

## 2. Représentat° des nombres, Erreurs d'arrondi

Sources d'erreurs : données, méthode de calcul, arrondi (précision finie)

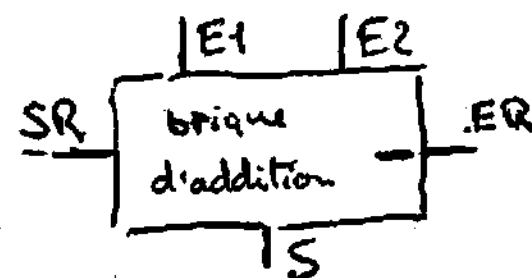
A. représentat° de nombres en base qq :  $\forall x \in \mathbb{R} : x = \epsilon B^N \sum_{v=1}^{\infty} x_v B^{-v}$   $\epsilon \in \{+1, -1\}$ ,  $x_v \in \{0, 1, \dots, B-1\}$

unique (si  $\forall m \in \mathbb{N} \exists v \geq m : x_v \neq B-1$ )

numérat°	base	B	digits	N : exposant	(x <sub>v</sub> ) : mantisse
binair	binair	2	0, 1		
octale	octale	8	0, 1, ..., 7		
décimale	décimale	10	0, 1, ..., 9		
hexadécimale (sic)	hexadécimale	16	0, 1, ..., 9, A, B, ..., F		

exemples  $3_{10} = 11_2$   
 $9_{10} = 9_{16} = 11_8 = 1001_2$

calcul binaire :  
 $0+0 = 0$   
 $0+1 = 1$   
 $1+1 = 10$   
 $1+0 = 1$



suffit : on cascade ↗

E1	E2	ER	S	SR
0	0	0	0	0
0	0	1	0	0
0	1	0	1	0
0	1	1	0	1
1	0	0	1	0
1	0	1	0	1
1	1	0	0	1
1	1	1	1	1

Cours ②

Cours ③

↓ 3/10/2005

139072  
 262144  
 524288  
 1048576  
 2097152  
 4194304  
 8388608  
 16777216  
 33554432  
 67108864  
 134217728  
 268435456  
 536870912  
 1073741824

2147483648 = 2<sup>31</sup>  
 4294967296 = 2<sup>32</sup>

B calcul à virgule fixe

util courants:  $N$  donné  $\rightarrow$  ne consomme pas de mémoire, on ne stocke que les chiffres (etc)  
 : entiers  $x = \sum_{v=0}^N x_v B^v$  «int» en C : 32bits +  $2^{31} \sim 2,1 \cdot 10^9$   
 inconvénient : inadapté à grandeurs sur plusieurs décades  $2^{31} \sim 2,1 \cdot 10^9$   
 $2^{31} \sim 2,1 \cdot 10^9$   
 $2^{31} \sim 2,1 \cdot 10^9$

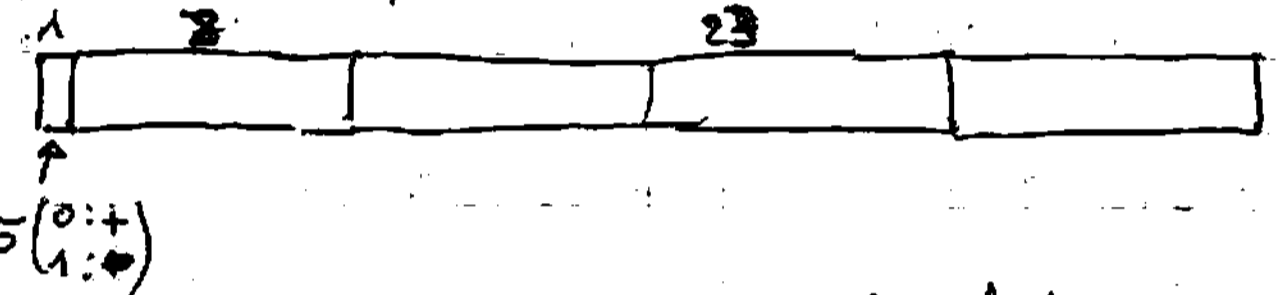
limites  $\frac{x_{MAX}}{x_{MIN}} = B^t$  ,  $x_{MAX} = B^{N-t}$ ,  $x_{MIN} = B^{N-t}$   
 $\sim 2^{32}$   $\sim 2^{32}-1$   $\sim 1$   $\rightarrow 2^8 = 256$

nombres négatifs :  $x \mapsto B^N + x$  ex: char:  $(N=8, B=2) \rightarrow \in \{-128, \dots, 127\}$   
 $-128_{10} = 10000000_2$   
 $-2_{10} = 256 - 2 \pmod{2^8} = 254 = 11111110_2$   
 $-1_{10} = 255 = 11111111_2$   
 $0 = 0$   
 $1_{10} = 00000001_2$   
 $127 = 01111111_2$

C nombres à virgule flottante

ou stocke  $N$  et  $(x_i)$

exemple : type «float» C sur Intel Pentium : 32bits



inconvenient : compromis de stockage, moins de place pour mantisse

avantage :  $|x|_{MIN} = B^{N_{min}-T}$  ,  $|x|_{MAX} = B^{N_{max}-T}$  ,  $\frac{|x|_{MAX}}{|x|_{MIN}} = B^{N_{max}-N_{min}}$   
 $\sim 2^{-151} \sim 10^{-45}$   $\sim 2^{127} \sim 10^{40}$  float 40  $\rightarrow N_{max}-N_{min} = 255, T=23$   
 $\sim 2^{278}$

D Règle d'arrondi

$x$  nombre réel,  $\tilde{x}$  approx :  $x - \tilde{x}$  erreur absolue  
 $-E = \frac{x - \tilde{x}}{\tilde{x}}$  (ou  $\frac{x - \tilde{x}}{x}$ ) erreur relative (prà  $\tilde{x}$ )

supposons qu'il n'y ait pas de dépassement pour l'exposant, B pair

$x = B^N \sum_{v=1}^t x_v B^{-v}$   $\rightarrow$  Arr(x) =  $\begin{cases} B^N \sum_{v=1}^t x_v B^{-v} & \text{si } x_{-t-1} < \frac{B}{2} \\ B^N \sum_{v=1}^t x_v B^{-v} + B^{-t} & \text{si } x_{-t-1} \geq \frac{B}{2} \end{cases}$

Arr : arrondi à t chiffres (t : longueur de mantisse)

$\frac{Arr_t(x) - x}{Arr_t(x)} \leq 0.5 B^{-t+1}$   $Arr_t(x) - x \leq 0.5 B^{N-t}$   
 $\leq 0.5 B^{-t+1}$  précision du calcul à virg. flottante, ex  $t=23 \rightarrow 2^{-24} \sim 10^{-7}$

$x = 5 \cdot 10^N \cdot m$ ,  $0.1 \leq m < 1$   $\wedge$   $\tilde{x} = 5 \cdot 10^N \cdot \tilde{m}$  une approx.

nombre de chiffres significatifs de  $\tilde{x} = \max \{t \in \mathbb{Z} \mid |m - \tilde{m}| \leq 0.5 \cdot 10^{-t+1}\}$

# $\epsilon$ calcul en virgule flottante

pb. : si  $a$  et  $b$  sont nombres à v.f. de long. de mantisse  $t$ ,

$a \square b$  ( $\square \in \{+, -\}$ ) n'en est pas forcément

(ex.  $0.123 \cdot 10^4 + 0.456 \cdot 10^{-3}$  (donc  $t=3$ ) =  $0.1230000456 \cdot 10^4$  ( $t=10$ ))

$\Rightarrow$  le résultat doit être arrondi à la fin  $\Rightarrow FL_t(a \square b)$

Habituellement les ordinateurs sont construits/progr. en sorte que

$$FL_t(a \square b) = Arr_t(a \square b) \approx$$

$$\Rightarrow FL_t(a \square b) = (a \square b)(1 + \epsilon) \quad , |\epsilon| < \tau$$

$$\begin{aligned} a &= \overset{FL}{Arr_t}(x) & \text{ex } 0.100 \cdot 10^1 &= Arr_2(0.9995 \cdot 10^0) \\ b &= Arr_t(y) & -0.998 \cdot 10^0 &= Arr_2(0.9984 \cdot 10^0) \end{aligned}$$

$$\Rightarrow |FL_t(FL_t(x) \square FL_t(y))| = |FL_t(x) \square FL_t(y)| (1 + \tau)$$

$$\begin{aligned} & \overset{0.200 \cdot 10^2}{=} (x(1 + \epsilon_x) \square y(1 + \epsilon_y))(1 + \epsilon) \\ & \overset{0.11 \cdot 10^{-2}}{=} x \pm y + x(\epsilon + \epsilon_x(1 + \epsilon)) \pm y(\epsilon + \epsilon_y(1 + \epsilon)) \end{aligned}$$

$$\square \epsilon \quad + - \quad \rightarrow \quad = x \pm y + x(\epsilon + \epsilon_x(1 + \epsilon)) \pm y(\epsilon + \epsilon_y(1 + \epsilon)) = (x \pm y)(1 + \delta) = (x \pm y) + F$$

ex  $F = 0.9995 \cdot 0.5003 \cdot 10^{-3} + 0.9984 \cdot 0.4006 \cdot 10^{-3} = 0.9000 \cdot 10^{-3}$

$$\delta = \frac{F}{x \pm y} = 0.82 \quad (!)$$

$$= x \pm y + (x \pm y) \epsilon + (\epsilon_x x + \epsilon_y y)(1 + \epsilon)$$

$$|\delta| = \frac{|x(\epsilon + \epsilon_x(1 + \epsilon)) \pm y(\epsilon + \epsilon_y(1 + \epsilon))|}{|x \pm y|} \leq \left| \epsilon + \frac{x\epsilon_x(1 + \epsilon) \pm y\epsilon_y(1 + \epsilon)}{x \pm y} \right|$$

$$\leq |\epsilon| + |1 + \epsilon| \frac{|\epsilon_x x| + |\epsilon_y y|}{|x \pm y|} \leq \tau + (1 + \tau)\tau \frac{|x| + |y|}{|x \pm y|}$$

$$\approx \tau \left(1 + \frac{|x| + |y|}{|x \pm y|}\right) \quad \sim 2 \text{ cas} \quad \delta \approx 2\tau \text{ et } \delta \approx \tau$$

$$\square \epsilon \quad \rightarrow \quad = x \square y + x \square y \cdot \epsilon + (1 + \epsilon)(x\epsilon_x \square y + \epsilon_y x \square y + \epsilon_x \epsilon_y x \square y)$$

$$= x \square y (1 + \epsilon + (1 + \epsilon)(\epsilon_x + \epsilon_y + \epsilon_x \epsilon_y))$$

$$\delta = \epsilon + (1 + \epsilon)(\epsilon_x + \epsilon_y + \epsilon_x \epsilon_y) \leq \tau + (1 + \tau)(2\tau + \tau^2) = 3(\tau + \tau^2) + \tau^3 \sim 3\tau$$

3/10  
cours  
↑  
↓ cours ⑦  
10/10

## F Evaluat° stable/instable

Toute formule compliquée se réduit in fine à une suite d'op. élem.

→ il faut s'assurer que chaque pas soit numériquement stable.

⇒ Importance du choix de l'algo!

ex.:  $ax^2 + bx + c = 0$ , soit  $|4ac| < b^2$

$$\rightarrow x_1 = \frac{1}{2a}(-b - \operatorname{sgn}(b)\sqrt{b^2 - 4ac}), \quad x_2 = \frac{1}{2a}(-b + \operatorname{sgn}(b)\sqrt{b^2 - 4ac})$$

$$= -\frac{\operatorname{sgn}(b)}{2a}(|b| - \sqrt{b^2 - 4ac})$$

→  $x_2$  numériquement instable

mieux vaut utiliser  $x_2 \approx x_1 = \frac{c}{a} \rightarrow x_2 = \frac{2c}{-b - \operatorname{sgn}(b)\sqrt{b^2 - 4ac}}$

## G Condition d'un pb

$$\varphi: D \rightarrow \mathbb{R}, \quad D \subset \mathbb{R}^n, \quad \varphi \in \mathcal{C}^2(D, \mathbb{R})$$

$$x \mapsto \varphi(x) = y$$

$$\delta y \approx \sum \frac{\partial \varphi}{\partial x_i}(x) (\tilde{x}_i - x_i)$$

err. relative  $\frac{\delta y}{y} = \sum \underbrace{\frac{x_i}{\varphi(x)} \frac{\partial \varphi}{\partial x_i}(x)}_{\text{coeff. de condition}} \frac{\tilde{x}_i - x_i}{x_i}$

$\leq 1$  bien conditionné  
 $> 1$  mal "

ex:  $\varphi(x) = ax \rightarrow \frac{x}{\varphi(x)} \frac{\partial \varphi}{\partial x} \Big|_x = 1$

$\varphi(x) = a+x \rightarrow \frac{x}{a+x} \begin{matrix} < 1 \\ > 1 \end{matrix} \quad \begin{matrix} x \cdot a > 0 \\ x \cdot a < 0 \end{matrix}$

→  $\cdot, \pm$  opérat° bien conditionnées

$+, -$  bien " si termes  $\hat{m}$  signe/signes différents  
mal " " signes diff /  $\hat{m}$  signe

ex  $x^2 + 2px - q = 0, \quad p, q > 0 \wedge p \gg q$

calc. plus grand zéro  $\varphi(p, q) = -p + \sqrt{p^2 + q}$

sol° :  $\bullet s = p^2$   
 $t = s + q$   
 $u = \sqrt{t}$

+ 2 méthodes : ①  $y = \varphi_1(u) = -p + u$   
ou ②  $v = -p - u$  et  $y = \varphi_2(v) = \frac{-q}{v}$

$$\frac{\delta y}{y} = \frac{p}{\varphi(p, q)} \frac{\partial \varphi}{\partial p} \varepsilon_p + \frac{q}{\varphi(p, q)} \frac{\partial \varphi}{\partial q} \varepsilon_q = -\frac{p}{(p^2 + q)^{1/2}} \varepsilon_p + \frac{p + (p^2 + q)^{1/2}}{2(p^2 + q)^{1/2}} \varepsilon_q \rightsquigarrow \text{b.c.}$$

mais méth. 1:  $\frac{\delta y}{y} = \frac{u}{-p+u} \varepsilon_u = \frac{1}{q(p(p^2+q)^{1/2} + p^2+q)} \varepsilon_u > \frac{2p^2}{q} \varepsilon_u \gg \varepsilon_u \rightsquigarrow \text{m.c.}!$

méth 2:  $\frac{\delta y}{y} = -\frac{(p^2+q)^{1/2}}{p+(p^2+q)^{1/2}} \varepsilon_u \rightsquigarrow \text{b.c.}$

suite cours ④ ph332  
10/10/2005

conclusion : coeffs de conditi<sup>o</sup> renseignent sur la stab. num.

Il peut y avoir amplif. et atténuat<sup>e</sup> des erreurs

coeffs  $< 1$  : pour un bon choix des étapes de calcul  $\rightarrow$  stabilité  
mauvais choix : quand même instab.

coeff  $> 1$  : ~~X~~ méthode stable,  $T_j$  instable

H autres techniques d'analyse d'erreur

calcul direct à priori  $\rightarrow$  cas le pire

\ analytiquement

\ numériquement = calcul d'intervalles

calcul à posteriori : quelles erreurs des données initiales auraient produit le résultat obtenu ?