

Quelle précision pour le raffinement itératif ?

Stef Graillat

LIP6 - Université Pierre et Marie Curie (Paris 6)

Rencontres arithmétique de l'informatique mathématique
22 janvier 2006, Montpellier



Raffinement itératif : améliorer la précision d'une solution approchée \hat{x} du système $Ax = b$.

① Calculer le résidu $r = b - A\hat{x}$.

② Résoudre $Ad = r$.

③ Mettre à jour $y = \hat{x} + d$.

(Recommencer à l'étape 1 sur nécessaire en remplaçant \hat{x} par y .)

→ méthode de Newton pour l'équation $F(x) = 0$ avec $F(x) = Ax - b$

Systèmes linéaire $Ax = b$ avec $A \in \mathbb{R}^{n \times n}$ inversible

On suppose que le solveur produit une solution approchée \hat{x} vérifiant

$$(A + \Delta A)\hat{x} = b, \quad |\Delta A| \leq \mathbf{u}W,$$

où W est une matrice positive dépendante de A , n et \mathbf{u} (mais pas de b)

C'est le cas pour l'élimination gaussienne avec pivotage partiel (GEPP) :

- matrice quelconque : $\mathbf{u}W = \gamma_{3n}|\hat{L}||\hat{U}|$ où \hat{L} et \hat{U} sont le résultat calculé de la décomposition LU de A
- matrice triangulaire : $\mathbf{u}W = \gamma_n|A|$

$$\gamma_n = \frac{n\mathbf{u}}{1 - n\mathbf{u}} \text{ avec } \mathbf{u} \text{ unité d'arrondi}$$

On utilise les **conditionnements** suivants :

- conditionnement du système linéaire

$$\text{cond}(A, x) := \frac{\|A^{-1}\| \|A\| \|x\|_{\infty}}{\|x\|_{\infty}}$$

- conditionnement « de la matrice »

$$\text{cond}(A) := \|A^{-1}\| \|A\|_{\infty}$$

- conditionnement de l'inversion

$$\kappa_{\infty} = \|A\|_{\infty} \|A^{-1}\|$$

On a toujours

$$\text{cond}(A, x) \leq \text{cond}(A) \leq \kappa_{\infty}.$$

- Wilkinson 1948 : premier programme de raffinement itératif
- Wilkinson 1963 : analyse pour une arithmétique en virgule fixe
- Moller 1967 : analyse pour une arithmétique en virgule flottante
- Jankowski et Woźniakowski 1977 : stabilité inverse du raffinement itératif
- Rump et Böhm 1983 : évaluation polynomiale précise par raffinement itératif
- Higham 1997 : raffinement itératif avec calcul du résidu en précision doublée
- Tisseur 2001 : algorithme de Newton en précision finie

Précision du raffinement itératif avec précision fixée

- 1 Calculer le résidu $r = b - A\hat{x}$.
- 2 Résoudre $Ad = r$.
- 3 Mettre à jour $y = \hat{x} + d$.
(Recommencer à l'étape 1 sur nécessaire en remplaçant \hat{x} par y .)

Théorème 1 (Higham)

Appliquons l'algorithme de raffinement itératif au système linéaire $Ax = b$ d'ordre n , en utilisant un solveur stable. Notons $\eta = 2n \text{cond}(A)$. Alors, en supposant que η est suffisamment plus petit que 1, le raffinement itératif réduit l'erreur directe d'un facteur proche de η à chaque étape, jusqu'à ce que $\|x - \hat{x}_i\|_\infty / \|x\|_\infty \lesssim 2n \text{cond}(A, x) \mathbf{u}$.

→ stabilité inverse du raffinement itératif

Précision du raffinement itératif avec précision mixte

- 1 Calculer le résidu $r = b - A\hat{x}$ en **précision doublée**.
- 2 Résoudre $Ad = r$.
- 3 Mettre à jour $y = \hat{x} + d$.
(Recommencer à l'étape 1 sur nécessaire en remplaçant \hat{x} par y .)

Théorème 2 (Higham)

Appliquons l'algorithme de raffinement itératif au système linéaire $Ax = b$ d'ordre n , en utilisant un solveur stable et en calculant le résidu avec une précision doublée. Notons $\eta = 2u \operatorname{cond}(A)$. alors, si η est suffisamment plus petit que 1, le raffinement itératif réduit l'erreur directe d'un facteur proche de η à chaque étape, jusqu'à ce que $\|x - \hat{x}_i\|_\infty / \|x\|_\infty \approx u$.

En fait, Higham montre que $\|x - \hat{x}_i\|_\infty / \|x\|_\infty \lesssim u + 2nu \operatorname{cond}(A, x)$.

→ résultat aussi précis que si on l'avait calculé avec une précision doublée

→ **stabilité directe** du raffinement itératif

Un bon critère d'arrêt ?

2 critères d'arrêt classiques :

- basé sur la norme du résidu $\|r\|_\infty$
- basé sur l'erreur relative direct estimée par $\|x_{i+1} - x_i\|_\infty / \|x_{i+1}\|_\infty$

Bornes d'erreur certifiées \Rightarrow évite le calcul par intervalles

Borne d'erreur absolue certifiée :

Trouver $\kappa \in \mathbb{F}$ tel que $\|x - \hat{x}_i\|_\infty \leq \kappa$

Bornes d'erreur relatives certifiées :

$$\frac{\|x - \hat{x}_i\|_\infty}{\|\hat{x}_i\|_\infty} \leq \text{fl}((\kappa / \|\hat{x}_i\|_\infty) / (1 - 2\mathbf{u}))$$

$$\frac{\|x - \hat{x}_i\|_\infty}{\|x\|_\infty} \leq \text{fl}((\kappa / (\|\hat{x}_i\|_\infty - \kappa)) / (1 - 3\mathbf{u}))$$

Algorithme 1 (Une étape de raffinement itératif)

```
function res = OneStepItRaf(A, b)
     $\hat{x} = \text{fl}(A \setminus b)$ 
     $\hat{r} = \text{fl}_e(b - A\hat{x})$     % calculé en précision doublée
     $\hat{d} = \text{fl}(A \setminus \hat{r})$ 
    res =  $\text{fl}(\hat{x} + \hat{d})$ 
```

Si $r = b - A\hat{x}$ et $d = A^{-1}r$ alors $x = \hat{x} + d$
→ transformation exacte entre le système (A, b) et (\hat{A}, r) .

→ on va approcher l'erreur $d = A^{-1}r$ par $A \setminus \hat{r}$.

$$\|\text{res} - x\|_\infty \leq \mathbf{u}\|x\|_\infty + (1 + \mathbf{u})(\mathbf{u}^2\| |A^{-1}|W|x\|_\infty + 2\gamma_{n+1}^2\| |A^{-1}||A|x\|_\infty + \mathbf{u}^2\| |A^{-1}|^2W^2|x\|_\infty) + \mathcal{O}(\mathbf{u}^3)$$

ou encore en relatif

$$\frac{\|\text{res} - x\|_\infty}{\|x\|_\infty} \leq \mathbf{u} + (1 + \mathbf{u}) \times (\mathbf{u}^2 \frac{\| |A^{-1}|W|x\|_\infty}{\|x\|_\infty} + 2\gamma_{n+1}^2 \frac{\| |A^{-1}||A|x\|_\infty}{\|x\|_\infty} + \mathbf{u}^2 \frac{\| |A^{-1}|^2W^2|x\|_\infty}{\|x\|_\infty}) + \mathcal{O}(\mathbf{u}^3)$$

avec $\mathbf{u}W = \gamma_{3n}|\widehat{L}||\widehat{U}|$ où \widehat{L} et \widehat{U} sont le résultat calculé de la décomposition LU de A

Une étape de raffinement itératif pour les systèmes triangulaire = méthode CENA de compensation pour les systèmes triangulaires

Pour les système triangulaires, $\mathbf{u}W = \gamma_n|A|$

$$\frac{\|\mathbf{res} - x\|_\infty}{\|x\|_\infty} \leq \mathbf{u} + (1 + \mathbf{u})(\mathbf{u}\gamma_n + 2\gamma_{n+1}^2 + \gamma_n^2 \text{cond}(A)) \text{cond}(A, x) + \mathcal{O}(\mathbf{u}^3)$$

→ résultat **pas** aussi précis que si on l'avait calculé avec une précision doublée à cause du $\text{cond}(A)$

$$\|\hat{x} - A^{-1}b\|_\infty \leq \frac{\|R(A\hat{x} - b)\|_\infty}{1 - \|I - RA\|_\infty} \quad (\text{Oishi, Rump})$$

Algorithme 2 (Ogita, Rump and Oishi)

Étant donné $A \in \mathbb{F}^{n \times n}$ et R un inverse approché, l'algorithme suivant calcule une borne supérieure α pour $\|RA - I\|_\infty$.

```
function  $\alpha = \text{Alpha.Std}(A, R)$   
  if  $(3n + 2)\mathbf{u} \geq 1$ , error("verification failed"), end  
   $\alpha_1 = \text{fl}(\|RA - I\|_\infty)$   
  if  $\alpha_1 \geq 1$ , error("verification failed"), end  
   $\alpha_1 = \text{fl}(\|R(|A|e)\|_\infty)$  %  $e = (1, \dots, 1)^T$   
   $\alpha = \text{fl}((\alpha_1 + \tilde{\gamma}_{3n+2}(\alpha_2 + 2))/(1 - 2\mathbf{u}))$ 
```

$$\|\hat{x} - A^{-1}b\|_{\infty} \leq \frac{\|R(A\hat{x} - b)\|_{\infty}}{1 - \|I - RA\|_{\infty}}$$

Algorithme 3 (Ogita, Rump and Oishi)

Étant donnés $A \in \mathbb{F}^{n \times n}$, $b \in \mathbb{F}^n$ et x une solution approchée de $Ax = b$ et R un inversé approché de A , l'algorithme suivant calcule une borne supérieure β pour $\|R(A\tilde{x} - b)\|_{\infty}$.

fonction $\beta = \text{Beta.Std}(A, b, \tilde{x}, R)$

$$r_{\text{mid}} = \text{fl}(A\tilde{x} - b)$$

$$r_{\text{rad}} = \text{fl}(\tilde{\gamma}_{2n+4}(|A||x| + |b|))$$

$$q = \text{fl}(|R|(\tilde{\gamma}_{n+1}|r_{\text{mid}}| + r_{\text{rad}})/(1 - (n+3)\mathbf{u}))$$

$$\beta = \text{fl}((\|Rr_{\text{mid}}\| + q)_{\infty}/(1 - 2\mathbf{u}))$$

$$\|\text{res} - x\|_\infty \leq \text{fl} \left(\left(\mathbf{u} \|\text{res}\|_\infty + \frac{\beta + \|R\|_\infty \delta}{1 - \alpha} \right) / (1 - (n + 4)\mathbf{u}) \right)$$

système $Ax = b$ avec $A \in \mathbb{F}^{10 \times 10}$ et b tel que $x = (1, \dots, 1)$ soit solution

conditionnement	erreur relative	erreur relative avec raffinement
10^8	10^{-9}	0
10^{18}	10^{-1}	0
10^{30}	1	10^{-1}
10^{40}	1	1

Raffinement itératif précis

Pour $r \in \mathbb{R}$, $\mathbb{F} \ni f = \square(r)$ est un **arrondi fidèle** de r (c-à-d l'arrondi vers $+\infty$ ou vers $-\infty$)

Algorithme 4

```
function resN = RafIt1(A, b, N)
```

```
   $\hat{x}_1 = \text{fl}(A \setminus b)$ 
```

```
  for  $i = 2 : N$ 
```

```
     $\hat{r}_i = \square(b - A(\sum_{k=1}^{i-1} \hat{x}_k))$ 
```

```
     $\hat{x}_i = \text{fl}(A \setminus \hat{r}_i)$ 
```

```
  resN =  $\square(\sum_{k=1}^N \hat{x}_k)$ 
```

\square est disponible via `accdot` de Rump, Ogita et Oishi

Important : trouver un bon inverse approché R , $\|I - RA\|_\infty \ll 1$

Si $\mathbf{u} \| |A^{-1}|W \|_\infty \leq 1/2$ alors $R = \text{inv}(A)$ convient

Théorème 3

Soit l'algorithme `RafIt1` appliqué au système linéaire régulier $Ax = b$ d'ordre n , en utilisant un solveur stable. Si $\mathbf{u} \| |A^{-1}|W \|_\infty \leq 1/2$ et $\mathbf{u} \text{cond}(A) < 1/8$ alors

$$|x - \text{res}_N| \leq 2\mathbf{u}|x| + (1 + 2\mathbf{u})G^{N-1}|\hat{x}_1 - x|$$

avec $\|G\|_\infty < 1$. Par conséquent,

$$\frac{\|x - \text{res}_N\|_\infty}{\|x\|_\infty} \leq 2\mathbf{u} + (1 + 2\mathbf{u})\|G\|_\infty^{N-1} \frac{\|\hat{x}_1 - x\|_\infty}{\|x\|_\infty}$$

Si $\text{inv}(A)$ ne convient pas, alors AccInv de Rump donne une expansion de A^{-1} comme somme de matrices flottantes

Algorithme 5

```
function resN = RafIt2(A, b, N)
    R = AccInv(A)           % R = R1 + ... + Rm avec ||I - RA||∞ < 1/2
    x̂1 = □(Rb)
    for i = 2 : N
        r̂i = □(b - A(∑k=1i-1 x̂k))
        x̂i = □(Rr̂i)
    end
    resN = □(∑k=1N x̂k)
```

Théorème 4

Soit l'algorithme RafIt2 appliqué au système linéaire régulier $Ax = b$ d'ordre n . Alors

$$|x - \text{res}_N| \leq 2\mathbf{u}|x| + (1 + 2\mathbf{u})G^{N-1}|\hat{x}_1 - x|$$

avec $\|G\|_\infty < 1$. Par conséquent,

$$\frac{\|x - \text{res}_N\|_\infty}{\|x\|_\infty} \leq 2\mathbf{u} + (1 + 2\mathbf{u})\|G\|_\infty^{N-1} \frac{\|\hat{x}_1 - x\|_\infty}{\|x\|_\infty}$$

$$\|\hat{x} - A^{-1}b\|_\infty \leq \frac{\|R(A\hat{x} - b)\|_\infty}{1 - \|I - RA\|_\infty}$$

ou

$$|\hat{x} - A^{-1}b| \leq |R(A\hat{x} - b)| + \frac{\|R(A\hat{x} - b)\|_\infty}{1 - \|G\|_\infty} \quad (\text{Yamamoto})$$

avec $G = RA - I$, $\|G\|_\infty < 1$ et

$$\left(\sum_{j=1}^n |G_{1,j}|, \dots, \sum_{j=1}^n |G_{n,j}| \right)^T \leq t$$

Conclusion :

- raffinement itératif = méthode de Newton = algorithme compensé
- raffinement itératif précis \rightarrow « pseudo-expansion » du résultat

Perspectives :

- Simulations numériques pour la précision et les performances
- Étude du raffinement itératif précis (c-à-d méthode de Newton précise) pour le calcul de zéros de polynômes (univariés plus multivariés) et de valeurs propres (généralisées, etc.)