

# **Programmation et Méthodes Numériques**

Représentation d'un nombre en machine, erreurs numériques, etc...

M. Casaletti ([massimiliano.casaletti@upmc.fr](mailto:massimiliano.casaletti@upmc.fr))

Laboratoire d'Electronique et Electromagnétisme (L2E)  
Université Pierre et Marie Curie

- **Déroulement** sur 9 semaines:
  - 8 Cours magistraux (4 M.N.+ 4 P.)
  - 7 TD (8h M.N. + 6h P.)
  - 6 TP (6h M.N. + 6h P.)
  - Mini projet (4h TD + 10h TP + soutenance)
  
- **Évaluation:**
  - 2 Examens (60%)
  - Mini projet (40%)
    - Simulation des circuits RLC
    - Rayonnement des antennes
    - Propagation de la chaleur
    - ...

## ● Points abordés:

- Représentation d'un nombre en machine, erreurs numériques, conditionnement d'un problème, etc...
- Interpolation polynomiale
- Intégration numérique
- Résolution numérique des équations différentielles ordinaires (EDO)

## ● Ouvrages:

- **Introduction à l'analyse numérique**, J. Bastien et J.N. Martin, DUNOD
- **Numerical Recipes : The Art of Scientific Computing**, Cambridge University Press

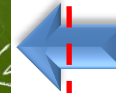
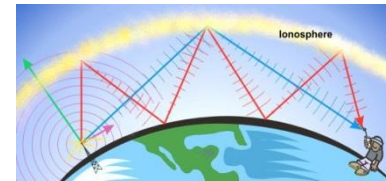
**Mathématique :**

- Dérivation
- Intégration
- Développement de Taylor
- Espaces vectoriels (base, produit scalaire, distance,...)
- Équations différentielles ordinaires

**Physique (mini projet) :**

- Mécanique (cinématique, lois de Newton,...)
- Théorie des circuits électroniques
- Champ électrique et magnétique
- .....

- **Représentation d'un nombre en machine, erreurs numériques, etc...**
  - **Introduction générale**
  - Représentation des nombres entiers et réels
  - Opérations élémentaires en virgule flottante
  - Conditionnement d'un problème
  - Opérations complexes
  - Conclusion: différentes sources d'erreur
- Interpolation polynomiale
- Intégration numérique
- Résolution numérique des équations différentielles ordinaires (EDO)

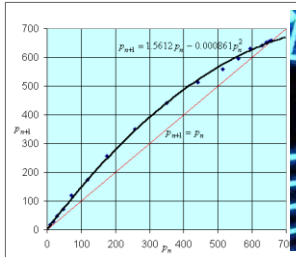
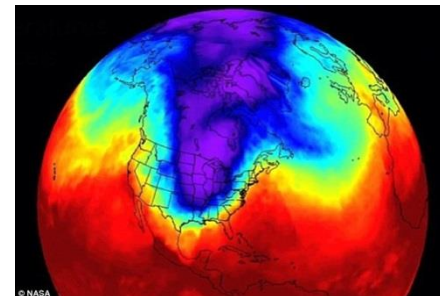


$$\frac{\ln x}{x}, f(e) = \frac{1}{e}$$

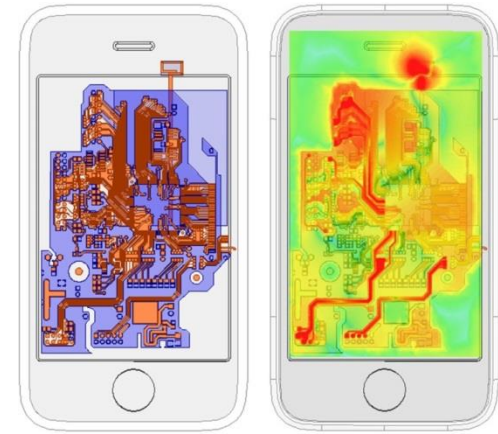
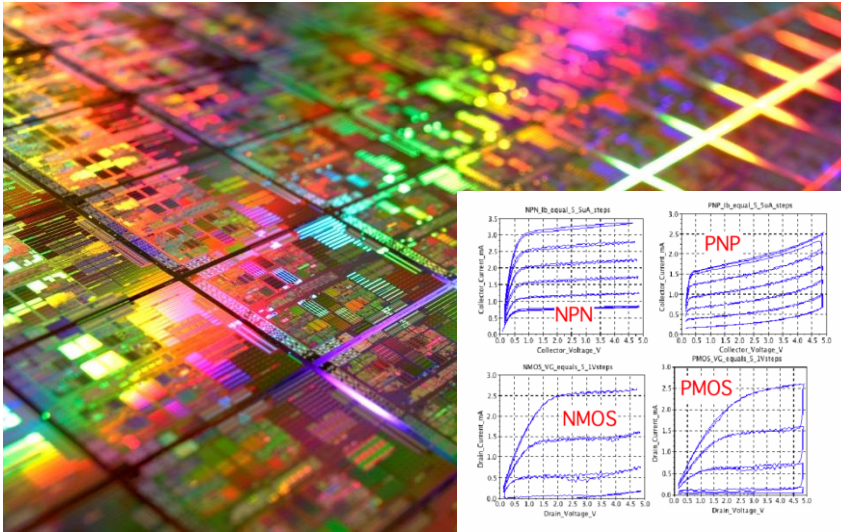
$$x = \frac{\ln x}{x \cdot x - \ln x} = \frac{1}{x^2} = 1 - \frac{1}{x}$$

$$x = 0 \Rightarrow \ln x = 1 \Rightarrow$$

$$f(x) = f(e)$$



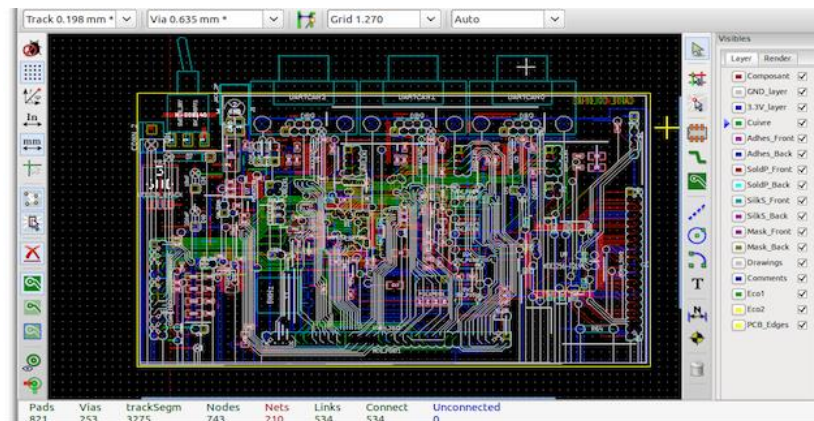
## Circuits électroniques



Dissipation de la chaleur

$$\frac{\partial u(\vec{r})}{\partial t} - \alpha \nabla^2 u(\vec{r}) = f(x, t)$$

Analyse des circuits complexes et non linéaires



Lois de Kirchhoff

$$\sum_{i \in \sigma} i_k(t) = 0$$

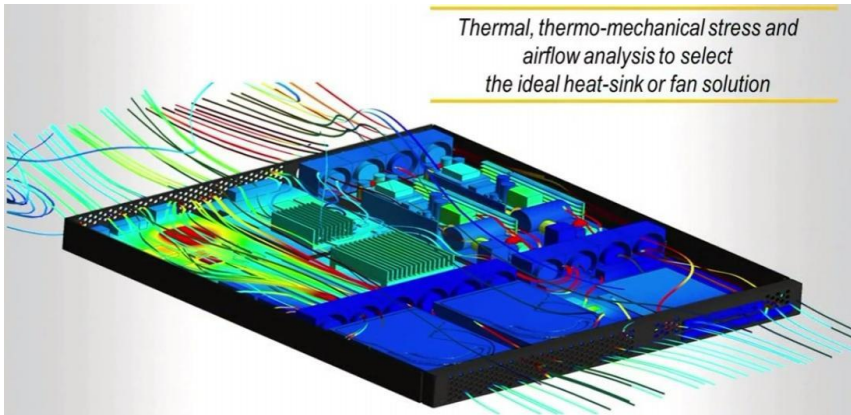
Optimisation non linéaire

$$f(x) = \sum_{i=1}^N a_i \varphi(\|x - x_i\|, c_i) + R(x)$$

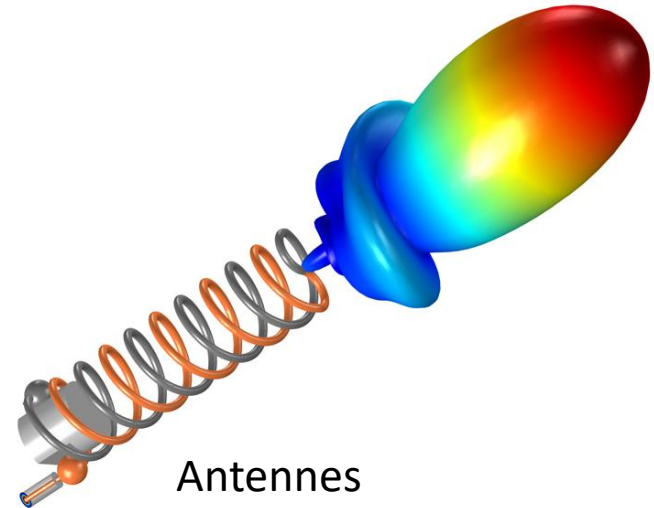
Optimisation des connexions circuits



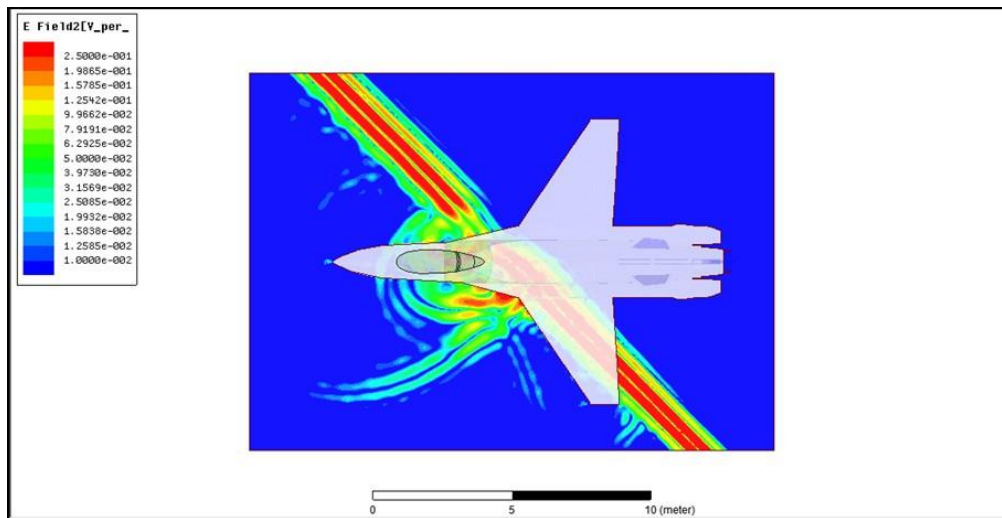
## Electronique haute fréquence et propagation électromagnétique



Circuits haute fréquence



Antennes



radar

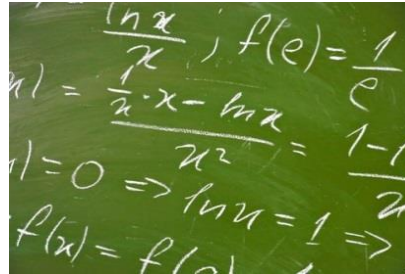
$$\begin{cases} \vec{\nabla} \times \vec{E}(\vec{r}) = -\mu_0 \mu_r(\vec{r}) \frac{\partial \vec{H}(\vec{r})}{\partial t} \\ \vec{\nabla} \times \vec{H}(\vec{r}) = \varepsilon_0 \varepsilon_r(\vec{r}) \frac{\partial \vec{E}(\vec{r})}{\partial t} + \vec{J}(\vec{r}) \\ \vec{\nabla} \cdot \vec{E}(\vec{r}) = -\frac{\rho(\vec{r})}{\varepsilon_0 \varepsilon_r(\vec{r})} \\ \vec{\nabla} \cdot \vec{H}(\vec{r}) = 0 \end{cases}$$

Equations de Maxwell





Bureautique, audio-vidéo,...



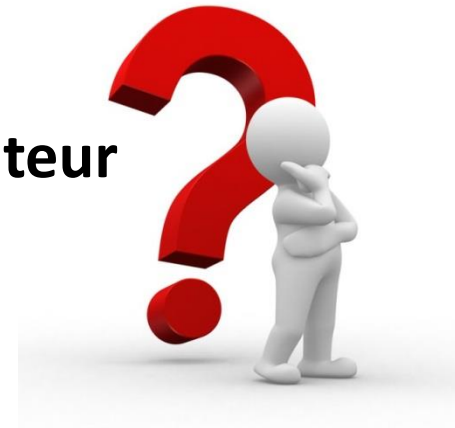
Calcul scientifique,  
Simulation, modélisation,...



Systèmes embarqués

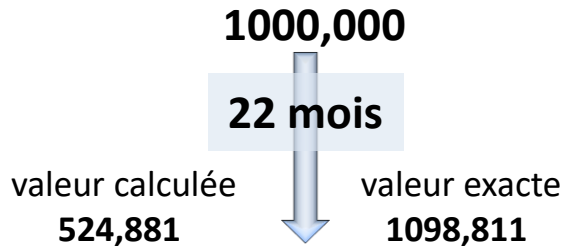
les ordinateurs effectuent des **calculs**

Puis-je faire confiance aux résultats de mon ordinateur



## Bourse de Vancouver (1982)

Création d'un nouvel indice de valeur initiale 1000  
recalculé après chaque transaction et **tronqué après le 3<sup>e</sup> chiffre**



### Source du problème :

les erreurs de troncature ont le même signe

$$\begin{aligned} 10,4562 &\rightarrow 10,456 + \Delta \\ 10,4568 &\rightarrow 10,456 + \Delta \end{aligned} \quad \Delta > 0$$



Source: Wikipedia

## Ariane V (1996)

**Source du problème :** représentation de l'accélération horizontale.

- L'accélération horizontale maximum d'Ariane 4 était d'environ 64.

↓  
la variable a été **codée sur 8 bits**

- Ariane 5 était plus rapide: son accélération pouvait atteindre la valeur 300 (qui vaut 100101100 en binaire et nécessite 9 bits).

↓  
**Overflow: 100000000 -> 00000000**

Le logiciel de contrôle face à des valeurs vraiment pas normales décida de l'autodestruction de la fusée.



04/06/1996 à Kourou  
Le lanceur fut détruit après  
37 secondes de vol.  
European Space Agency (ESA)  
(coût 500 million de dollars)

## ● Dysfonctionnement dans l'unité de calcul en virgule flottante du *Pentium 5* (1994)

valeur correcte

$$4\,195\,835,0 / 3\,145\,727,0 = 1,333\,820\,449\,136\,241\,002$$

valeur erronée retournée par le processeur:

$$4\,195\,835,0 / 3\,145\,727,0 = 1,333\,739\,068\,902\,037\,589$$



Source: Wikipedia

## ● Bug dans Excel 2007

**Source du problème :** Conversion des nombres binaires en caractères pour la visualisation

“That's because **0.1** has *no exact representation in binary*... it's a repeating binary number. It's sort of like how  $1/3$  has no representation in decimal.  $1/3$  is  $0.33333333$  and you have to keep writing 3's forever”

	A	B	C	D	E	F
1						
2		850.00	x	77.10	=	100,000.00
3		1700.00	x	38.55	=	100,000.00
4		3400.00	x	19.28	=	100,000.00
5		6800.00	x	9.64	=	100,000.00
6		13600.00	x	4.82	=	100,000.00
7		27200.00	x	2.41	=	100,000.00
8		425.00	x	154.20	=	100,000.00
9		212.50	x	308.40	=	100,000.00
10		106.25	x	616.80	=	100,000.00
11		53.13	x	1233.60	=	100,000.00
12		26.56	x	2467.20	=	100,000.00
13		13.28	x	4934.40	=	100,000.00
14						

Source: [www.joelonsoftware.com](http://www.joelonsoftware.com)

- **Représentation d'un nombre en machine, erreurs numériques, etc...**
  - Introduction générale
  - **Représentation des nombres entiers et réels**
  - Opérations élémentaires en virgule flottante
  - Conditionnement d'un problème
  - Opérations complexes
  - Conclusion: différentes sources d'erreur
- Interpolation polynomiale
- Intégration numérique
- Résolution numérique des équations différentielles ordinaires (EDO)

- On appelle **base** un entier  $\beta$  supérieur ou égal à 2
- Un **chiffre** sera un entier (symbole) compris entre 0 et  $\beta - 1$

Exemples:  $\beta=2$       0,1 (bit)  
               $\beta=10$      0,1,2,3,4,5,6,7,8,9  
               $\beta=16$      0,1,2,3,4,5,6,7,8,9,A,B,C,D,E,F

- Un **nombre x de n chiffres** = une séquence  $x_{n-1}$   $x_{n-2}$   $x_{n-3}$   $x_1$   $x_0$  telle que

$$x = \sum_{i=0}^{n-1} x_i \beta^i$$

- Proposition:** Pour tout entier  $y \geq 1$ , il existe un entier unique  $n$  et des entiers  $0 \leq x_i \leq \beta - 1$  avec  $x_{n-1} \neq 0$  tels que

$$y = \sum_{i=0}^{n-1} x_i \beta^i$$

# Exemples

<b>1</b>	<b>1</b>	<b>0</b>	<b>1</b>	<b>0</b>
$x_4$	$x_3$	$x_2$	$x_1$	$x_0$



$$x = \sum_{i=0}^{n-1} x_i \beta^i$$

$$\beta = 2, \quad x = 0 \cdot 2^0 + 1 \cdot 2 + 0 \cdot 2^2 + 1 \cdot 2^3 + 1 \cdot 2^4 = 2 + 8 + 16 = 26$$

$$\beta = 10, \quad x = 0 \cdot 10^0 + 1 \cdot 10 + 0 \cdot 10^2 + 1 \cdot 10^3 + 1 \cdot 10^4 = 11101$$

$$\beta = 16, \quad x = 0 \cdot 16^0 + 1 \cdot 16 + 0 \cdot 16^2 + 1 \cdot 16^3 + 1 \cdot 16^4 = 69648$$

$$x = 101 \quad \Rightarrow \quad \{x_i\} = ?$$

$$\begin{array}{ll} 101 : 2 = 50 & r = 1 = x_0 \\ 50 : 2 = 25 & r = 0 = x_1 \\ 25 : 2 = 12 & r = 1 = x_2 \\ \beta = 2 \quad 12 : 2 = 6 & r = 0 = x_3 \\ 6 : 2 = 3 & r = 0 = x_4 \\ 3 : 2 = 1 & r = 1 = x_5 \\ 1 : 2 = 0 & r = 1 = x_6 \end{array}$$

<b>1</b>	<b>1</b>	<b>0</b>	<b>0</b>	<b>1</b>	<b>0</b>	<b>1</b>
$x_6$	$x_5$	$x_4$	$x_3$	$x_2$	$x_1$	$x_0$

$$\beta = 16$$

$$\begin{array}{ll} 101 : 16 = 6 & r = 5 = x_0 \\ 6 : 16 = 0 & r = 6 = x_1 \end{array}$$

<b>6</b>	<b>5</b>
$x_1$	$x_0$

## Divisions successives module $\beta$

$$x = \sum_{i=0}^{n-1} x_i \beta^i = x_n \beta^n + \dots + x_1 \beta^1 + x_0$$

$$\Rightarrow \frac{x}{\beta} = \sum_{i=0}^{n-1} x_i \beta^{i-1} = x_n \beta^{n-1} + \dots + x_1 + \frac{x_0}{\beta}$$

$x_0$  reste de la division

$$x = -26 \quad \Rightarrow \quad \{x_i\} = ?$$

$\beta = 2$

<b>s</b>	<b>1</b>	<b>1</b>	<b>0</b>	<b>1</b>	<b>0</b>
----------	----------	----------	----------	----------	----------

$+\rightarrow s = 0$   
 $-\rightarrow s = 1$

$|x| = 26$

$n = 6$

<b>1</b>	<b>1</b>	<b>1</b>	<b>0</b>	<b>1</b>	<b>1</b>
----------	----------	----------	----------	----------	----------

$x = -26$

$x=0$  est représenté par 2 séquences: 

0	0	0	0	0	0
---	---	---	---	---	---

 et 

1	0	0	0	0	0
---	---	---	---	---	---

$$-2^{n-1} + 1 \geq x \geq 2^{n-1} - 1$$

**Complément à deux**  $n = 6$

*si*

$x > 0$	$x \rightarrow \{x_i\}$	<table border="1" style="display: inline-table; text-align: center;"> <tr> <td>0</td><td>0</td><td>0</td><td>1</td><td>0</td><td>1</td> </tr> </table>	0	0	0	1	0	1	$\rightarrow x = 5$
0	0	0	1	0	1				
$x < 0$	$2^{n-1} - x > 0 \rightarrow \{x_i\}$	<table border="1" style="display: inline-table; text-align: center;"> <tr> <td>1</td><td>0</td><td>0</td><td>1</td><td>0</td><td>1</td> </tr> </table>	1	0	0	1	0	1	$\rightarrow -(2^5 - 5) = -27$
1	0	0	1	0	1				

$$-2^{n-1} \geq x \geq 2^{n-1} - 1$$



$$x = -131 \rightarrow \{x_i\} = ?$$



$$|x| = 131$$

$$\beta = 2$$

$$131 : 2 = 65 \quad r = 1 = x_0$$

$$65 : 2 = 32 \quad r = 1 = x_1$$

$$32 : 2 = 16 \quad r = 0 = x_2$$

$$16 : 2 = 8 \quad r = 0 = x_3$$

$$8 : 2 = 4 \quad r = 0 = x_4$$

$$4 : 2 = 2 \quad r = 0 = x_5$$

$$2 : 2 = 1 \quad r = 0 = x_6$$

$$1 : 2 = 0 \quad r = 1 = x_7$$

$$|x| = 131$$

1	0	0	0	0	0	1	1
$x_7$	$x_6$	$x_5$	$x_4$	$x_3$	$x_2$	$x_1$	$x_0$



$$x = -131$$

$$= \begin{array}{|c|c|c|c|c|c|c|c|c|} \hline 0 & 1 & 0 & 0 & 0 & 0 & 0 & 1 & 1 \\ \hline \end{array} + \begin{array}{|c|c|c|c|c|c|c|c|c|} \hline 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 1 \\ \hline \end{array}$$

$$= \begin{array}{|c|c|c|c|c|c|c|c|c|} \hline 1 & 0 & 1 & 1 & 1 & 1 & 1 & 0 & 0 \\ \hline \end{array} + \begin{array}{|c|c|c|c|c|c|c|c|c|} \hline 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 1 \\ \hline \end{array}$$

$$= \begin{array}{|c|c|c|c|c|c|c|c|c|} \hline 1 & 0 & 1 & 1 & 1 & 1 & 1 & 0 & 1 \\ \hline \end{array}$$

$x_8 \quad x_7 \quad x_6 \quad x_5 \quad x_4 \quad x_3 \quad x_2 \quad x_1 \quad x_0$

$$x = s \left( \underbrace{\sum_{i=0}^{n-1} x_i \beta^i}_{x_e \text{ Partie entière}} + \underbrace{\sum_{j=-1}^{-q} x_j \beta^j}_{x_f \text{ Partie fractionnelle}} \right)$$

$\downarrow$   
 Signe

**Remarque:** possibilité que  $q = \infty$

Exemple:  $\beta=10$   $1/3=0,3333333...$

Un ordinateur possède  
une **mémoire limitée**



Nombre limité de chiffres non nuls!!!



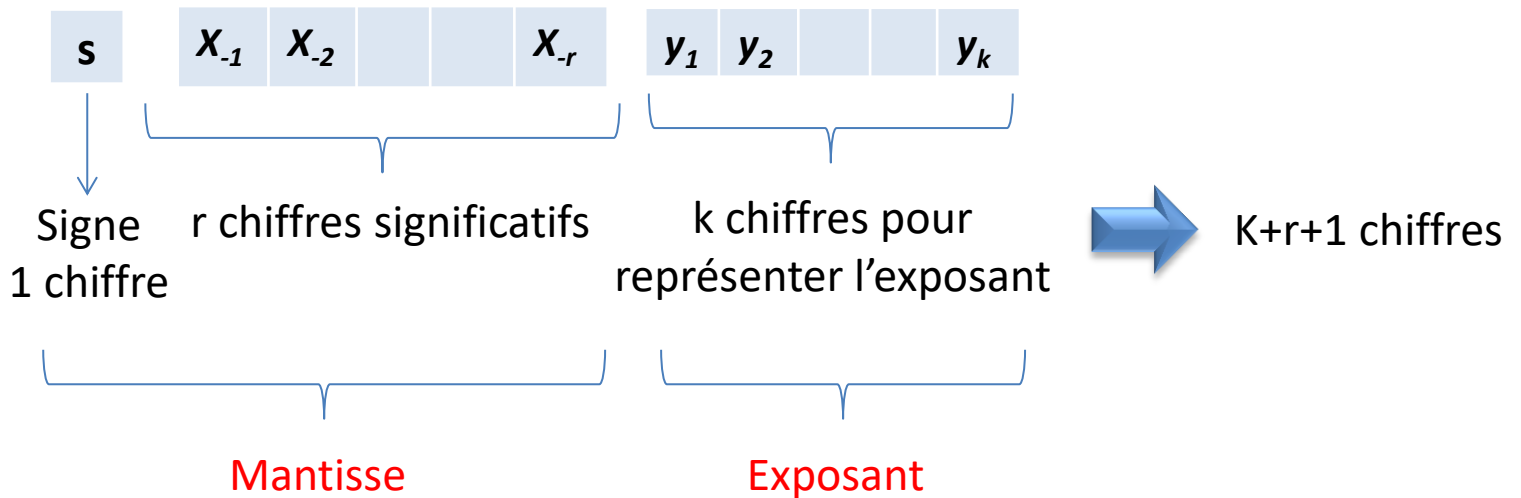
Solution

## Ecriture à virgule flottante normalisée

On fixe un nombre  $r$  de **chiffre significatifs**

$$x = s \left( \sum_{i=-1}^{-r} x_i \beta^{-i} \right) \cdot \beta^k \rightarrow s \ 0, x_{-1} x_{-2} \dots x_{-r} \cdot \beta^k$$

On place la virgule juste avant le premier chiffre non nul du développement  
(virgule flottante)



Développement illimité (exact) de  $x$

$$x = 0, x_{-1}x_{-2}\dots x_{-r}x_{-r-1}\dots\beta^j$$



$$m' = \begin{cases} 0, x_{-1}x_{-2}\dots x_{-r} & \text{si } x_{-r-1} < \beta / 2 \\ 0, x_{-1}x_{-2}\dots x_{-r} + \beta^{-r} & \text{si } x_{-r-1} \geq \beta / 2 \end{cases}$$

Hypothèse:  $j$  représentable avec  $k$  chiffres

On définit **l'arrondi de  $x$** , noté  $fl(x) = m' \beta^j$

**Les  $r-1$  premiers chiffres sont conservés et peuvent être considérés comme exacts !!!**

$$e_r \triangleq \frac{|fl(x) - x|}{|x|} = \frac{|m' \beta^j - 0, x_{-1} x_{-2} \dots x_{-r} x_{-r-1} \dots \beta^j|}{|0, x_{-1} x_{-2} \dots x_{-r} x_{-r-1} \dots \beta^j|} = \frac{|0, 0 \dots 0 (a_{-r} - x_{-r}) (-x_{-r-1}) \dots|}{|0, x_{-1} x_{-2} \dots x_{-r} x_{-r-1} \dots|}$$

Par la définition d'arrondi on a  $(a_{-r} - x_{-r})(-x_{-r-1}) \dots \leq \frac{\beta}{2} \beta^{-r-1}$

$$\frac{|fl(x) - x|}{|x|} \leq \frac{\frac{\beta}{2} \beta^{-r-1}}{|0, x_{-1} x_{-2} \dots x_{-r} x_{-r-1} \dots|} \leq \frac{\frac{\beta}{2} \beta^{-r-1}}{\beta^{-1}} = \frac{\beta}{2} \beta^{-r}$$

$$x_{-1} \neq 0 \Rightarrow |0, x_{-1} x_{-2} \dots x_{-r} x_{-r-1} \dots| \geq \beta^{-1}$$

●  $\varepsilon = \frac{\beta}{2} \beta^{-r}$  est appelée **précision machine**.



$$fl(x) = x(1 + \alpha) \quad |\alpha| \leq \varepsilon$$

**Exemple 1:**  $x = -131.23411$   $\beta=10$   $r=6$  (chiffres mantisse)

$$x = -131.23411 \quad x = -0.13123411 \cdot 10^3 \quad fl(x) = -0.131234 \cdot 10^3 \quad \varepsilon = \frac{\beta}{2} \beta^{-r} = 5 \cdot 10^{-6}$$

$$e_r = \frac{|-0.13123411 \cdot 10^3 + 0.131234 \cdot 10^3|}{|-0.13123411 \cdot 10^3|} = 8.382 \cdot 10^{-7} < \varepsilon$$

**Exemple 2:**  $x = 12.23$   $\beta=2$   $r=6$  (chiffres mantisse)  $\varepsilon = \frac{\beta}{2} \beta^{-r} = 2^{-6} = 1.56 \cdot 10^{-2}$

Partie entière ( $x_e$ )

$$12 : 2 = 6 \quad r = 0 = x_0$$

$$6 : 2 = 3 \quad r = 0 = x_1$$

$$3 : 2 = 1 \quad r = 1 = x_2$$

$$1 : 2 = 0 \quad r = 1 = x_3$$

Partie fractionnelle ( $x_f$ )

Multiplications successives

$$x_f = \sum_{i=-1}^{-q} x_i \beta^i = x_{-1} \beta^{-1} + \dots + x_{-q+1} \beta^{-q+1} + x_{-q} \beta^{-q}$$

$$\Rightarrow \beta x_f = x_{-1} + x_{-2} \beta^{-1} + \dots + x_{-q} \beta^{-q+1}$$

$x_{-1}$  partie entière  $\beta x_f$

$$0,23 \cdot 2 = 0,46 \Rightarrow x_{-1} = 0$$

$$0,46 \cdot 2 = 0,92 \Rightarrow x_{-2} = 0$$

$$0,92 \cdot 2 = 1,84 \Rightarrow x_{-3} = 1$$

$$(1,84 - 1) \cdot 2 = 1,68 \Rightarrow x_{-4} = 1$$

$$0,68 \cdot 2 = 1,36 \Rightarrow x_{-5} = 1$$

$$0,36 \cdot 2 = 0,72 \Rightarrow x_{-6} = 0$$

$$0,72 \cdot 2 = 1,44 \Rightarrow x_{-7} = 1$$

$$0,44 \cdot 2 = 0,88 \Rightarrow x_{-8} = 0$$



$$x = 1100,001110101\dots = 0,1100001110101\dots \cdot 2^4$$

$x_3 x_2 x_1 x_0$   $x_{-1} x_{-2} x_{-3} x_{-4} x_{-5} x_{-6} x_{-7} x_{-8}$

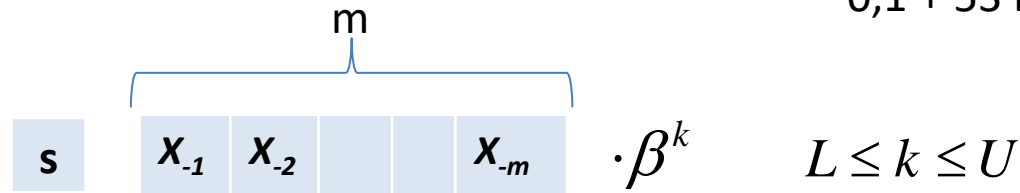


$$fl(x) = 0,11000 \cdot 2^4$$



Elle est la norme la plus employée actuellement pour le calcul des nombres à virgule flottante dans le domaine informatique (PC, système embarqué, tablette, smartphone, ...).

- 1) simple precision (32 bits) : 1 bit de signe, 8 bits d'exposant, 23 bits de mantisse  
0,1 + 23 bits  $\rightarrow$  m=24
- 2) double precision (64 bits) : 1 bit de signe, 11 bits d'exposant, 52 bits de mantisse  
0,1 + 53 bits  $\rightarrow$  m=53

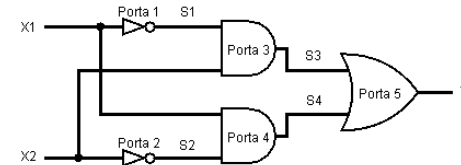
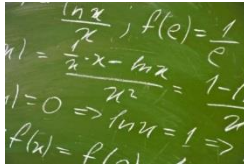


- 1) en base 2,  $L = -126$ ,  $U = 127$ ,  $\text{eps} = 2^{-24}$ ;  
en base 10,  $L \approx -38$ ,  $U \approx 38$ ,  $\text{eps} \approx 10^{-8}$
- 2) en base 2,  $L = -1023$ ,  $U = 1024$ ,  $\text{eps} \approx 2^{-53}$ ;  
en base 10,  $L \approx -308$ ,  $U \approx 308$ ,  $\text{eps} \approx 10^{-16}$



- **Représentation d'un nombre en machine, erreurs numériques, etc...**
  - Introduction générale
  - Représentation des nombres entiers et réels
  - **Opérations élémentaires en virgule flottante**
  - Conditionnement d'un problème
  - Opérations complexes
  - Conclusion: différentes sources d'erreur
- Interpolation polynomiale
- Intégration numérique
- Résolution numérique des équations différentielles ordinaires (EDO)





## Opérations théoriques

- ❑ Addition  $x + y$
- ❑ Soustraction  $x - y$
- ❑ Multiplication  $x \cdot y$
- ❑ Division  $x / y$



## Opérations en virgule flottante

$$x \oplus y = fl(fl(x) + fl(y))$$

$$x \ominus y = fl(fl(x) - fl(y))$$

$$x \odot y = fl(fl(x) \cdot fl(y))$$

$$x \oslash y = fl(fl(x) / fl(y))$$

## Exemple: propagation des erreurs

$$fl(x) = x(1 + \alpha), \quad |\alpha| < \varepsilon$$

$$fl(y) = y(1 + \beta), \quad |\beta| < \varepsilon$$



$$x \oplus y = fl(x(1 + \alpha) + y(1 + \beta))$$

$$= (x(1 + \alpha) + y(1 + \beta))(1 + \delta), \quad |\delta| < \varepsilon$$

$$= (x + x\alpha + y + y\beta)(1 + \delta)$$

$$= (x + y)(1 + \gamma)$$

$$\gamma = x\alpha + y\beta + (x + x\alpha + y + y\beta)\delta \Rightarrow |\gamma| \leq |x||\alpha| + |y||\beta| + |x||\delta| + |x||\alpha||\delta| + |y||\delta| + |y||\beta||\delta|$$

$$|\gamma| \leq |x|\varepsilon + |y|\varepsilon + |x|\varepsilon + |x|\varepsilon^2 + |y|\varepsilon + |y|\varepsilon^2 \simeq (2|x| + 2|y|)\varepsilon$$

## ● propriété de commutativité



$$a \oplus b = b \oplus a \quad a \odot b = b \odot a$$

## ● propriétés d'associativité et de distributivité



$$a \oplus (b \oplus c) \neq (a \oplus b) \oplus c$$

$$a \odot (b \odot c) \neq (a \odot b) \odot c$$

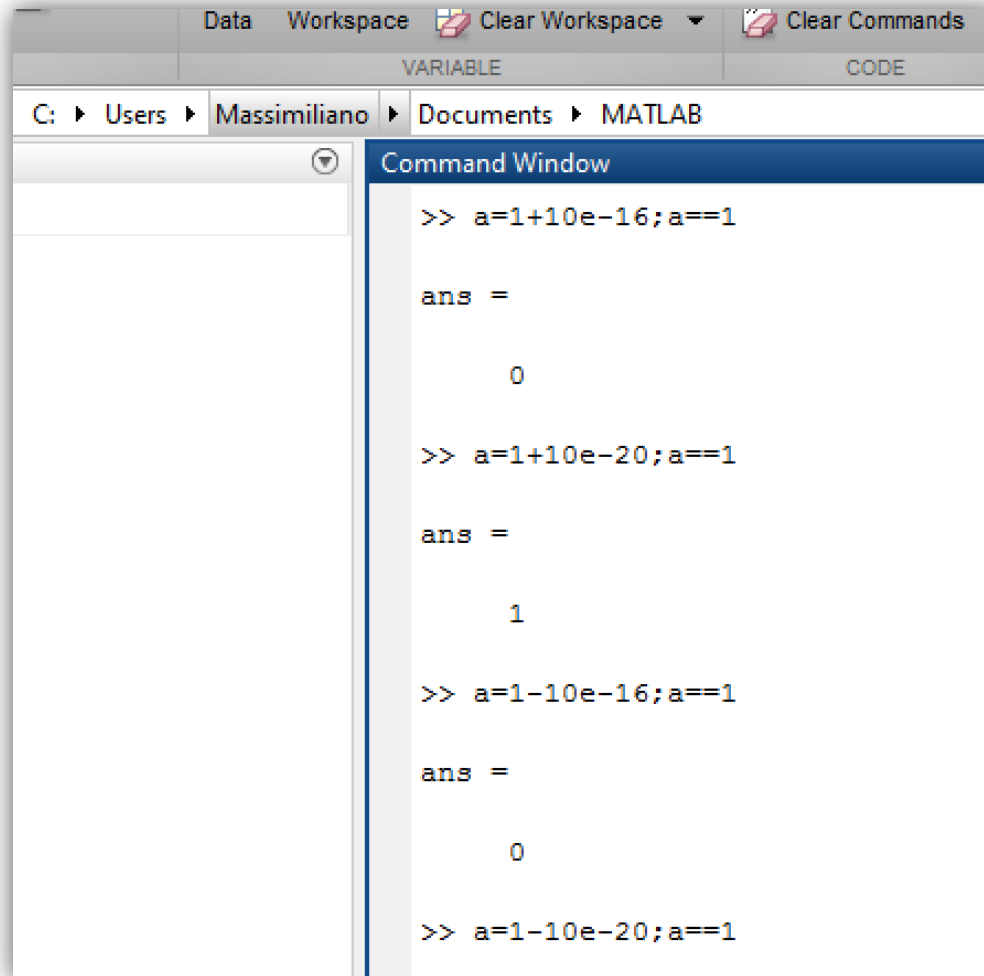
$$a \odot (b \oplus c) \neq a \odot b \oplus a \odot c$$

$$a \odot (b / c) \neq a$$

$$(a / b) \odot b \neq a$$

$$(a \odot b) / c \neq (a / c) \odot b$$

$$1 \pm b = 1 \quad \text{si } b \in \mathbb{R} \wedge |b| < \varepsilon$$



The image shows a screenshot of the MATLAB Command Window. The window title is 'Command Window'. The path bar shows 'C: \> Users \> Massimiliano \> Documents \> MATLAB'. The window contains the following commands and outputs:

```
>> a=1+10e-16;a==1  
  
ans =  
  
0  
  
>> a=1+10e-20;a==1  
  
ans =  
  
1  
  
>> a=1-10e-16;a==1  
  
ans =  
  
0  
  
>> a=1-10e-20;a==1
```

**MATLAB**  
(double précision)

## Exemple 2: somme de deux nombres voisins

$$x_1 = 0,191019721 \cdot 10^3$$

$$x_2 = 0,191017083 \cdot 10^3$$

$r=6$  chiffres significatifs  $\beta = 10$ ,  $x_1 \oplus x_2 = ?$

$$fl(x_1) \quad \begin{array}{|c|c|c|c|c|c|} \hline 1 & 9 & 1 & 0 & 2 & 0 \\ \hline \end{array} 10^3$$

$$\varepsilon = \frac{\beta}{2} \beta^{-r} = 5 \cdot 10^{-6}$$

$$fl(x_2) \quad \begin{array}{|c|c|c|c|c|c|} \hline 1 & 9 & 1 & 0 & 1 & 7 \\ \hline \end{array} 10^3$$

$$e_{x_1} = \frac{|fl(x_1) - x_1|}{|x_1|} = 1.4606 \cdot 10^{-6} < \varepsilon$$

$$fl(x_1) + fl(x_2)$$

$$\begin{array}{|c|c|c|c|c|c|} \hline 3 & 8 & 2 & 0 & 3 & 7 \\ \hline \end{array} 10^3$$

$$e_{x_2} = \frac{|fl(x_2) - x_2|}{|x_2|} = 4.3452 \cdot 10^{-7} < \varepsilon$$

**6 chiffres significatifs**

$$x_1 + x_2 = 3.82036 \cdot 10^2 \quad x_1 \oplus x_2 = 3.82037 \cdot 10^2 \quad e_{x_1+x_2} = \frac{|3.82036 - 3.82037|}{|3.82036|} = 5.13 \cdot 10^{-7}$$

**Même ordre de précision ( $< \varepsilon$ )**

## Exemple 3: différence de deux nombres voisins

$$x_1 = 0,191019721 \cdot 10^3$$

$$x_2 = 0,191017083 \cdot 10^3$$

6 chiffres significatifs  $\beta = 10$ ,  $x_1 \ominus x_2 = ?$

$$fl(x_1) = \begin{array}{|c|c|c|c|c|c|} \hline 1 & 9 & 1 & 0 & 2 & 0 \\ \hline \end{array} 10^3$$

$$\varepsilon = \frac{\beta}{2} \beta^{-r} = 5 \cdot 10^{-6}$$

$$fl(x_2) = \begin{array}{|c|c|c|c|c|c|} \hline 1 & 9 & 1 & 0 & 1 & 7 \\ \hline \end{array} 10^3$$

$$e_{x_1} = \frac{|fl(x_1) - x_1|}{|x_1|} = 1.4606 \cdot 10^{-6} < \varepsilon$$

$$fl(x_1) - fl(x_2) = \begin{array}{|c|c|c|c|c|c|} \hline 0 & 0 & 0 & 0 & 0 & 3 \\ \hline \end{array} 10^3$$

$$e_{x_2} = \frac{|fl(x_2) - x_2|}{|x_2|} = 4.3452 \cdot 10^{-7} < \varepsilon$$

1 chiffre significatif

$$x_1 - x_2 = 2.63799 \cdot 10^{-3} \quad x_1 \ominus x_2 = 3 \cdot 10^{-3}$$

$$e_{x_1 - x_2} = \frac{|2.6380 - 3|}{|2.6380|} = 0.1372 \approx 13.7\%$$

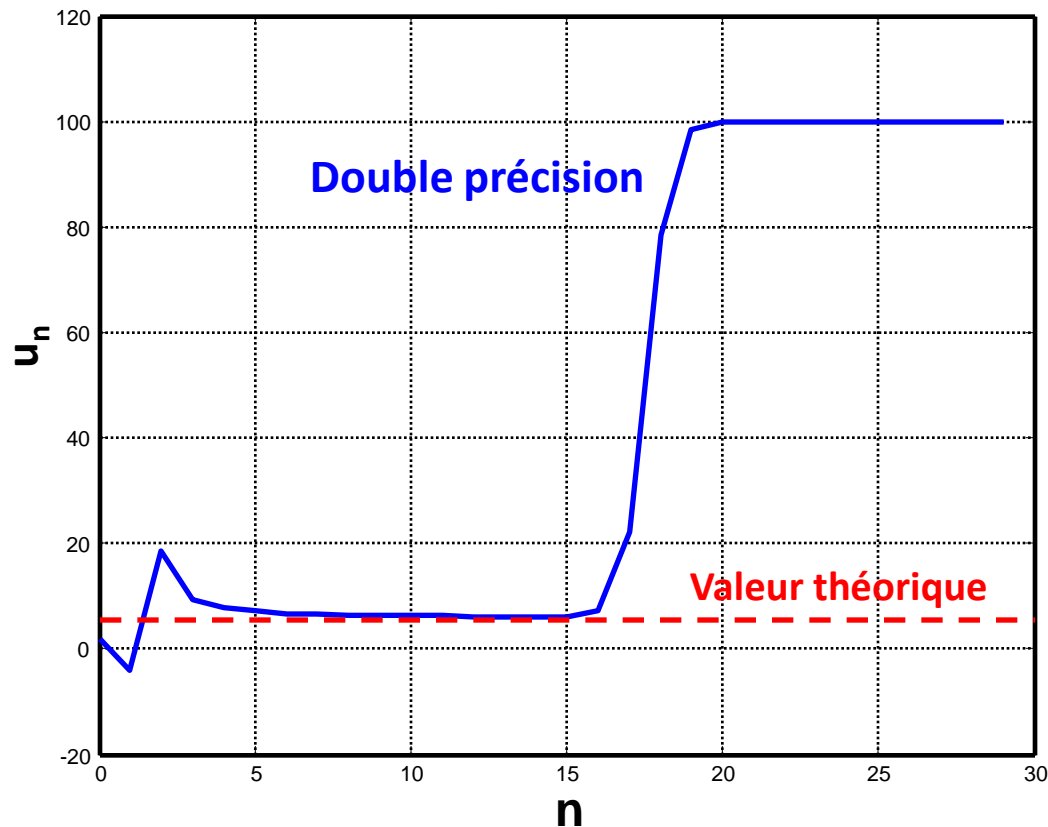
!!!

**Phénomène de cancellation**

$$e_{x_1 - x_2} \gg \varepsilon$$

## Exemple 3: Suite numérique de Muller

$$u_n = \begin{cases} u_0 = 2 \\ u_1 = -4 \\ u_{n+1} = 111 - \frac{1130}{u_n} + \frac{3000}{u_n \cdot u_{n-1}} \end{cases}$$

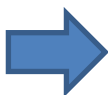


$$x^2 - 2ax + \varepsilon = 0 \quad \begin{cases} x_1 = a + \sqrt{a^2 - \varepsilon} \\ x_2 = a - \sqrt{a^2 - \varepsilon} \end{cases}$$

Si  $\varepsilon \ll a \Rightarrow \sqrt{a^2 - \varepsilon} \simeq |a| \Rightarrow \begin{cases} x_1 = a + |a| \\ x_2 = a - |a| \end{cases} \Rightarrow \begin{cases} x_2 = 0 & \text{si } a > 0 \\ x_1 = 0 & \text{si } a < 0 \end{cases}$  **Phénomène de cancellation**

**Solution:**

$$\varepsilon = x_1 x_2$$



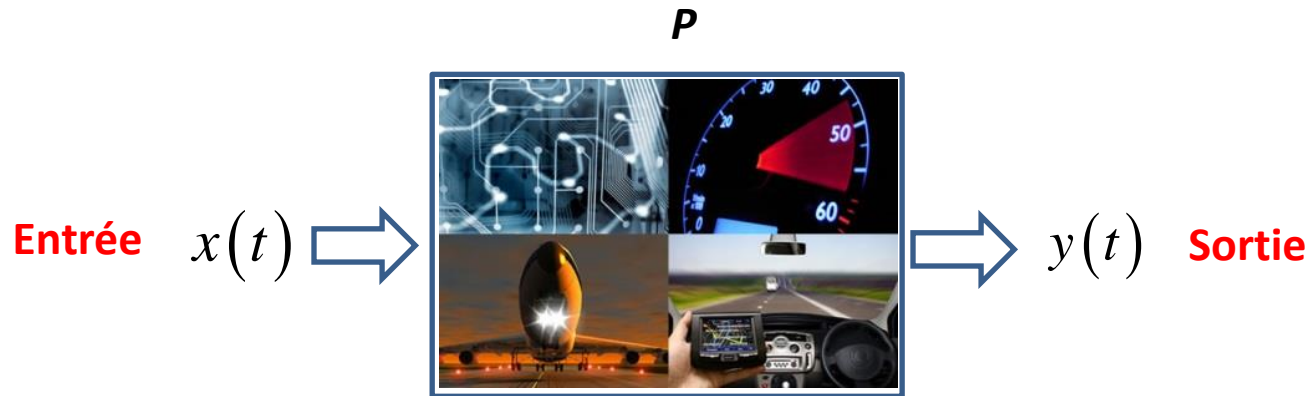
$$\begin{cases} x_1 = a + \operatorname{sgn}(a) \sqrt{a^2 - \varepsilon} \\ x_2 = \frac{\varepsilon}{x_1} \end{cases}$$

Si  $\varepsilon \simeq a^2 \Rightarrow \sqrt{a^2 - \varepsilon} \simeq 0 \Rightarrow \begin{cases} x_1 = a \\ x_2 = a \end{cases}$

**Phénomène de cancellation**  
**aucune solution possible**



- **Représentation d'un nombre en machine, erreurs numériques, etc...**
  - Introduction générale
  - Représentation des nombres entiers et réels
  - Opérations élémentaires en virgule flottante
  - **Conditionnement d'un problème**
  - Opérations complexes
  - Conclusion: différentes sources d'erreur
- Interpolation polynomiale
- Intégration numérique
- Résolution numérique des équations différentielles ordinaires (EDO)

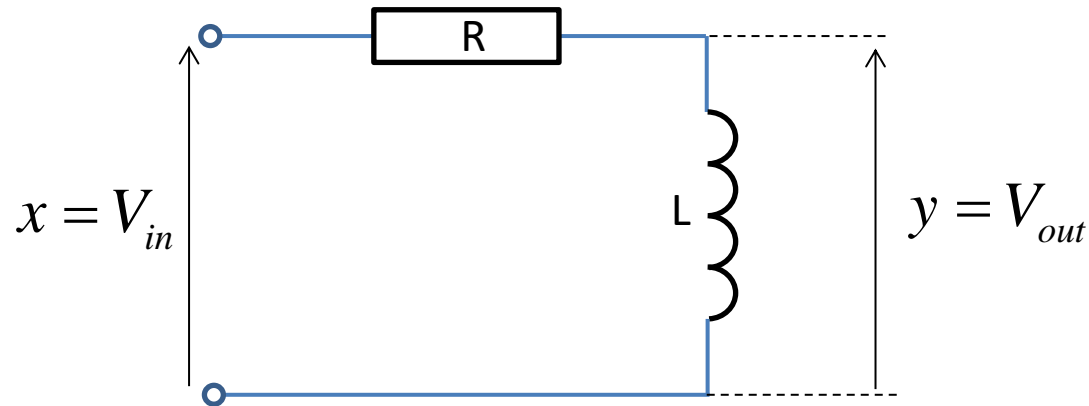


**Forme explicite**

$$y(t) = f[x(t)]$$

**Forme implicite**

$$g[x(t), y(t)] = 0$$

Exemple

Forme explicite

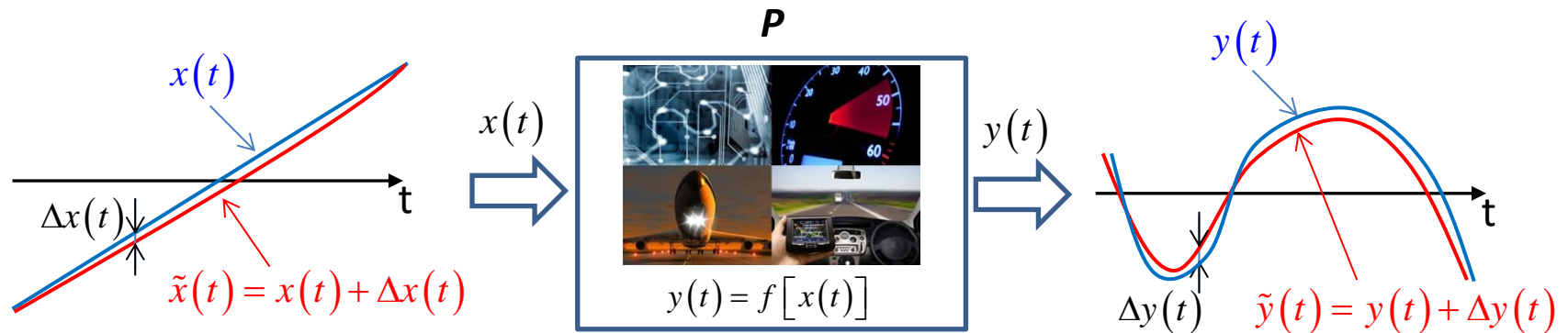
$$V_{out} = V_{in} \frac{j\omega L}{R + j\omega L} \quad f(x) = \left[ \frac{j\omega L}{R + j\omega L} \right] x$$

Forme implicite

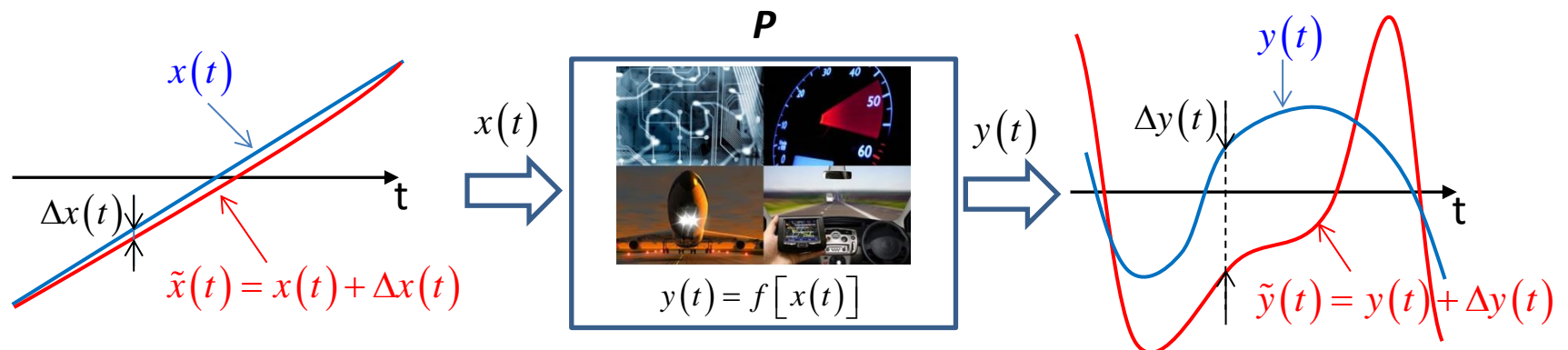
Loi des mailles

$$V_{in} + \frac{R}{R + j\omega L} V_{out} = 0 \quad g(x, y) = x + \frac{R}{R + j\omega L} y$$

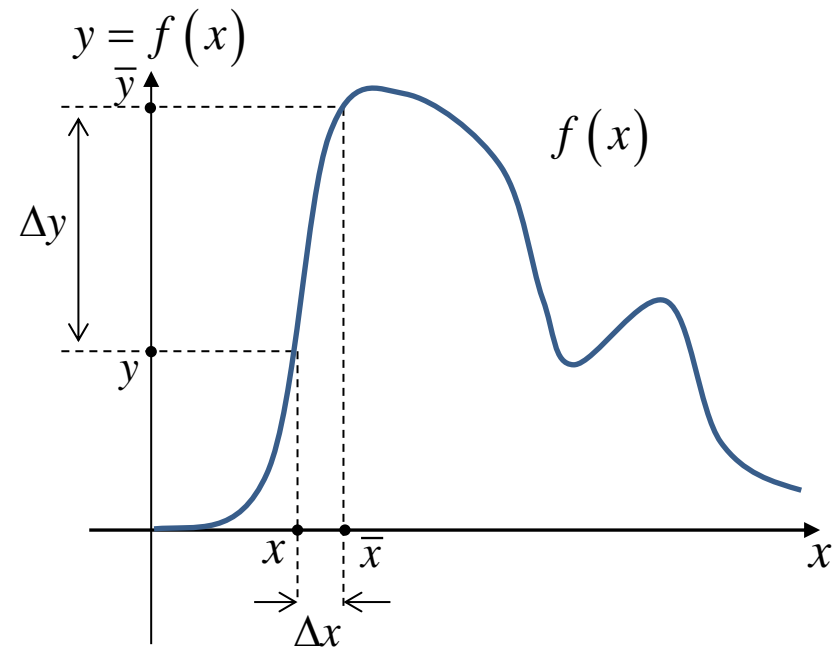
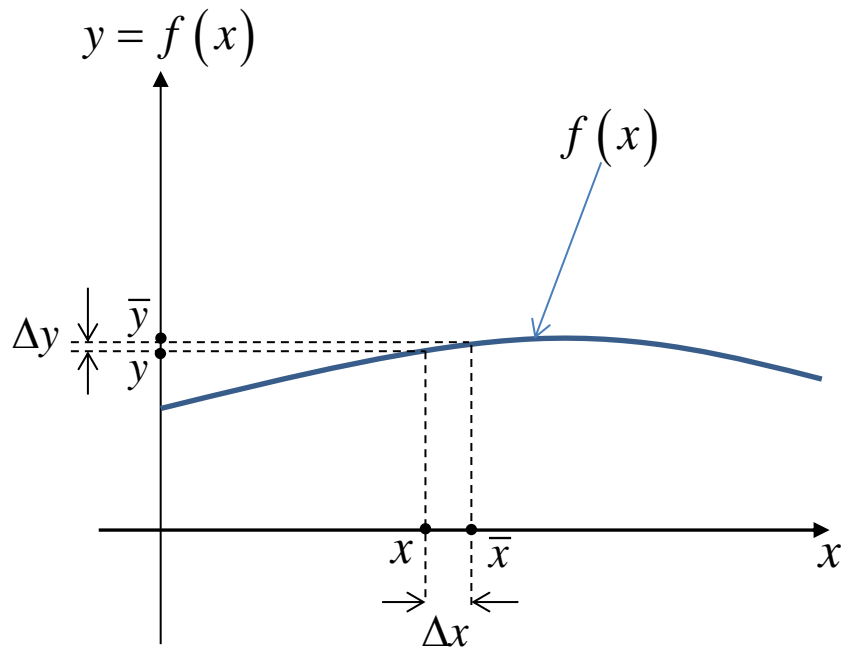
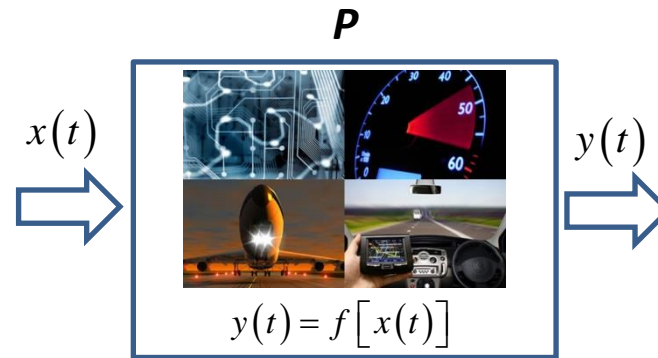
- le **conditionnement** mesure la dépendance de la solution d'un **problème numérique**  $P$  par rapport aux données du problème (mesure de la difficulté du calcul numérique)



On dit que  $P$  est **bien conditionné** si un petit changement  $\Delta x$  de  $x$  entraîne un petit changement  $\Delta Y$  de  $Y$ .



On dit que  $P$  est **mal conditionné** si un petit changement  $\Delta x$  de  $x$  entraîne un grand changement  $\Delta Y$  de  $Y$ .



- Le **conditionnement d'un problème par rapport à l'erreur absolue** est donnée par le nombre de condition absolu  $k$

$$k = \frac{\|\Delta y\|}{\|\Delta x\|}$$

rapport d'amplification des erreurs

- Le **conditionnement d'un problème par rapport à l'erreur relative** est donnée par le nombre de condition relatif  $k_r$

$$k_r = \frac{\frac{\|\Delta y\|}{\|y\|}}{\frac{\|\Delta x\|}{\|x\|}} = \frac{\|\Delta y\| \|x\|}{\|\Delta x\| \|y\|}$$

rapport d'amplification des erreurs relatives

Le problème  $P$  est bien conditionné si  $k$  n'est pas *très grand*. Sinon, ce problème  $P$  est mal conditionné.

## Normes vectorielles

$$\mathbf{x}, \mathbf{y} \in \mathbb{R}^N$$

$$\|\mathbf{x}\| \geq 0 \quad \|\mathbf{x}\| = 0 \Leftrightarrow \mathbf{x} = 0 \quad (\text{séparation})$$

$$\|\alpha \mathbf{x}\| = |\alpha| \|\mathbf{x}\| \quad (\text{homogénéité})$$

$$\|\mathbf{x} + \mathbf{y}\| \leq \|\mathbf{x}\| + \|\mathbf{y}\| \quad (\text{inégalité triangulaire})$$

### Exemples

$$\|\mathbf{x}\|_2 = \sqrt{\sum_{i=1}^N x_i^2}$$

$$\|\mathbf{x}\|_1 = \sum_{i=1}^N |x_i|$$

$$\|\mathbf{x}\|_\infty = \max_{i=1, \dots, N} |x_i|$$

## Normes matricielles

(subordonnée à une norme vectorielle)

$$\underline{\underline{\mathbf{A}}}, \underline{\underline{\mathbf{B}}} \in \mathbb{R}^{N \times N}$$

$$\|\underline{\underline{\mathbf{A}}}\| \geq 0 \quad \|\underline{\underline{\mathbf{A}}}\| = 0 \Leftrightarrow \underline{\underline{\mathbf{A}}} = 0$$

$$\|\alpha \underline{\underline{\mathbf{A}}}\| = |\alpha| \|\underline{\underline{\mathbf{A}}}\|$$

$$\|\underline{\underline{\mathbf{A}}} + \underline{\underline{\mathbf{B}}}\| \leq \|\underline{\underline{\mathbf{A}}}\| + \|\underline{\underline{\mathbf{B}}}\|$$

$$\|\underline{\underline{\mathbf{A}\mathbf{B}}}\| = \|\underline{\underline{\mathbf{A}}}\| \|\underline{\underline{\mathbf{B}}}\|$$

### Exemples

$$\|\underline{\underline{\mathbf{A}}}\|_p = \max_{\|\mathbf{x}\|_p \neq 0} \frac{\|\underline{\underline{\mathbf{A}}}\mathbf{x}\|_p}{\|\mathbf{x}\|_p}$$

$$\|\underline{\underline{\mathbf{A}}}\|_2 = \sqrt{\rho\{\underline{\underline{\mathbf{A}}}^T \underline{\underline{\mathbf{A}}}\}}$$

**En dimension finie, toutes les normes sont équivalentes**



Soit  $P$  définit par  $y=x_1+x_2$

$$\tilde{x}_1 = x_1 + \Delta x_1 \quad \tilde{x}_2 = x_2 + \Delta x_2$$

$$\tilde{y} = f(\tilde{x}_1, \tilde{x}_2) = \tilde{x}_1 + \tilde{x}_2 = (x_1 + \Delta x_1) + (x_2 + \Delta x_2) = (x_1 + x_2) + (\Delta x_1 + \Delta x_2) = y + \Delta y$$

$$\|y\| = |x_1 + x_2| \quad \|\Delta y\| = |\Delta x_1 + \Delta x_2|$$

$$\Delta \mathbf{x} = (\Delta x_1, \Delta x_2)$$

$$|\Delta x_1|, |\Delta x_2| \leq \varepsilon$$



$$\|\mathbf{x}\|_1 = \|(x_1, x_2)\|_1 = |x_1| + |x_2|$$

$$\|\Delta \mathbf{x}\|_1 = \|(\Delta x_1, \Delta x_2)\|_1 = |\Delta x_1| + |\Delta x_2| \leq 2\varepsilon$$



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

$$k_r = \frac{\|\Delta y\| \|\mathbf{x}\|_1}{\|\Delta \mathbf{x}\|_1 \|y\|} = \frac{|\Delta x_1 + \Delta x_2| (|x_1| + |x_2|)}{(|\Delta x_1| + |\Delta x_2|) |x_1 + x_2|} \leq \frac{(|\Delta x_1| + |\Delta x_2|) (|x_1| + |x_2|)}{(|\Delta x_1| + |\Delta x_2|) |x_1 + x_2|} = \frac{|x_1| + |x_2|}{|x_1 + x_2|}$$

Soit  $P$  définit par  $y=x_1-x_2$

$$\tilde{y} = f(\tilde{x}_1, \tilde{x}_2) = \tilde{x}_1 - \tilde{x}_2 = (x_1 + \Delta x_1) - (x_2 + \Delta x_2) = (x_1 - x_2) + (\Delta x_1 - \Delta x_2) = y + \Delta y$$

$$\|y\| = |x_1 - x_2|$$

$$\|\Delta y\| = |\Delta x_1 - \Delta x_2| \leq 2 \max \{|\Delta x_1|, |\Delta x_2|\}$$

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

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



$$k = \frac{\|\Delta y\|}{\|\Delta \mathbf{x}\|_\infty} = \frac{|\Delta x_1 - \Delta x_2|}{\max \{|\Delta x_1|, |\Delta x_2|\}} \leq \frac{2 \max \{|\Delta x_1|, |\Delta x_2|\}}{\max \{|\Delta x_1|, |\Delta x_2|\}} = 2$$

$$k_r = \frac{\|\Delta y\| \|\mathbf{x}\|_\infty}{\|\Delta \mathbf{x}\|_\infty \|y\|} = \frac{|\Delta x_1 - \Delta x_2| \max \{|x_1|, |x_2|\}}{\max \{|\Delta x_1|, |\Delta x_2|\} |x_1 - x_2|} \leq \frac{2 \max \{|\Delta x_1|, |\Delta x_2|\} \max \{|x_1|, |x_2|\}}{\max \{|\Delta x_1|, |\Delta x_2|\} |x_1 - x_2|} = \frac{2 \max \{|x_1|, |x_2|\}}{|x_1 - x_2|}$$

Si  $x_1 \approx x_2$  le conditionnement relatif est très grand. La soustraction est donc mal conditionnée. (phénomène d'annulation).

Soit  $\mathbf{P}$  défini par  $y = x_1 * x_2$

$$\tilde{y} = \tilde{x}_1 \tilde{x}_2 = (x_1 + \Delta x_1)(x_2 + \Delta x_2) = (x_1 x_2) + (x_1 \Delta x_2 + x_2 \Delta x_1 + \Delta x_1 \Delta x_2) = y + \Delta y$$

$$\|\mathbf{x}\|_2 = \|(x_1, x_2)\|_2$$

$$\|\Delta \mathbf{x}\|_2 = \|(\Delta x_1, \Delta x_2)\|_2$$

$$\|y\| = |x_1 x_2|$$

$$\begin{aligned} \|\Delta y\| &= |x_1 \Delta x_2 + x_2 \Delta x_1 + \Delta x_1 \Delta x_2| \approx |x_1 \Delta x_2 + x_2 \Delta x_1| = |(x_1, x_2) \cdot (\Delta x_2, \Delta x_1)| \\ &= |\mathbf{x} \cdot \Delta \mathbf{x}| \leq \|\mathbf{x}\|_2 \|\Delta \mathbf{x}\|_2 = \|x_1, x_2\|_2 \|\Delta x_1, \Delta x_2\|_2 \end{aligned}$$



$$k = \frac{\|\Delta y\|}{\|\Delta \mathbf{x}\|_2} \leq \frac{\|x_1, x_2\|_2 \|\Delta x_1, \Delta x_2\|_2}{\|\Delta x_1, \Delta x_2\|_2} = \|x_1, x_2\|_2 = \|\mathbf{x}\|_2$$

$$k_r = \frac{\|\Delta y\| \|\mathbf{x}\|_2}{\|\Delta \mathbf{x}\|_2 \|y\|} \leq \|\mathbf{x}\|_2 \frac{\|\mathbf{x}\|_2}{\|y\|} = \frac{(\|\mathbf{x}\|_2)^2}{|x_1 x_2|} = \frac{x_1^2}{|x_1 x_2|} + \frac{x_2^2}{|x_1 x_2|} = \frac{|x_1| |x_1|}{|x_1| |x_2|} + \frac{|x_2| |x_2|}{|x_1| |x_2|} = \frac{|x_1|}{|x_2|} + \frac{|x_2|}{|x_1|}$$

Soit  $P$  défini par le système linéaire  $\underline{\underline{A}}\mathbf{y} = \mathbf{x}$  de solution exacte et unique.

Forme implicite

$$\underline{\underline{A}}\mathbf{y} - \mathbf{x} = 0$$



Forme explicite

$$\mathbf{y} = \underline{\underline{A}}^{-1}\mathbf{x} = f(\mathbf{x})$$

$$\|\mathbf{x}\| = \|\underline{\underline{A}}\mathbf{y}\| \leq \|\underline{\underline{A}}\| \|\mathbf{y}\| \Rightarrow \|\mathbf{y}\| \geq \frac{\|\mathbf{x}\|}{\|\underline{\underline{A}}\|}$$

$$\tilde{\mathbf{x}} = \mathbf{x} + \Delta\mathbf{x}$$

$$\underline{\underline{A}}\tilde{\mathbf{y}} = \tilde{\mathbf{x}} \Leftrightarrow \underline{\underline{A}}(\mathbf{y} + \Delta\mathbf{y}) = \mathbf{x} + \Delta\mathbf{x} \Rightarrow \underline{\underline{A}}\Delta\mathbf{y} = \Delta\mathbf{x} \Rightarrow \Delta\mathbf{y} = \underline{\underline{A}}^{-1}\Delta\mathbf{x}$$

$$\|\Delta\mathbf{y}\| = \|\underline{\underline{A}}^{-1}\Delta\mathbf{x}\| \leq \|\underline{\underline{A}}^{-1}\| \|\Delta\mathbf{x}\|$$

$$\frac{\|\Delta\mathbf{y}\|}{\|\mathbf{y}\|} \leq \frac{\|\underline{\underline{A}}^{-1}\| \|\Delta\mathbf{x}\|}{\frac{\|\mathbf{x}\|}{\|\underline{\underline{A}}\|}} = \|\underline{\underline{A}}^{-1}\| \|\underline{\underline{A}}\| \frac{\|\Delta\mathbf{x}\|}{\|\mathbf{x}\|}$$

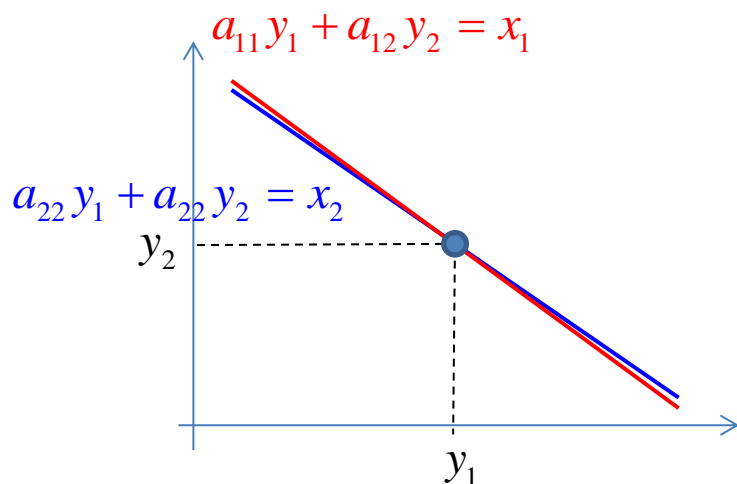
$$k_r \leq \|\underline{\underline{A}}^{-1}\| \|\underline{\underline{A}}\|$$

$$\underline{\underline{\mathbf{A}}}\mathbf{y} = \mathbf{x} \quad \mathbf{y} = \underline{\underline{\mathbf{A}}}^{-1}\mathbf{x} = f(\mathbf{x})$$

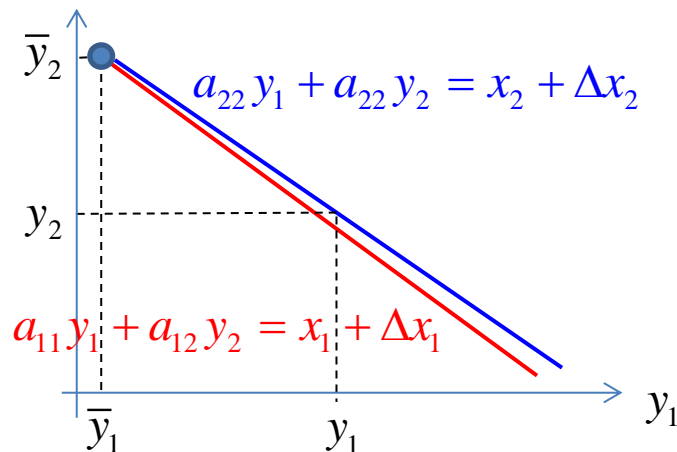
$$\begin{bmatrix} a_{11} & a_{12} \\ a_{21} & a_{22} \end{bmatrix} \begin{bmatrix} y_1 \\ y_2 \end{bmatrix} = \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} \Leftrightarrow \begin{cases} a_{11}y_1 + a_{12}y_2 = x_1 \\ a_{21}y_1 + a_{22}y_2 = x_2 \end{cases}$$

Système linéaire mal conditionnée: **droites presque parallèles**

**Solution exacte**  $\underline{\underline{\mathbf{A}}}\mathbf{y} = \mathbf{x}$



**Solution système perturbé**  $\underline{\underline{\mathbf{A}}}\mathbf{y} = \tilde{\mathbf{x}} = \mathbf{x} + \Delta\mathbf{x}$

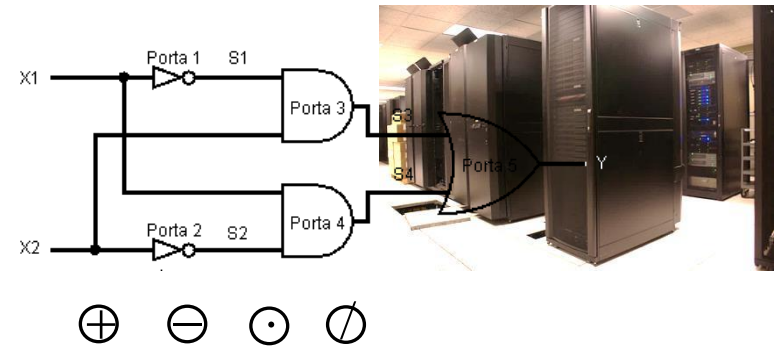


**Exemple:** 
$$\begin{cases} 2y_1 + 2y_2 = 3 \\ 0.499y_1 + 1.001y_2 = 1.5 \end{cases}$$

- **Représentation d'un nombre en machine, erreurs numériques, etc...**
  - Introduction générale
  - Représentation des nombres entiers et réels
  - Opérations élémentaires en virgule flottante
  - Conditionnement d'un problème
  - **Opérations complexes**
  - Conclusion: différentes sources d'erreur
- Interpolation polynomiale
- Intégration numérique
- Résolution numérique des équations différentielles ordinaires (EDO)

Handwritten mathematical formulas on a green chalkboard, including  $\ln x$ ,  $f(e) = \frac{1}{e}$ ,  $\ln x = \frac{1}{x} \cdot x - \ln x$ ,  $\frac{1}{x^2} = \frac{1}{x^2}$ ,  $\ln u = 1 \Rightarrow$ , and  $f(x) = f(x)$ .

$$f(x)$$



Solution

Représentation de  $f$  comme combinaison des fonctions élémentaires

$$f(x) = \sum_{n=0}^{\infty} f_n(x) \quad f(x) = \lim_{n \rightarrow \infty} u_n$$

$f$  limite d'une série ou d'une suite



Valeur exacte pour un nombre d'itérations  $n \rightarrow \infty$

## Série de Taylor

$$f(x) = \sum_{k=0}^n \frac{f^{(k)}(x_0)}{k!} (x - x_0)^k + e(x, x_0), \quad e(x, x_0) = o\left((x - x_0)^k\right)$$

Si  $f^{(n)}(x_0)$  est représentable comme une combinaison des opérations de machine élémentaires



OK, mais de combien de termes a-t-on besoin



## Formule de Taylor-Lagrange

Si  $f$  est de classe  $C^n$  sur  $[x_0, x]$  et admet une dérivée d'ordre  $n+1$ , alors il existe  $\xi \in ]x_0, x[$  tel que

$$e(x, x_0) = \frac{f^{(n+1)}(\xi)}{(n+1)!} (x - x_0)^{n+1}$$

L'erreur dépend de la "*variabilité*" de la fonction et de la distance entre  $x$  et le centre de la série  $x_0$



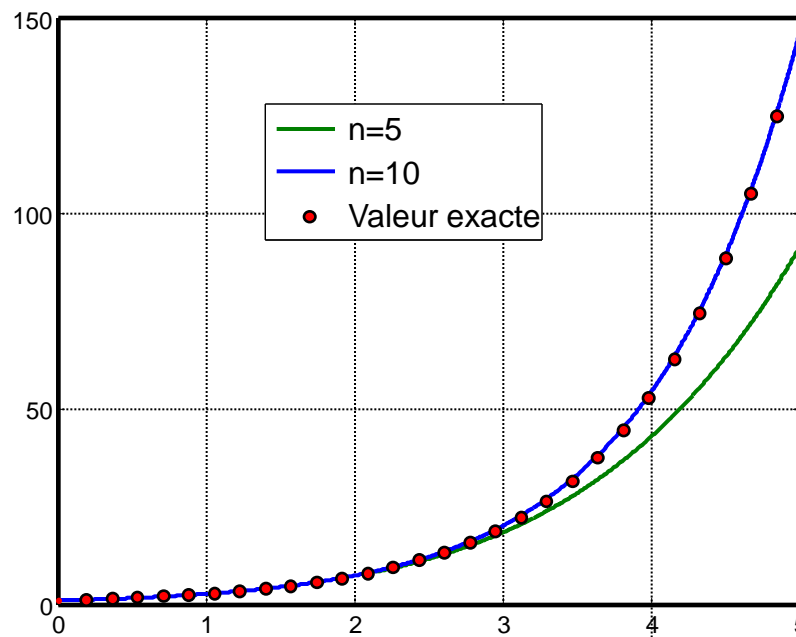
## Exemple 1: fonction exponentielle

$$f^{(k)}(x) = e^x \quad \forall k \geq 0 \quad \text{Si } x_0 = 0 \text{ on a } f(x) = f^{(k)}(x) = 1$$



$$f(x) = \sum_{k=0}^n \frac{x^k}{k!} + \text{err}(x, x_0) \quad \text{err}(x, x_0) = \frac{e^\zeta x^{n+1}}{(n+1)!} \leq \frac{e^{x_{\max}} x^{n+1}}{(n+1)!} \quad \lim_{n \rightarrow \infty} \text{err}(x, x_0) = 0$$

On appelle  $\alpha$  l'erreur relative maximale désirée, on détermine  $n$  par la relation:  $\frac{e^x x^{n+1}}{(n+1)!} \leq \alpha e^x$



$$\alpha \geq \frac{x^{n+1}}{(n+1)!}$$

# Exemple 1: fonction exponentielle

$$f(x) = e^x$$

$$x = 12.3451 = 12 + 0.3451 = x_I + x_F$$

$$f(x) = e^{x_I + x_F} = e^{x_I} \cdot e^{x_F}$$

$$e^{x_I} = \prod_{i=1}^{x_I} e$$

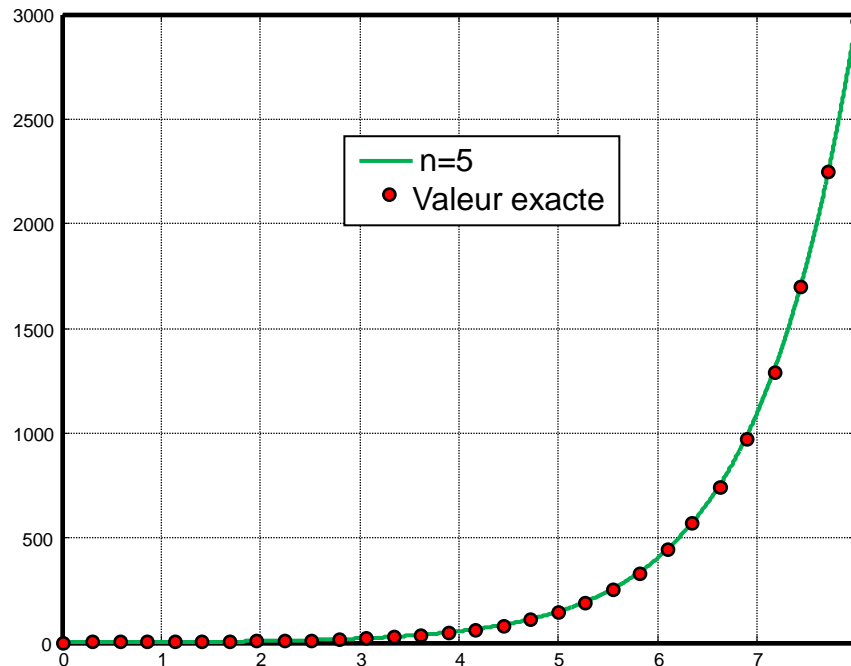
$$e^{x_F} = \sum_{k=0}^n \frac{x_F^k}{k!} + \text{er}(x_F), \quad \text{avec } x_F \leq 1$$

$$\text{err}(x_F) = \frac{e^\zeta x_F^{n+1}}{(n+1)!} \leq \frac{1}{(n+1)!} = \text{err}(1)$$



$$\alpha \geq \frac{1}{(n+1)!}$$

$$n = 5 \Rightarrow \alpha = 0.0083$$



## Exemple 2: fonction sinus

Si  $x_0 = 0$  on a

$$f(x) = \sin(x)$$

$$f^{(n)}(x) = \begin{cases} (-1)^{\frac{n-1}{2}} \cos(x) & n \text{ impair} \\ (-1)^{\frac{n}{2}} \sin(x) & n \text{ pair} \end{cases}$$

$$f(a) = 0$$

$$f^{(n)}(a) = \begin{cases} (-1)^{\frac{n-1}{2}} & n \text{ impair} \\ 0 & n \text{ pair} \end{cases}$$



$$f(x) = \sin(x) = x - \frac{x^3}{6} + \frac{x^5}{5!} \dots + \frac{(-1)^n x^{2n+1}}{(2n+1)!} + o(x^{2n+2}) = \sum_{n=0}^{\infty} \frac{(-1)^n x^{2n+1}}{(2n+1)!} + err(x, x_0)$$

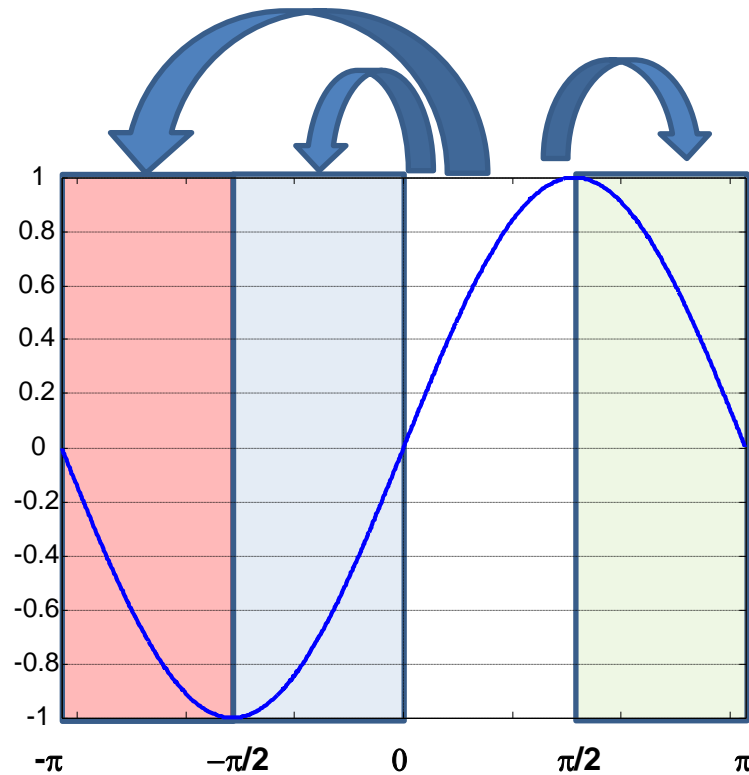
$$|e(x, 0)| = \left| \frac{f^{(n+1)}(\zeta)}{(n+1)!} x^{n+1} \right| \leq \frac{|f^{(n+1)}(\zeta)|}{(n+1)!} |x^{n+1}| \leq \frac{|x^{n+1}|}{(n+1)!}$$

## Exemple 2: fonction sinus

$f(x) = \sin(x) = \sin(x + 2\pi) = f(x + 2\pi)$  fonction périodique de période  $2\pi$

$f(-x) = \sin(-x) = -\sin(x) = -f(x)$  fonction impaire

$f(x) = f(\pi - x)$  pour  $0 < x \leq \frac{\pi}{2}$



## Rappels de notation

$x$  Donnée exacte

$\bar{x} = x + \Delta x$  Donnée perturbée

$$\bar{y} = f(\bar{x}) = f(x + \Delta x)$$

Fonction exacte calculée en précision infinie

$$y_1 = \hat{f}(\bar{x}) = \hat{f}(x + \Delta x)$$

Fonction approchée calculée en précision infinie

$$y^* = fl(\hat{f}(\bar{x})) = fl(\hat{f}(x + \Delta x))$$

Fonction approchée calculée en précision finie

Soit  $f$  une fonction et  $\hat{f}$  sa version approchée.

- Un algorithme est dit **stable** si pour tout

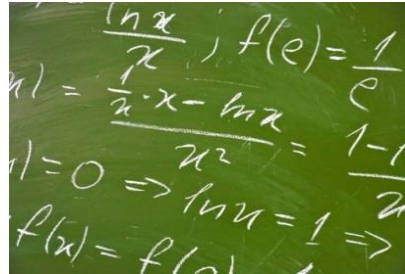
$$\frac{\|\bar{y} - y^*\|}{\|\bar{y}\|} = \frac{\|\hat{f}(\bar{x}) - f(\bar{x})\|}{\|f(\bar{x})\|} = O(\varepsilon) \quad \text{avec} \quad \frac{\|x - \bar{x}\|}{\|x\|} = O(\varepsilon)$$

Un problème est stable, si pour une donnée  $\tilde{x}$  pas très loin de  $x$  on obtient une solution  $\hat{f}(\bar{x})$  pas très loin de  $f(\bar{x})$ .

- **Représentation d'un nombre en machine, erreurs numériques, etc...**
  - **Introduction générale**
  - **Représentation des nombres entiers et réels**
  - **Opérations élémentaires en virgule flottante**
  - **Conditionnement d'un problème**
  - **Opérations complexes**
  - **Conclusion: différentes sources d'erreur**
- Interpolation polynomiale
- Intégration numérique
- Résolution numérique des équations différentielles ordinaires (EDO)



Bureautique, audio-vidéo,...



Calcul scientifique,  
Simulation, modélisation,...



Systèmes embarqués

les ordinateurs effectuent des **calculs**

Puis-je faire confiance aux résultats de mon ordinateur

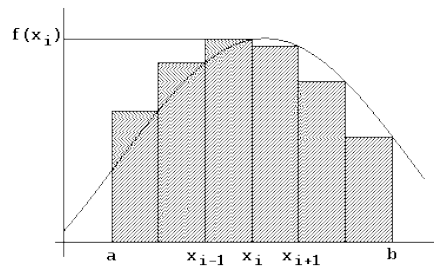


- Erreurs sur les données:** imprécision des mesures physiques, etc.



On peut étudier l'influence de ces erreurs sur le résultat final (conditionnement, etc.)

- Erreurs de méthode:** elles sont dues à l'algorithme utilisé.



(approximation d'une intégrale, d'une somme infinie, ...)

- Erreurs de calcul en machine:** elles sont liées à l'arrondi de calcul pour les nombre flottants.





## Rappels de notation

$x$  Donnée exacte

$\bar{x} = x + \Delta x$  Donnée perturbée

$$y = f(x)$$

$$\bar{y} = f(\bar{x}) = f(x + \Delta x)$$

Fonction exacte calculée en précision infinie

$$y_1 = \hat{f}(\bar{x}) = \hat{f}(x + \Delta x)$$

Fonction approchée calculée en précision infinie

$$y^* = fl(\hat{f}(\bar{x})) = fl(\hat{f}(x + \Delta x))$$

Fonction approchée calculée en précision finie

**Erreur absolue**

$$|y - y^*| = |y - y^* + \bar{y}_1 - \bar{y}_1 + \bar{y} - \bar{y}|$$

$$\leq \underbrace{|y - \bar{y}|}_{\text{Conditionnement du problème}} + \underbrace{|\bar{y} - \bar{y}_1|}_{\text{Erreur d'approximation de la fonction}} + \underbrace{|\bar{y}_1 - y^*|}_{\text{erreur de l'algorithme (arithmétique en précision finie)}}$$

Conditionnement  
du problème

Erreur d'approximation  
de la fonction

erreur de l'algorithme  
(arithmétique en  
précision finie)