

(1.31) Definition (predefined functions).

Let  $M = F(\beta, m)$  be the set of machine numbers of the computer under consideration, and  $rd$  be the rounding function in  $M$ . The set  $FP$  of *predefined functions*, i.e. the functions that the computer can calculate by operating with the elements of  $M$ , is made up of three classes.

- The set of predefined functions corresponding to *arithmetic operations*. If  $\cdot$  is one of the arithmetic operations between real numbers  $+$ ,  $-$ ,  $\times$ ,  $/$  then the corresponding predefined function is indicated by the symbol  $\odot$  (a small circle containing the symbol of the operation considered) and is defined, for each pair  $\xi, \vartheta$  of elements of  $F(\beta, m)$  belonging to the domain of the operation  $\cdot$ , by

$$\xi \odot \vartheta = rd(\xi \cdot \vartheta)$$

- The set of predefined functions corresponding to the usual *elementary functions* ( $\sin$ ,  $\cos$ ,  $\arcsin$ ,  $\arccos$ ,  $\ln$ ,  $\exp$  ...). If  $f: A \rightarrow R$  is one of the elementary functions then the corresponding predefined function is indicated by the symbol  $F$  and is defined, for each element  $\xi$  of  $F(\beta, m)$  belonging to the domain  $A$  of the elementary function  $f$ , by

$$F(\xi) = rd(f(\xi))$$

- The set of predefined functions corresponding to *comparisons* between real numbers ( $<$ ,  $\leq$ ,  $=$ ,  $\neq$ ,  $\geq$ ,  $>$ ). In this case, since the elements of  $F(\beta, m)$  are real numbers, they are compared as such. Thus, the predefined functions corresponding to comparisons are simply the restrictions to  $F(\beta, m) \times F(\beta, m)$  of comparisons between real numbers (and it is not necessary to introduce new symbols to indicate them).

(1.32) Definition (algorithm, naive algorithm).

Let  $f_1, \dots, f_k$  be elementary functions or arithmetic operations and let  $f: A \rightarrow R$ , where  $A$  is a suitable subset of  $R$ , be the function obtained by *composing*  $f_1, \dots, f_k$ :

$$f(x) = f_1 \circ \dots \circ f_k(x)$$

(for example:  $f(x) = \sin(x) + \cos(x)$ , where  $f_3(x) = \sin(x)$ ,  $f_2(x) = \cos(x)$  and  $f_1(x_1, x_2) = x_1 + x_2$ ). If we ask *Scilab* to evaluate the function  $f$  with the instruction

```
> f(x)
```

the returned value will be

$$F_1 \circ \dots \circ F_k(rd(x))$$

where  $F_1, \dots, F_k(x)$  are, respectively, the predefined functions corresponding to  $f_1, \dots, f_k(x)$ .

The expression  $F_1 \circ \dots \circ F_k(rd(x))$  defines a function  $\varphi: A \rightarrow M$  called the *naive algorithm* for  $f$  (for the function in the example:  $\varphi(x) = \text{SEN}(rd(x)) \oplus \text{COS}(rd(x))$ , defined for every  $x$

in  $\mathbb{R}$ ). The term *algorithm* generally refers to a *finite* sequence of operations for calculating predefined functions.

Except that very special cases, there will be values of  $x$  for which  $f(x) \neq \varphi(x)$ . In these cases, we use  $\varphi(x)$  to approximate  $f(x)$ , and it is interesting to have *information on the error committed*.

To obtain this information we introduce the notions of *accurate algorithm*, *stable algorithm* and *well-conditioned computation of the value of a function*.

(1.33) Definition (accurate algorithm).

Let  $f:A \rightarrow \mathbb{R}$  be a function,  $\varphi:A \rightarrow \mathbb{M}$  the algorithm used to approximate the values of  $f$  and  $x \in A$ .

The algorithm  $\varphi$  is said to be *accurate* (when used to approximate the value of  $f$  at  $x$ ) if there exists a real number  $\varepsilon$  such that:

- (1)  $\varphi(x) = (1 + \varepsilon) f(x)$
- (2)  $\varepsilon$  'small'

If the algorithm is accurate for every  $x \in B \subset A$ , the algorithm is said to be accurate on  $B$ . In that case  $\varepsilon$  will depend on  $x$ .

(1.34) Remark.

- Let  $f$  and  $x$  be such that  $f(x) \neq 0$ . Condition (1) of the previous Definition is *equivalent* to the following:

$$\varepsilon = \frac{\varphi(x) - f(x)}{f(x)}$$

In this case, then, the algorithm is accurate is equivalent to saying that the relative error committed by approximating  $f(x)$  with  $\varphi(x)$  is 'small'.

- If the algorithm is accurate we have:  $f(x) = 0 \Leftrightarrow \varphi(x) = 0$ .
- The definition of accurate algorithm is *qualitative* because the term 'small' relative to  $\varepsilon$  is not quantified. The concrete meaning of the term 'small' depends on the individual case. For example, if, as in the case of the bisection method, it only matters that  $\varphi(x)$  and  $f(x)$  have the same sign,  $\varepsilon$  'small' means  $\varepsilon > -1$ .

Exercise: We use  $\lambda$  to approximate  $L > 0$ . What relative error  $\varepsilon$  is made using  $\lambda = 0$ ? What value of  $\lambda$  should be used to obtain a relative error  $\varepsilon = 1$ ?

(1.35) Definition (stable algorithm).

Let  $f:A \rightarrow \mathbb{R}$  be a function,  $\varphi:A \rightarrow \mathbb{M}$  be the algorithm used to approximate the values of  $f$  and  $x \in A$ .

The algorithm  $\varphi$  is said to be *stable* (when used to approximate the value of  $f$  at  $x$ ) if

there exist real numbers  $\varepsilon_a, \varepsilon_v$  such that:

$$(1) \varphi(x) = (1 + \varepsilon_v) f((1 + \varepsilon_a)x)$$

$$(2) \varepsilon_a, \varepsilon_v \text{ 'small'}$$

If the algorithm is stable for every  $x \in B \subset A$ , we say that the algorithm is stable on B.

In this case  $\varepsilon_a, \varepsilon_v$  will depend on x.

(1.36) Remark.

- If an algorithm is accurate then it is stable ( $\varepsilon_a = 0, \varepsilon_v = \varepsilon$ ); a stable algorithm *may not* be accurate.
- Informally: a stable algorithm returns a *good approximation* ( $\varepsilon_v$  'small') of the value of f at a point *close* to x ( $\varepsilon_a$  'small').

(1.37) Remark ('good' algorithm).

The notion of stability formalizes the idea of a 'good' algorithm for approximating the values of a given f. For example, if f is an elementary function and  $\varphi$  is the naive algorithm for f, then, calling F the predefined function corresponding to f, we have:

$$\varphi(x) = F(\text{rd}(x)) = \text{rd}(f(\text{rd}(x)))$$

(1.38) Theorem (relative error and perturbation).

Recalling the Definition of relative error committed by approximating a real number t with its rounded  $\text{rd}(t)$  and Theorem (1.28) of Lesson 5 on the limitation of the relative error, we obtain:

Let x be a real number and rd be the rounding function in  $F(\beta, m)$ . There exists a real number  $\varepsilon$  such that:

$$\text{rd}(x) = (1 + \varepsilon)x \quad \text{and} \quad |\varepsilon| < u$$

The equality expresses the rounded of x as a (small) *multiplicative perturbation* of x.

(Proof: if  $x \neq 0$  then  $\varepsilon$  is the relative error committed by approximating x with  $\text{rd}(x)$ ; if  $x = 0$  (and therefore  $\text{rd}(x) = 0$ ) the equality holds, for example, with  $\varepsilon = 0$ .)

(1.39) Remark (continuation of the previous one).

Using the previous Theorem twice we finally get:

$$\varphi(x) = (1 + \varepsilon_2)f((1 + \varepsilon_1)x) \quad \text{with} \quad |\varepsilon_1| < u \quad \text{e} \quad |\varepsilon_2| < u$$

The algorithm  $\varphi$  returns *the best possible approximation* of the value of f at the point *closest* to x. In this sense,  $\varphi$  is the 'best possible' algorithm that the computer can use to approximate  $f(x)$ . Hence, generalizing, the idea that a 'good' algorithm for

approximating the value of a function at a given point is an algorithm that returns a good approximation of the value of the function at a point close to the one where we wanted to compute it.

(1.40) Definition (well-conditioned computation of the value of a function).

Let  $f:A \rightarrow \mathbb{R}$  be a function and  $x \in A$ . The computation of the value of  $f$  at  $x$  is *well-conditioned* if: for every 'small' real number  $\alpha$  there exists a 'small' real number  $\varepsilon_v$  such that

$$f((1 + \alpha)x) = (1 + \varepsilon_v)f(x)$$

Informally: the computation of the value of  $f$  at  $x$  is well conditioned if the value of  $f$  at any point 'near'  $x$  is a 'good' approximation of the value of  $f$  at  $x$ .

(1.41) Remark.

- The property that the computation of the value of  $f$  at  $x$  is well-conditioned concerns *only* the function  $f$ . In particular, it is not related to the *algorithm* chosen to approximate the values of  $f$ .
- If  $f(x) \neq 0$ , the value of  $\varepsilon_v$ , once  $\alpha$  is assigned, is *determined*. Specifically,  $\varepsilon_v$  is:

$$\varepsilon_v = \frac{f((1 + \alpha)x) - f(x)}{f(x)}$$

(1.42) Theorem (stability + well-conditioning  $\Rightarrow$  accuracy).

Let  $f:A \rightarrow \mathbb{R}$  be a function,  $x \in A$ , and  $\varphi$  be the algorithm used to approximate  $f(x)$ . If the algorithm is *stable* and the computation of  $f$  at  $x$  is *well-conditioned*, then the algorithm is *accurate*.

Proof. By the stability of the algorithm there exist  $\varepsilon_1$  and  $\varepsilon_2$  such that:

$$\varphi(x) = (1 + \varepsilon_2)f((1 + \varepsilon_1)x) \quad \text{and} \quad \varepsilon_1, \varepsilon_2 \text{ 'small'}$$

By the well-conditioning of the calculation of  $f$  at  $x$  there exists  $\varepsilon_3$  such that:

$$f((1 + \varepsilon_1)x) = (1 + \varepsilon_3)f(x) \quad \text{and} \quad \varepsilon_3 \text{ 'small'}$$

Then we can rewrite:

$$\varphi(x) = (1 + \varepsilon_2)(1 + \varepsilon_3)f(x)$$

and, introducing  $(1 + \varepsilon_2)(1 + \varepsilon_3) = 1 + t$ , i.e.  $t = \varepsilon_2 + \varepsilon_3 + \varepsilon_2\varepsilon_3$ , we get:

$$\varphi(x) = (1 + t)f(x) \quad \text{and} \quad t \text{ 'small'}$$

so the algorithm is accurate.

(1.43) Remark (stability of naive algorithms in elementary cases).

- From what we deduced in Remarks (1.37) and (1.39), if  $f:A \rightarrow R$  is an elementary function and  $\varphi$  is the naive algorithm for  $f$ ,  $\varphi$  is stable on  $A$ : *for each elementary function the naive algorithm is stable.*
- Let  $f(x_1, x_2) = x_1 + x_2$ . The naive algorithm for  $f$  is:

$$\varphi(x_1, x_2) = \text{rd}(x_1) \oplus \text{rd}(x_2)$$

Recalling the definition of  $\oplus$  (see Definition (1.31)) and using Theorem (1.38) three times we obtain:

$$\varphi(x_1, x_2) = (1 + \varepsilon_3) \left( (1 + \varepsilon_1)x + (1 + \varepsilon_2)x \right) \quad , \quad \text{with } |\varepsilon_j| \leq u \quad , \quad j = 1, 2, 3$$

So, the naive algorithm for the sum is stable.

Similarly, *the naive algorithm for each of the arithmetic operations is shown to be stable.*