

Algorithme de Horner compensé en précision finie et applications

Stef Graillat

LIP6/PEQUAN - Université Pierre et Marie Curie (Paris 6)

Séminaire SPIRAI/SALSA

1 juin 2007, Paris



Plan de l'exposé

- 1 Motivations
- 2 Évaluation précise de polynômes
- 3 Applications

- 1 Motivations
- 2 Évaluation précise de polynômes
- 3 Applications

Proposer des algorithmes et des logiciels :

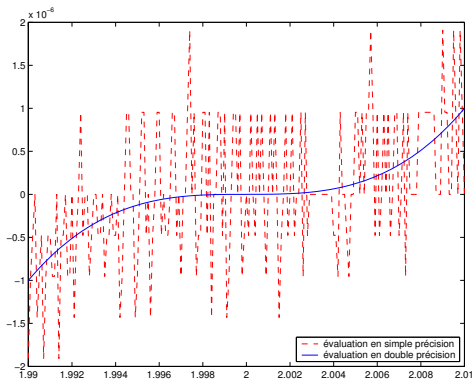
- plus précis que ceux utilisant les précisions de la norme IEEE 754
 - ▷ on veut que la précision vérifie une «version améliorée» de l'estimation empirique classique
- efficace en terme de performance sans sacrifier la portabilité
 - ▷ on utilise seulement la précision IEEE 754 simple ou double
- avec une borne d'erreur pour contrôler la précision du résultat
 - ▷ borne d'erreur dynamique et certifiée calculable en précision finie

Un exemple : le schéma de Horner pour l'évaluation polynomiale
→ le schéma de Horner compensé¹

¹SG, N. Louvet, Ph. Langlois. Compensated Horner Scheme. Research Report, 2005

Imprecision dans l'évaluation polynomiale

Évaluation du polynôme $p(x) = (x - 2)^3 = x^3 - 6x^2 + 12x - 8$ pour environ 200 points au voisinage de $x = 2$ en **simple** et **double** précision



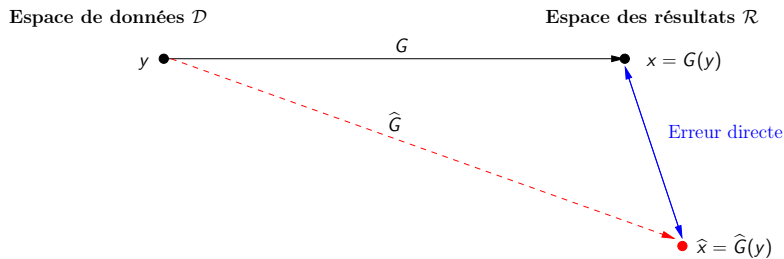
But : Résoudre des problèmes numériques en étant **précis** et **fiable**

- **Comprendre** l'influence de la précision finie sur la qualité numérique du logiciel scientifique pour **contrôler et limiter ses effets néfastes**
 - résultat imprécis ;
 - instabilité numérique.
- **Améliorer** la précision des résultats

Comment être plus précis à faible coût

Comprendre l'influence de la précision finie sur la qualité numérique du logiciel scientifique pour contrôler et limiter ses effets néfastes :

- Contrôler les effets de la précision finie :
 - Comment mesurer la **difficulté de résolution** d'un problème ?
 - Comment apprécier la **fiabilité de l'algorithme** de résolution ?
 - Comment estimer la **précision de la solution** calculée ?
- Limiter les effets de la précision finie
 - Comment **améliorer la précision de la solution** ?



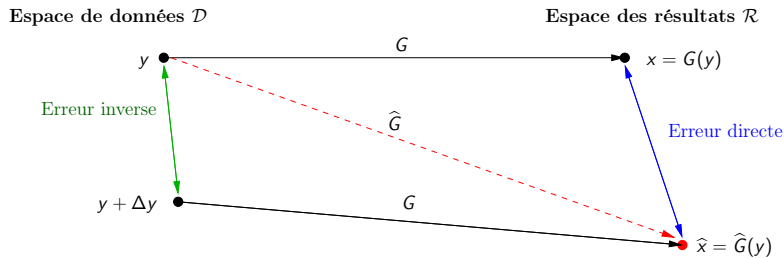
- Analyse directe

- Analyse inverse

consiste à identifier \hat{x} à la solution d'un problème perturbé :

$$\hat{x} = G(y + \Delta y).$$

Analyse d'erreur



- Analyse directe

- Analyse inverse

consiste à identifier \hat{x} à la solution d'un problème perturbé :

$$\hat{x} = G(y + \Delta y).$$

- **Comment estimer la précision de la solution calculée ?**

Au premier ordre, on a l'estimation empirique :

$$\text{erreur directe} \lesssim \text{conditionnement} \times \text{erreur inverse.}$$

- **Comment mesurer la difficulté de résolution d'un problème ?**

Le conditionnement caractérise la sensibilité de la solution d'un problème à des perturbations sur les données.

$$\text{Conditionnement} : K(P, y) := \lim_{\varepsilon \rightarrow 0} \sup_{\Delta y \in \mathcal{P}(\varepsilon)} \left\{ \frac{\|\Delta x\|_{\mathcal{R}}}{\|\Delta y\|_{\mathcal{D}}} \right\}$$

- **Comment apprécier la fiabilité de l'algorithme de résolution ?**

L'erreur inverse mesure la distance entre le problème que l'on a effectivement résolu et le problème initial.

$$\text{Erreur inverse} : \eta(\hat{x}) = \min_{\Delta y \in \mathcal{D}} \{ \|\Delta y\|_{\mathcal{D}} : \hat{x} = G(y + \Delta y) \}$$

- **Comment estimer la précision de la solution calculée ?**

Au premier ordre, on a l'estimation empirique :

$$\text{erreur directe} \lesssim \text{conditionnement} \times \text{erreur inverse.}$$

- **Comment mesurer la difficulté de résolution d'un problème ?**

Le conditionnement caractérise la sensibilité de la solution d'un problème à des perturbations sur les données.

$$\text{Conditionnement} : K(P, y) := \lim_{\varepsilon \rightarrow 0} \sup_{\Delta y \in \mathcal{P}(\varepsilon)} \left\{ \frac{\|\Delta x\|_{\mathcal{R}}}{\|\Delta y\|_{\mathcal{D}}} \right\}$$

- **Comment apprécier la fiabilité de l'algorithme de résolution ?**

L'erreur inverse mesure la distance entre le problème que l'on a effectivement résolu et le problème initial.

$$\text{Erreur inverse} : \eta(\hat{x}) = \min_{\Delta y \in \mathcal{D}} \{ \|\Delta y\|_{\mathcal{D}} : \hat{x} = G(y + \Delta y) \}$$

- **Comment estimer la précision de la solution calculée ?**

Au premier ordre, on a l'estimation empirique :

$$\text{erreur directe} \lesssim \text{conditionnement} \times \text{erreur inverse.}$$

- **Comment mesurer la difficulté de résolution d'un problème ?**

Le conditionnement caractérise la sensibilité de la solution d'un problème à des perturbations sur les données.

$$\text{Conditionnement} : K(P, y) := \lim_{\varepsilon \rightarrow 0} \sup_{\Delta y \in \mathcal{P}(\varepsilon)} \left\{ \frac{\|\Delta x\|_{\mathcal{R}}}{\|\Delta y\|_{\mathcal{D}}} \right\}$$

- **Comment apprécier la fiabilité de l'algorithme de résolution ?**

L'erreur inverse mesure la distance entre le problème que l'on a effectivement résolu et le problème initial.

$$\text{Erreur inverse} : \eta(\hat{x}) = \min_{\Delta y \in \mathcal{D}} \{ \|\Delta y\|_{\mathcal{D}} : \hat{x} = G(y + \Delta y) \}$$

Plan de l'exposé

- 1 Motivations
- 2 Évaluation précise de polynômes
- 3 Applications

Nombres à virgule flottante

Un nombre flottant normalisé $x \in \mathbb{F}$ est un nombre qui s'écrit sous la forme

$$x = \pm \underbrace{x_0.x_1 \dots x_{p-1}}_{\text{mantisse}} \times b^e, \quad 0 \leq x_i \leq b-1, \quad x_0 \neq 0$$

b : la base, p : précision, e : exposant vérifiant $e_{\min} \leq e \leq e_{\max}$

Précision machine $\epsilon = b^{1-p}$, $|1^+ - 1| = \epsilon$

Approximation de \mathbb{R} par \mathbb{F} , arrondi fl : $\mathbb{R} \rightarrow \mathbb{F}$

Soit $x \in \mathbb{R}$ alors

$$\text{fl}(x) = x(1 + \delta), \quad |\delta| \leq u.$$

L'unité d'arrondi u vaut $u = \epsilon/2$ pour l'arrondi au plus près.

Modèle standard de l'arithmétique à virgule flottante

Soient $x, y \in \mathbb{F}$,

$$\text{fl}(x \circ y) = (x \circ y)(1 + \delta), \quad |\delta| \leq \mathbf{u}, \quad \circ \in \{+, -, \cdot, /\}$$

Norme IEEE 754 (1985)

Type	Taille	Mantisse	Exposant	Unité d'arrondi	Intervalle
Simple	32 bits	23+1 bits	8 bits	$\mathbf{u} = 2^{-24} \approx 5,96 \times 10^{-8}$	$\approx 10^{\pm 38}$
Double	64 bits	52+1 bits	11 bits	$\mathbf{u} = 2^{-53} \approx 1,11 \times 10^{-16}$	$\approx 10^{\pm 308}$

- Évaluation de $p(x)$ plus précise : le **schéma de Horner compensé** et l'**estimation empirique compensée**
- Une borne améliorée et **certifiée** de l'erreur
- Les résultats théoriques et expérimentaux montrent que
 - précision : identique à celle obtenue si les calculs avaient été effectués avec **deux fois la précision courante**,
 - vitesse : **deux fois plus rapide** que l'implémentation double-double de référence

Plus de précision, comment ?

Augmenter la précision interne des calculs :

- de façon matérielle
 - précision étendue sur l'architecture x86
- de façon logicielle
 - expansions de longueur finie : double-double (Briggs, Bailey, Hida, Li), quad-double (Bailey, Hida, Li)
 - expansions de longueur variable : Priest, Shewchuk
 - précision arbitraire : MP, MPFUN/ARPREC, MPFR

Corriger les erreurs d'arrondis :

- sommation compensée (Kahan,1965) et doublement compensée (Priest,1991), etc.
- sommation et produit scalaire : Ogita, Rump et Oishi (2005)
→ résultat avec une précision identique à celle obtenue si on faisait les calculs avec 2 fois la précision courante

Plus de précision, comment ?

Augmenter la précision interne des calculs :

- de façon matérielle
 - précision étendue sur l'architecture x86
- de façon logicielle
 - expansions de longueur finie : double-double (Briggs, Bailey, Hida, Li), quad-double (Bailey, Hida, Li)
 - expansions de longueur variable : Priest, Shewchuk
 - précision arbitraire : MP, MPFUN/ARPREC, MPFR

Corriger les erreurs d'arrondis :

- sommation compensée (Kahan,1965) et doublement compensée (Priest,1991), etc.
- sommation et produit scalaire : Ogita, Rump et Oishi (2005)
→ résultat avec une précision identique à celle obtenue si on faisait les calculs avec 2 fois la précision courante

Estimation empirique pour les algorithmes inverses-stables :

précision du résultat \approx conditionnement \times précision de calcul

- 1 précision IEEE-754 : double ($\mathbf{u} = 2^{-53} \approx 10^{-16}$)
- 2 Conditionnement pour l'évaluation de $p(x) = \sum_{i=0}^n a_i x^i$:

$$\text{cond}(p, x) = \frac{\sum_{i=0}^n |a_i| |x|^i}{|\sum_{i=0}^n a_i x^i|} = \frac{\tilde{p}(|x|)}{|p(x)|}, \text{ toujours } \geq 1.$$

- 3 Précision de la solution $\hat{p}(x)$:

$$\frac{|p(x) - \hat{p}(x)|}{|p(x)|} \leq \alpha(n) \times \text{cond}(p, x) \times \mathbf{u}$$

avec $\alpha(n) \approx 2n$

Que signifie doubler la précision courante ?

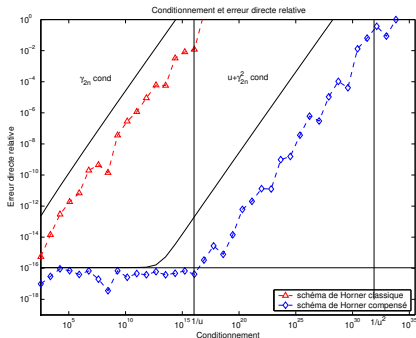
Estimation empirique compensée :

$$\text{précision du résultat} \lesssim \text{précision} + \text{conditionnement} \times \text{précision}^2$$

Trois régimes pour la précision de l'évaluation de $\hat{p}(x)$:

1) conditionnement $\leq 1/u$: la précision de $\hat{p}(x)$ est optimale

$$\frac{|\hat{p}(x) - p(x)|}{|p(x)|} \approx u$$



Que signifie doubler la précision courante ?

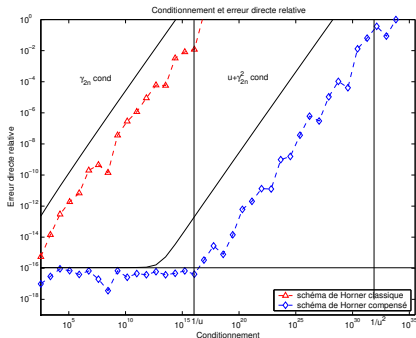
Estimation empirique compensée :

$$\text{précision du résultat} \lesssim \text{précision} + \text{conditionnement} \times \text{précision}^2$$

Trois régimes pour la précision de l'évaluation de $\hat{p}(x)$:

2) $1/u \leq \text{conditionnement} \leq 1/u^2$: la précision de $\hat{p}(x)$

$$\frac{|\hat{p}(x) - p(x)|}{|p(x)|} \approx \text{cond} \times u^2$$



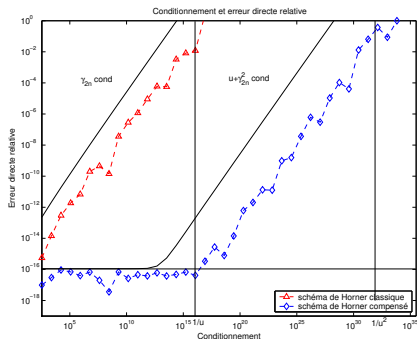
Que signifie doubler la précision courante ?

Estimation empirique compensée :

$$\text{précision du résultat} \lesssim \text{précision} + \text{conditionnement} \times \text{précision}^2$$

Trois régimes pour la précision de l'évaluation de $\hat{p}(x)$:

3) perte totale de précision quand conditionnement $> 1/u^2$.



Algorithme 1 (Algorithme de Horner classique)

```
function res = Horner(p, x)
    sn = an
    for i = n - 1 : -1 : 0
        pi = fl(si+1 · x)           % erreur d'arrondi πi
        si = fl(pi + ai)         % erreur d'arrondi σi
    end
    res = s0
```

$$\gamma_n = nu / (1 - nu) \approx nu$$

$$\frac{|p(x) - \text{Horner}(p, x)|}{|p(x)|} \leq \gamma_{2n} \text{cond}(p, x)$$

Transformations exactes pour la somme

$$x = \text{fl}(a \pm b) \Rightarrow a \pm b = x + y \quad \text{avec } y \in \mathbb{F},$$

Algorithmes de Dekker (1971) et Knuth (1974)

Algorithme 2 (Transformation exacte pour la somme de 2 flottants, nécessite $|a| \geq |b|$)

```
function [x, y] = FastTwoSum(a, b)
    x = fl(a + b)
    y = fl((a - x) + b)
```

Algorithme 3 (Transformation exacte pour la somme de 2 flottants)

```
function [x, y] = TwoSum(a, b)
    x = fl(a + b)
    z = fl(x - a)
    y = fl((a - (x - z)) + (b - z))
```


Transformations exactes pour le produit (1/2)

$$x = \text{fl}(a \cdot b) \Rightarrow a \cdot b = x + y \quad \text{avec } y \in \mathbb{F},$$

Algorithme TwoProduct de Veltkamp et Dekker (1971)

$$a = x + y \quad \text{et} \quad x \text{ et } y \text{ ne se chevauchent pas avec } |y| \leq |x|.$$

Algorithme 4 (Séparation exacte d'un flottant en deux parties)

```
function [x, y] = Split(a, b)
    factor = fl(2s + 1)           % u = 2-p, s = [p/2]
    c = fl(factor · a)
    x = fl(c - (c - a))
    y = fl(a - x)
```

Algorithme 5 (Transformation exacte pour le produit de 2 flottants)

```
function  $[x, y] = \text{TwoProduct}(a, b)$ 
```

```
   $x = \text{fl}(a \cdot b)$ 
```

```
   $[a_1, a_2] = \text{Split}(a)$ 
```

```
   $[b_1, b_2] = \text{Split}(b)$ 
```

```
   $y = \text{fl}(a_2 \cdot b_2 - (((x - a_1 \cdot b_1) - a_2 \cdot b_1) - a_1 \cdot b_2))$ 
```

Transformations exactes pour la somme et le produit

Étant donnés $a, b, c \in \mathbb{F}$,

- $\text{FMA}(a, b, c)$ est l'arrondi au plus près de $a \cdot b + c \in \mathbb{F}$
- $\text{ADD3}(a, b, c)$ est l'arrondi au plus près de $a + b + c \in \mathbb{F}$

Algorithme 6 (Transformation exacte pour la somme de 2 flottants)

```
function [x, y] = TwoSumADD3(a, b)
    x = fl(a + b)
    y = ADD3(a, b, -x)
```

Algorithme 7 (Transformation exacte pour le produit de 2 flottants)

```
function [x, y] = TwoProductFMA(a, b)
    x = fl(a · b)
    y = FMA(a, b, -x)
```

Algorithme 8 (Transformation exacte pour le schéma de Horner)

fonction $[\text{Horner}(p, x), p_\pi, p_\sigma] = \text{EFTHorner}(p, x)$

$s_n = a_n$

for $i = n - 1 : -1 : 0$

$[p_i, \pi_i] = \text{TwoProduct}(s_{i+1}, x)$

$[s_i, \sigma_i] = \text{TwoSum}(p_i, a_i)$

Soit π_i le coefficient de degré i de p_π

Soit σ_i le coefficient de degré i de p_σ

end

$\text{Horner}(p, x) = s_0$

$$p(x) = \text{Horner}(p, x) + (p_\pi + p_\sigma)(x)$$

avec $p_\pi(x) = \sum_{i=0}^{n-1} \pi_i x^i$ et $p_\sigma(x) = \sum_{i=0}^{n-1} \sigma_i x^i$ avec π_i et σ_i dans \mathbb{F}

Algorithme 9 (Algorithme de Horner compensé)

```
function res = CompHorner( $p, x$ )  
[ $h, p_\pi, p_\sigma$ ] = EFTHorner( $p, x$ )  
 $c$  = Horner( $p_\pi + p_\sigma, x$ )  
res = fl( $h + c$ )
```

Théorème 1

Soient p un polynôme de degré n à coefficients flottants et x un nombre flottant. Alors en absence d'underflow, on a

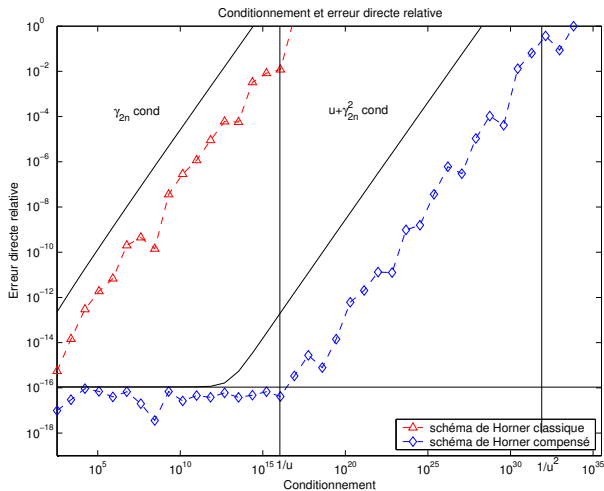
$$\frac{|\text{CompHorner}(p, x) - p(x)|}{|p(x)|} \leq \mathbf{u} + \gamma_{2n}^2 \text{cond}(p, x).$$

- argument clé de la preuve :

$$(\widetilde{p}_\pi + \widetilde{p}_\sigma)(|x|) \leq \gamma_{2n} \widetilde{p}(|x|)$$

- on a une borne similaire en présence d'underflow

Évaluation de $p_n(x) = (x - 1)^n$ pour $x = \text{fl}(1.333)$ et $n = 3, \dots, 42$



Nous comparons

- **Horner** : algorithme de Horner en double précision IEEE 754
- **CompHorner** : algorithme de Horner compensé
- **DDHorner** : algorithme de Horner avec les double-double

Tous les tests de performances ont été fait en langage C avec la double précision IEEE 754

Pour chaque polynôme p_n avec n variant 3 à 42 :

- on effectue 100 tests mesurant le nombre de cycles,
- on garde la valeur moyenne, la valeur minimum et la valeur maximum des 10 plus petits nombres de cycles.

Pentium 4 : 3.0GHz, 1024kB cache L2 - GCC 3.4.1				
ratio	minimum	moyen	maximum	théorique
CompHorner/Horner	1.5	2.9	3.2	13
DDHorner/Horner	2.3	8.4	9.4	17

Intel Celeron : 2.4GHz, 256kB cache L2 - GCC 3.4.1				
ratio	minimum	moyen	maximum	théorique
CompHorner/Horner	1.4	3.1	3.4	13
DDHorner/Horner	2.3	8.4	9.4	17

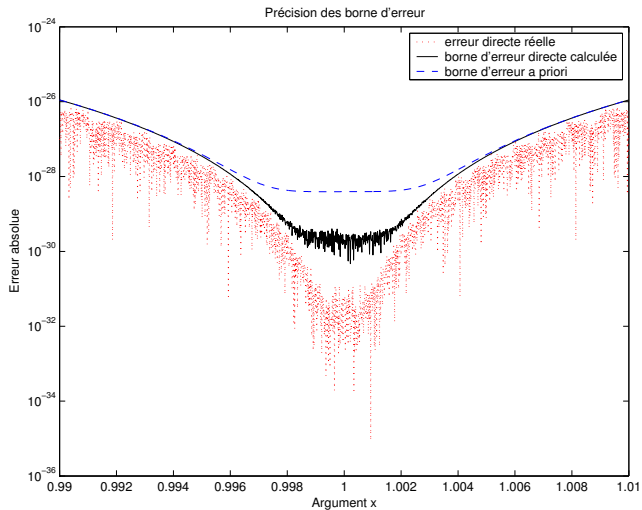
→ algorithme de Horner compensé = algorithme de Horner avec les double-double **sans renormalisation**

Théorème 2

Soient p un polynôme de degré n à coefficients flottants, x un nombre flottant et $\text{res} = \text{CompHorner}(p, x)$ l'évaluation précise par l'algorithme de Horner compensé de p en x . L'erreur directe absolue commise lors de l'évaluation vérifie

$$|\text{CompHorner}(p, x) - p(x)| \leq \text{fl}((\mathbf{u}|\text{res}| + (\gamma_{4n+2}\text{HornerSum}(|p_\pi|, |p_\sigma|, |x|) + 2\mathbf{u}^2|\text{res}|))).$$

Précision de la borne certifiée pour $p_5(x) = (x - 1)^5$



Plan de l'exposé

- 1 Motivations
- 2 Évaluation précise de polynômes
- 3 Applications**

L'évaluation polynomiale précise est utilisée dans

- l'évaluation des fonctions élémentaires
- l'évaluation par interpolation polynomiale
- la géométrie algorithmique (intersection de surface, prédicats géométriques)
- le calcul scientifique (calcul des racines de polynômes)
- robotique (polynômes à plusieurs variables)

Le prédicat **Orient3D** : détermine si un point D est du côté droit ou gauche du plan orienté défini par A , B et C . Le résultat dépend du signe du déterminant

$$\text{Orient3D}(a, b, c, d) = \text{signe} \begin{vmatrix} a_x & a_y & a_z & 1 \\ b_x & b_y & b_z & 1 \\ c_x & c_y & c_z & 1 \\ d_x & d_y & d_z & 1 \end{vmatrix}$$

→ revient à évaluer un polynôme à plusieurs variables.

Approximation de racines simples par la méthode de Horner (1/2)

Polynôme $p(z) = \sum_{i=0}^n a_i z^i$ avec une racine simple x

Définition 1

Étant donné un polynôme p et une racine simple x , on définit le conditionnement de la racine par

$$\text{cond}(p, x) = \limsup_{\varepsilon \rightarrow 0} \left\{ \frac{|\Delta x|}{\varepsilon |x|} : |\Delta a_i| \leq \varepsilon |a_i| \right\}.$$

Théorème 3

On a la relation suivante

$$\text{cond}(p, x) = \frac{\tilde{p}(|x|)}{|x| |p'(x)|}.$$

Approximation de racines simples par la méthode de Horner (2/2)

Algorithme de Newton pour les racines simples :

Algorithme 10

$$x_0 = \xi$$

$$x_{i+1} = x_i - \frac{p(x_i)}{p'(x_i)} ; \quad x_{i+1} = x_i - \frac{\text{CompHorner}(p, x_i)}{p'(x_i)}$$

Théorème 4

Si ξ est suffisamment proche de x alors l'erreur relative décroît jusqu'au premier i pour lequel

$$\frac{|\hat{x}_i - x|}{|x|} \leq \mathbf{u} + \text{cond}(p, x)\mathbf{u}^2$$

Et les polynômes à coefficients complexes ?

$$p(z) = \sum_{j=0}^n a_j z^j, \quad a_j \in \mathbb{C}, z = x + iy \in \mathbb{C}$$

→ Écrire $p(z) = p_r(x, y) + iq_i(x, y)$ avec p_r et q_r à coefficients réels puis évaluer p_r et q_r avec l'algorithme de Horner précédent

Pb : nécessite des manipulations formelles

⇒ il faut de nouvelles transformations exactes pour les complexes

Transformations exactes complexes (1/2)

Pour $x, y \in \mathbb{F} + i\mathbb{F}$,

$$\text{fl}(x \circ y) = (x \circ y)(1 + \varepsilon_1), \text{ pour } \circ \in \{+, -\} \text{ et } |\varepsilon_\nu| \leq \mathbf{u},$$

et

$$\text{fl}(x \cdot y) = (x \cdot y)(1 + \varepsilon_1), |\varepsilon_1| \leq \sqrt{2}\gamma_2.$$

Algorithme 11 (Transformation exacte pour la somme de 2 flottants $x = a + ib$ and $y = c + id$)

```
function [s, e] = TwoSumCplx(x, y)
```

```
    [s1, e1] = TwoSum(a, c)
```

```
    [s2, e2] = TwoSum(b, d)
```

```
    s = s1 + is2
```

```
    e = e1 + ie2
```

$$x + y = s + e, \quad s = \text{fl}(x + y), \quad |e| \leq \mathbf{u}|s|, \quad |e| \leq \mathbf{u}|x + y|$$

Transformations exactes complexes (2/2)

Algorithme 12 (Transformation exacte pour le produit de 2 flottants $x = a + ib$ and $y = c + id$)

```
function [p, e, f, g] = TwoProductCplx(x, y)
```

$$[z_1, h_1] = \text{TwoProduct}(a, c)$$

$$[z_2, h_2] = \text{TwoProduct}(b, d)$$

$$[z_3, h_3] = \text{TwoProduct}(a, d)$$

$$[z_4, h_4] = \text{TwoProduct}(b, c)$$

$$[z_5, h_5] = \text{TwoSum}(z_1, -z_2)$$

$$[z_6, h_6] = \text{TwoSum}(z_3, z_4)$$

$$p = z_5 + iz_6$$

$$e = h_1 + ih_3$$

$$f = -h_2 + ih_4$$

$$g = h_5 + ih_6$$

$$x \cdot y = p + e + f + g \quad p = \text{fl}(x \cdot y), \quad |e + f + g| \leq \sqrt{2}\gamma_2|x \cdot y|$$

- Le schéma de Horner compensé fournit
 - une **précision double de la précision courante**,
 - un algorithme **deux fois plus rapide** que la version double-double,
 - une borne d'erreur **dynamique et certifiée**.
- Développements passés, en cours et futurs
 - schéma de Horner compensé : underflow, avec FMA, pour le FMA
 - même technique avec la **méthode de Newton** pour d'autres problèmes

La future révision de la norme IEEE 754 devrait inclure les fonctions `tailadd`, `tailsubtract` et `tailmultiply` qui calculent l'erreur d'arrondi pour l'addition, la soustraction et la multiplication.



Nicholas J. Higham.

Accuracy and stability of numerical algorithms.

Society for Industrial and Applied Mathematics (SIAM), Philadelphia, PA, second edition, 2002.



Takeshi Ogita, Siegfried M. Rump, and Shin'ichi Oishi.

Accurate sum and dot product.

SIAM J. Sci. Comput., 26(6) :1955–1988, 2005.



Douglas M. Priest.

On Properties of Floating Point Arithmetics : Numerical Stability and the Cost of Accurate Computations.

PhD thesis, Mathematics Department, University of California, Berkeley, CA, USA, November 1992.



Jonathan Richard Shewchuk.

Adaptive precision floating-point arithmetic and fast robust geometric predicates.

Discrete Comput. Geom., 18(3) :305–363, 1997.



Françoise Tisseur.

Newton's method in floating point arithmetic and iterative refinement of generalized eigenvalue problems.

SIAM J. Matrix Anal. Appl., 22(4) :1038–1057, 2001.