

EXERCICE 1: ENCODAGE DES NOMBRES ENTIERS

On considère la base $\beta = 2$.

1. Écrire sur 8 bits, en numération simple à position, les valeurs décimales suivantes : 31, 247
2. Écrire sur 8 bits, en complément à 2, les valeurs décimales suivantes : 44, -9, 0, -31.
3. Définir la plage des valeurs entières représentables avec un variable du type :
 - `short` : 2 octets
 - `int` : 4 octets
 - `long` : 8 octets

EXERCICE 2: ENCODAGE DES NOMBRES FLOTTANTS

1. Rappeler la définition d'un nombre flottant ? Qu'est-ce qu'un nombre flottant normalisé ?
2. On considère la base $\beta = 2$. Convertir en décimal les représentations (signe, mantisse, exposant) suivantes :

$$x = (-1)^s \left(\sum_{i=1}^r m_i \beta^{-i} \right) \beta^e$$

	s	m	e
a =	1	10 1101	3
b =	0	1100 0001	-5
c =	0	1 0000 0001	8

3. On considère une base β et soit $x \in \mathbb{R}$. Définir l'arrondi de x , noté $A(x)$, en utilisant r chiffres significatifs pour la mantisse. Définir l'erreur de représentation et la précision machine ε .
4. On considère la base $\beta = 10$. Représenter le réel 231,345 en virgule flottante avec 4 chiffres significatifs. Calculer les erreurs absolue et relative de représentation.

EXERCICE 3: UN FAUX PARADOXE (FACULTATIF)

1. Montrer qu'en base 10, on a : $1 = 0,9999\dots$
2. Montrer que, pour toute base β , si $b = \beta - 1$, on a $1 = 0,bbb\dots$

EXERCICE 4: (FACULTATIF)

1. Soit ε un nombre dont la taille est de l'ordre de la précision machine. Justifier l'approximation :

$$1 - \frac{1}{\varepsilon} = -\frac{1}{\varepsilon}$$

EXERCICE 5: CONDITIONNEMENT D'UN PROBLÈME

Calculer le nombre de conditions absolues et relatives des problèmes suivants :

1. $f : x \rightarrow \beta x, \quad \beta \in \mathbb{R}$
2. $f : (x_1, x_2) \rightarrow x_1 + x_2$
3. $f : x \rightarrow \sqrt{x}$
4. $f : (x_1, x_2) \rightarrow x_1 - x_2$

EXERCICE 6: ÉVALUATION D'UNE FONCTION POLYNÔME PAR SCHÉMA DE HORNER

Soit une fonction polynôme de degré n définie par $p(x) = \sum_{i=0}^n a_i x^i$. La fonction $p(x)$ peut être écrite sous la forme :

$$p(x) = \sum_{i=0}^n a_i x^i = a_0 + x \left\{ a_1 + x \left[a_2 + x + \left(\dots (a_{n-1} + x a_n) \right) \right] \right\} \quad (\text{Schéma de Horner})$$

Exo 1.

1) 31:

$$31/2 \quad 15 \dots 1$$

$$15/2 \quad 7 \dots 1$$

$$7/2 \quad 3 \dots 1$$

$$3/2 \quad 1 \dots 1$$

$$1/2 \quad 0 \dots 1$$

$$31_d = 00011111_b$$

247:

$$247/2 \quad 123 \dots 1$$

$$123/2 \quad 61 \dots 1$$

$$61/2 \quad 30 \dots 1$$

$$30/2 \quad 15 \dots 0$$

$$15/2 \quad 7 \dots 1$$

$$7/2 \quad 3 \dots 1$$

$$3/2 \quad 1 \dots 1$$

$$1/2 \quad 0 \dots 1$$

$$247_d = 11110111_b$$

$$256 = 100000000$$

$$247 = 256 - 9$$

$$\begin{array}{r} 100000000 \\ - 00001001 \\ \hline 011110111 \end{array}$$

2) 44:

$$44/2 \quad 22 \dots 0$$

$$22/2 \quad 11 \dots 0$$

$$11/2 \quad 5 \dots 1$$

$$5/2 \quad 2 \dots 1$$

$$2/2 \quad 1 \dots 0$$

$$1/2 \quad 0 \dots 1$$

$$00101100$$

-9:

$$9 \Rightarrow 9/2 \quad 4 \dots 1$$

$$4/2 \quad 2 \dots 0$$

$$2/2 \quad 1 \dots 0$$

$$1/2 \quad 0 \dots 1$$

$$9_{10} = 00001001_b \Rightarrow 11110110$$

$$\begin{array}{r} 11110110 \\ + 1 \\ \hline 11110111_b = (-9)_{10} \end{array}$$

$$(0)_{10} = 0000\ 0000_2.$$

$$\begin{array}{rcl}
 -31: & 31 \Rightarrow & 31/2 \quad 15 \dots 1 \quad 00011111 \\
 & & 15/2 \quad 7 \dots 1 \quad \downarrow \\
 & & 7/2 \quad 3 \dots 1 \quad 11100000 \\
 & & 3/2 \quad 1 \dots 1 \quad \downarrow \\
 & & 1/2 \quad 0 \dots 1 \quad 11100001_2 = (-31)_{10}.
 \end{array}$$

$$-256 \sim 255$$

3) short 2 octets \Rightarrow 16 bits.

$$\text{unsigned short} : [0, 2^{16}-1] = [0, 65535]$$

$$\text{short} : [-32768, 32767]$$

int 4 octets \Rightarrow 32 bits

$$\text{unsigned} : [0, 2^{32}-1]$$

$$\text{signed} : [-2^{31}, 2^{31}-1]$$

Exo 2.

$$\begin{array}{c}
 \text{signe} \quad \text{mantisse} \\
 1 \quad 1 \\
 1) \text{ Nb flottant} = \underbrace{(s, m, e)}_{\text{triplet}}^{\text{exposant}} = (-1)^s \cdot m \cdot \beta^e
 \end{array}$$

$$\begin{array}{lcl}
 1337 \Rightarrow & s=0 & \\
 & m=1337 & 0,1337 \cdot 10^4 \\
 & e=4, \beta_{10} &
 \end{array}$$

pour garantir l'unicité,

$$0, \boxed{1}337 \cdot 10^4$$

ce nombre doit être $\neq 0$.

$$\begin{array}{lcl}
 \leftarrow & s=0 & \\
 & m=0,1337 & \\
 & e=5, \beta_{10} &
 \end{array}$$

Normalisé: \Rightarrow 1^{er} chiffre de la mantisse est \neq de 0.
 \Rightarrow Garantir l'unicité de la représentation.

$$2) a = \left\{ \overset{3}{1} \overset{m}{101101} \overset{e}{3} \right\} \rightarrow \text{convertir en d cimal}$$

$$= (-1)^1 0,101101 \cdot 2^3$$

$$\downarrow 1 \cdot 2^{-1} + 0 \cdot 2^{-2} + 1 \cdot 2^{-3} \dots$$

$$101101 = (2^{-1} + 2^{-3} + 2^{-4} + 2^{-6}) \cdot 2^3$$

$$= 2^2 + 2^0 + 2^{-1} + 2^{-3} = 4 + 1 + 0,5 + 0,125$$

$$a = -5,625$$

$$b = \left\{ 0 \ 11000001 \ -5 \right\}$$

$$b = (2^{-1} + 2^{-2} + 2^{-8}) \cdot 2^{-5}$$

$$0,11000001 \ 2^{-5}$$

$$= 2^{-6} + 2^{-7} + 2^{-13}$$

$$= +0,02355957$$

$$c = \left\{ 0 \ 100000001 \ 8 \right\}$$

$$c = (2^{-1} + 2^{-9}) \cdot 2^8 = 2^7 + 2^{-1} = 128,5$$

3) Soit $x \in \mathbb{R}$

Sa repr sentation en flottant normalis 

$$S \ m \ e$$

$$\text{avec } m = m_{-1} m_{-2} m_{-3} m_{-4} \dots m_{-r} m_{-r-1}$$

$A(x)$ = arrondi   r chiffres.

$$\hookrightarrow A(x) = (S \ m' \ e)$$

$$m' = m_{-1} m_{-2} \dots m_{-r+2} m_{-r+1} m_{-r}' \text{ et } m_{-r}' =$$

$$\begin{cases} m_{-r} & \text{si } m_{-r-1} < \frac{\beta}{2} \\ m_{-r} + 1 & \text{si } m_{-r-1} \geq \frac{\beta}{2} \end{cases}$$

$$\begin{array}{r} 0,1337 \\ 0,13 \quad | \quad 7 \\ \hline 0,134 \end{array}$$



$$\beta = 3 \quad 0,122112$$

$$A(x)_6 = 0,1221$$

$$A(x)_5 = 0,12212$$

$$1 < \frac{\beta}{2} = \frac{3}{2}$$

$$2 > \frac{\beta}{2}$$

Erreur de représentation.

$$|A(x) - x| = \text{erreur absolue.}$$

$$\frac{|A(x) - x|}{|x|} = \text{erreur relative}$$

Précision machine $\frac{\beta}{2} \cdot \beta^{-r}$ lié à m'_r

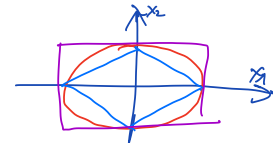
$$4) \beta = 10 \quad x = 231,345$$

$$x = (-1)^0 \cdot 0,231345 \cdot 10^3$$

$$A(x)_4 = (-1)^0 \cdot 0,2313 \cdot 10^3$$

$$E_a = |A(x)_4 - x| = 0,045$$

$$E_r = \frac{|A(x)_4 - x|}{|x|} = \frac{0,045}{231,345} = 2 \cdot 10^{-4}$$



Exo 5. Conditionnement.

$$1) f: x \rightarrow \beta x, \beta \in \mathbb{R}$$

$$\text{Soit } x' = x + \Delta x$$

$$y' = f(x') = \beta(x + \Delta x) = \beta x + \beta \Delta x$$

$$y' = y + \Delta y$$

Comme $x \in \mathbb{R}$ (1 dimension) $\Rightarrow L_1 = L_\infty$

$$\|x\|_1 = \|x\|_\infty = |x|$$

$$\|\Delta x\|_1 = \|\Delta x\|_\infty = |\Delta x|$$

$$\|y\| = \|\beta x\| = |\beta| |x|$$

$$\text{norme } L_1: \|x\|_1 = \sum_{i=1}^n |x_i| \quad /$$

$$L_2: \|x\|_2 = \sqrt{\sum_{i=1}^n x_i^2} \quad /$$

$$L_\infty: \|x\|_\infty = \max_i (|x_i|) \quad /$$

erreur de troncature

x + une petite erreur. capteur bruité

...

$$K_a = \frac{\|\Delta y\|}{\|\Delta x\|} = |\beta|$$

$$K_r = \frac{\frac{\|\Delta y\|}{\|y\|}}{\frac{\|\Delta x\|}{\|x\|}} = \frac{|\beta \Delta x| |x|}{|\beta x| |\Delta x|} = \frac{1}{1}$$

$$\|\Delta y\| = \|\beta \Delta x\| = |\beta| \|\Delta x\|$$

erreur restée indolore
en terme de %.

$$2) y: (x_1, x_2) \mapsto x_1 + x_2$$

$$\Delta x = (\Delta x_1, \Delta x_2)$$

$$y' = x_1 + \Delta x_1 + x_2 + \Delta x_2 = x_1 + x_2 + \Delta x_1 + \Delta x_2$$

$$y' = y + \Delta y$$

$$\|x\|_1 = |x_1| + |x_2|$$

$$\|x\|_\infty = \max(|x_1|, |x_2|)$$

$$\|\Delta x\|_1 = |\Delta x_1| + |\Delta x_2|$$

$$\|\Delta x\|_\infty = \max(|\Delta x_1|, |\Delta x_2|)$$

$$\|y\|_1 = \|y\|_\infty = |x_1 + x_2|$$

$$\|\Delta y\|_1 = \|\Delta y\|_\infty = |\Delta x_1 + \Delta x_2|$$

en norme 1:

$$k_a = \frac{\|\Delta y\|_1}{\|\Delta x\|_1} = \frac{|\Delta x_1 + \Delta x_2|}{|\Delta x_1| + |\Delta x_2|} \leq \frac{|\Delta x_1| + |\Delta x_2|}{|\Delta x_1| + |\Delta x_2|} \Rightarrow k_a \leq 1$$

erreur diminuée en absolue

$$k_r = k_a \frac{\|x\|_1}{\|y\|_1} \leq \frac{|x_1| + |x_2|}{|x_1 + x_2|}$$

en norme ∞ :

$$k_a = \frac{\|\Delta y\|_\infty}{\|\Delta x\|_\infty} = \frac{|\Delta x_1 + \Delta x_2|}{\max(|\Delta x_1|, |\Delta x_2|)} \leq \frac{2 \max(|\Delta x_1|, |\Delta x_2|)}{\max(|\Delta x_1|, |\Delta x_2|)} \Rightarrow k_a \leq 2$$

$|\Delta x_1| \leq \max(|\Delta x_1|, |\Delta x_2|)$
 $|\Delta x_2| \leq \max(|\Delta x_1|, |\Delta x_2|)$

$$k_r = k_a \frac{\|x\|_\infty}{\|y\|_\infty} \leq \frac{2 \max(|x_1|, |x_2|)}{|x_1 + x_2|}$$

$$x_1 = x_2 = 10^{-5} \quad y' = x_1 + \Delta x_1 - (x_2 + \Delta x_2)$$

$$\Delta x_1 = ? \quad y' = \Delta x_1 - \Delta x_2 \quad y = x_1 - x_2 = 0 = 10^{-5}$$

$$\Delta x_2 = ?$$

$$e_a = |y - y'| = |\Delta x_1 - \Delta x_2|$$

$$e_r = \frac{|y - y'|}{|y|} = \frac{|\Delta x_1 - \Delta x_2|}{10^{-5}}$$

phénomène d'accumulation

→ erreur énorme

$$3) y = \sqrt{x}$$

$$x' = x + \Delta x$$

$$y' = \sqrt{x + \Delta x} \text{ petit devant } x \Rightarrow \text{DL ordre 1.}$$

$$y' = \sqrt{x} + \frac{d}{dx}(\sqrt{x}) \Delta x = \underbrace{\sqrt{x}}_y + \underbrace{\frac{1}{2\sqrt{x}} \Delta x}_{\Delta y}.$$

$x, \Delta x, y, \Delta y$ en dimension 1 \Rightarrow toutes les normes équivalentes

$$\|x\| = |x| \quad \|y\| = |\sqrt{x}| = \sqrt{x}$$

$$\|\Delta x\| = |\Delta x| \quad \|\Delta y\| = \left| \frac{\Delta x}{2\sqrt{x}} \right| = \frac{1}{2\sqrt{x}} |\Delta x|$$

$$k_a = \frac{\|\Delta y\|}{\|\Delta x\|} = \frac{\frac{1}{2\sqrt{x}} |\Delta x|}{|\Delta x|} = \frac{1}{2\sqrt{x}}$$

$$k_r = k_a \frac{\|x\|}{\|y\|} = \frac{1}{2\sqrt{x}} \cdot \frac{|x|}{\sqrt{x}} = \frac{|x|}{2|x|} = \frac{1}{2} \Rightarrow \text{diviser l'erreur par 2.}$$

$$\|x\|_p^n = \sqrt[p]{\sum_{i=1}^n x_i^p}$$

$$f(x) = f(a) + \frac{f'(a)}{1!} (x-a) + \frac{f''(a)}{2!} (x-a)^2 + \dots$$

1. Écrire l'algorithme qui implémente le schéma de Horner
2. Démontrer que l'algorithme d'évaluation de Horner est linéaire par rapport au degré du polynôme évalué
3. Généraliser au cas :

$$p(x) = \sum_{i=0}^n a_i \prod_{j=1}^i (x - c_{j-1}) = a_0 + a_1(x - c_0) + a_2(x - c_0)(x - c_1) + \dots + a_n(x - c_0)(x - c_1) \dots (x - c_{n-1})$$

EXERCICE 7: ÉVALUATION DES FONCTIONS COMPLEXES

Proposer une méthode pour calculer les fonctions suivantes en utilisant les fonctions élémentaires (+, -, *, /)

1. \sqrt{x} , $x > 0$
2. $\sin(x)$, $\cos(x)$
3. e^x
4. $\log(x)$, $x > 0$

EXERCICE 8: CALCUL DES FONCTIONS ET ERREURS DE REPRÉSENTATION

Proposer une méthode pour éviter une perte de précision dans les calculs suivants :

1. $e^x - \sin(x) - \cos(x)$
2. $\log(x) - 1$
3. $\log(x) + \log(1/x)$