

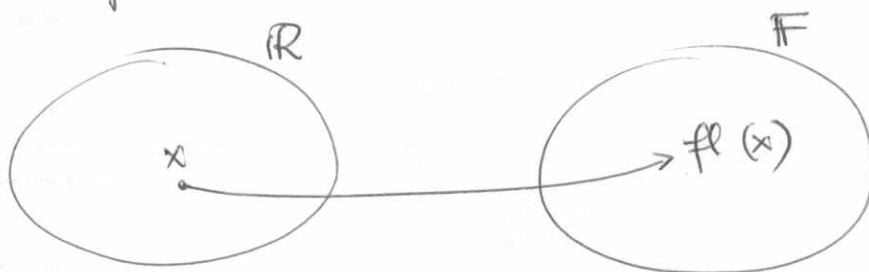
# THE FLOATING POINT SYSTEM

$\mathbb{R}$  = set of real numbers  $\mathbb{R} = (-\infty, +\infty)$

$\mathbb{F}$  = set of machine numbers  $\mathbb{F} \subset \mathbb{R}$

$\mathbb{F}$  is a finite (i.e. its cardinality is finite) and bounded set

Not any  $x \in \mathbb{R}$  can be exactly represented on the computer, thus we denote by  $fl(x)$  the representative number of  $x$  in the set  $\mathbb{F}$ .



In order to describe  $\mathbb{F}$  we start by the scientific (or exponential) form of a real number.

let us consider  $x \in \mathbb{R}$  and we write it as

$$x = (-1)^s 0.a_1 a_2 \dots a_t \dots \beta^e = (-1)^s \underbrace{a_1 a_2 \dots a_t}_{m} \beta^{e-t}$$

where  $s = 0$  or  $1$

$\beta$  = basis

$e$  = exponent of the basis.

$a_i$  are integer numbers :  $0 \leq a_i < \beta - 1, i \geq 1$   
with  $a_1 \neq 0$

~~0.9999999999~~  $m$  is said mantissa

We cannot store infinite digits  $a_1 a_2 \dots a_t \dots a_\infty$  but only a finite number, say  $t$

$$\Rightarrow fl(x) = (-1)^s a_1 a_2 \dots \tilde{a}_t \beta^{e-t} = (-1)^s 0.a_1 a_2 \dots \tilde{a}_t \beta^e$$

$$\text{The last digit } \tilde{a}_t = \begin{cases} a_t & \text{if } a_{t+1} < \beta/2 \\ a_{t+1} & \text{if } a_{t+1} \geq \beta/2 \end{cases}$$

ex 1:  $x = 123.45678 = (-1)^0 0.12345678 \cdot 10^3$

$t = 6$   $fl(x) = (-1)^0 0.123457 \cdot 10^3 \neq x$

ex 2:  $x = 123.45678$

$t = 8$   $fl(x) = (-1)^0 0.12345678 \cdot 10^3 = x$

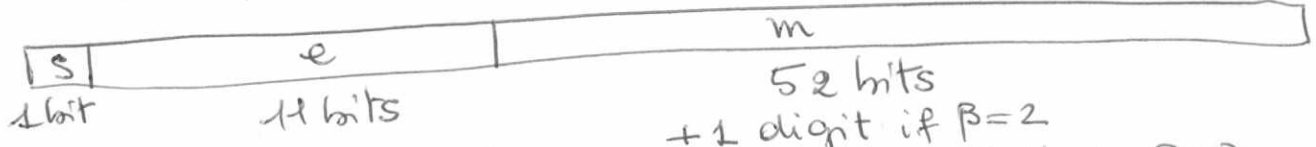
Thus  $fl(x)$  can be stored in a finite number of Bytes. Why?

Single precision format = 4 Bytes = 32 bits



$-126 \leq e \leq 127$   
 $23 \text{ bits} = 23 \text{ digits for } \beta=2$   
 $\cong 6,9 \text{ digits for } \beta=10$   
 $(2^{23} = 10^x \Rightarrow x = \log_{10} 2^{23})$

Double precision format = 8 Bytes = 64 bits



$-1022 \leq e \leq 1023$   
 $+1 \text{ digit if } \beta=2$   
 $53 \text{ bit} = 53 \text{ digits in } \beta=2$   
 $\sim 16 \text{ digits in } \beta=10$

We denote by:  $t = n^{\circ}$  of digits with respect to  $\beta$

$L, U =$  lower and upper bounds  
for  $e$  ( $L < 0 < U$ )

$\Rightarrow$  the set of numbers the machine can represent  
is denoted by  $\mathbb{F} = \mathbb{F}(\beta, t, L, U) :$

$$\mathbb{F}(\beta, t, L, U) = \{ x = (-1)^s 0.a_1 \dots a_t \cdot \beta^e, \text{ with } s=0,1 \\ 0 \leq a_i \leq \beta-1, a_1 \neq 0, L \leq e \leq U \} \cup \{0\}$$

Remark: zero cannot be represented by the  
floating point form, since it has all  
 $a_i = 0$ .

$\mathbb{F}$  is said floating point system

"floating point" means that the point is not fixed,  
but it moves to match the form  $x = (-1)^s 0.a_1 \dots a_t \beta^e$

Properties of  $\mathbb{F}$ :

1.  $\text{card}(\mathbb{F}) = 2(\beta-1)\beta^{t-1}(U-L+1)+1 < +\infty$   
( $\mathbb{F}$  is finite)

2. the smallest positive number is  $x_{\min} = \beta^{L-1}$

3. the largest positive number is  $x_{\max} = (1 - \beta^{-t})\beta^U$

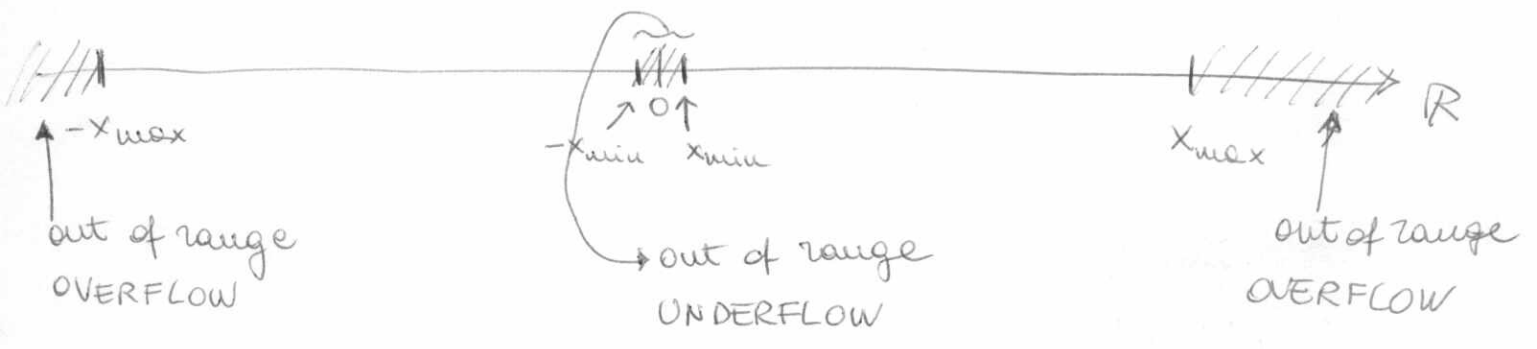
4.  $\mathbb{F}$  is symmetric w.r.t. zero

5. floating point numbers are more dense close to 0  
and less dense far from zero

Example  $\mathbb{F}(2, 3, -2, 3)$   
 $\beta \leq L \leq U$

$x_{\min} = 0.100 \cdot \beta^{-2} = [1 \cdot \beta^{-1} + 0 \cdot \beta^{-2} + 0 \cdot \beta^{-3}] \cdot \beta^{-2} = \beta^{-3} = 2^{-3} = \frac{1}{8}$   
read these numbers following the positional notation

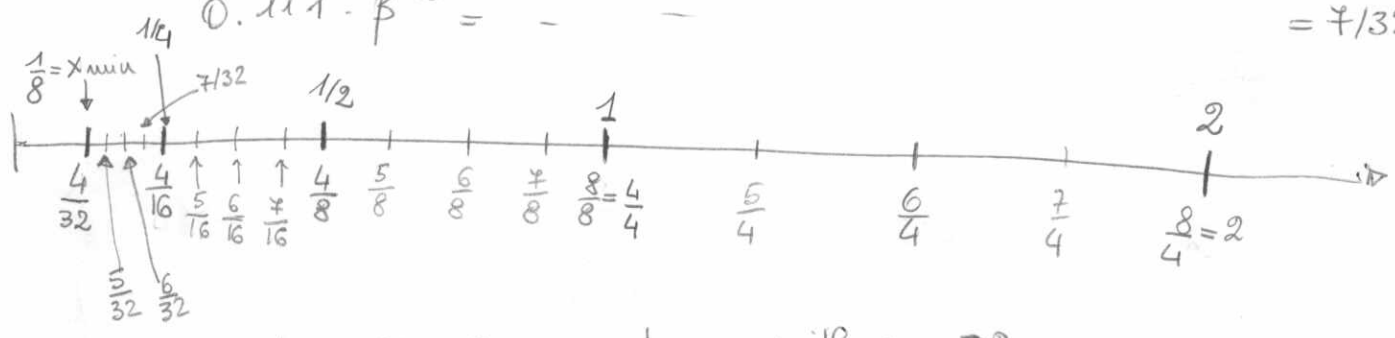
$x_{\max} = 0.111 \cdot \beta^{+3} = [1 \cdot \beta^1 + 1 \cdot \beta^2 + 1 \cdot \beta^3] \beta^3 = [\frac{1}{2} + \frac{1}{4} + \frac{1}{8}] \cdot 8 = 7$



$\mathbb{F}$  is bounded

Now we compute the ~~max~~ f.p. numbers between  $x_{\min}$  and  $x_{\max}$

$x_{\min} = 0.100 \cdot \beta^{-2} = [1 \cdot \beta^{-1} + 0 \cdot \beta^{-2} + 0 \cdot \beta^{-3}] \beta^{-2} = \frac{1}{8} = \frac{4}{32}$   
 $0.101 \cdot \beta^{-2} = [1 \cdot \beta^{-1} + 0 \cdot \beta^{-2} + 1 \cdot \beta^{-3}] \beta^{-2} = (\frac{1}{2} + \frac{1}{8}) \cdot \frac{1}{4} = \frac{5}{8} \cdot \frac{1}{4} = \frac{5}{32}$   
 $0.110 \cdot \beta^{-2} = \dots = \frac{6}{32}$   
 $0.111 \cdot \beta^{-2} = \dots = \frac{7}{32}$



We wrote all fp numbers with  $e = -2$ , now we

first  $e = -1$

$0.100 \cdot \beta^{-1} = 1/4 = 4/16$	$e = 0$
$0.101 \cdot \beta^{-1} = 5/16$	
$0.110 \cdot \beta^{-1} = 6/16$	
$0.111 \cdot \beta^{-1} = 7/16$	

then with  $e = 0, e = 1, \dots$

$e=1: 0.100 \cdot \beta^1 = 1 = 4/4$

$0.101 \cdot \beta^1 = 5/4$

$0.110 \cdot \beta^1 = 6/4$

$0.111 \cdot \beta^1 = 7/4$

-----

In Matlab  $\beta=2, t=53 (=52+1)$

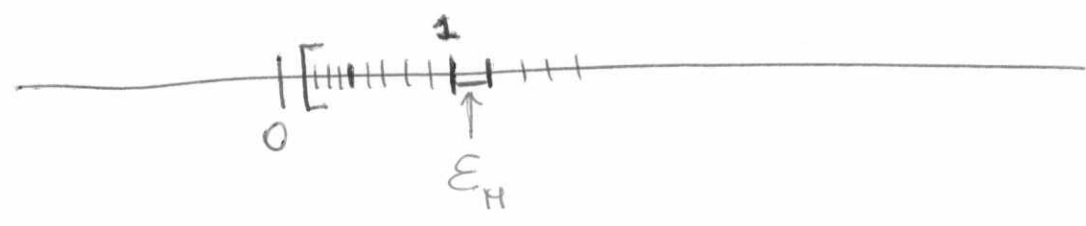
$U = -1024, L = 1024$

$x_{min} \approx 2.2251 \cdot 10^{-308} \gg \text{realmin}$

$x_{max} \approx 1.7977 \cdot 10^{308} \gg \text{realmax}$

MACHINE PRECISION  $\epsilon_H$

It is the distance between 1 and the smallest floating point number greater than 1.



$$\epsilon_H = 0.\underbrace{10\dots01}_t \cdot \beta^1 - 0.\underbrace{100\dots0}_t \cdot \beta^1 =$$

$$= 0.10\dots01 \cdot \beta^1 - 0.10\dots00 \cdot \beta^1 =$$

$$\frac{0.00\dots\underbrace{1}_t \cdot \beta^1}{\beta^t} = \beta^{1-t}$$

$$\Rightarrow \boxed{\epsilon_H = \beta^{1-t}}$$

Remark :  $\epsilon_H$  is not the smallest positive f.p. but it is the distance between 1 and the smallest fp number greater than 1.

In  $\mathbb{F}(2, 3, -2, 3)$  we have

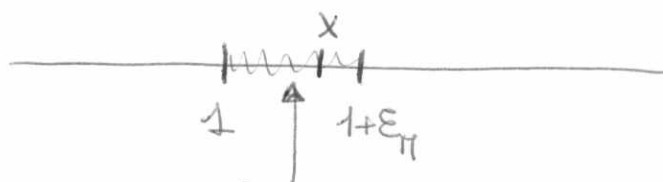
$$x_{\min} = \frac{1}{8}$$

while  $\epsilon_H = \beta^{1-t} = 2^{1-3} = \frac{1}{4}$ , In matlab  $\epsilon_H = 2^{1-53} = 2^{-52} \approx 2.22 \cdot 10^{-16}$

Why is  $\epsilon_H$  so important?

Not every  $x \in \mathbb{R} \cap [x_{\min}, x_{\max}]$  is a f.p. number

let us consider  $x$  between 1 and  $1 + \epsilon_H$

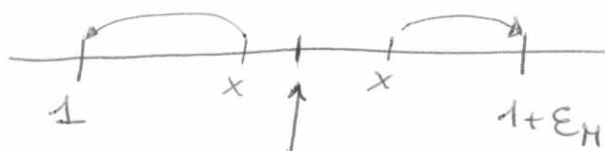


here no other f.p. number is present

This  $x \notin \mathbb{F} \Rightarrow \text{fl}(x)$  is the floating point number closest to  $x$  and  $\text{fl}(x) \neq x$ .

$\text{fl}(x) = 1$  or  $\text{fl}(x) = 1 + \epsilon_H$ ?

It depends if  $x$  is on the left or on the right of the middle point between 1 and  $(1 + \epsilon_H)$



$$\text{middle point} = \frac{1 + (1 + \epsilon_H)}{2} = 1 + \frac{\epsilon_H}{2}$$

if  $x \in [1, 1 + \frac{\epsilon_H}{2}] \Rightarrow \text{fl}(x) = 1$

if  $x \in (1 + \frac{\epsilon_H}{2}, 1 + \epsilon_H] \Rightarrow \text{fl}(x) = 1 + \epsilon_H$

We are interested in evaluating which error we generate ~~ex~~ when we approximate  $x$  with  $fl(x)$ .

We have:

$$\boxed{\frac{|x - fl(x)|}{|x|} \leq \frac{1}{2} \epsilon_H} \quad \forall x \in [-x_{max}, -x_{min}] \cup [x_{min}, x_{max}]$$

$\frac{1}{2} \epsilon_H$  is the semi-amplitude of the interval  $[1, 1 + \epsilon_H]$

and  $\frac{|x - fl(x)|}{|x|}$  is the relative rounding error

generated by the floating point system, or equivalently, by the machine arithmetics.

We measure relative (and not absolute) errors, since relative errors are more meaningful and they are independent of the order of magnitude of the measure.

We set  $\| \underline{u} = \frac{1}{2} \epsilon_H$  roundoff unit  $\|$ .

Obviously, if  $x \in \mathbb{R} \cap \mathbb{F} \Rightarrow fl(x) = x$  and the error is null.

Remark  $x = 0.1$  (in decimal expansion) =  $0.000110011001100 \dots$  ( $\beta = 2$ )  
has infinite digits in binary expansion.

$$\Rightarrow fl(0.1) \neq 0.1 \quad \text{and} \quad 0 < \frac{|0.1 - fl(0.1)|}{|0.1|} \leq \frac{1}{2} \epsilon_H = \frac{1}{2} \beta^{1-t}$$

The larger  $t$  is and the smaller the rounding errors are.





(9)

Remark: Why did I write  $\text{fl}(\bar{x} + \bar{y})$  if  $\bar{x}, \bar{y} \in \mathbb{F}$ ?

Because, even if  $\bar{x}, \bar{y} \in \mathbb{F}$ , it is not guaranteed that  $\bar{x} + \bar{y} \in \mathbb{F}$ , thus we have to approximate  $(\bar{x} + \bar{y})$  by its float.

---

We want to compute  $x \cdot y$ , the machine computes  $\text{fl}(\bar{x} \cdot \bar{y})$  and the error is

$$e_{\cdot} = \frac{|\text{fl}(\bar{x} \cdot \bar{y}) - xy|}{|xy|} \leq 3u + O(u^2)$$

independently of  $x$  and  $y$ , thus when we multiply two fp numbers the error on the result is of the same order as the errors on the data.

---

We say that the sum is potentially UNSTABLE; even if the errors on data are small, we can have large errors on the solution.

The product is a STABLE operation; when the errors on the data are small, the error on the result is small, too.

---

Ex:  $x = 0.1234$      $y = -0.1231$   
 $x + y = 0.1234 - 0.1231 = 3 \cdot 10^{-4}$

$$\bar{x} = 0.123 \quad \bar{y} = -0.123$$

$$\beta = 10 \quad t = 3$$

$$E_H = 10^{-2}$$

$$u = \frac{1}{2} \cdot 10^{-2}$$

$$\Rightarrow \bar{x} + \bar{y} = 0 \quad e_{+} = \frac{|0 - 3 \cdot 10^{-4}|}{3 \cdot 10^{-4}} = 1 = 100\% \text{ di errore}$$

$$C_{xy} = \frac{|x| + |y|}{|x + y|} = \frac{0.2465}{3 \cdot 10^{-4}} \approx 821,6 \quad u = \frac{1}{2} \beta^{-t} = \frac{1}{2} \cdot 10^{-2}; \quad C_{xy} u \approx 4,108$$

$\rightarrow e_{+} \leq C_{xy} \cdot u$

Associativity is violated whenever a situation of overflow or underflow occurs.

$$a = 1.e+308$$

$$b = 1.1e+308$$

$$c = -1.001e+308$$

$$a + (b+c) = 1.0990e+308$$

$$(a+b)+c = \text{Inf.}$$

Unicity of zero is violated

$$1.e+10 + 1.e-7 = 1.e+10$$

i.e.  $1.e-7$  is neutral in this addition

$$x+y = x \quad \forall y < |x|u \quad (u = \text{roundoff unity})$$

---

Reference : Scientific computing (or Calculo Científico)  
chapter 1.