chosen in order to illustrate some basic use of the error terms  $\delta$  and  $\theta$  which play a key role in proofs of more complicated results<sup>5</sup>.

It is also worth noting that this is a rather simple result and that proofs of more complicated results can rapidly involve more error terms. The use of a proof assistant then becomes a good way to get some confidence that no term is accidentally dropped or mixed with another which can easily happen with pen and paper proofs.

## 2 Specification of Floating-Point Arithmetic

We now detail the formal specification of floating-point arithmetic used in our Coq development and already informally introduced in Section 1.3.

As seen in Property 1.5 and Theorem 1.8, both definitions and proofs make intensive use of real numbers with a bounded absolute value. To ease the manipulation of these error terms, we use a dependent record bounded  $r$  packing a real number with a proof that its absolute value is less than a non-negative real number  $r$ :

```
Record bounded (r : ℝ) :=
{ bounded_val :> ℝ; bounded_prop : |bounded_val| ≤ r }.
```

Since the set of floating-point values  $\mathbb{F}$  is a subset of  $\mathbb{R}$ , we will similarly define floating point values as a value in  $\mathbb{R}$  along with a proof that it lies in  $\mathbb{F}$ :

```
Record Ff format :=
{ F_val :> ℝ; F_prop : format F_val }.
```

<sup>5</sup> Moreover, this choice improves modularity of the Coq code by sharing lemmas with other results.

where `format` is a predicate over  $\mathbb{R}$  identifying real numbers that are in  $\mathbb{F}$ .

The floating-point arithmetic specification is then given by the following large record which will be used as parameter of all our subsequent developments.

```
Record Float_spec := {
  (** format x means that  $x \in \mathbb{R}$  is a floating-point value *)
  format :  $\mathbb{R} \rightarrow \text{Prop}$ ;
  (** The type of floating-point values (coercible to  $\mathbb{R}$ ). *)
   $\mathbb{F}$  := Ff format;
  (** 0 and 1 must be floating-point numbers. *)
  format0 : format 0; format1 : format 1;
  (** The opposite of a floating point number is a floating point number. *)
  format_opp x : format x -> format (- x);
  (** Bound on the relative error (normalized numbers, no underflow). *)
  eps :  $\mathbb{R}$ ; eps_pos :  $0 \leq \text{eps}$ ; eps_lt_1 :  $\text{eps} < 1$ ;
  (** Bound on the absolute error (denormalized, when underflow occurs). *)
  eta :  $\mathbb{R}$ ; eta_pos :  $0 < \text{eta}$ ;
  (** Some rounding. *)
  frnd :  $\mathbb{R} \rightarrow \mathbb{F}$ ; frnd_spec x :
   $\exists (d : \text{bounded eps}) (e : \text{bounded eta}), \text{frnd } x = (1+d)x+e :> \mathbb{R}$ ;
  (** Addition. *)
  fplus (x y :  $\mathbb{F}$ ) :  $\mathbb{F}$  := frnd (x+y);
  fplus_spec x y :  $\exists d : \text{bounded eps}, \text{fplus } x \ y = (1+d)(x+y) :> \mathbb{R}$ ;
  fplus_spec2 (x y :  $\mathbb{F}$ ) :  $y \leq 0 \rightarrow \text{fplus } x \ y \leq x$ ;
  (** Multiplication. *)
  fmult (x y :  $\mathbb{F}$ ) :  $\mathbb{F}$  := frnd (x*y);
  fmult_spec (x y :  $\mathbb{F}$ ) := frnd_spec (x*y);
  fmult_spec2 x :  $0 \leq \text{fmult } x \ x$ ;
  (** Square root. *)
  fsqrt (x :  $\mathbb{F}$ ) :  $\mathbb{F}$  := frnd  $\sqrt{x}$ ;
  fsqrt_spec x :  $\exists d : \text{bounded eps}, \text{fsqrt } x = (1+d)\sqrt{x} :> \mathbb{R}$ ;
}.
```

The unary minus `fopp` is then derived from `format_opp` and the subtraction and division of  $x$  and  $y$  are defined respectively as `fplus x (fopp y)` and `frnd (x/y)`.

Having performed our proofs with a proof assistant, we are guaranteed that the above record contains all the hypotheses about floating-point arithmetic used in these proofs. It is interesting to notice that this specification of floating-point arithmetic is really broad. In particular, it encompasses floating-point formats with gradual or abrupt underflow (i.e., denormalized numbers are respectively used or not) and any rounding mode. Fixed-point arithmetic can even be handled by just setting `eps` to 0, i.e., no relative, only absolute error occur. However, most of our developments are carried on with floating-point arithmetic in mind and the proved bounds might be pretty poor in a fixed-point arithmetic setting.



It is common practice in numerical analysis to ignore underflows [10]. Although this gives good indications on the numerical behavior of algorithms, underflows can appear with any practical implementation of floating-point arithmetic, potentially breaking such results. In our development, they are taken into account, thanks to the `eta` constant.

In our Coq development, the above specification of floating-point arithmetic is proved to hold for the floating point format with gradual underflow and rounding to nearest with any tie-break rule modeled in the Flocq library [4] with parameters corresponding to the binary64 format<sup>4</sup> (albeit without NaNs nor overflows, c.f., record `binary64` in our development). Other formats such as binary32<sup>6</sup> could be obtained by just modifying two constants defining size of the mantissa and minimal exponent.

In contrary to underflows, not handling NaNs and overflows does not constitute an actual issue. In fact, results considering those special values can easily be derived from results in our model with only finite values as will be shown in Section 5.

### 3 Combining Error Terms

#### 3.1 Bounded Error Terms

Values of the type `bounded` defined at beginning of Section 2 are coercible to  $\mathbb{R}$  and we developed a few helpful lemmas about them. The two following lemmas can be used to create such values.

**Lemma 3.1** (`bounded_le_1`)  $\forall r_1, r_2 \in \mathbb{R}, |r_1| \leq r_2 \Rightarrow \exists b : \text{bounded } 1, r_1 = b r_2$

**Lemma 3.2** (`bounded_scale`)

$\forall r_1, r_2 \in \mathbb{R}, b_1 : \text{bounded } r_1, 0 < r_2 \Rightarrow \exists b_2 : \text{bounded } r_2, b_1 = b_2 \frac{r_1}{r_2}$

It is often needed to say that a value of type `bounded b` is also of type `bounded b'` for any  $b' \geq b$  (for instance, in proof of Theorem 1.8, to state that the  $\theta_{n-i} : \text{bounded } \gamma_{n-i}$  are all of type `bounded  $\gamma_{n-1}$` ):

**Lemma 3.3** (`widen_bounded`)

$\forall r, r' \in \mathbb{R}, \forall b : \text{bounded } r, r \leq r' \Rightarrow \exists b' : \text{bounded } r', b = b'$

It is also common to get bounds of the form  $b e$  with  $b : \text{bounded } r$  and  $e$  a complicated expression we want to replace by a simpler expression  $e' \geq e$ :

**Lemma 3.4** (`bounded_larger_factor`)

$\forall r, r_1, r_2 \in \mathbb{R}, \forall b : \text{bounded } r, |r_1| \leq |r_2| \Rightarrow \exists b' : \text{bounded } r, b r_1 = b' r_2$

Error terms are compatible with basic arithmetic operations:

**Lemma 3.5** (`bounded_opp`, `bounded_plus`, `bounded_mult`)

$\forall r : \mathbb{R}, \forall b : \text{bounded } r, \exists b' : \text{bounded } r, b' = -b$

$\forall r_1, r_2 : \mathbb{R}, \forall b_1 : \text{bounded } r_1, \forall b_2 : \text{bounded } r_2, \exists b' : \text{bounded } (r_1 + r_2), b' = b_1 + b_2$

$\forall r_1, r_2 : \mathbb{R}, \forall b_1 : \text{bounded } r_1, \forall b_2 : \text{bounded } r_2, \exists b' : \text{bounded } (r_1 r_2), b' = b_1 b_2$

<sup>6</sup> Usual implementation of type `float` in C.

Finally, probably the most important lemma about error terms allows to factor them and was exemplified between (2) and (3) in the proof of Theorem 1.8:

**Lemma 3.6** (`bounded_distr1`, `big_bounded_distr1`)

$$\begin{aligned} \forall r, r_1, r_2 \in \mathbb{R}, \forall b_1, b_2 : \text{bounded } r, \exists b' : \text{bounded } r, b_1 r_1 + b_2 r_2 &= b' (|r_1| + |r_2|) \\ \forall r \in \mathbb{R}, \forall n : \mathbb{N}, \forall a \in \mathbb{R}^n, \forall b : (\text{bounded } r)^n, \exists b' : \text{bounded } r, \sum_i b_i a_i &= b' \left( \sum_i |a_i| \right) \end{aligned}$$

It is worth noting that this last property involves tuples ( $a$  and  $b$ ) and sums of an arbitrary number of terms ( $\sum_i b_i a_i$  for instance). Those are efficiently handled thanks to the `bigop` operator from the `SSReflect` library [2].

### 3.2 Accumulating Relative Errors

As already seen, for instance in the proof of Theorem 1.8 (between (1) and (2)), error terms of the form  $(1 + \delta_1) \dots (1 + \delta_n)$ , with  $|\delta_i| \leq \text{eps}$ , easily occur when relative errors accumulate. The terms  $\gamma_n := \frac{n \text{eps}}{1 - n \text{eps}}$  nicely enable to compact them into  $1 + \theta_n$ ,  $|\theta_n| \leq \gamma_n$  as will be exposed in this section.

Most of the following lemmas require hypothesis of the form  $n \text{eps} < 1$  for various values of  $n$ . In our Coq code, a set of small lemmas<sup>7</sup> allow to easily manipulate these hypothesis so that they do not constitute an annoying burden in practice. First, a few very basic properties of the  $\gamma_n$  are proved. Namely, that they are non negative (provided  $n \text{eps} < 1$ ), strictly less than 1 (provided  $2n \text{eps} < 1$ ) and constitute a monotone sequence: for all  $n \leq m$ ,  $\gamma_n \leq \gamma_m$  (provided  $m \text{eps} < 1$ ).

We then get some more interesting properties.

**Lemma 3.7** ([10, Lemma 3.3], `gamma_{mult{,_nat},plus{,_mult,_,_eps}}`)

For all  $n, m \in \mathbb{N}$

$$\begin{aligned} n \leq m &\Rightarrow 2m \text{eps} < 1 \Rightarrow \gamma_n \gamma_m \leq \gamma_n \\ (nm) \text{eps} < 1 &\Rightarrow n \gamma_m \leq \gamma_{nm} \\ (n+m) \text{eps} < 1 &\Rightarrow \gamma_n + \gamma_m + \gamma_n \gamma_m \leq \gamma_{n+m} \\ (n+m) \text{eps} < 1 &\Rightarrow \gamma_n + \gamma_m \leq \gamma_{n+m} \\ (n+1) \text{eps} < 1 &\Rightarrow \gamma_n + \text{eps} \leq \gamma_{n+1} \end{aligned}$$

Allowing to prove properties about the  $\theta_n : \text{bounded } \gamma_n$ .

<sup>7</sup> For instance:  $(n+1) \text{eps} < 1 \Rightarrow n \text{eps} < 1$ .

**Lemma 3.8** (`[10, Lemma 3.3]`, `gammap1_{mult{,_eps}1,div{,_le}}`)

For all  $n, m \in \mathbb{N}$ , for all  $\theta_n : \text{bounded } \gamma_n$ ,  $\theta_m : \text{bounded } \gamma_m$  and  $\delta : \text{bounded } \text{eps}$

$$\begin{aligned}
2(n+m)\text{eps} < 1 &\Rightarrow \\
&\exists \theta_{n+m} : \text{bounded } \gamma_{n+m}, (1 + \theta_n)(1 + \theta_m) = 1 + \theta_{n+m} \\
2(n+1)\text{eps} < 1 &\Rightarrow \\
&\exists \theta_{n+1} : \text{bounded } \gamma_{n+1}, (1 + \theta_n)(1 + \delta) = 1 + \theta_{n+1} \\
m \leq n \Rightarrow 2(n+m)\text{eps} < 1 &\Rightarrow \\
&\exists \theta_{n+m} : \text{bounded } \gamma_{n+m}, (1 + \theta_n)/(1 + \theta_m) = 1 + \theta_{n+m} \\
2(n+2m)\text{eps} < 1 &\Rightarrow \\
&\exists \theta_{n+2m} : \text{bounded } \gamma_{n+2m}, (1 + \theta_n)/(1 + \theta_m) = 1 + \theta_{n+2m}
\end{aligned}$$

**Lemma 3.9** (`phi_gamma,inv_phi_gamma`)

For all  $i, j, n \in \mathbb{N}$ ,  $\delta : (\text{bounded } \text{eps})^n$ , if  $0 \leq i \leq j \leq n$  then

$$\begin{aligned}
2(j-i)\text{eps} < 1 &\Rightarrow \exists \theta_{j-i} : \text{bounded } \gamma_{j-i}, \prod_{k=i}^{j-1} (1 + \delta_k) = 1 + \theta_{j-i} \\
(j-i)\text{eps} < 1 &\Rightarrow \exists \theta_{j-i} : \text{bounded } \gamma_{j-i}, \frac{1}{\prod_{k=i}^{j-1} (1 + \delta_k)} = 1 + \theta_{j-i}
\end{aligned}$$

It should be noted that, unlike the terms  $\gamma_n$ , the subscripts  $n$  in the notation  $\theta_n$  are only present to improve readability but don't hold any semantic value.

It is also interesting to notice about the division [10, §3.4] that, although (according to Lemma 3.8)  $(1 + \theta_n)/(1 + \theta_m) = 1 + \theta_{n+m}$  only holds for  $m \leq n$ , according to the above lemma  $\prod_n (1 + \delta_i)/\prod_m (1 + \delta_i) = 1 + \theta_{n+m}$  even when  $m > n$ .

The notations  $\delta$  and  $\theta_n$ , with  $|\delta| \leq \text{eps}$  and  $|\theta_n| \leq \gamma_n$ , used in the above lemmas are particularly convenient and popular to carry proofs about error bounds. In fact, like the Landau big O notation, they greatly simplify proofs by enabling the use of simple equalities<sup>8</sup> instead of a collection of inequalities, or limits. However, it is easy to misuse them or forget hypotheses, such as  $m \leq n$  in Lemma 3.8. The use of a proof assistant ensures that this does not happen.

### 3.3 First Applications

Thanks to all the above lemmas, the Theorem 1.8, bounding the rounding error of a sum, is easily proven (c.f., `fsum_l2r_err_abs` in our Coq development). Similar results are also proved for the dotproduct of two vectors of floating point values.

**Lemma 3.10** (`fdotprod_l2r_err_abs`)

For all  $n \in \mathbb{N}$  and  $a, b \in \mathbb{F}^n$ , if  $2n \text{eps} < 1$ , then

$$\left| \text{fl} \left( \sum_{i=0}^{n-1} a_i b_i \right) - \sum_{i=0}^{n-1} a_i b_i \right| \leq \gamma_n \left( \sum_{i=0}^{n-1} |a_i b_i| \right) + 2n \text{eps} a.$$

<sup>8</sup> See for instance the proof of Theorem 1.8, page 6.

Another example, in case the first operand  $a$  is constituted of real numbers, which have to be rounded to floating-point values prior to computation of the dotproduct. This will be used in the last application of this paper.

**Lemma 3.11** (`fdotprod_l2r_fstr_err`)

For all  $n \in \mathbb{N}$ ,  $a \in \mathbb{R}^n$  and  $b \in \mathbb{F}^n$ , if  $2(n+1)\text{eps} < 1$ , then

$$\left| \text{fl} \left( \sum_{i=0}^{n-1} a_i b_i \right) - \sum_{i=0}^{n-1} a_i b_i \right| \leq \gamma_{n+1} \left( \sum_{i=0}^{n-1} |a_i b_i| \right) + 2 \left( n + \sum_{i=0}^{n-1} |b_i| \right) \text{e.t.a.}$$

Due to the presence of existential quantifiers in most intermediate lemmas, proof style in the Coq proof assistant heavily relies on forward proving. This does not appear to add much burden to the proof writing process, as long as proofs are well structured into lemmas of reasonable size<sup>9</sup>. Otherwise, one can first provide some dummy terms and later step back to replace them by the, then more obvious, adequate terms. Using the `evvar` mechanism of Coq might also be a solution. In the “big enough numbers” for  $\varepsilon, \eta$  proofs about limits [5], constraints of the form  $\eta \leq \eta_k$  are first accumulated as a list of the  $\eta_k$  and then simplified as  $\eta \leq \min_k \eta_k$ . Unfortunately, we are manipulating equalities of arithmetic expressions involving error terms rather than inequalities and we do not see such an easy simplification in this setting.

## 4 Errors on Matrix Operations

### 4.1 Real Numbers Matrices

As stated in the introduction (Section 1.1), we intend to prove numerical analysis results on algorithms involving matrices. In our Coq development, we borrow matrix algebra from the `SSReflect` library [8]. But we also need some results which are specific to matrices of real numbers. We therefore introduce some basic definitions and lemmas about pointwise orders and absolute values, dotproducts and quadratic norms.

First, the pointwise extensions of the order  $\leq$  and  $<$  as well as the absolute value  $|\cdot|$  are defined and a group of lemmas are proved about them. Most of these lemmas are just lifting of the existing results on the real field  $\mathbb{R}$ : reflexivity and transitivity of the order  $\leq$ , compatibility of this order with matrix addition or scaling, triangular inequality of the absolute value ( $\forall A, B \in \mathbb{R}^{n \times m}, |A + B| \leq |A| + |B|$ ) and so on.

In order to deal with quadratic norms, we first define *positive (semi-)definite* matrices. A matrix  $P \in \mathbb{R}^{n \times n}$  is said to be positive semi-definite, written  $P \succeq 0$ , when for all  $x \in \mathbb{R}^n$ ,  $x^T P x \geq 0$  and it is said to be positive definite, written  $P \succ 0$ , when for all  $x \neq 0$ ,  $x^T P x > 0$ . Thus for a symmetric ( $P^T = P$ ) positive definite matrix  $P$ , we define the *dotproduct* of two vectors  $x, y \in \mathbb{R}^n$  as  $x^T P y$  and the quadratic norm  $\|x\|_P$  as  $\sqrt{x^T P x}$ . The dotproduct is then proven to actually be a dotproduct (bilinear, symmetric, definite and non negative). It follows that the quadratic norm is definite

<sup>9</sup> Which is just good programming practice.

non-negative and satisfies the scaling property ( $\|\lambda x\|_P = |\lambda| \|x\|_P$ ) which eventually enables to prove two usual inequalities: the Cauchy-Schwarz inequality:

$$\forall x, y \in \mathbb{R}^n, |x^T P y| \leq \|x\|_P \|y\|_P$$

and the triangular inequality:

$$\forall x, y \in \mathbb{R}^n, \|x + y\|_P \leq \|x\|_P + \|y\|_P.$$

In the particular case when  $P := I$ , the quadratic norm will be written  $\|\cdot\|_2$  and a few additional properties are proved:  $\forall x, y \in \mathbb{R}^n$ ,

$$\begin{aligned} |x| \leq |y| &\Rightarrow \|x\|_2 \leq \|y\|_2 \\ \| |x| \|_2 &= \|x\|_2 \\ \|[1 \dots 1]^T\|_2 &= \sqrt{n} \\ \|x\|_1 &\leq \sqrt{n} \|x\|_2 \end{aligned}$$

where  $\|x\|_1 := \sum_i |x_i|$ .

We eventually needed 110 lemmas. Thanks to the nice SSReflect matrices [8], they are proved using only 394 lines of tactics (hence an average of 3.6 lines of tactic per lemma, the longest proof being the Cauchy-Schwarz inequality with 36 lines of tactic).

## 4.2 Main Application: Cholesky Decomposition

As explained in Section 1.1, given a matrix  $A$  we want to check its positive definiteness, that is to prove  $A \succ 0$ . This will be done by proving that there exists a constant  $c \in \mathbb{R}$  such that when the Cholesky decomposition (c.f., Figure 1) of  $A - cI$ , performed with floating-point arithmetic, runs to completion without error (square root of negative value or division by zero), then  $A \succ 0$ . We follow the proof in [13]<sup>10</sup>.

The first lemmas proved deal with the two “basic blocks” of the Cholesky decomposition: the assignments performed in the inner then the outer loop (c.f., Figure 1, page 3). The two following lemmas are proved with tools similar to the one required for Lemmas 3.10 and 3.11 about floating-point sums and dotproducts.

**Lemma 4.1** ([13, Lemma 2.1], lemma\_2\_1) *For all  $n \in \mathbb{N}$ ,  $a, b \in \mathbb{R}^n$ ,  $c, d \in \mathbb{F}$ , if  $d \neq 0$  and  $2(n+1)\text{eps} < 1$ , then*

$$\left| c - \sum_i a_i b_i - d \tilde{y} \right| \leq \gamma_{n+1} \left( \sum_i |a_i b_i| + |d \tilde{y}| \right) + 2\text{eta}(k+1+|d|).$$

where  $\tilde{y} := \text{fl} \left( \frac{c - \sum_i a_i b_i}{d} \right)$ .

<sup>10</sup> Actually, part of it. We only consider matrices of real numbers whereas [13] also handles complex numbers. [13] also offers improved bounds for sparse matrices and a non-positive-definiteness check.

**Lemma 4.2** ([13, Lemma 2.2], lemma\_2\_2\_{1,2}) *For all  $n \in \mathbb{N}$ ,  $a \in \mathbb{F}^n$ ,  $c \in \mathbb{F}$ , if  $2(n+2)\text{eps} < 1$  and  $\text{fl}(c - \sum_i a_i^2) \geq 0$ , then*

$$\left| c - \sum_i a_i^2 - \tilde{y}^2 \right| < \gamma_{n+2} \left( \sum_i a_i^2 + \tilde{y}^2 \right) + 2\text{eta}(k+1).$$

and

$$\tilde{y}^2 + \sum_i a_i^2 \leq \frac{c + 2\text{eta}k}{1 - \gamma_{n+2}}$$

where  $\tilde{y} := \text{fl}\left(\sqrt{c - \sum_i a_i^2}\right)$ .

Then, given two matrices  $A, \tilde{R} \in \mathbb{F}^{n \times n}$  the proposition `cholesky_spec A  $\tilde{R}$`  expresses that  $\tilde{R}$  is the floating-point Cholesky factor of  $A$  and is defined as follows

$$\begin{aligned} \forall i, j \in \llbracket 1, n \rrbracket, i < j \Rightarrow \tilde{R}_{i,j} &= \text{fl}\left(\frac{A_{i,j} - \sum_{k=1}^{i-1} \tilde{R}_{k,i} \tilde{R}_{k,j}}{\tilde{R}_{i,i}}\right) \\ \wedge \forall j \in \llbracket 1, n \rrbracket, \tilde{R}_{j,j} &= \text{fl}\left(\sqrt{A_{j,j} - \sum_{k=1}^{j-1} \tilde{R}_{k,j}^2}\right). \end{aligned} \quad (4)$$

This corresponds to the algorithm presented in Figure 1, page 3.

Eventually, the proposition `cholesky_success A  $\tilde{R}$`  states that the floating-point Cholesky decomposition of  $A$  runs to completion without error (and returns  $\tilde{R}$ ):

$$\text{cholesky\_spec } A \tilde{R} \wedge \forall i \in \llbracket 1, n \rrbracket, \tilde{R}_{i,i} > 0.$$

The constraint  $\tilde{R}_{i,i} > 0$  obviously prevents divisions by zero. The way it prevents taking square roots of negative values is a bit more subtle. According to the specification of  $\text{fl}(\sqrt{\cdot})$  (c.f., `fsqrt_spec`, page 8),  $\text{fl}(\sqrt{x}) = (1+d)\sqrt{x}$  and  $\sqrt{x}$  is defined as 0 for all  $x \leq 0$  in the Coq standard library for reals. Thus,  $\text{fl}(\sqrt{x}) > 0$  guarantees that  $x > 0$ .

With this specification of the floating-point Cholesky decomposition, the following main theorem can be proved about it.

**Theorem 4.3** ([13, Theorem 2.3]) *For all  $n \in \mathbb{N}$  ( $n \geq 1$ ), for all  $A, \tilde{R} \in \mathbb{F}^{n \times n}$ ,  $m \in \mathbb{R}$ , if  $2(n+2)\text{eps} < 1$ ,  $A^T = A$ , for all  $i$ ,  $A_{i,i} \leq m$  and `cholesky_success A  $\tilde{R}$` , then*

$$\forall x \in \mathbb{R}^n, x \neq 0 \Rightarrow -|x|^T \Delta_A |x| < x^T A x$$

where  $\Delta_{A_{i,j}} := \alpha_{i,j} d_i d_j + 4\text{eta}(n+2+m)$  with  $\alpha_{i,j} := \gamma_{\min(i,j)+2}$  and  $d_i := \sqrt{\frac{A_{i,i} + 2\text{eta}}{1 - \alpha_{i,i}}}$ .

This theorem is proved thanks to the above Lemmas 4.1 and 4.2 and the lemmas about matrices of real numbers described in Section 4.1.

Here is an idea of the kind of tiny mistakes that can be hard to spot in pen and paper proofs and may only be unveiled when trying to mechanically check them. In the original paper [13], the proof accidentally made the assumption that  $\tilde{R}_{i,i} \leq A_{i,i}$  where  $\tilde{R}_{i,i} := \text{fl}\left(\sqrt{A_{i,i} - \sum_{k=1}^{i-1} \tilde{R}_{k,i}^2}\right)$ . This may be slightly false since  $\sqrt{x}$  can be larger

than  $x \in (0, 1)$  but this was completely overlooked by the author of this paper<sup>11</sup> when first checking the proof by hand. Although very minor, the issue only appeared when translating the proof in Coq (and was easily fixed since  $\sqrt{x} \leq x + 1$ ).

The previous theorem is pretty useless by itself but the following corollary looks closer from what we are looking for.

**Corollary 4.4** ([13, Corollary 2.4]) *For all  $n \in \mathbb{N}$  ( $n \geq 1$ ), for all  $A, \tilde{A} \in \mathbb{F}^{n \times n}$ ,  $c \in \mathbb{R}$ , if  $2(n+2)\text{eps} < 1$ ,  $A^T = A$  and for all  $i$ ,  $0 \leq A_{i,i}$  and*

$$\forall x \in \mathbb{R}^n, \|x\|_2 = 1 \Rightarrow |x|^T \Delta_A |x| \leq c$$

and  $\tilde{A}^T = \tilde{A}$  and

$$\forall i, j \in \llbracket 1, n \rrbracket, i < j \Rightarrow \tilde{A}_{i,j} = A_{i,j}$$

$$\forall i \in \llbracket 1, n \rrbracket, \tilde{A}_{i,i} \leq A_{i,i} - c$$

then, if there exists  $\tilde{R} \in \mathbb{F}^{n \times n}$  such that `cholesky_success`  $\tilde{A} \tilde{R}$ , we have

$$A \succ 0.$$

*Proof (sketch)* By definition of  $\tilde{A}$ , for any  $x \in \mathbb{R}^n$ , if  $\|x\| = 1$ , then  $x^T A x \geq x^T \tilde{A} x + c$ . From the theorem, we get  $x^T \tilde{A} x > -|x|^T \Delta_{\tilde{A}} |x|$ . Since  $\Delta_{\tilde{A}} \leq \Delta_A$  (elementwise inequality, c.f., `Delta_At_1e_Delta_A` in the code),  $x^T A x > -|x|^T \Delta_A |x| + c \geq 0$ .

Finally, it is proven that any value larger than

$$\frac{\gamma_{2n+2}}{2} \text{tr}(A) + 4\text{eta}(n+1) \left( 2(n+2) + \max_i A_{i,i} \right) \quad (5)$$

will work as constant  $c$  in the above corollary as long as  $4(n+2)\text{eps} < 1$ . Thus, an appropriate constant  $c$  can easily be computed, for instance with floating-point arithmetic with rounding toward  $+\infty$ . Then  $\tilde{A}$  is computed by subtracting  $cI$  to  $A$  and the floating point Cholesky decomposition is performed. If it runs to completion without error, this rigorously proves that  $A \succ 0$ . This automatic positive definiteness check is efficient as it is performed with  $O(n^3)$  floating-point operations for a matrix  $A$  of size  $n \times n$ .

The bound  $4(n+2)\text{eps} < 1$  is not a practical issue since it only implies  $n < \frac{1}{4\text{eps}} - 2$  which is huge for practical values of  $\text{eps}$ . For instance, for the binary64 format,  $\text{eps} = 2^{-53}$  which leads to the constraint  $n < 2^{51} - 2$  ( $\simeq 10^{15}$ ). In practice,  $\text{eta}$  being very small compared to  $\text{eps}$  ( $\text{eps} = 2^{-53}$  and  $\text{eta} = 2^{-1075}$  for the binary64 format), the second term of the bound (5) is negligible with respect to the first term  $\frac{\gamma_{2n+2}}{2} \text{tr}(A)$ . This means that, when the diagonal coefficients of the matrix  $A$  are of order of magnitude<sup>12</sup> 1, the bound (5) is mostly  $n^2 \text{eps}$ .

Thanks to the previous corollary, positive definiteness check can be performed on matrices  $A$  of floating-point values ( $A \in \mathbb{F}^{n \times n}$ ). However, if the matrix  $X$  we want

<sup>11</sup> As well as the author and reviewers of the original paper.

<sup>12</sup> A preprocessing can fit all diagonal coefficients in the interval  $[0.25, 1]$ . Indeed,  $A \succ 0$  when  $A' := D^T A D \succ 0$  for any non-singular matrix  $D$ . By choosing  $D = \text{diag}(d_0, \dots, d_{n-1})$ , we have  $A'_{i,j} = d_i d_j A_{i,j}$ . By choosing the  $d_i$  as powers of 2, these multiplications are exact in floating point arithmetic (if no underflow (nor overflow) occurs) and we can guarantee  $A'_{i,i} \in [0.25, 1]$  [13].

to check has coefficients in the real field ( $X \in \mathbb{R}^{n \times n}$ ), we first have to round them to floating-point values in  $\mathbb{F}$  and we will end up checking some matrix  $A \in \mathbb{F}^{n \times n}$  such that  $|X - A|_{i,j} \leq R_{i,j} := \text{eps} |A|_{i,j} + \text{eta}$ . Such interval matrices are easily handled thanks to the following corollary.

**Corollary 4.5** ([13, Corollary 2.7]) *For all  $n \in \mathbb{N}$  ( $n \geq 1$ ), for all  $A, \tilde{A} \in \mathbb{F}^{n \times n}$ ,  $R \in \mathbb{R}^{n \times n}$ ,  $c, r \in \mathbb{R}$ , if  $2(n+2)\text{eps} < 1$ ,  $A^T = A$ ,  $R \geq 0$  and for all  $i$ ,  $0 \leq A_{i,i}$  and*

$$\forall x \in \mathbb{R}^n, \|x\|_2 = 1 \Rightarrow |x|^T \Delta_A |x| \leq c$$

and

$$\forall x \in \mathbb{R}^n, \|x\|_2 = 1 \Rightarrow |x|^T R |x| \leq r$$

and  $\tilde{A}^T = \tilde{A}$  and

$$\forall i, j \in \llbracket 1, n \rrbracket, i < j \Rightarrow \tilde{A}_{i,j} = A_{i,j}$$

$$\forall i \in \llbracket 1, n \rrbracket, \tilde{A}_{i,i} \leq A_{i,i} - c - r$$

then, if there exists  $\tilde{R} \in \mathbb{F}^{n \times n}$  such that `cholesky_success`  $\tilde{A} \tilde{R}$ , we have

$$\forall X \in \mathbb{R}^{n \times n}, X^T = X \Rightarrow |X - A| \leq R \Rightarrow X \succ 0.$$

Since  $n \max \{R_{i,j} \mid i, j \in \llbracket 1, n \rrbracket\}$  is a suitable value for  $r$ , this gives an effective criterion for positive definiteness of a matrix  $X$  with coefficients in the real field  $\mathbb{R}$ .

The whole Coq development eventually counts 3.5 kloc. Among them, 0.4 are devoted to the specification of floating-point arithmetic (described in Section 2), 0.2 to bounded error terms (Section 3.1), 0.6 to the  $\gamma_n$  terms and their properties (Section 3.2), 0.4 to basic lemmas about sums and dotproducts (Sections 1.4 and 3.3) and 0.8 to matrices of real numbers (Section 4). Finally, the main theorem and corollaries (this section) take 1.0 kloc and the remainder (0.1 kloc) is constituted of miscellaneous lemmas. This appears particularly reasonable, considering the original result is a far from trivial 6 page long paper proof [13].

#### 4.3 Another Application: Impact of Rounding Errors on Ellipsoidal Invariants

To assert the reusability of our developments for numerical analysis results involving matrices, we targeted another property. We will see that three quarters of the previous development can be directly reused.

As introduced in Section 1.2, we are given an invariant “ $x^T P x \leq \lambda$ ” for the program of Figure 2 executed with real arithmetic and we want to prove, under some conditions, that it is also an invariant for the same program executed with floating-point arithmetic.

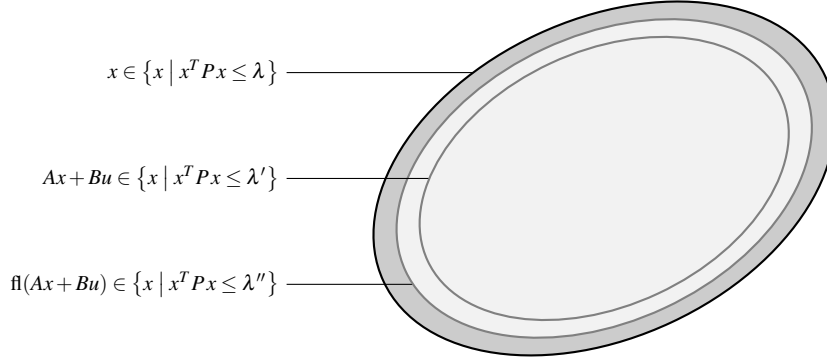
Here is the basic idea. The fact that the invariant is inductive for real arithmetic can be expressed as the following property

$$\forall x \in \mathbb{R}^n, \forall u \in \mathbb{R}^p, x^T P x \leq \lambda \Rightarrow \|u\|_\infty \leq 1 \Rightarrow (Ax + Bu)^T P (Ax + Bu) \leq \lambda.$$

But in practice, there will be some margin and we will be given the property

$$\forall x \in \mathbb{R}^n, \forall u \in \mathbb{R}^p, x^T P x \leq \lambda \Rightarrow \|u\|_\infty \leq 1 \Rightarrow (Ax + Bu)^T P (Ax + Bu) \leq \lambda'.$$





**Fig. 4** When  $x$  is in an ellipsoid of radius  $\lambda$ ,  $Ax + Bu$  is in a smaller ellipsoid of radius  $\lambda' < \lambda$  and  $\text{fl}(Ax + Bu)$  is then in a slightly larger ellipsoid of radius  $\lambda'' > \lambda'$  (and hopefully  $\lambda'' \leq \lambda$ ). The figure is not to scale, the difference between  $\lambda'$  and  $\lambda''$  being usually orders of magnitude smaller than with  $\lambda$ .

for some  $\lambda' < \lambda$ . That is if  $x$  is in the ellipsoid  $\{x \mid x^T P x \leq \lambda\}$ , then  $Ax + Bu$  is in the smaller ellipsoid  $\{x \mid x^T P x \leq \lambda'\}$ . Computed with floating-point arithmetic,  $\text{fl}(Ax + Bu)$  will be equal to  $Ax + Bu$  plus a small rounding error and will then lie in a slightly larger ellipsoid  $\{x \mid x^T P x \leq \lambda''\}$  with  $\lambda'' > \lambda'$ . If  $\lambda'' \leq \lambda$ , then

$$\forall x \in \mathbb{R}^n, \forall u \in \mathbb{R}^p, x^T P x \leq \lambda \Rightarrow \|u\|_\infty \leq 1 \Rightarrow \text{fl}(Ax + Bu)^T P \text{fl}(Ax + Bu) \leq \lambda,$$

which means that the invariant is also inductive for floating-point arithmetic. This is illustrated on Figure 4.

The following theorem shows that  $\lambda'' \leq (\sqrt{\lambda'} + \sqrt{\lambda} a + b)^2$  with  $a$  and  $b$  very small constants depending on  $A, B, P$  and the floating-point format used. That is  $\lambda''$  is only slightly larger than  $\lambda'$  (and then hopefully less than  $\lambda$ ).

**Theorem 4.6** For all  $n, p \in \mathbb{N}$ ,  $A, P \in \mathbb{R}^{n \times n}$ ,  $B \in \mathbb{R}^{n \times p}$ ,  $s, s', \lambda, \lambda' \in \mathbb{R}$ , if  $2(n + p + 1)\text{eps} < 1$ ,  $P^T = P$ ,  $P \succ 0$ ,  $sP - I \succeq 0$ ,  $s'I - P \succeq 0$ , then for all  $x \in \mathbb{F}^n$ ,  $u \in \mathbb{F}^p$  if

$$x^T P x \leq \lambda, \|u\|_\infty \leq 1 \text{ and } (Ax + Bu)^T P (Ax + Bu) \leq \lambda'$$

then

$$\text{fl}(Ax + Bu)^T P \text{fl}(Ax + Bu) \leq (\sqrt{\lambda'} + \sqrt{\lambda} a + b)^2$$

where  $a := \gamma_{n+p+1} \sqrt{ss'} \sqrt{n} \|A\|_F + 2\sqrt{ss'} n \sqrt{n} \text{eta}$

and  $b := \gamma_{n+p+1} \sqrt{s'} \sqrt{p} \|B\|_F + 2\sqrt{s'} (n + 2p) \sqrt{n} \text{eta}$

where  $\|M\|_F$  denotes the Frobenius norm of the matrix  $M$  (i.e.,  $\|M\|_F := \sqrt{\sum_{i,j} M_{i,j}^2}$ ).

For typical values ( $n \leq 10$ ,  $p \leq 10$ , coefficients of  $A$  and  $B$  of order of magnitude 1,  $ss' \leq 10^4$ ,  $s' \leq \lambda$  and  $\text{eta} \ll \text{eps} \simeq 10^{-16}$  for binary64),  $\sqrt{\lambda} a + b \leq 10^{-10} \sqrt{\lambda}$  which is very small with respect to a typical relative difference of  $10^{-4}$  between  $\lambda$  and  $\lambda'$  [12]. Thus  $\lambda'' \leq \lambda$  always holds in practical cases.

Here are the main lines of the proof.  $\text{fl}(Ax + Bu)$  being defined as  $\text{fl}([A \ B][x^T \ u^T]^T)$  (recall that  $\text{fl}(expr)$  means  $expr$  computed from left to right), for all  $i$ ,  $\text{fl}(Ax + Bu)_i$  is a dotproduct and Lemma 3.11 gives for all  $i$ ,  $\text{fl}(Ax + Bu)_i = (Ax + Bu)_i + e$  with  $|e| \leq b(x, u)$  (c.f., Lemma 3.11 for the details of  $b$ ).  $\|u\|_\infty \leq 1$  is an hypothesis. Then, the hypothesis  $sP - I \succeq 0$  and  $x^T P x \leq \lambda$  allow to prove  $\|x\|_\infty \leq \sqrt{s\lambda}$  (lemma\_2 in ellipsoid\_error.v, proved using the lemmas on matrices of Section 4.1). These are the main ingredients of the following lemma.

**Lemma 4.7** (lemma\_3 in ellipsoid\_error.v) *For all  $n, p \in \mathbb{N}$ ,  $A, P \in \mathbb{R}^{n \times n}$ ,  $B \in \mathbb{R}^{n \times p}$ ,  $s, \lambda \in \mathbb{R}$ , if  $2(n + p + 1)\text{eps} < 1$ ,  $P \succ 0$ ,  $sP - I \succeq 0$ , then for all  $x \in \mathbb{F}^n$ ,  $u \in \mathbb{F}^p$  if*

$$x^T P x \leq \lambda \text{ and } \|u\|_\infty \leq 1$$

*then there exists  $d : (\text{bounded } 1)^n$  such that*

$$\text{fl}(Ax + Bu) = Ax + Bu + \text{diag}(d) \left( \sqrt{s\lambda} (\gamma_{n+p+1} c_A + 2n c_{\text{eta}}) + \gamma_{n+p+1} c_B + 2(n + 2p) c_{\text{eta}} \right)$$

where  $c_A := [\sum_j A_{0,j}, \dots, \sum_j A_{n-1,j}]^T$ ,  $c_B := [\sum_j A_{0,j}, \dots, \sum_j A_{n-1,j}]^T$  and  $c_{\text{eta}} := [\text{eta}, \dots, \text{eta}]^T$ .

From the hypothesis  $s'I - P \succeq 0$ , we get for all  $x \in \mathbb{R}^n$ ,  $\|\text{diag}(d)x\|_P \leq \sqrt{s'}\|x\|_2$  (lemma\_4) and the proof of the theorem finally follows thanks to matrix manipulations and lemmas of Section 4.1.

Among the 3.1 kloc needed to prove this theorem, 2.3 are shared with the development performed for the previous Section 4.2. Again, 0.8 kloc of Coq is a reasonably small amount of code for translating such a non trivial 4 page long paper proof.

## 5 Considering Potential Overflows

Until now, the possibility of overflows has been ignored. This section shows how previous results can be extended taking them into account. The basic idea is to use the fact that, as long as no overflow occurs, results are the same than in the previous model without overflows. An extension to overflows of the floating-point model of Section 2 is first given before exemplifying it on the Cholesky decomposition result of Section 4.2.

### 5.1 Specification of Floating-Point Arithmetic with Overflows

To the normalized and denormalized numbers already considered, the IEEE754 standard adds special values  $-\infty$  and  $+\infty$  to handle overflows and NaN (Not a Number) to handle undefined operations (for instance  $0/0$ ). The new following specification handles this additional values by basically stating that when the result of an operation is not one of them, then the operation behaves as already defined in Section 2.

```

Record Float_infnan_spec := {
  (** Type of floating-point values (either finite, infinite or NaN). *)
  FI : Set;
  (** finite f means that the floating-point number f is finite. *)
  finite : FI → Prop;
  (** Underlying unbounded floating-point format. FI and F fis match when finite
  holds. *)
  fis : Float_spec;
  (** Associates the corresponding value in F fis for finite values or 0 for infinities
  and NaN. *)
  FI2F : FI → F fis;
  FI2F_spec x : (FI2F x ≠ 0 :> ℝ) → finite x;
  (** Opposite *)
  fiopp : FI → FI;
  fiopp_spec x : finite (fiopp x) →
  FI2F (fiopp x) = fopp (FI2F x) :> ℝ;
  fiopp_spec_f x : finite (fiopp x) → finite x;
  (** Addition *)
  fiplus : FI → FI → FI;
  fiplus_spec x y : finite (fiplus x y) →
  FI2F (fiplus x y) = fplus (FI2F x) (FI2F y) :> ℝ;
  fiplus_spec_fl x y : finite (fiplus x y) → finite x;
  fiplus_spec_fr x y : finite (fiplus x y) → finite y;

```

Multiplication and square root are very similar (and omitted for brevity). Only the division slightly differs because  $f/\pm\infty = 0$  for any finite  $f$ .

```

fidiv : FI → FI → FI;
fidiv_spec x y : finite (fidiv x y) → finite y →
FI2F (fidiv x y) = fdiv (FI2F x) (FI2F y) :> ℝ;
fidiv_spec_fl x y : finite (fidiv x y) → finite x;
}.

```

In our Coq development, the above specification is proved to hold for the bit-level model of binary64 in the Flocq library [4] (c.f., record `binary64_infnan`) with the `binary64` of Section 2 as underlying model `fis`.

## 5.2 Example: Cholesky Decomposition

The specification `cholesky_spec_infnan` of the Cholesky decomposition is the same as (4), page 14, with the arithmetic operations replaced by the one of the new specification above. `cholesky_success_infnan` is then defined as

$$\text{cholesky\_spec\_infnan } A \tilde{R} \wedge \forall i \in \llbracket 1, n \rrbracket, \text{FI2F } \tilde{R}_{i,i} > 0.$$

It is worth noting that  $\text{FI2F } \tilde{R}_{i,i} > 0$  implicitly states that  $\tilde{R}_{i,i}$  is finite. This enables to prove the following lemma

**Lemma 5.1** (`cholesky_success_infnan_cholesky_success`) *For  $A, \tilde{R} \in \text{FI}^{n \times n}$ , if `cholesky_success_infnan A`  $\tilde{R}$ , then `cholesky_success (MFI2F A) (MFI2F  $\tilde{R}$ )` with `MFI2F M` the matrix with coefficients `FI2F  $M_{i,j}$` .*

Thus, if the other hypothesis of Corollary 4.4 hold,  $\text{MFI2F } A \succ 0$  (i.e.,  $A \succ 0$  if all element of  $A$  are finite).

The Coq development for this extension to overflows counts 0.5 kloc. Among them, 0.1 are devoted to the specification described in Section 5.1, 0.2 to prove that the bit-level model of binary64 from the Flocq library satisfy it and 0.2 to the Cholesky decomposition example.

## 6 Conclusion

We formally proved, using the proof assistant Coq [1,6], two results bounding rounding errors of numerical computations and involving matrices and common numerical analysis tools [10]. Our Coq development is available at <http://cavale.enseeiht.fr/formalbounds2014/>. It indicates that performing such proofs within proof assistants is tractable and that a large part of the proof effort could be reused for similar results. Our development is based on a broad high-level floating-point specification and we have proved that this specification is satisfied by the model of the binary64 format from the Flocq library of floating-point arithmetic for Coq [4].

The fact that we were able to translate, far from trivial, multiple pages paper proofs in about 1 kloc of Coq is a very encouraging achievement. It is also worth noting that performing mechanically checked proofs gave the opportunity to fix a few small mistakes in the proofs, thus asserting the interest of formalized proofs.

We eventually hope that a large part of our code can be reused in future developments, for instance about numerical integration of ODEs.

It would also be interesting to study a recent work offering sharper bounds [15].

**Acknowledgements** The author wants to express its deepest thanks to Sylvie Boldo and Guillaume Melquiond as well as to Érik Martin-Dorel and Pierre-Marie Pédrot for their help regarding this work.

## References

1. Yves Bertot, Pierre Castéran, Gérard Huet, and Christine Paulin-Mohring. *Interactive theorem proving and program development : Coq'Art : the calculus of inductive constructions*. Texts in theoretical computer science. Springer, Berlin, New York, 2004. Données complémentaires <http://coq.inria.fr>.
2. Yves Bertot, Georges Gonthier, Sidi Ould Biha, and Ioana Pasca. Canonical big operators. In Otmane Aït Mohamed, César A. Muñoz, and Sofiène Tahar, editors, *TPHOLS*, volume 5170 of *Lecture Notes in Computer Science*, pages 86–101. Springer, 2008.
3. Sylvie Boldo, François Clément, Jean-Christophe Filliâtre, Micaela Mayero, Guillaume Melquiond, and Pierre Weis. Formal proof of a wave equation resolution scheme: The method error. In Matt Kaufmann and Lawrence C. Paulson, editors, *Interactive Theorem Proving, First International Conference, ITP 2010, Edinburgh, UK, July 11-14, 2010. Proceedings*, volume 6172 of *Lecture Notes in Computer Science*, pages 147–162. Springer, 2010.

4. Sylvie Boldo and Guillaume Melquiond. Flocq: A Unified Library for Proving Floating-point Algorithms in Coq. In *Proceedings of the 20th IEEE Symposium on Computer Arithmetic*, pages 243–252, Tübingen, Germany, July 2011.
5. Cyril Cohen. Construction of real algebraic numbers in coq. In Lennart Beringer and Amy P. Felty, editors, *ITP*, volume 7406 of *Lecture Notes in Computer Science*, pages 67–82. Springer, 2012.
6. The Coq development team. *The Coq proof assistant reference manual*, 2012. Version 8.4.
7. Florent de Dinechin, Christoph Quirin Lauter, and Guillaume Melquiond. Assisted verification of elementary functions using Gappa. In Hisham Haddad, editor, *Proceedings of the 2006 ACM Symposium on Applied Computing (SAC), Dijon, France, April 23-27, 2006*, pages 1318–1322. ACM, 2006.
8. Georges Gonthier, Assia Mahboubi, and Enrico Tassi. A Small Scale Reflection Extension for the Coq system. Research Report RR-6455, INRIA, 2008.
9. John Harrison. Floating point verification in HOL. In E. Thomas Schubert, Phillip J. Windley, and Jim Alves-Foss, editors, *Higher Order Logic Theorem Proving and Its Applications, 8th International Workshop, Aspen Grove, UT, USA, September 11-14, 1995, Proceedings*, volume 971 of *Lecture Notes in Computer Science*, pages 186–199. Springer, 1995.
10. Nicholas Higham. *Accuracy and Stability of Numerical Algorithms*. Society for Industrial and Applied Mathematics, Philadelphia, PA, USA, 1996.
11. IEEE Computer Society. IEEE Standard for Floating-Point Arithmetic. *IEEE Standard 754-2008*, 2008.
12. Pierre Roux and Pierre-Loïc Garoche. Computing quadratic invariants with min- and max-policy iterations: A practical comparison. In Cliff B. Jones, Pekka Pihlajasaari, and Jun Sun, editors, *FM 2014: Formal Methods - 19th International Symposium, Singapore, May 12-16, 2014. Proceedings*, volume 8442 of *Lecture Notes in Computer Science*, pages 563–578. Springer, 2014.
13. Siegfried Rump. Verification of positive definiteness. *BIT Numerical Mathematics*, 46:433–452, 2006.
14. Siegfried Rump. Verification methods: Rigorous results using floating-point arithmetic. *Acta Numerica*, 19:287–449, May 2010.
15. Siegfried Rump and Claude Pierre Jeannerod. Improved backward error bounds for lu and cholesky factorizations. *SIAM Journal on Matrix Analysis and Applications*, 35(2):684–698, 2014.