

Επιστημονικός Υπολογισμός

Ε.Γαλλόπουλος

ΤΜΗΥΠ, Π. Πατρών

Διάλεξη 7: 15 Νοεμβρίου 2017

- 1 Μελέτη σφαλμάτων (στην α.κ.υ).
- 2 Πίσω ανάλυση σφάλματος
- 3 Δείκτες κατάστασης
- 4 Παραδείγματα πίσω ανάλυσης σφάλματος
- 5 Εντολή FMA
- 6 Το πρόβλημα της άθροισης
- 7 Προς μια ακριβέστερη αριθμητική

8 Υλοποιήσεις

Ας θεωρήσουμε ότι το πρόβλημα αντιστοιχεί στον υπολογισμό της απεικόνισης

$f: \mathcal{U} \subset \mathbb{R}^m \rightarrow \mathbb{R}^n$ για ορισμένα δεδομένα.

Προσοχή Αν τρέξουμε τον αλγόριθμο διακρίνουμε:

$x \in \mathcal{U}$ τα m στοιχεία στο πεδίο ορισμού της f ;

$f(x)$ τα n στοιχεία της τιμής της συνάρτησης στο x , χωρίς λάθη υπολογισμών;

x^* τα στοιχεία εισόδου που χρησιμοποιούνται στην υλοποίηση, οπότε $x^* \in F$. Αν τα μόνα λάθη που υπάρχουν στο x^* είναι αυτά που προέρχονται από τη στρογγύλευση του x , τότε $x^* = \text{fl}(x)$.

$f(x^*)$ η τιμή της $f(x^*)$, υπολογισμένη χωρίς σφάλματα (δηλ. με την αριθμητική του \mathbb{R}).

$f_{\text{prog}}(x^*)$ η υλοποίηση του υπολογισμού της $f(x^*)$ με πεπερασμένο πλήθος διακριτών πράξεων α.κ.υ. στο \mathcal{F} .

Να εκτιμήσουμε το **απόλυτο** ή **σχετικό σφάλμα**

$$\|f_{\text{prog}}(x^*) - f(x)\| \text{ ή αν } f(x) \neq 0 \frac{\|f_{\text{prog}}(x^*) - f(x)\|}{\|f(x)\|}.$$

Αν το σφάλμα είναι μικρό (π.χ. της τάξης του ϵ της μηχανής) ο αλγόριθμος, f , θεωρείται ακριβής.

Εμπρός Σφάλμα (forward error).

Το παραπάνω αποκαλείται (απόλυτο ή σχετικό) εμπρός^α σφάλμα.

^αΟ χαρακτηρισμός ((εμπρός)) θα εξηγηθεί στη συνέχεια των διαλέξεων.

ΠΡΟΣΟΧΗ Γενικά είναι ανέφικτο να υπολογιστεί. **Στόχος είναι η εκτίμησή του, π.χ. βρίσκοντας σφικτά φράγματα.**

- Το $\| \cdot \|$ συμβολίζει κάποια μετρική - π.χ.
- ... απόλυτη τιμή, στους βαθμωτούς,
- ... μία από τις γνωστές νόρμες, για διανύσματα και μητρώα
- ... πίνακας απολύτων τιμών (απόσταση κατά συνιστώσες).
- Οι έννοιες του απόλυτου και σχετικού σφάλματος είναι γενικές και ξεπερνούν την α.κ.υ.
- Αν το x^* χρησιμοποιείται ως προσέγγιση του x , τότε

$$\text{ΑΣ: } \|x^* - x\| \text{ ΣΣ αν } x \neq 0 \quad \frac{\|x^* - x\|}{\|x\|}$$

- Μερικές φορές θέτουμε ΣΣ το $\frac{\|x^* - x\|}{\|x^*\|}$
- Αν $x = 0$ τότε αυτό το ΣΣ είναι πάντα 1 όσο καλή και να είναι η προσέγγιση x^* !!! (... το άλλο ΣΣ δεν ορίζεται)
- Για προσεγγίσεις πολύ μικρών τιμών χρησιμοποιείται το απόλυτο σφάλμα.
- Χονδρικά, αν $\frac{|x^* - x|}{|x|} \leq 5 \cdot 10^{-d}$ τότε λέμε ότι οι (αριθμοί) x^* και x συμφωνούν σε d δεκαδικά ψηφία.

Από το ϵ_M στο u – Μονάδα στρογγύλευσης

- Στη MATLAB για οποιοδήποτε α.κ.υ. x , η εντολή $\text{eps}(x)$ επιστρέφει την απόστασή του από τον διαδοχικό του. Στη βιβλιογραφία αυτό αναφέρεται και ως $\text{ulp}(x)$ (units in the last place).
- Αν χρησιμοποιούμε t bits για την αναπαράσταση της ουράς (το 1ο κρυμμένο) και $x = m \times \beta^e$ τότε η απόσταση του x από τον αμέσως επόμενο α.κ.υ. είναι $\text{ulp}(x) = (m + 2^{-(t-1)})2^e - m2^e = 2^{e+1-t}$.
- Η μέγιστη σχετική απόσταση είναι $\frac{2^{e+1-t}}{2^e} = 2^{1-t}$.
- ... προφανώς ίση με ϵ_M .

Θα εκφράζουμε τα σφάλματα ως πολλαπλάσια της **μονάδας στρογγύλευσης** που θα χρησιμεύσει ως μονάδα μέτρησης των σφαλμάτων.

Μονάδα στρογγύλευσης (unit roundoff)

Η μονάδα στρογγύλευσης είναι το μέγιστο δυνατό σχετικό σφάλμα για τον επιλεγμένο τρόπο στρογγύλευσης. Για στρογγύλευση προς το πλησιέστερο,

$$u = \max_{z \neq 0} \frac{|z - \text{fl}(z)|}{|z|} = \frac{2^{1-t}}{2}.$$

IEEE double: $u = 1.1102e - 016$; IEEE single: $u = 5.9605e - 008$.

Αρχή ακριβούς στρογγύλευσης (υπενθύμιση)

Αν $\tilde{\odot}$ είναι η υλοποίηση της αριθμητικής πράξης \odot , τότε αν $x, y \in F$ ισχύει ότι $x\tilde{\odot}y = \mathbf{fl}(x \odot y) \in F$.

Το υπολογισμένο αποτέλεσμα είναι ακριβώς ίδιο με το να εκτελούνταν η πράξη με ((θεϊκή)) αριθμητική και μετά να εφαρμοζόταν στρογγύλευση.

ΕΠΟΜΕΝΩΣ

$$|x\tilde{\odot}y - x \odot y| = |\mathbf{fl}(x \odot y) - x \odot y| \leq \mathbf{u}|x \odot y|$$

άρα

$$-\mathbf{u}(x \odot y) \leq \mathbf{fl}(x \odot y) - x \odot y \leq \mathbf{u}(x \odot y)$$

$$(1 - \mathbf{u})(x \odot y) \leq x\tilde{\odot}y \leq (1 + \mathbf{u})(x \odot y)$$

Μοντέλο διάδοσης

Μετά από κάθε αριθμητική πράξη \odot επί δεδομένων α.κ.υ. $x, y \in F$ και δεν υπάρχει υπερ- ή υποχείλιση, ισχύει:

$$\text{fl}(x \odot y) = (1 + \delta)(x \odot y), \text{ για } \delta \text{ τ.ώ. } |\delta| \leq \mathbf{u}$$

Σχετικά με το δ :

- φράσσεται σε απόλυτη τιμή! Επομένως $-\mathbf{u} \leq \delta \leq \mathbf{u}$ (με πιο ενδελεχή ανάλυση $|\delta| < \mathbf{u}$.)
- διαφέρει για κάθε πράξη και στοιχεία x, y
- δεν απαιτείται να το γνωρίζουμε ακριβώς - χρησιμοποιούμε το μοντέλο γνωρίζοντας μόνον πώς φράσσεται το δ !
- το μοντέλο δεν προβλέπει πότε $\delta = 0$ (αυτή η πληροφορία μπορεί να είναι χρήσιμη, αλλά χάνεται).

Ιδιότητες αριθμητικών πράξεων στο πεδίο των πραγματικών \mathbb{R}

Αξιώματα πρόσθεσης

- A0 Το άθροισμα $x + y \in \mathbb{R}$.
- A1 Η πρόσθεση είναι αντιμεταθετική: $x + y = y + x$.
- A2 Η πρόσθεση είναι προσεταιριστική
 $x + (y + z) = (x + y) + z$.
- A3 Υπάρχει στοιχείο 0 ώστε $x + 0 = x$ για κάθε $x \in \mathbb{R}$.
- A4 Για κάθε $x \in \mathbb{R}$ υπάρχει αντίστροφο στοιχείο $-x \in \mathbb{R}$ ως προς την πρόσθεση, δηλ. $x + (-x) = 0$.

Αξιώματα πολλαπλασιασμού

- P0 Το γινόμενο $x \times y \in \mathbb{R}$.
- P1 Αντιμεταθετική ιδιότητα πολλαπλασιασμού: $x \times y = y \times x$.
- P2 Προσεταιριστική ιδιότητα πολλαπλασιασμού:
 $x \times (y \times z) = (x \times y) \times z$.
- P3 Υπάρχει στοιχείο 1 ώστε $x \times 1 = x$ για κάθε $x \in \mathbb{R}$.
- P4 Για κάθε μη μηδενικό x υπάρχει αντίστροφο ως προς τον πολλαπλασιασμό $\frac{1}{x} \in \mathbb{R}$ ώστε $x \times (\frac{1}{x}) = 1$.

Ε Επιμεριστική ιδιότητα: $x \times (y + z) = x \times y + x \times z$.

Παράδειγμα : Δεν ισχύει πάντα το A2

$$t_1 = \text{fl}(x + y) \quad s_1 = \text{fl}(y + z)$$

$$t_2 = \text{fl}(t_1 + z) \quad s_2 = \text{fl}(x + s_1)$$

και δεν υπάρχει λόγος να ισχύει πάντα (αν και πολλές φορές ισχύει!)

$$\text{fl}(\text{fl}(x + y) + z) = \text{fl}(x + \text{fl}(y + z)).$$

Παράδειγμα : Δεν ισχύει πάντα το A2

Listing 1: Παράδειγμα

```
(1+eps/2)+eps/2 % = 1  
1+(eps/2+eps/2) % = 1.0000000000000000  
(-10^20+10^20)+1 % = 1  
-10^20+(10^20+1) % = 0
```

Listing 2: Παράδειγμα

```
a = 0.1234567000000000;  
b = 4.711325195312500e+008; c = -b;  
(a+b)+c % = 0.123456716537476  
a+(b+c) % = 0.1234567000000000
```

Παράδειγμα : Δεν ισχύει πάντα το Π4

Γενικά για $x \in F$, $\text{fl}(x \cdot \text{fl}(\frac{1}{x})) \neq 1$.

Listing 3: Παράδειγμα

```
index = []; for i=1:170
    if ((1/i)*i ≈ 1)
        index = [index i];
    end;
end;
index
    49      98     103     107     161
```

Προειδοποιήσεις: Προσοχή για τους συγγραφείς μεταφραστών!

Επειδή στο \mathcal{F} δεν ισχύουν όλες οι ιδιότητες πεδίου, δεν μπορούμε να κάνουμε τις ίδιες απλοποιήσεις και μετατροπές των αριθμητικών εκφράσεων που επιτρέπονται στο \mathbb{R} .

Π.χ. αν με δύο προσθετές που υπολογίζουν ταυτόχρονα το άθροισμα των α.κ.υ. x_1, x_2 και των x_3, x_4 το άθροισμα $x_1 + x_2 + x_3 + x_4$ μπορεί να υπολογισθεί ταχύτερα από το κλασικό $((x_1 \dot{+} x_2) \dot{+} x_3) \dot{+} x_4$ χρησιμοποιώντας $(x_1 \dot{+} x_2) \dot{+} (x_3 \dot{+} x_4)$.

Το σφάλμα και τα αποτελέσματα μπορεί να είναι διαφορετικά!

Προειδοποίηση Οι επεξεργαστές Intel εκτελούν τις πράξεις σε εκτεταμένη διπλή ακρίβεια (80 bits) ... Αυτό ορισμένες φορές οδηγεί σε απρόσμενα αποτελέσματα (double rounding).

Προσοχή Όταν αφαιρούνται δύο αριθμοί που είναι σχεδόν ίσοι και περιέχουν μικρά σφάλματα “αναδύονται σκουπίδια” ακόμα και αν χρησιμοποιούσαμε αριθμητική άπειρης ακρίβειας!

Αν $A = A_1 + \text{θόρυβος}$, $A' = A_2 + \text{θόρυβος}$ και οι τιμές A_1, A_2 είναι σχεδόν ίσες και θέλουμε να υπολογίσουμε το (μικρό) $A_1 - A_2$ αφαιρώντας $A - A'$, τότε:

$$(A_1 + \text{θόρυβος}) - (A_2 + \text{θόρυβος}) = (A_1 - A_2) + \text{θόρυβος}$$

- **καταστροφική απαλοιφή** όταν $|A_1 - A_2| = O(\text{θόρυβος})$ ή λιγότερο.
- Η μόλυνση από τα σκουπίδια επηρεάζει και έχει καταστροφικά αποτελέσματα αν, για παράδειγμα, τα ((ασήμαντα σκουπίδια)) πολλαπλασιαστούν με μεγάλους αριθμούς και χρησιμοποιηθούν περαιτέρω.
- Διαβάστε την ιστορία Catastrophic cancellation in the high seas της A. Langville (Lan01).

Στρογγύλευση: Από $x \in \mathbb{R}$ στο $\mathbf{fl}(x) \in F$

Αν $x \in G$ τότε μπορούμε να γράψουμε $\mathbf{fl}(x) = x(1 + \delta)$ για κάποιο $|\delta| \leq \mathbf{u}$.

Σφάλμα πράξεων

Έστω ότι $x, y \in F$ και $x \odot y \in G$ και ότι $\odot \in \{\pm, \times, /, \sqrt{\cdot}\}$. Τότε

$$\mathbf{fl}(x \odot y) = (x \odot y)(1 + \delta), \text{ για κάποιο } |\delta| \leq \mathbf{u}$$

και

$$\mathbf{fl}(x \odot y) = \frac{x \odot y}{1 + \delta}, \text{ για κάποιο } |\delta| \leq \mathbf{u}$$

Πώς δικαιολογείται σε α.κ.υ. (IEEE 64 bits, $\mathbf{u} = 2^{-53} \approx 10^{-16}$) ότι $10^{20} \tilde{+} 10 = 10^{20}$
(θυμηθείτε παράδειγμα εισαγωγής)

$$10^{20} \tilde{+} 10 = (10^{20} + 10)(1 + \delta), \quad |\delta| \leq \mathbf{u},$$

άρα για να ισχύει είναι ισοδύναμο να υπάρχει αποδεκτό δ για οποίο να ισχύει η ισότητα

$$10^{20} = (10^{20} + 10)(1 + \delta)$$

Λύνοντας ως προς δ ,

$$\begin{aligned} 10^{20} &= (10^{20} + 10)(1 + \delta), \\ 10^{20} &= 10^{20} + 10 + 10^{20}\delta + 10\delta \\ \delta &= \frac{-10}{10^{20} + 10} \approx -10^{-19} \end{aligned}$$

επομένως το υπολογισμένο αποτέλεσμα εξηγείται καθώς το παραπάνω $|\delta| \leq \mathbf{u} \approx 10^{-16}$.

((Προς τα εμπρός ανάλυση)) και εκτίμηση σφάλματος

Παράδειγμα

Θεωρούμε ότι οι μεταβλητές x_j περιέχουν α.κ.υ.

$$\begin{aligned}(x_1 \tilde{+} x_2) \tilde{+} x_3 &= ((x_1 + x_2)(1 + \delta_1) + x_3)(1 + \delta_2) \\ &= x_1(1 + \delta_1)(1 + \delta_2) + x_2(1 + \delta_1)(1 + \delta_2) + x_3(1 + \delta_2)\end{aligned}$$

επομένως αν θέσουμε

$$E := (x_1 \tilde{+} x_2) \tilde{+} x_3 - (x_1 + x_2 + x_3)$$

$$E = (x_1 + x_2)(\delta_1 + \delta_2 + \delta_1\delta_2) + x_3\delta_2$$

επομένως μπορούμε να φράξουμε ως εξής:

$$|E| \leq (|x_1| + |x_2|)(2\mathbf{u} + \mathbf{u}^2) + |x_3|\mathbf{u}.$$

Με τον ίδιο τρόπο καταλήγουμε και σε άνω φράγμα για το σφάλμα στον υπολογισμό του $x_1 + (x_2 + x_3)$. Θέτουμε για συντομία $\hat{E} := x_1 \tilde{+} (x_2 \tilde{+} x_3) - (x_1 + x_2 + x_3)$. Εντέλει έχουμε τα παρακάτω άνω φράγματα για τους δύο εναλλακτικούς τρόπους υπολογισμού:

$$\begin{aligned} |E| &\leq (|x_1| + |x_2|)(2\mathbf{u} + \mathbf{u}^2) + |x_3|\mathbf{u} \\ |\hat{E}| &\leq |x_1|\mathbf{u} + (2\mathbf{u} + \mathbf{u}^2)(|x_2| + |x_3|). \end{aligned}$$

- Αναδεικνύεται λεπτομερώς η (διαφορετική) συνεισφορά του κάθε όρου στο άνω φράγμα.
- Διαφαίνεται ότι με βάση τα x_i , θα μπορούσαμε να επιλέξουμε σειρά άθροισης που να ελαχιστοποιεί το άνω φράγμα (αλλά όχι κατ' ανάγκη το σφάλμα!).

Πώς (με ποια σειρά) αθροίζουμε 3 α.κ.υ? Η παρακάτω συζήτηση θα γίνει αποκλειστικά με βάση τα παραπάνω. Περισσότερα σε επόμενη διάλεξη αφιερωμένη στην άθροιση.

Υπάρχουν **3 μη ισοδύναμοι αριθμητικά** τρόποι άθροισης (**δεν ισχύει** προσεταιριστικότητα):

$$(x_1 + x_2) + x_3, x_1 + (x_2 + x_3), (x_1 + x_3) + x_2$$

Οι υπόλοιποι (9) τρόποι είναι αριθμητικά ισοδύναμοι με έναν απο τους παραπάνω (γιατί **ισχύει** αντιμεταθετικότητα), π.χ.

$$\begin{aligned} & x_3 \dot{+} (x_1 \dot{+} x_2), x_3 \dot{+} (x_2 \dot{+} x_1), (x_2 \dot{+} x_3) \dot{+} x_1, \\ & (x_3 \dot{+} x_2) \dot{+} x_1, x_2 \dot{+} (x_1 \dot{+} x_3), x_2 \dot{+} (x_3 \dot{+} x_1), \\ & x_1 \dot{+} (x_3 \dot{+} x_2), (x_3 \dot{+} x_1) \dot{+} x_2, (x_2 \dot{+} x_1) \dot{+} x_3, \end{aligned}$$

Με ποιόν τρόπο προκύπτει το μικρότερο άνω φράγμα? Με βάση τα παραπάνω, αυτός που αφήνει το μεγαλύτερο στοιχείο τελευταίο (εκτός παρένθεσης):

- Αν $|x_3| = \max(|x_1|, |x_2|, |x_3|)$ τότε $\max_{x_j \in F} |E| \leq \max_{x_j \in F} |\hat{E}|$
- Αν $|x_1| = \max(|x_1|, |x_2|, |x_3|)$ τότε $\max_{x_j \in F} |\hat{E}| \leq \max_{x_j \in F} |E|$
- Αν $|x_2| = \max(|x_1|, |x_2|, |x_3|)$ τότε επιλέγουμε $(x_1 \tilde{+} x_3) \tilde{+} x_2$.
- Όμως: Το άνω φράγμα δείχνει την χειρότερη περίπτωση!
- Τα πράγματα μπορεί να είναι πολύ καλύτερα!!
- ... για παράδειγμα αν γνωρίζουμε ότι $x_1 = -x_2$, ή/και το αντίστοιχο δ να είναι 0.
- Αυτά δεν προβλέπονται αν εφαρμόσουμε το μοντέλο χωρίς ειδικές τροποποιήσεις.

$$\begin{aligned} |E| &\leq (|x_1| + |x_2|)(2\mathbf{u} + \mathbf{u}^2) + |x_3|\mathbf{u} \\ &\leq (|x_1| + |x_2| + |x_3|)2\mathbf{u} + \underbrace{(|x_1| + |x_2|)\mathbf{u}^2}_{\text{αμελητέο}} \end{aligned}$$

αν αγνοήσουμε όρους με παράγοντα \mathbf{u}^2 (δηλ. όρους 2ης τάξης) μπορούμε να φράξουμε αμφότερους όρους $||, |\hat{E}|$

$$||, |\hat{E}| \leq (|x_1| + |x_2| + |x_3|)2\mathbf{u}.$$

ΠΡΟΣΟΧΗ Από τα παραπάνω και μόνον το σχετικό σφάλμα φράσσεται ως εξής (ομοίως και για το $|\hat{E}|$):

$$\frac{|E|}{|x_1 + x_2 + x_3|} \leq \frac{|x_1| + |x_2| + |x_3|}{|x_1 + x_2 + x_3|} 2\mathbf{u}$$

ΕΡΩΤΗΣΗ: Ποιο είναι το μειονέκτημα αυτού του φράγματος;

Ποιο είναι το μειονέκτημα αυτού του φράγματος;

- Το φράγμα εξαρτάται από τις τιμές των $x_1, x_2, x_3 \dots$
- Ο όρος $\frac{|x_1|+|x_2|+|x_3|}{|x_1+x_2+x_3|}$ μπορεί να γίνει πολύ μεγάλος.

Αν υπολογίσουμε (μη μηδενικό) $x_1 \tilde{\times} x_2 \tilde{\times} x_3$, τότε

$$\begin{aligned}\frac{|(x_1 \tilde{\times} x_2) \tilde{\times} x_3 - x_1 x_2 x_3|}{|x_1 x_2 x_3|} &= \frac{|(x_1 x_2)(1 + \delta_1)x_3(1 + \delta_2) - x_1 x_2 x_3|}{|x_1 x_2 x_3|} \\ &= \frac{|x_1 x_2 x_3(\delta_1 + \delta_2 + \delta_1 \delta_2)|}{|x_1 x_2 x_3|} \\ \frac{|(x_1 \tilde{\times} x_2) \tilde{\times} x_3 - x_1 x_2 x_3|}{|x_1 x_2 x_3|} &\leq 2\mathbf{u} + \mathbf{u}^2\end{aligned}$$

Όταν οι όροι x_j είναι ομόσημοι, $|x_1| + |x_2| + |x_3| = |x_1 + x_2 + x_3|$, επομένως

$$\frac{|E|}{|x_1 + x_2 + x_3|} \leq 2\mathbf{u}.$$

Παραδείγματα:

- άθροισμα ομόσημων όρων (π.χ. μη αρνητικών στατιστικών μεγεθών, μετρήσεων, ...)
- υπολογισμός νόρμας πραγματικών διανυσμάτων

Εύρεση άνω φράγματος του εμπρός σφάλματος για τον υπολογισμό των

- ... $(x_1 + x_2) + x_3$ όταν $x_i \in G$ και όχι κατ' ανάγκη α.κ.υ.
- ... αθροίσματος n αριθμών (για εναλλακτικούς τρόπους άθροισης)
- ... εσωτερικού γινομένου διανυσμάτων ...

Π.χ. στην πρώτη περίπτωση

$$\begin{aligned}\text{fl}((x_1 + x_2) + x_3) &= \text{fl}(\text{fl}(\text{fl}(x_1) + \text{fl}(x_2)) + \text{fl}(x_3)) \\ &= ((x_1(1 + \delta_1) + x_2(1 + \delta_2))(1 + \delta_3) + (1 + \delta_4)x_3)(1 + \delta_5)\end{aligned}$$

Το παραπάνω μπορεί να ξαναγραφτεί ως :

$$x_1(1 + \delta_1)(1 + \delta_3)(1 + \delta_5) + x_2(1 + \delta_2)(1 + \delta_3)(1 + \delta_5) + x_3(1 + \delta_4)(1 + \delta_5)$$

Οι εκφράσεις γίνονται πολύπλοκες ακόμα και στην απλή αυτή μελέτη!

Παράδειγμα - το γνωστό κουίζ

$$\begin{aligned}10^{20} - 10 - 10^{20} + 20 &= 20 \\10^{20} + 20 - 10^{20} - 10 &= -10 \\-10 + 20 - 10^{20} + 10^{20} &= 0 \\10^{20} - 10^{20} + 20 - 10 &= 10\end{aligned}$$

Για παράδειγμα, $10^{20} - 10 - 10^{20} + 20 = 20$?

Προσοχή: Θεωρούμε ότι οι πράξεις εκτελούνται από αριστερά προς τα δεξιά.

$$\begin{aligned}&(((10^{20} - 10)(1 + \delta_1) - 10^{20})(1 + \delta_2) + 20)(1 + \delta_3) \\10^{20}(1 + \delta_1)(1 + \delta_2)(1 + \delta_3) - 10(1 + \delta_1)(1 + \delta_2)(1 + \delta_3) - 10^{20}(1 + \delta_2)(1 + \delta_3) + 20(1 + \delta_3)\end{aligned}$$

Παράδειγμα - το γνωστό κουίζ

$$\begin{aligned}10^{20} - 10 - 10^{20} + 20 &= 20 \\10^{20} + 20 - 10^{20} - 10 &= -10 \\-10 + 20 - 10^{20} + 10^{20} &= 0 \\10^{20} - 10^{20} + 20 - 10 &= 10\end{aligned}$$

Για παράδειγμα, $10^{20} - 10 - 10^{20} + 20 = 20$?

Προσοχή: Θεωρούμε ότι οι πράξεις εκτελούνται από αριστερά προς τα δεξιά.

$$\begin{aligned}&(((10^{20} - 10)(1 + \delta_1) - 10^{20})(1 + \delta_2) + 20)(1 + \delta_3) \\10^{20}(1 + \delta_1)(1 + \delta_2)(1 + \delta_3) - 10(1 + \delta_1)(1 + \delta_2)(1 + \delta_3) - 10^{20}(1 + \delta_2)(1 + \delta_3) + 20(1 + \delta_3)\end{aligned}$$

Ερώτημα: Είναι ίσο με 20 για επιτρεπτές τιμές των δ ?

Παράδειγμα - το γνωστό κουίζ

$$\begin{aligned}10^{20} - 10 - 10^{20} + 20 &= 20 \\10^{20} + 20 - 10^{20} - 10 &= -10 \\-10 + 20 - 10^{20} + 10^{20} &= 0 \\10^{20} - 10^{20} + 20 - 10 &= 10\end{aligned}$$

Για παράδειγμα, $10^{20} - 10 - 10^{20} + 20 = 20$?

Προσοχή: Θεωρούμε ότι οι πράξεις εκτελούνται από αριστερά προς τα δεξιά.

$$\begin{aligned}&(((10^{20} - 10)(1 + \delta_1) - 10^{20})(1 + \delta_2) + 20)(1 + \delta_3) \\10^{20}(1 + \delta_1)(1 + \delta_2)(1 + \delta_3) - 10(1 + \delta_1)(1 + \delta_2)(1 + \delta_3) - 10^{20}(1 + \delta_2)(1 + \delta_3) + 20(1 + \delta_3)\end{aligned}$$

Ερώτημα: Είναι ίσο με 20 για επιτρεπτές τιμές των δ ?

Αν $\delta_1 = 1/(10^{19} - 1)$, $\delta_2 = \delta_3 = 0$, τότε ισχύει ισότητα.

Αφού $|\delta_j| \leq \mathbf{u}$, οι τιμές είναι επιτρεπτές, άρα το αποτέλεσμα είναι αιτιολογημένο.

Για την απλοποίηση όρων όπως $p_n = \prod_{i=1}^n (1 + \delta_i)$ όταν γνωρίζουμε ότι $|\delta_i| \leq u$. Αμέσως βλέπουμε ότι

$$(1 - u)^n \leq p_n \leq (1 + u)^n.$$

και ότι $p_n = 1 + nu + O(u^2)$.

Ακόμα καλύτερα:

Λήμμα

Αν $|\delta_i| \leq u$ και $\rho_i = \pm 1$ για $i = 1 : n$ και $nu < 1$ τότε υπάρχει κάποια τιμή θ_n τ.ώ.

$$\prod_{i=1}^n (1 + \delta_i)^{\rho_i} = 1 + \theta_n,$$

όπου

$$|\theta_n| \leq \frac{nu}{1 - nu} := \gamma_n.$$

Εύρεση άνω φράγματος του εμπρός σφάλματος για τον υπολογισμό του

- $(x_1 + x_2) + x_3$ όταν $x_j \in \mathcal{G}$ και όχι κατ' ανάγκη α.κ.υ.

Με βάση το Λήμμα, υπάρχουν $\theta_3, \zeta_3, \theta_2$ τέτοια ώστε

$$|\theta_3| \leq \gamma_3, |\zeta_3| \leq \gamma_3, |\theta_2| \leq \gamma_2$$

και μπορούμε να γράψουμε

$$\mathbf{fl}((x_1 + x_2) + x_3) = x_1(1 + \theta_3) + x_2(1 + \zeta_3) + x_3(1 + \theta_2)$$

επομένως ένα (χαλαρό) άνω φράγμα για το σφάλμα θα είναι

$$\begin{aligned} |\mathbf{fl}((x_1 + x_2) + x_3) - (x_1 + x_2 + x_3)| &\leq |x_1|\gamma_3 + |x_2|\gamma_3 + |x_3|\gamma_2 \\ &\leq (|x_1| + |x_2| + |x_3|)\gamma_3. \end{aligned}$$

Το σχετικό σφάλμα φράσσεται άμεσα από

$$\frac{|\mathbf{fl}((x_1 + x_2) + x_3) - (x_1 + x_2 + x_3)|}{|x_1 + x_2 + x_3|} \leq \frac{(|x_1| + |x_2| + |x_3|)}{|x_1 + x_2 + x_3|} \gamma_3.$$

Για ομόσημα το άνω φράγμα είναι $\gamma_3 = \gamma_3 = \frac{3\mathbf{u}}{1-3\mathbf{u}}$.

Κριτική

Εκκινώντας από τα στοιχεία εισόδου και παρακολουθώντας το σφάλμα σε κάθε πράξη προσπαθούμε να φράξουμε το μέγιστο απόλυτο ή σχετικό σφάλμα που θα μπορούσε να προκύψει στο τελικό αποτέλεσμα.

Η ιδέα είναι απλή

- ... η εφαρμογή της μπορεί να είναι περίπλοκη
- ... σκληρή άσκηση σε ανισότητες
- ... τεράστιες εκφράσεις, κ.λπ.

Έχουν γίνει πολλές προσπάθειες αυτοματοποίησης της ανάλυσης των σφαλμάτων που υπεισέρχονται στις υπολογιστικές διαδικασίες με μέτρια ή μεγαλύτερη επιτυχία.

Εμπρός ανάλυση σφάλματος: Θεωρούμε τον αλγόριθμο ως μιά σειρά **στοιχειωδών πράξεων**.

Σε κάθε βήμα, υπολογίζεται μία τιμή α_{k+1} με βάση προηγούμενες τιμές και στοιχεία εισόδου, π.χ. $\alpha_{k+1} = g_k(\alpha_1, \dots, \alpha_k)$. Μερικές από τις τιμές μπορεί να είναι δεδομένα εισόδου. Στη συνέχεια, υπολογίζουμε φράγματα για τα σφάλματα στα τελικά αποτελέσματα.

Έχουν γίνει πολλές προσπάθειες αυτοματοποίησης της ανάλυσης των σφαλμάτων που υπεισέρχονται στις υπολογιστικές διαδικασίες με μέτρια ή μεγαλύτερη επιτυχία.

Εμπρός ανάλυση σφάλματος: Θεωρούμε τον αλγόριθμο ως μία σειρά **στοιχειωδών πράξεων**.

Σε κάθε βήμα, υπολογίζεται μία τιμή α_{k+1} με βάση προηγούμενες τιμές και στοιχεία εισόδου, π.χ. $\alpha_{k+1} = g_k(\alpha_1, \dots, \alpha_k)$. Μερικές από τις τιμές μπορεί να είναι δεδομένα εισόδου. Στη συνέχεια, υπολογίζουμε φράγματα για τα σφάλματα στα τελικά αποτελέσματα.

Ανάλυση διαστημάτων: Θεωρούμε ότι κάθε δεδομένο x εγκλείεται σε κάποιο διάστημα $[x_L, x_U]$ οπότε οι πράξεις επί των δεδομένων εκτελούνται χρησιμοποιώντας ((αριθμητική διαστημάτων)) (interval arithmetic).

Έχουν γίνει πολλές προσπάθειες αυτοματοποίησης της ανάλυσης των σφαλμάτων που υπεισέρχονται στις υπολογιστικές διαδικασίες με μέτρια ή μεγαλύτερη επιτυχία.

Εμπρός ανάλυση σφάλματος: Θεωρούμε τον αλγόριθμο ως μιά σειρά **στοιχειωδών πράξεων**.

Σε κάθε βήμα, υπολογίζεται μία τιμή α_{k+1} με βάση προηγούμενες τιμές και στοιχεία εισόδου, π.χ. $\alpha_{k+1} = g_k(\alpha_1, \dots, \alpha_k)$. Μερικές από τις τιμές μπορεί να είναι δεδομένα εισόδου. Στη συνέχεια, υπολογίζουμε φράγματα για τα σφάλματα στα τελικά αποτελέσματα.

Ανάλυση διαστημάτων: Θεωρούμε ότι κάθε δεδομένο x εγκλείεται σε κάποιο διάστημα $[x_L, x_U]$ οπότε οι πράξεις επί των δεδομένων εκτελούνται χρησιμοποιώντας ((αριθμητική διαστημάτων)) (interval arithmetic).

Πίσω ανάλυση σφάλματος: Στη συνέχεια!

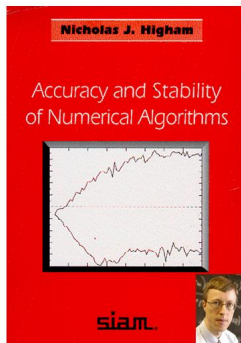
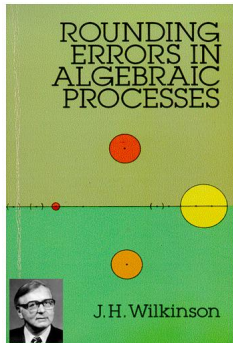
Η ιδέα της εμπρός ανάλυσης είναι απλή, αλλά

- ... η εφαρμογή της μπορεί να είναι περίπλοκη
- ... σκληρή άσκηση σε ανισότητες
- ... τεράστιες εκφράσεις, κ.λπ.
- ΠΡΟΣΟΧΗ ... δεν λέει κάτι ξεκάθαρο για τον αλγόριθμο.

Η ιδέα της εμπρός ανάλυσης είναι απλή, αλλά

- ... η εφαρμογή της μπορεί να είναι περίπλοκη
- ... σκληρή άσκηση σε ανισότητες
- ... τεράστιες εκφράσεις, κ.λπ.
- ΠΡΟΣΟΧΗ ... δεν λέει κάτι ξεκάθαρο για τον αλγόριθμο.

Vel Kahan: "In the 1950's almost no programmers knew how to distinguish numerical instability from ill condition. In other words, the distinction between an algorithm's hypersensitivity to its own internal rounding errors, and a true solution's hypersensitivity to roundoff-like errors in its problem's data, as contributors to wrongly computed solutions, was obscure. By the 1960s the distinction had been clarified through "Backward Error Analyses""



J.H. Wilkinson: "It is a matter of some surprise that one of the simplest methods of solution generally leads to an error, the expected value of which is precisely that resulting from random perturbations ... This means that when the original elements are not exactly representable ... the errors resulting from any initial rounding that may be necessary are as serious as those arising from any initial rounding from all the steps in the solution..."

Καταστάσεις και δείκτες τους (**conditioning**)

Έννοιες που χρησιμεύουν για τη μέτρηση και ποσοτικοποίηση της επίδρασης διαταραχών και σφαλμάτων στα δεδομένα του προβλήματος και στις πράξεις α.κ.υ. επί των υπολογισμένων αποτελεσμάτων.

Κατάσταση προβλήματος

Αφορά στην ευαισθησία/μεταβλητότητα των αποτελεσμάτων της f που οφείλονται αποκλειστικά σε (μικρές) διαταραχές των στοιχείων εισόδου.

Κατάσταση (υλοποίησης) αλγορίθμου

Αφορά στην επίδραση της αριθμητικής πεπερασμένης ακρίβειας στην υλοποίηση του αλγορίθμου (μέσω προγράμματος.)

Ορισμός

Με τον γνωστό συμβολισμό, έστω ότι υπάρχει x_{prog} κοντά στο x τέτοιο ώστε $f_{\text{prog}}(x) = f(x_{\text{prog}})$. Τότε ο αλγόριθμος χαρακτηρίζεται (αριθμητικά) ((προς τα πίσω ευσταθής)) στο x .

Αν αυτό συμβαίνει για κάθε x στο πεδίο ορισμού της f ο αλγόριθμος αποκαλείται ((προς τα πίσω ευσταθής))

((Ένας πίσω ευσταθής αλγόριθμος απαντά με ακρίβεια σε ένα παραπλήσιο ερώτημα.))

Κρίσιμη παρατήρηση: Αν κατασκευάσουμε x_{prog} ((κοντά)) στο x ώστε αν το αποτέλεσμα που υπολογίστηκε με πράξεις α.κ.υ. είναι $z_{\text{prog}} = f_{\text{prog}}(x)$ τότε $z_{\text{prog}} = f(x_{\text{prog}})$. Τότε

$$\|z_{\text{prog}} - z\| = \|f_{\text{prog}}(x) - f(x)\| = \|f(x_{\text{prog}}) - f(x)\|.$$

Αντί να εκτιμούμε το προς τα εμπρός σφάλμα απευθείας, εκτιμούμε πρώτα τη διαφορά $\|f(x_{\text{prog}}) - f(x)\|$ που αφορά την ευαισθησία της f σε αλλαγές.

Αν υπάρχει x_{prog} ώστε το $\|x_{\text{prog}} - x\|$ να είναι μικρό, τότε το πρόβλημα ανάγεται στο μαθηματικό (και όχι αριθμητικό) πρόβλημα της ευαισθησίας της f σε μικρές διαταράξεις (= perturbations) των στοιχείων εισόδου.

((Ρίξαμε το σφάλμα προς τα πίσω))¹ αναγόντάς το σε (εικονικές) διαταραχές των στοιχείων εισόδου. Από εκεί και πέρα δεν χρειάζεται να ασχοληθούμε με τα σφάλματα των ενδιάμεσων πράξεων!

Η f_{prog} αποτελείται από όλες τις στοιχειώδεις πράξεις (ό,τι επηρεάζει δεδομένα α.κ.υ. του προβλήματος) και με τη σειρά που εκτελούνται από το πρόγραμμα που υλοποιεί την f .

- Δεν περιμένουμε η υλοποίηση f_{prog} να λύνει το πρόβλημα με πιστότητα αν τα δεδομένα x^* είναι κοντά σε στοιχεία όπου η f είναι ((ευαίσθητη)).
- Αυτό που μας ενδιαφέρει είναι η υλοποίηση της f να μην εισάγει από μόνη της ((μεγάλα)) σφάλματα.

¹ Θυμηθείτε το συμπέρασμα του Wilkinson!

Έχουμε ήδη συναντήσει την πίσω ευστάθεια!

Θυμηθείτε την Αρχή Ακριβούς Στρογγύλευσης:

Αν $\tilde{\odot}$ είναι η υλοποίηση της αριθμητικής πράξης \odot , τότε αν $x, y \in F$ ισχύει ότι

$$x \tilde{\odot} y = \text{fl}(x \odot y) = (x \odot y)(1 + \delta), \text{ για κάποιο } |\delta| \leq \mathbf{u}.$$

Παράδειγμα: από άμεση επαλήθευση,

$$f_{\text{prog}}(x) = f(x_{\text{prog}})$$

όπου

$$\begin{aligned} x &= [\xi_1, \xi_2]^\top \\ x_{\text{prog}} &= [\xi_1(1 + \delta), \xi_2(1 + \delta)] \end{aligned}$$

και

$$\|x_{\text{prog}} - x\| = \|[\xi_1\delta, \xi_2\delta]^\top\| \leq \|x\|\mathbf{u}$$

Συμπέρασμα Όλες οι απλές αριθμητικές πράξεις ικανοποιούν την αρχή ακριβούς στρογγύλευσης επομένως είναι πίσω ευσταθείς.

Ασυμπτωτικός δείκτης κατάστασης προβλήματος (Rice' 66)

Δίνονται δύο νορμισμένοι γραμμικοί χώροι \mathcal{X}, \mathcal{Y} και έστω $f: \Omega \subset X \rightarrow Y$, όπου Ω ανοικτό χωρίο. Έστω x^* σταθερή τιμή και ότι $y^* := f(x^*)$. Υποθέτουμε ότι τα x^*, y^* δεν είναι μηδέν. Ο **ασυμπτωτικός δείκτης κατάστασης της απεικόνισης f στο x^*** για σχετικά απειροελάχιστες αλλαγές του x^* ορίζεται ως

$$\text{cond}(f, x^*) := \lim_{\delta \rightarrow 0} \sup_{\|h\|=\delta} \left\{ \frac{\frac{\|f(x^*+h) - f(x^*)\|}{\|f(x^*)\|}}{\frac{\|h\|}{\|x^*\|}} \right\}$$

εφόσον το όριο υπάρχει. **Με άλλα λόγια: Ποιά είναι η σχετική αλλαγή στην τιμή μιας συνάρτησης αν αλλάξουμε σχετικά απειροελάχιστα τις τιμές των μεταβλητών της?**

Παράδειγμα: Αν $X = Y = \mathbb{R}$ και η f διαφορίσιμη τότε

$$\text{cond}(f, x^*) = |f'(x^*)| \frac{|x^*|}{|f(x^*)|}$$

Προσοχή:

- Αν για κάποιο σύνολο τιμών $x^* \in X^* \subseteq \mathcal{X}$ ισχύει ότι $\text{cond}(f, x^*)$ είναι μη φραγμένο ή πάρα πολύ μεγάλο, τότε ο υπολογισμός της f καθίσταται προβληματικός επί των τιμών του \mathcal{X} ανεξαρτήτως του αλγορίθμου που χρησιμοποιείται.
- Για τις τιμές αυτές, η συνάρτηση f είναι “δύσκολη” στη διαχείρισή της γιατί μικρές αλλαγές στις τιμές της μεταβλητής της x^* οδηγεί σε πολύ μεγάλες αλλαγές στην τιμή της $f(x^*)$.
- Το πρόβλημα της εύρεσης του $f(x^*)$ όταν ο ασυμπτωτικός δείκτης κατάστασης είναι πολύ μεγάλος αποκαλείται **κακώς τοποθετημένο**.
- Τα κακώς τοποθετημένα προβλήματα και στρατηγικές για την αντιμετώπισή τους είναι παρουσιάζουν εξαιρετικό ενδιαφέρον γιατί είναι ιδιαίτερα δύσκολα αλλά απαντώνται πολύ συχνά στην πραγματικότητα.

Καλά και κακά τοποθετημένα προβλήματα I

Θυμηθείτε τις κατηγορίες των προβλημάτων (άμεσα, αντίστροφα, ταυτοποίησης). Η (κακή) τοποθέτηση ως ζήτημα συνήθως δεν αφορά τα άμεσα αλλά τις άλλες δύο κατηγορίες. Η συζήτηση συνήθως διεξάγεται για τα **αντίστροφα προβλήματα**.

Καλά τοποθετημένο (well-posed) πρόβλημα

Ως **καλά τοποθετημένο** εννοούμε ένα πρόβλημα π.χ. τη λύση του $\phi(y) = x$ όπου $\phi : Y \rightarrow X$ είναι γραμμικοί χώροι και 1) για κάθε $\hat{x} \in X$ υπάρχει λύση $\hat{y} \in Y$. 2) Η λύση είναι μοναδική, δηλαδή η συνάρτηση ϕ είναι αντιστρέψιμη, έστω $\phi^{-1} = f$. 3) Για κάθε \bar{y} πλησίον της λύσης \hat{y} υπάρχει \bar{x} πλησίον του \hat{x} τέτοιο ώστε $\phi(\bar{y}) = \bar{x}$. Ισοδύναμα, η $f = \phi^{-1}$ είναι συνεχής.

Κακά τοποθετημένο (ill-posed) πρόβλημα

Κακά τοποθετημένο λέγεται ένα πρόβλημα για το οποίο δεν ισχύει κάποιες από τις 3 παραπάνω συνθήκες.

- Οι παραπάνω ορισμοί είναι κλασικοί και παλαιοί (Hadamard' 1902).

Καλά και κακά τοποθετημένα προβλήματα II

- Κακώς τοποθετημένα είναι πολλά προβλήματα σε πραγματικές εφαρμογές (δείτε παρακάτω).
- Η επίλυσή τους είναι μεγάλη πρόκληση σήμερα και απαιτεί **διεπιστημονικές** προσεγγίσεις που συνδυάζουν μαθηματικά, πιθανότητες, στατιστική, αριθμητική ανάλυση, HPC (high performance computing) καθώς γνώσεις της εφαρμογής όπου προκύπτει το κάθε πρόβλημα.
- Τεχνικές επίλυσης κακά τοποθετημένων προβλημάτων αποτελούν σημαντική πρόκληση για τη σημερινή έρευνα σε πολλές περιοχές.
- Το Gene Golub SIAM Summer School 2018 είναι αφιερωμένο στο ζήτημα.
- Δείτε πάλι σελ. 33-37 από το σύγγραμμα Quarteroni et al. (QSS00, Ch. 2).
- (προαιρετικά) διαβάστε όσο μπορείτε από την αναφορά (Kab08).
- (προαιρετικά, ακολουθήστε τους συνδέσμους) για συζητήσεις σχετικά με ζητήματα τοποθέτησης στην επεξεργασία σήματος, στη μηχανική μάθηση.

Διαβάστε:

(Kab08) Attempting to understand a substantially complex phenomenon and solve a problem such that the probability of error is high, we usually arrive at an unstable (ill-posed) problem. Ill-posed problems are ubiquitous in our daily lives. Indeed, everyone realizes how easy it is to make a mistake when reconstructing the events of the past from a number of facts of the present (for example, to reconstruct a crime scene based on the existing direct and indirect evidence, determine the cause of a disease based on the results of a medical examination, and so on). The same is true for tasks that involve predicting the future (predicting a natural disaster or simply producing a one week weather forecast) or “reaching into” inaccessible zones to explore their structure (subsurface exploration in geophysics or examining a patient’s brain using NMR tomography).

Definitions and examples of inverse and ill-posed problems

S. I. Kabanikhin

Survey paper

Abstract. The terms “inverse problems” and “ill-posed problems” have been steadily and surely gaining popularity in modern science since the middle of the 20th century. A little more than fifty years of studying problems of this kind have shown that a great number of problems from various branches of classical mathematics (computational algebra, differential and integral equations, partial differential equations, functional analysis) can be classified as inverse or ill-posed, and they are among the most complicated ones (since they are unstable and usually nonlinear). At the same time, inverse and ill-posed problems began to be studied and applied systematically in physics, geophysics, medicine, astronomy, and all other areas of knowledge where mathematical methods are used. The reason is that solutions to inverse problems describe important properties of media under study, such as density and velocity of wave propagation, elasticity parameters, conductivity, dielectric permittivity and magnetic permeability, and properties and location of inhomogeneities in inaccessible areas, etc.

In this paper we consider definitions and classification of inverse and ill-posed problems and describe some approaches which have been proposed by outstanding Russian mathematicians A. N. Tikhonov, V. K. Ivanov and M. M. Lavrentiev.

Key words. Inverse and ill-posed problems, regularization.

AMS classification. 65J20, 65J10, 65M32.

1. Introduction.
2. Classification of inverse problems.
3. Examples of inverse and ill-posed problems.
4. Some first results in ill-posed problems theory.
5. Brief survey of definitions of inverse problems.

Well-posed problems	Ill-posed problems
Arithmetic	
Multiplication by a small number A $Aq = f$	Division by a small number $q = A^{-1}f \quad (A \ll 1)$
Algebra	
Multiplication by a matrix $Aq = f$	$q = A^{-1}f$, A is an ill-conditioned, degenerate or rectangular $m \times n$ -matrix
Calculus	
Integration $f(x) = f(0) + \int_0^x q(\xi) d\xi$	Differentiation $q(x) = f'(x)$
Differential equations	
The Sturm-Liouville problem $u''(x) - q(x)u(x) = \lambda u(x)$, $u(0) - hu'(0) = 0$, $u(1) - Hu'(1) = 0$	The inverse Sturm-Liouville problem. Find $q(x)$ using spectral data $\{\lambda_n, \ u_n\ \}$
Integral geometry	
Find integrals $\int_{\Gamma(\xi, \eta)} q(x, y) ds$	Find q from $\int_{\Gamma(\xi, \eta)} q(x, y) ds = f(\xi, \eta)$
Integral equations	
Volterra equations and Fredholm equations of the second kind $q(x) + \int_0^x K(x, \xi)q(\xi) d\xi = f(x)$ $q(x) + \int_a^b K(x, \xi)q(\xi) d\xi = f(x)$	Volterra equations and Fredholm equations of the first kind $\int_0^x K(x, \xi)q(\xi) d\xi = f(x)$ $\int_a^b K(x, \xi)q(\xi) d\xi = f(x)$
Operator equations $Aq = f$	
$\exists m > 0: \forall q \in Q$ $m(q, q) \leq \langle Aq, q \rangle$	$A: D(A) \subset Q \rightarrow R(A) \subset F$ A is a compact linear operator with singular values $\sigma_n \searrow 0, n \rightarrow \infty$

Λεπτομερής μέτρηση ευαισθησίας συνάρτησης

Ιακωβιανό μητρώο

Το μητρώο $J = \left[\frac{\partial f_i}{\partial \xi_j} \right] \in \mathbb{R}^{m \times n}$ αποκαλείται Ιακωβιανό^{σ'} της f στο x .

^{σ'} γενίκευση της 1ης παραγώγου για συναρτήσεις πολλών μεταβλητών.

Δείκτες ευαισθησίας

Έστω $x = (\xi_1, \dots, \xi_m)$ και $y = (\eta_1, \dots, \eta_n) = f(x)$ και ομοίως $y^* = f(x^*)$ για κάποιο x^* . Τότε αν η f είναι διαφορίσιμη στο x , οι mn παράγοντες

$$\kappa_{ij} = \frac{\xi_j \frac{\partial f_i}{\partial \xi_j}}{f_i(x)}$$

δείχνουν την ευαισθησία της f ως προς διαταραχές των στοιχείων x .

Ειδικές περιπτώσεις: Αν $x = 0, f(x) = 0$ χρησιμοποιούμε δείκτη, $|f'(x)|$.

Δείκτες κατάστασης κοινών (βαθμωτών) συναρτήσεων

Άθροισμα n αριθμών

$$\frac{\sum_{j=1}^n |\xi_j|}{|\sum_{j=1}^n \xi_j|}$$

Εσωτερικό γινόμενο

$$\frac{2 \sum_{j=1}^n |\xi_j \psi_j|}{|\sum_{j=1}^n \xi_j \psi_j|}$$

Τιμή πολωνύμου σε μορφή δύναμης

$$\frac{\sum_{j=1}^n |\alpha_j| |\chi^j|}{|\sum_{j=1}^n \alpha_j \chi^j|}$$

Απλή ρίζα πολωνύμου σε μορφή δύναμης

$$\frac{\sum_{j=1}^n |\alpha_j| |\xi^j|}{|\xi \rho'(\xi)|}$$

Έστω $f(\xi_1, \xi_2) = \xi_1 - \xi_2$. Τότε $f : \mathbb{R}^2 \rightarrow \mathbb{R}$ και $\frac{\partial f}{\partial x} = [1, -1]$

$$\begin{aligned} \frac{|f(x + \Delta x) - f(x)|}{|f(x)|} &= \frac{|((\xi_1 + \Delta\xi_1) - (\xi_2 + \Delta\xi_2)) - (\xi_1 - \xi_2)|}{|\xi_1 - \xi_2|} \\ &= \frac{|\Delta\xi_1 - \Delta\xi_2|}{|\xi_1 - \xi_2|} \end{aligned}$$

Το σχετικό σφάλμα γίνεται μεγάλο όταν η πραγματική τιμή $|f(\xi_1, \xi_2)| = |\xi_1 - \xi_2|$ είναι πολύ μικρή σε σύγκριση με τη διαφορά $|\Delta\xi_1 - \Delta\xi_2|$.

Είναι η περίπτωση της **καταστροφικής απαλοιφής** και αφορά στο πρόβλημα (αφαίρεση) για τα συγκεκριμένα δεδομένα (ομόσημα και σε μικρή απόσταση μεταξύ τους) και όχι στην υλοποίηση (αλγόριθμος αφαίρεσης) στην ALU.

Ορισμός

Αν γνωρίζουμε ποιό είναι το x_{prog} που ελαχιστοποιεί το $\frac{\|x - x_{\text{prog}}\|}{\|x\|}$ για κάθε $x \in \mathcal{U}$ τότε μπορούμε να ορίσουμε το δείκτη κατάστασης του αλγορίθμου (στην πραγματικότητα - της υλοποίησής του). Θα είναι η μικρότερη τιμή, που συμβολίζουμε με $\text{cond}(f_{\text{prog}})$, για την οποία ισχύει η ανισότητα

$$\frac{\|x - x_{\text{prog}}\|}{\|x\|} \leq \text{cond}(f_{\text{prog}}) \mathbf{u}.$$

για όλα τα έγκυρα x_{prog} .

- Ο δείκτης $\text{cond}(f_{\text{prog}}) \mathbf{u}$ δηλώνει το πίσω σφάλμα ως προς όλο το πεδίο ορισμού \mathcal{U} .
- Η τιμή δεν είναι ενιαία: μπορεί να υπάρχουν περιοχές του \mathcal{U} που είναι μεγάλη, ενώ σε άλλες να είναι μικρή.
- Στις περιπτώσεις αυτές, μπορεί να προτιμάμε κάποιον αλγόριθμο για την “εύκολη” περιοχή και άλλον για την “δύσκολη”.
- Συχνά αδυνατούμε να βρούμε την καλύτερη δυνατή τιμή του άνω φράγματος και

Συμπεριφορά της LU για την επίλυση του $Ax = b$ και επιλογή αλγορίθμου:

- α) Αν διαγώνια κυρίαρχο, η LU χωρίς οδήγηση είναι πίσω ευσταθής.
- β) Για τα περισσότερα μητρώα, η LU συμπεριφέρεται ως πίσω ευσταθής.
- γ) Υπάρχουν ειδικές περιπτώσεις (π.χ. μητρώο gfrrp που η LU δεν είναι πίσω ευσταθής.)

```
function A = gfpp(T, c)
%      Reference:
%      N. J. Higham and D. J. Higham, Large growth factors in
%      Gaussian elimination with pivoting, SIAM J. Matrix Analysis and
%      Appl., 10 (1989), pp. 155-164.
%      N. J. Higham, Accuracy and Stability of Numerical Algorithms,
%      SIAM, 2nd ed.,
%      Philadelphia, PA, 2002; sec. 9.4.
```

```
function A = gfpp(T, c)
%GFPP (Author: Nick Higham)
% Matrix giving maximal growth factor for Gaussian elim. with pivoting.
% GFPP(T) is a matrix of order N for which Gaussian elimination
% with partial pivoting yields a growth factor  $2^{N-1}$ .
% T is an arbitrary nonsingular upper triangular matrix of order N-1.
% GFPP(T, C) sets all the multipliers to C ( $0 \leq C \leq 1$ )
% and gives growth factor  $(1+C)^{N-1}$  - but note that for  $T \sim \text{EYE}$ 
% ...
% GFPP(N, C) (a special case) is the same as GFPP(EYE(N-1), C) and
% generates the well-known example of Wilkinson.

if ~isequal(T, triu(T)) | any(~diag(T))
error('First argument must be a nonsingular upper triangular matrix.')
end
if nargin == 1, c = 1; end
if c < 0 | c > 1
error('Second parameter must be a scalar between 0 and 1 inclusive.')
end
m = length(T);
if m == 1 % Handle the special case T = scalar
n = T; m = n-1; T = eye(n-1);
else
n = m+1;
end
A = zeros(n); L = eye(n) - c*tril(ones(n), -1);
A(:,1:n-1) = L*(T; zeros(1,n-1)); theta = max(abs(A(:)));
A(:,n) = theta * ones(n,1); A = A/theta;
```

$$\begin{aligned} \frac{\|f_{\text{prog}}(x^*) - f(x)\|}{\|f(x)\|} &= \frac{\|f(x_{\text{prog}}^*) - f(x)\|}{\|f(x)\|} \\ &\leq \frac{\|f(x_{\text{prog}}^*) - f(x^*)\|}{\|f(x)\|} + \frac{\|f(x^*) - f(x)\|}{\|f(x)\|}. \end{aligned}$$

Φράσσουμε κάθε όρο ξεχωριστά:

$$\begin{aligned} \frac{\|f(x_{\text{prog}}^*) - f(x^*)\|}{\|f(x)\|} &\leq \frac{\|f(x_{\text{prog}}^*) - f(x^*)\|}{\|f(x^*)\|} \frac{\|f(x^*)\|}{\|f(x)\|} \\ &\leq \text{cond}(f, x^*) \frac{\|x^* - x_{\text{prog}}^*\|}{\|x^*\|} \frac{\|f(x^*)\|}{\|f(x)\|} \\ &\leq \text{cond}(f, x^*) \text{cond}(f_{\text{prog}}) \mathbf{u} \frac{\|f(x^*)\|}{\|f(x)\|}. \end{aligned}$$

και εξ ορισμού

$$\frac{\|f(x^*) - f(x)\|}{\|f(x)\|} \leq \text{cond}(f, x) \frac{\|x^* - x\|}{\|x\|} \leq \text{cond}(f, x) \mathcal{E}$$

όπου

$$\frac{\|x^* - x\|}{\|x\|} \leq \mathcal{E}$$

Επομένως

$$\frac{\|f_{\text{prog}}(x^*) - f(x)\|}{\|f(x)\|} \leq \text{cond}(f, x^*) \text{cond}(f_{\text{prog}}) \mathbf{u} \frac{\|f(x^*)\|}{\|f(x)\|} + \text{cond}(f, x) \mathcal{E}$$

Χονδρικά, όταν ο παράγοντας \mathcal{E} είναι μικρός και $\frac{\|f(x^*)\|}{\|f(x)\|} \approx 1$, ισχύει η ανισότητα:

προς τα εμπρός σφάλμα < δείκτης κατάστασης προβλήματος × πίσω σφάλμα

Ορισμός Ονομάζουμε ((προς τα πίσω ανάλυση σφάλματος)) τη διαδικασία εύρεσης φράγματος για το πίσω σφάλμα.

Στόχοι:

- 1 Να ερμηνεύσουμε τα σφάλματα στρογγύλευσης στους υπολογισμούς ως ισοδύναμες διαταραχές στα δεδομένα.
- 2 Να αναχθεί το πρόβλημα της εύρεσης φράγματος ή της εκτίμησης του εμπρός σφάλματος σε πρόβλημα της μαθηματικής θεωρίας διαταραχών.

Ισχύει ότι

$$\frac{\|f_{\text{prog}}(x^*) - f(x)\|}{\|f(x)\|} \leq \text{cond}(f; x^*) \text{cond}(f_{\text{prog}}) \mathbf{u} \frac{\|f(x^*)\|}{\|f(x)\|} + \text{cond}(f; x) \mathcal{E}$$

Ισχύει ότι

$$\frac{\|f_{\text{prog}}(x^*) - f(x)\|}{\|f(x)\|} \leq \text{cond}(f, x^*) \text{cond}(f_{\text{prog}}) \mathbf{u} \frac{\|f(x^*)\|}{\|f(x)\|} + \text{cond}(f, x) \mathcal{E}$$

Βήματα για να φράξουμε το εμπρός σφάλμα $\frac{\|f_{\text{prog}}(x^*) - f(x)\|}{\|f(x)\|}$ χρησιμοποιώντας την πίσω ανάλυση σφάλματος:

Υπολογισμός κατάστασης αλγορίθμου $\text{cond}(f_{\text{prog}})$, από την οποία προκύπτει το πίσω σφάλμα.

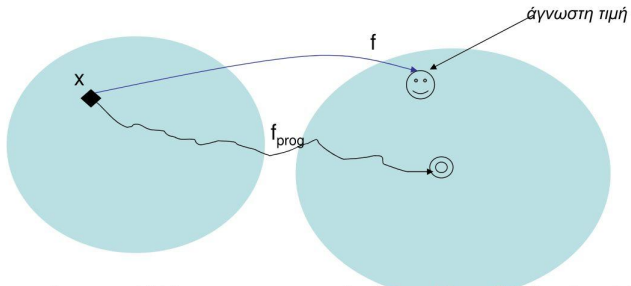
Υπολογισμός κατάστασης προβλήματος $\text{cond}(f, x)$

Εκτίμηση του \mathcal{E} αν χρειάζεται.

Όταν ο παράγοντας \mathcal{E} είναι μικρός και $\frac{\|f(x^*)\|}{\|f(x)\|} \approx 1$:

προς τα εμπρός σφάλμα $<$ δείκτης κατάστασης \times πίσω σφάλμα

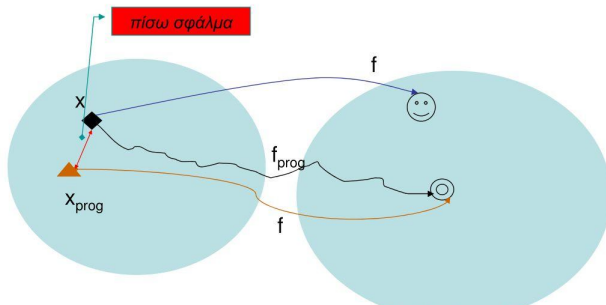
Το πρόβλημα: Να εκτιμήσουμε φράγμα για το εμπρός σφάλμα



Πόσο πρέπει να αλλάξει το x για να υπολογίσει η f (ακριβώς) το $f_{\text{prog}}(x)$?

$$f(x + ?) = f_{\text{prog}}(x)$$

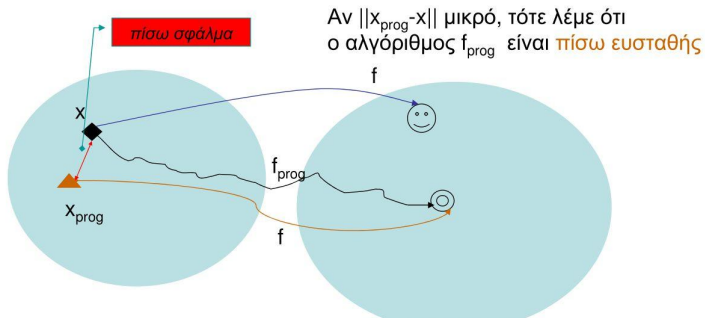
\swarrow
 $x_{\text{prog}} - x$



Πόσο πρέπει να αλλάξει το x ώστε η f να υπολογίζει (ακριβώς) το $f_{\text{prog}}(x)$?

$$f(x + ?) = f_{\text{prog}}(x)$$

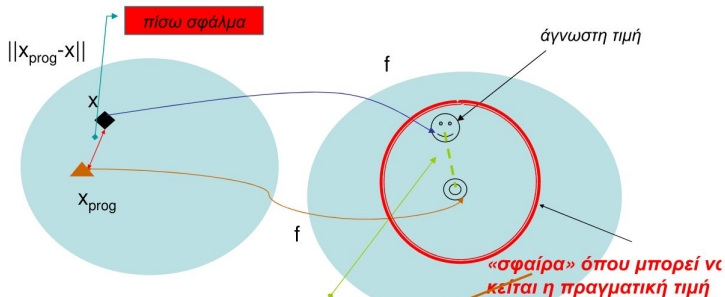
\swarrow
 $x_{\text{prog}} - x$



Πόσο πρέπει να αλλάξει το x ώστε η f να υπολογίζει (ακριβώς) το $f_{\text{prog}}(x)$?

$$f(x + ?) = f_{\text{prog}}(x)$$

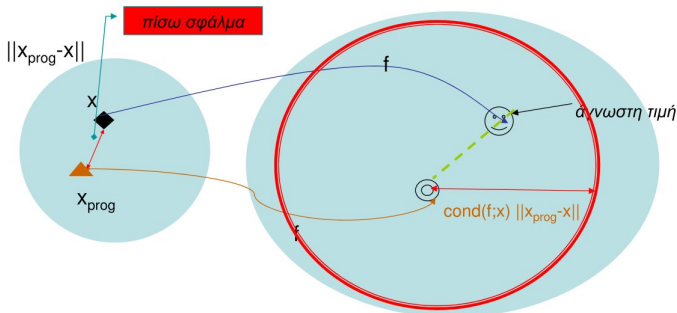
\nwarrow
 $x_{\text{prog}} - x$



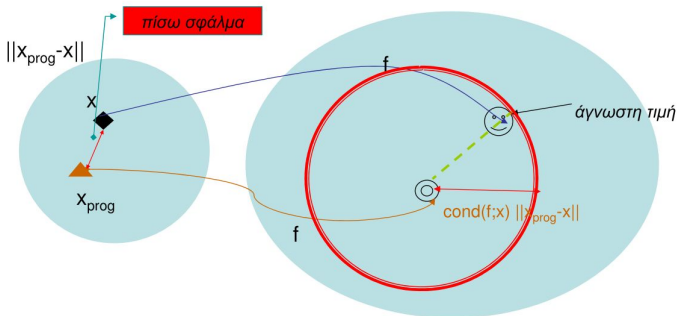
Η εκτίμηση του **εμπρός σφάλματος** γίνεται «μαθηματικά» βάσει

- του «δείκτη κατάστασης του προβλήματος» (δηλ. της f)
- χρησιμοποιώντας το πίσω σφάλμα.

$$\|f_{\text{prog}}(x) - f(x)\| \leq \text{cond}(f; x) \|x_{\text{prog}} - x\|$$



Αν η συνάρτηση είναι πολύ ευαίσθητη, το $\text{cond}(f;x)$ είναι μεγάλο (π.χ. αν $\frac{df}{dx}(x)$ μεγάλο) επομένως το $f(x)$ μπορεί να απέχει πολύ από το $f(x_{\text{prog}})$



Για δεδομένο $\text{cond}(f;x)$, η $\|f(x_{\text{prog}}) - f(x)\|$ θα κυμαίνεται ανάλογα με το πίσω σφάλμα $\|x_{\text{prog}} - x\|$

Αν $s_n = x^\top y$ τότε

$$\begin{aligned}\tilde{s}_1 &= \mathbf{fl}(\xi_1 \psi_1) = \xi_1 \psi_1 (1 + \delta_1) \\ \tilde{s}_2 &= \mathbf{fl}(\tilde{s}_1 + \mathbf{fl}(\xi_2 \psi_2)) \\ &= (\xi_1 \psi_1 (1 + \delta_1) + \xi_2 \psi_2 (1 + \delta_2))(1 + \delta_3) \\ &= \xi_1 \psi_1 (1 + \delta_1)(1 + \delta_3) + \xi_2 \psi_2 (1 + \delta_2)(1 + \delta_3)\end{aligned}$$

όπου $|\delta_i| \leq u$.

Έχουμε

$$\begin{aligned}\tilde{s}_n = & \xi_1 \psi_1 \prod_{\substack{j=1 \\ j \neq 2}}^{n+1} (1 + \delta_j) + \xi_2 \psi_2 \prod_{j=2}^{n+1} (1 + \delta_j) + \\ & \dots \xi_3 \psi_3 \prod_{j=3}^{n+1} (1 + \delta_j) + \dots + \xi_n \psi_n \prod_{j=n}^{n+1} (1 + \delta_j).\end{aligned}$$

Από το Λήμμα:

$$\begin{aligned}\tilde{s}_n = & \xi_1 \psi_1 (1 + \theta_n) + \xi_2 \psi_2 (1 + \hat{\theta}_n) + \\ & \dots \xi_3 \psi_3 (1 + \theta_{n-1}) + \dots + \xi_n \psi_n (1 + \theta_2).\end{aligned}$$

Το υπολογισμένο DOT είναι το ακριβές εσωτερικό γινόμενο για τα στοιχεία $\xi_1, \dots, \xi_n, \psi_1(1 + \theta_n), \psi_2(1 + \hat{\theta}_n), \dots, \psi_n(1 + \theta_2)$, όπου για τα θ ισχύει το φράγμα

$$|\theta_j| \leq \frac{j\mu}{1 - j\mu} = \gamma_j.$$

$$\begin{aligned} \mathbf{fl}(x^\top y) &= (x + \Delta x)^\top y = x^\top (y + \Delta y), \\ \text{όπου } |\Delta x| &\leq \gamma_n |x|, |\Delta y| \leq \gamma_n |y|. \end{aligned}$$

Δηλαδή το υπολογισμένο DOT είναι το ίδιο με το ακριβές DOT για στοιχεία εισόδου $x + \Delta x, y$. Επίσης

$$|\Delta x| \leq \gamma_n |x| \Rightarrow \|\Delta x\|_1 \leq \gamma_n \|x\|_1$$

άρα το σχετικό πίσω σφάλμα είναι φραγμένο ως εξής

$$\frac{\|\Delta x\|_1}{\|x\|_1} \leq \gamma_n$$

και δείκτης κατάστασης του αλγορίθμου το $\text{cond}(f_{\text{prog}}) = \frac{\gamma_n}{\mathbf{u}}$.

ο παραπάνω αλγόριθμος υπολογισμού του DOT είναι προς τα πίσω ευσταθής

ΠΡΟΣΟΧΗ

- ❶ Το ακριβές φράγμα εξαρτάται από τη σειρά υπολογισμού. Το παραπάνω φράγμα αντιστοιχεί στην άθροιση “από τα αριστερά προς τα δεξιά”.
- ❷ Στις συζητήσεις μας, οι τιμές των φραγμάτων είναι θεωρητικές. Για αξιόπιστες τιμές, πρέπει να ληφθεί υπόψη και το **σφάλμα στον υπολογισμό του γ_n** .
- ❸ π.χ. αντικαθιστώντας το γ_n με $\mathbf{fl}(\gamma_{n+1}) \geq \gamma_n$.

$$f([x; y]) := \sum_{j=1}^n \xi_j \psi_j = x^\top y, \quad x, y \in \mathbb{R}^n$$

Θέτουμε $X := [x; y] \in \mathbb{R}^{2n}$ και υπολογίζουμε το δείκτη κατάστασης με βάση την ((ευαισθησία)) του αποτελέσματος ως προς κάθε στοιχείο:

$$\begin{aligned} K &= \frac{1}{|x^\top y|} \left[\left| \xi_1 \frac{\partial f}{\partial \xi_1} \right|, \dots, \left| \xi_n \frac{\partial f}{\partial \xi_n} \right|, \left| \psi_1 \frac{\partial f}{\partial \psi_1} \right|, \dots, \left| \psi_n \frac{\partial f}{\partial \psi_n} \right| \right] \\ &= \frac{1}{|x^\top y|} \left[|\xi_1 \psi_1|, \dots, |\xi_n \psi_n|, |\xi_1 \psi_1|, \dots, |\xi_n \psi_n| \right] \end{aligned}$$

επομένως

$$\|K\|_1 = \frac{2 \sum_{j=1}^n |\xi_j| |\psi_j|}{\left| \sum_{j=1}^n \xi_j \psi_j \right|}$$

Η κατάσταση εξαρτάται αρ'ο το $|x^\top y|$, δηλ. από την τιμή $\cos(x, y)$. Υπάρχει πρόβλημα όταν $\|x\|, \|y\| \gg 0$ και $\cos(x, y) \approx 0$.

Προς τα πίσω ανάλυση $\frac{|x^\top y - \mathbf{fl}(x^\top y)|}{|x^\top y|} \leq \gamma_n \frac{2 \sum_{j=1}^n |\xi_j| |\psi_j|}{|\sum_{j=1}^n \xi_j \psi_j|}$

Προς τα εμπρός ανάλυση Χρησιμοποιώντας προς τα εμπρός ανάλυση συμπεραίνουμε ότι

$$\begin{aligned} |x^\top y - \mathbf{fl}(x^\top y)| &\leq \gamma_n \sum_{i=1}^n |\xi_i| |\psi_i| = \gamma_n |x|^\top |y| \\ &\leq nu |x|^\top |y| + O(u^2) \end{aligned}$$

Αν $|x|^\top y \ll |x|^\top |y|$ το άνω φράγμα μπορεί να είναι μεγάλο, επομένως δεν έχουμε καμμία εξασφάλιση όσο αφορά το μέγεθος του σχετικού σφάλματος.

Προσοχή: Όταν οι αντίστοιχες τιμές είναι ομόσημες, τότε το σχετικό σφάλμα είναι μικρό.

Η πίσω ανάλυση σφάλματος:

- Μπορεί να είναι δύσκολη
- Μπορεί να είναι **ανέφικτη**.

Διαστατικό επιχείρημα: Γενικά αν η διάσταση του διανύσματος εξόδου είναι μεγαλύτερη από τη διάσταση της εισόδου, τότε θα έχουμε δυσκολία να αποδείξουμε πίσω ευστάθεια.

Η πίσω ανάλυση σφάλματος:

- Μπορεί να είναι δύσκολη
- Μπορεί να είναι **ανέφικτη**.

Διαστατικό επιχείρημα: Γενικά αν η διάσταση του διανύσματος εξόδου είναι μεγαλύτερη από τη διάσταση της εισόδου, τότε θα έχουμε δυσκολία να αποδείξουμε πίσω ευστάθεια.

Παράδειγμα Ο υπολογισμός

$$\mathbf{fl}(C) = \mathbf{fl}(ab^T), a, b \in \mathbb{R}^n$$

άρα

$$\mathbf{fl}(\gamma_{ij}) = \alpha_i \beta_j (1 + \delta_{ij})$$

που σημαίνει 'οτι

$$\begin{aligned} \mathbf{fl}(C) &= ab^T + E \Rightarrow \\ |C - \mathbf{fl}(C)| &= |E| \leq \mathbf{u} \mathbf{e} \mathbf{e}^T \end{aligned}$$

όπου $[E]_{ij} = \alpha_i \beta_j \delta_{ij}$ και $\mathbf{e} \mathbf{e}^T$ είναι ίσο με το μητρώο με όλα τα στοιχεία του 1. Δηλ. το καλύτερο δυνατό άνω φράγμα για το εμπρός σφάλμα!

Για πίσω ευστάθεια θάπρεπε να υπάρχουν \tilde{a}, \tilde{b} κοντά στα a, b τ.ώ.

$$\mathbf{fl}(c) = \tilde{a}\tilde{b}^\top = (a + \Delta a)(b + \Delta b)^\top$$

Τότε θα ίσχυε

$$E = \overbrace{a\Delta b^\top + \Delta a(b^\top + \Delta b^\top)}^{\text{τάξη} \leq 2}$$

Αν ο αλγόριθμος υπολογισμού εξωτερικού γινομένου ήταν πίσω ευσταθής θα ίσχυε ότι

$$n = \mathbf{rank}(E) \leq 2$$

πράγμα που είναι γενικά αδύνατο (δεν υπάρχουν αρκετά δεδομένα εισόδου στα οποία να ((αναθέσουμε)) τα σφάλματα στα αποτελέσματα)

Το εξωτερικό γινόμενο διανυσμάτων δεν είναι πίσω ευσταθής πράξη

ΠΡΟΣΟΧΗ: Αυτό δεν σημαίνει ότι ο αλγόριθμος υπολογισμού έχει πρόβλημα! Μόνον ότι δεν μπορούμε να εφαρμόσουμε πίσω ανάλυση σφάλματος για να φράξουμε το εμπρός σφάλμα

....

Παράδειγμα

Μέθοδος Horner

$$s_n = \alpha_n$$

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

$$s_k = xs_{k+1} + \alpha_k$$

end

Χρησιμοποιώντας το γνωστό λήμμα:

$$\hat{s}_{n-1} = (xs_n(1 + \delta_1) + \alpha_{n-1})(1 + \delta_2) = x\alpha_n(1 + \theta_2) + \alpha_{n-1}(1 + \delta_2)$$

$$\hat{s}_{n-2} = (x\hat{s}_{n-1}(1 + \delta_3) + \alpha_{n-2})(1 + \delta_4)$$

$$\begin{aligned}\hat{s}_0 &= \alpha_0(1 + \delta) + \alpha_1x(1 + \theta_3) + \cdots \alpha_{n-1}x^{n-1}(1 + \theta_{2n-1}) + \alpha_nx^n(1 + \theta_{2n}) \\ &= f_{\text{prog}}(\alpha_0, \dots, \alpha_n, x) = f(\alpha_0(1 + \theta_1), \dots, \alpha_n(1 + \theta_{2n}), x)\end{aligned}$$

για κάποια θ_j φραγμένα σε απόλυτη τιμή από τα αντίστοιχα γ_j .

Επομένως, το προς τα πίσω σφάλμα της μεθόδου Horner είναι μικρό.

Ειδικότερα:

Αφού θ_{2n} είναι ο παράγοντας που χαρακτηρίζει τη μέγιστη ασάφεια στους συντελεστές, είναι ασφαλές να συμπεράνουμε ότι αν κάθε συντελεστής δεν είναι ακριβώς γνωστός αλλά γνωρίζουμε ότι $|\alpha_j - \hat{\alpha}_j| \leq \theta_{2n} |\alpha_j|$ για τον καθένα, τότε το σφάλμα που προκύπτει από την ασάφεια αυτή θα είναι σίγουρα μεγαλύτερο από οποιοδήποτε σφάλμα οφείλεται στις πράξεις α.κ.υ. που εκτελούνται από τον αλγόριθμο.

Αυτό από μόνο του δεν είναι αρκετό για να εγγυηθεί το σφάλμα των αποτελεσμάτων αλλά μας επιτρέπει να εξετάσουμε το προς τα εμπρός σφάλμα ανεξάρτητα από τις λεπτομέρειες του αλγορίθμου.

Εύκολα αποδεικνύεται η σχέση

$$\frac{|p(x) - \hat{s}_0|}{|p(x)|} \leq \gamma_{2n} \frac{\sum_{k=0}^n |\alpha_k| |x|^k}{|p(x)|},$$

από την οποία προκύπτει ότι δεν μπορούμε να εγγυηθούμε μικρό προς τα εμπρός σφάλμα.

Γιατί τόσο ενδιαφέρον για πολυώνυμα

Φαίνονται ((εύκολα)) αλλά επιπλέον :

- Τα πολυώνυμα είναι από τις πιο σημαντικά ((αλγεβρικά αντικείμενα)) και η ορθή χρήση τους είναι σημαντικό θέμα.
- Το πρόβλημα του υπολογισμού τιμών πολυωνύμου αυτό καθαυτό εμφανίζεται σε πολλές εφαρμογές.
- Ο αλγόριθμος αποτελεί υπολογιστικό πυρήνα για την εύρεση των ριζών.
- Δυστυχώς, φαίνονται εύκολα αλλά είναι ((ύπουλα)); βλ. το άρθρο (Wil84)!

- Ορισμένοι επεξεργαστές περιέχουν εντολή FMA (Fused Multiply and Add).
- Πλεονέκτημα στο χρόνο εκτέλεσης: Εκτελεί την πράξη $z + x * y$ στον ίδιο χρόνο περίπου που απαιτεί μια απλή $x * y$ ή $x + y$.
- Πλεονέκτημα στην ακρίβεια: προκύπτει **ένα μόνον σφάλμα στρογγύλευσης**.

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

αντί

$$\text{fl}(z + \text{fl}(xy)) = (z + xy(1 + \delta_1))(1 + \delta_2)$$

- Περιέχεται στο πρότυπο IEEE-754-2008.
- FMA3: $(a,b,c) \leftarrow a + b * c$; FMA4: $(a,b,c,\mathbf{d}) \leftarrow a + b * c$
- Δείτε (M^+10), Wikipedia
- Αρχικά (1990) στο IBM RS-6000. Σήμερα στα Intel Haswell, AMD Piledriver, Bulldozer

Εντολή Fused Multiply-Add (FMA) II

<https://software.intel.com/en-us/blogs/2011/06/13/haswell-new-instruction-descriptions-now-available/>

COO Τμήμα Μηχανικών Google Calendar CEID Webmail: Κ ΕΥΣΤΡΑΤΙΟΣ ΓΑΛ upatras eclass Ομάδες Google Ιδιώτες | Τράπεζα diary meteo.gr: Ο Κ

Floating Point Multiply Accumulate – Our floating-point multiply accumulate significantly increases peak flops and provides improved precision to further improve transcendental mathematics. They are broadly usable in high performance computing, professional quality imaging, and face detection. They operate on scalar, 128-bit packed single and double precision data types, and 256-bit packed single and double-precision data types. [These instructions were described previously, in the initial Intel® AVX specification].

	Double Precision Packed FP	Single Precision Packed FP	Double Precision Scalar FP	Single Precision Scalar FP
Fused Multiply-Add $A = A \times B + C$ $C \Leftarrow A \times B$	VFMADD132PD VFMADD213PD VFMADD231PD _mm_fmadd_pd _mm256_fmadd_pd	VFMADD132PS VFMADD213PS VFMADD231PS _mm_fmadd_ps _mm256_fmadd_ps	VFMADD132SD VFMADD213SD VFMADD231SD _mm_fmadd_sd _mm256_fmadd_sd	VFMADD132SS VFMADD213SS VFMADD231SS _mm_fmadd_ss _mm256_fmadd_ss
Fused Multiply-Alternating Add/Subtract $A = A \times B + C$ $A = A \times B - C$ $C \Leftarrow A \times B$ $C \Leftarrow A \times B$	VFMACDSUB132PD VFMACDSUB213PD VFMACDSUB231PD _mm_fmaddsub_pd _mm256_fmaddsub_pd	VFMACDSUB132PS VFMACDSUB213PS VFMACDSUB231PS _mm_fmaddsub_ps _mm256_fmaddsub_ps		
Fused Multiply-Alternating Subtract/Add $A = A \times B - C$ $A = A \times B + C$ $C \Leftarrow A \times B$ $C \Leftarrow A \times B$	VFMSUBADD132PD VFMSUBADD213PD VFMSUBADD231PD _mm_fmsubadd_pd _mm256_fmsubadd_pd	VFMSUBADD132PS VFMSUBADD213PS VFMSUBADD231PS _mm_fmsubadd_ps _mm256_fmsubadd_ps		
Fused Multiply-Subtract $A = A \times B - C$ $C \Leftarrow A \times B$	VFMSUB132PD VFMSUB213PD VFMSUB231PD _mm_fmsub_pd _mm256_fmsub_pd	VFMSUB132PS VFMSUB213PS VFMSUB231PS _mm_fmsub_ps _mm256_fmsub_ps	VFMSUB132SD VFMSUB213SD VFMSUB231SD _mm_fmsub_sd _mm256_fmsub_sd	VFMSUB132SS VFMSUB213SS VFMSUB231SS _mm_fmsub_ss _mm256_fmsub_ss
Fused Negative Multiply-Add $A = -A \times B + C$ $C \Leftarrow -A \times B$	VFNMADD132PD VFNMADD213PD VFNMADD231PD _mm_fnmadd_pd _mm256_fnmadd_pd	VFNMADD132PS VFNMADD213PS VFNMADD231PS _mm_fnmadd_ps _mm256_fnmadd_ps	VFNMADD132SD VFNMADD213SD VFNMADD231SD _mm_fnmadd_sd _mm256_fnmadd_sd	VFNMADD132SS VFNMADD213SS VFNMADD231SS _mm_fnmadd_ss _mm256_fnmadd_ss
Fused Negative Multiply-Subtract $A = -A \times B - C$ $C \Leftarrow -A \times B$	VFNMSUB132PD VFNMSUB213PD VFNMSUB231PD _mm_fnmsub_pd _mm256_fnmsub_pd	VFNMSUB132PS VFNMSUB213PS VFNMSUB231PS _mm_fnmsub_ps _mm256_fnmsub_ps	VFNMSUB132SD VFNMSUB213SD VFNMSUB231SD _mm_fnmsub_sd _mm256_fnmsub_sd	VFNMSUB132SS VFNMSUB213SS VFNMSUB231SS _mm_fnmsub_ss _mm256_fnmsub_ss

The vector instructions build upon the expanded (256-bit) register state added in Intel® AVX, and as such as supported by any operating system that supports Intel® AVX.

For developers, please note that the instructions span multiple CPUID leaves. You should be careful to check all applicable bits before using these instructions.

Please check out the specification and stay tuned for supporting tools over the next couple of months.

Mark Buxton
Software Engineer
Intel Corporation



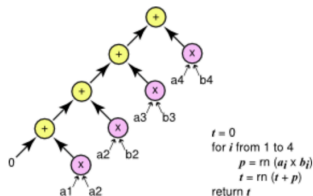
PRECISION AND PERFORMANCE: FLOATING POINT AND IEEE 754 COMPLIANCE FOR NVIDIA GPUS

TB-06711-001_v5.5 | July 2013

White paper

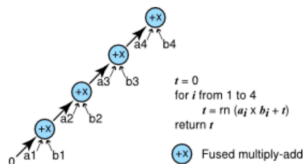
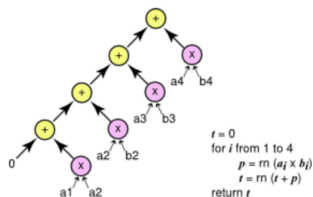
Εσωτερικό γινόμενο

$\alpha = (1.907607, -0.7862027, 1.148311, .9604002); \quad b = \dots$
 $(-0.9355000, -0.6915108, 1.724470, -0.7097529)$



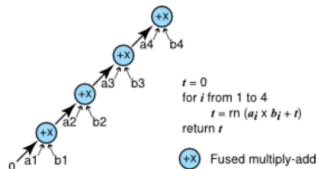
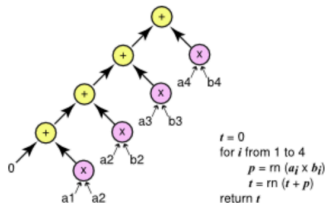
Εσωτερικό γινόμενο

$\mathbf{a} = (1.907607, -0.7862027, 1.148311, 0.9604002); \quad \mathbf{b} = \dots$
 $(-0.9355000, -0.6915108, 1.724470, -0.7097529)$



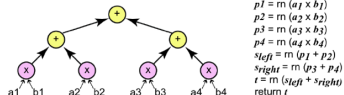
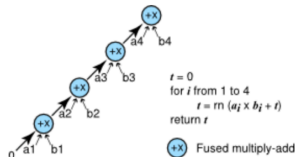
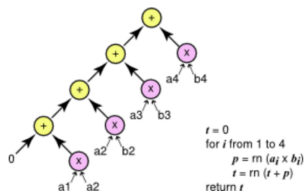
Εσωτερικό γινόμενο

$\alpha = (1.907607, -0.7862027, 1.148311, 0.9604002); \quad b = \dots$
 $(-0.9355000, -0.6915108, 1.724470, -0.7097529)$



Εσωτερικό γινόμενο

$\alpha = (1.907607, -.7862027, 1.148311, .9604002); \quad b = \dots$
 $(-.9355000, -.6915108, 1.724470, -.7097529)$

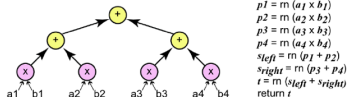
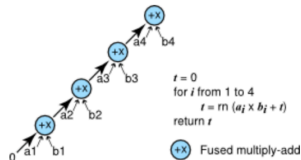
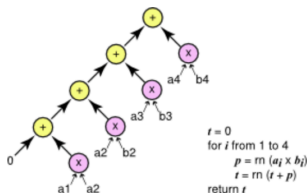


method	result	float value
exact	.0559587528435...	0x3D65350158...
serial	.0559588074	0x3D653510
FMA	.0559587515	0x3D653501
parallel	.0559587478	0x3D653500

Figure 5 Algorithms Results vs. the Correct Mathematical Dot Product

Εσωτερικό γινόμενο

$\alpha = (1.907607, -0.7862027, 1.148311, 0.9604002); \quad b = \dots$
 $(-0.9355000, -0.6915108, 1.724470, -0.7097529)$



method	result	float value
exact	.059587528435...	0x3D65350158...
serial	.059587515	0x3D653510
FMA	.059587515	0x3D653501
parallel	.059587478	0x3D653500

Figure 5 Algorithms Results vs. the Correct Mathematical Dot Product

Το πρόβλημα της άθροισης (Higham'02)

Δίδεται σύνολο αριθμών $\mathcal{S} = \{\xi_1, \dots, \xi_n\}$ και ζητούμε το άθροισμα $\text{sum}(\mathcal{S})$.

- Μια από τις πιο συχνές πράξεις σε επιστημονικούς και εμπορικούς υπολογισμούς
- π.χ. συνοπτικά στατιστικά χαρακτηριστικά, π.χ. μέσος όρος, απόκλιση, νόρμα, ...
- Αξίζει τον κόπο μια πιο προσεκτική θεώρηση
- ... σε σχέση με την *ταχύτητα* και την *ακρίβεια*

Αναδρομική άθροιση

$s = 0$

for $i = 1 : n$

$s = s + \xi_i$

end

Σχετικά με την αναδρομική άθροιση

- Πολύ απλή υλοποίηση
- Ιδιαίτερα χρήσιμη για ομόσημα στοιχεία μετά από ταξινόμηση σε αύξουσα σειρά κατ' απόλυτη τιμή
- Αποφυγή προβλημάτων ((απορρόφησης)) από μεγάλα στοιχεία,
- ... πρέπει να είναι γνωστά τα στοιχεία και να ταξινομήσουμε πριν την άθροιση.

Εύκολα βλέπουμε ότι

$$\begin{aligned}\mathbf{fl}(s) &= (\cdots ((\xi_1 + \xi_2)(1 + \delta_1) + \xi_3)(1 + \delta_2) \cdots + \xi_n)(1 + \delta_{n-1}) \\ &= \xi_1(1 + \theta_{n-1}) + \xi_2(1 + \hat{\theta}_{n-1}) + \xi_3(1 + \hat{\theta}_{n-2}) + \cdots + \xi_n(1 + \theta_1)\end{aligned}$$

όπου ως συνήθως

$$|\theta_j| \leq \gamma_j = \frac{\mu}{1 - \mu}$$

Επομένως, το πίσω σφάλμα εξαρτάται άμεσα από το n , όπως θα περιμέναμε...

- Αναγωγική άθροιση ανά ζεύγη² πολύ χρήσιμη για παράλληλη άθροιση
- Ενθετική άθροιση
- Αντισταθμισμένη άθροιση και νεώτερες παραλλαγές

Χρήσιμα χαρακτηριστικά

- είναι οι αριθμοί ομόσημοι?
- είναι η ακολουθία διαθέσιμη από την αρχή?
- ποιό είναι το εύρος (μέγιστο, ελάχιστο)?
- είναι η ακολουθία διατεταγμένη/ταξινομημένη?

² pairwise / cascade / fan-in

- 1 ταξινόμηση του \mathcal{S} κατ' απόλυτη τιμή σε αύξουσα σειρά
 $\mathcal{L} := \xi_1 \leq \xi_2 \leq \xi_3 \leq \dots \leq \xi_n$
- 2 Διαγραφή του $\xi_1 + \xi_2$ απο την \mathcal{L} και ένθεση στην \mathcal{L} στην κατάλληλη θέση ώστε να παραμείνει μονοτονική. Επαναρίθμηση του \mathcal{L} και αν περιέχει 2 ή περισσότερα στοιχεία, επιστροφή στο (2).
- 3 Η σειρά των προσθέσεων δημιουργεί ένα δένδρο με φύλλα τους αρχικούς αριθμούς $\{\xi_1, \dots, \xi_n\}$ υπό πρόσθεση και ρίζα την άθροιση από την οποία προκύπτει το τελικό άθροισμα. Οι υπόλοιποι κόμβοι αντιστοιχούν στις ενδιάμεσες αθροίσεις.

Ενδιαφέρον Θυμίζει την κωδικοποίηση Huffman.

(KW00)

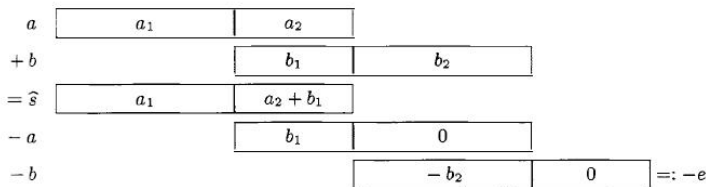
Αν τα δεδομένα έχουν μικτό πρόσημο (θετικά και αρνητικά στοιχεία) η εύρεση του αθροίσματος που ελαχιστοποιεί το μέγιστο σφάλμα στρογγύλευσης για όλους τους δυνατούς τρόπους άθροισης είναι NP-hard.

Εφόσον το πρόβλημα κατασκευής του βέλτιστου δένδρου άθροισης είναι πολύ δύσκολο:

- *χαλαρώνουμε* τις απαιτήσεις και αναζητούμε κάτι λιγότερο φιλόδοξο,
- *αξιοποιούμε* ό,τι πληροφορία υπάρχει για το \mathcal{S} .
- Αν χρειάζεται μεγάλη ακρίβεια, σκεφτείτε τη χρήση αντισταθμισμένης άθροισης ή άθροιση με εκτεταμένη ακρίβεια
- Για τις περισσότερες μεθόδους, το σφάλμα είναι στη χειρότερη περίπτωση ανάλογο με το n . Για πολύ μεγάλο n καθίσταται ενδιαφέρουσα η χρήση της αντισταθμισμένης άθροισης ή της αναγωγικής άθροισης
- Αν υπάρχει πιθανότητα σημαντικής (ως καταστροφικής) απαλοيفής λόγω μεικτών προσήμων, η αναδρομική άθροιση με φθίνουσα διάταξη είναι συνήθως καλύτερη.
- Αν το σύνολο ομόσημο, όλες οι μέθοδοι προσφέρουν σχετικό σφάλμα το πολύ n^{-1} με καλύτερη την **επανορθωμένη άθροιση**³ (compensated summation).

³Στις επόμενες διαφάνειες.

Ακριβής άθροιση \rightarrow επανορθωμένη άθροιση



Σχήμα: Ιδιοφυής ιδέα (Gill'51, Kahan'65) από (Higham'02)

Για κάθε α.κ.υ. $|a| \geq |b|$ μπορούμε να υπολογίσουμε

$$\hat{s} = \text{fl}(a + b), \quad \hat{e} = \text{fl}((a - \hat{s}) + b)$$

χρησιμοποιώντας στρογγύλευση προς το πλησιέστερο. Το υπολογισμένο \hat{e} είναι το **ακριβές σφάλμα** της άθροισης! Ειδικότερα,

$$a + b = \hat{s} + \hat{e}$$

Παρόλα αυτά, $\text{fl}(\hat{s} + \hat{e}) = \hat{s}$ γιατί το αποτέλεσμα $\hat{s} = \text{fl}(a + b)$ είναι το καλύτερο δυνατό (λόγω αρχής ακριβούς στρογγύλευσης)

Ιδέα Να το κρατήσουμε για την επόμενη πρόσθεση (αν υπάρχει). Αν πρέπει να προσθέσουμε και άλλον αριθμό, c , τότε (ίσως) μπορούμε να ((επανορθώσουμε)) με το άθροισμα των κρατούμενων. Πλεονεκτήματα επανορθωμένης άθροισης Γενικά είναι πιο ακριβής από την συνηθισμένη άθροιση.

Εμπρός σφάλμα (Knutth, Kahan)

φράσσεται ως

$$\left| \sum_{j=1}^n \xi_j - \text{fl} \left(\sum_{j=1}^n \xi_j \right) \right| \leq (2\mathbf{u} + o(n\mathbf{u}^2)) \sum_{j=1}^n |\xi_j|.$$

Πίσω σφάλμα

$$\text{fl}(s) = \sum_{j=1}^n (1 + \mu_j) \xi_j, \quad \text{όπου } |\mu_j| \leq 2\mathbf{u} + o(n\mathbf{u}^2)$$

δηλ. όσον αφορά στους όρους 1ης τάξης (που εξαρτώνται από το \mathbf{u}) το πίσω σφάλμα είναι ανεξάρτητο του n .

A Floating-Point Technique for Extending the Available Precision

T. J. DEKKER*

Received July 26, 1971

Abstract. A technique is described for expressing multilength floating-point arithmetic in terms of singlelength floating point arithmetic, i.e. the arithmetic for an available (say: single or double precision) floating-point number system. The basic algorithms are exact addition and multiplication of two singlelength floating-point numbers, delivering the result as a doublelength floating-point number. A straightforward application of the technique yields a set of algorithms for doublelength arithmetic which are given as ALGOL 60 procedures.

Let x and y be singlelength floating-point numbers and let

$$z = fl(x + y);$$

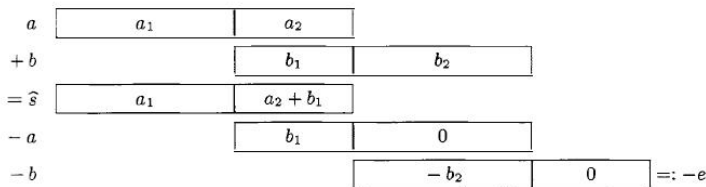
i.e. z is the result of a singlelength floating-point addition of x and y . Let zz be the correction term exactly satisfying

$$z + zz = x + y.$$

It will be shown that, under various conditions, zz can be obtained by the formula

$$zz = fl((x - z) + y).$$

Ακριβής άθροιση \rightarrow επανορθωμένη άθροιση



Σχήμα: Ιδιοφυής ιδέα (Gill'51, Kahan'65) από (Higham'02)

Για κάθε α.κ.υ. $|a| \geq |b|$ μπορούμε να υπολογίσουμε

$$\hat{s} = \text{fl}(a + b), \quad \hat{e} = \text{fl}((a - \hat{s}) + b)$$

χρησιμοποιώντας στρογγύλευση προς το πλησιέστερο. Το υπολογισμένο \hat{e} είναι το **ακριβές σφάλμα** τής άθροισης! Ειδικότερα,

$$a + b = \hat{s} + \hat{e}$$

Παρόλα αυτά, $\text{fl}(\hat{s} + \hat{e}) = \hat{s}$ γιατί το αποτέλεσμα $\hat{s} = \text{fl}(a + b)$ είναι το καλύτερο δυνατό (λόγω αρχής ακριβούς στρογγύλευσης)

Ιδέα Να το κρατήσουμε για την επόμενη πρόσθεση (αν υπάρχει). Αν πρέπει να προσθέσουμε και άλλον αριθμό, c , τότε (ίσως) μπορούμε να ((επανορθώσουμε)) με το άθροισμα των κρατούμενων. Πλεονεκτήματα επανορθωμένης άθροισης Γενικά είναι πιο ακριβής από την συνηθισμένη άθροιση.

Εμπρός σφάλμα (Knutth, Kahan)

φράσσεται ως

$$\left| \sum_{j=1}^n \xi_j - \text{fl} \left(\sum_{j=1}^n \xi_j \right) \right| \leq (2\mathbf{u} + o(n\mathbf{u}^2)) \sum_{j=1}^n |\xi_j|.$$

Πίσω σφάλμα

$$\text{fl}(s) = \sum_{j=1}^n (1 + \mu_j) \xi_j, \quad \text{όπου } |\mu_j| \leq 2\mathbf{u} + o(n\mathbf{u}^2)$$

δηλ. όσον αφορά στους όρους 1ης τάξης (που εξαρτώνται από το \mathbf{u}) το πίσω σφάλμα είναι ανεξάρτητο του n .

Θα δοκιμάσουμε βήμα-βήμα την επανόρθωση στον γνωστό υπολογισμό

$$((10^{20} - 10) - 10^{20}) + 20$$

```
a1=10^20; b1=-10; s1 = a1+b1;  
e1 = -((a1 - s1)+b1); % επιστρέφει 10  
a1=-10^20; b1=s1; s2=a1+b1;  
e2 = -((a1 - s2)+b1); % επιστρέφει 0  
a1=s2; b1=20; s3 = a1+b1;  
e3 = -((a1 - s2)+b1); επιστρέφει 20  
s = s3+e1+e2+e3  
= 10
```

- Η επανόρθωση “διόρθωσε” το αποτέλεσμα και το “επανέφερε” στην ακριβή του τιμή.
- Πρόκειται για το best case scenario, δεν συμβαίνει πάντα!!!

Listing 4: Επανόρθωση στην άθροιση 2 αριθμών

```
function [s,e] = kahan_cs(a,b);  
if (abs(a)< abs(b))  
    temp =a; a=b; b=temp;  
end  
s = a+b;  
e = (a-s)+b;
```

Listing 5: Επανόρθωση χωρίς ελέγχους στην άθροιση 2 αριθμών

```
function [s,e] = nocomp2sum(x,y);  
% following Muller et al.  
s = a+b;  
a_dot = s-b; b_dot = s-a_dot;  
da = a-a_dot; db = b-b_dot;  
e = da + db;
```

Άθροιση $n \geq 3$ στοιχείων I

Listing 6: Επανόρθωση στην άθροιση

```
function [s] = casc_sum(x);  
s = x(1);  
e = 0;  
for i=2:length(x);  
    [s,c(i)] = kahan_cs(s,x(i));  
    e = e + c(i);  
end  
s = s+e;
```

Παρατήρηση: Εναλλακτικά μπορούμε να εφαρμόσουμε αναδρομικά τον ίδιο αλγόριθμο για την άθροιση των στοιχείων του C .

ΠΑΡΑΔΕΙΓΜΑ: (όπως πριν)

$x = [1e20, -10, -1e20, 20]$

ΑΠ	$\text{sum}(x)$	$\text{casc_sum}(x)$
10	20	10

Listing 7: Επανόρθωση στην άθροιση

```
function [s] = casc_sum(x);  
s = x(1);  
e = 0;  
for i=2:length(x);  
    [s,c(i)] = kahan_cs(s,x(i));  
    e = e + c(i);  
end  
s = s+e;
```

ΠΑΡΑΔΕΙΓΜΑ: $x = 1/n * \text{ones}(n,1)$ Θεωρητικά $\text{sum}(x) = 1$
 $x = 1/n * \text{ones}(n,1)$

Έστω $x = 1/n * \text{ones}(n, 1)$ Θεωρητικά πρέπει να ισχύει ότι $\text{sum}(x) = 1$

n	sum(x)	casc_sum(x)	abs(1-sum(x))
10	9.999999999999999e-01	1	1.1102e-016
100	1.0000000000000001e+00	1	6.6613e-016
10^4	9.999999999980838e-01	1	1.9162e-012

- 1 Διάταξη \mathcal{S} σε αύξουσα απόλυτη τιμή $\mathcal{L} := \xi_1 \leq \xi_2 \leq \xi_3 \leq \dots \leq \xi_n$
- 2 Διαγραφή των ξ_1, ξ_2 από το σύνολο \mathcal{L} και ένθεση του $\xi_1 + \xi_2$ στο \mathcal{L} στη σωστή θέση ώστε η ακολουθία να παραμείνει διατεταγμένη κατ' αύξουσα τιμή.
- 3 Επαναρίθμηση των στοιχείων του \mathcal{L} και επιστροφή στο (2) αν έχει τουλάχιστον 3 στοιχεία.

Summation algorithms can be represented as binary tree \mathcal{T} , where the leaves are the inputs, the root the output and the internal nodes are the partial sums computed in the course of the algorithm.

Δείτε την τελευταία παράγραφο:



matteo-frigo commented on Mar 31, 2014

Just to make it more interesting, here are a couple of issues with FMA, in addition to those pointed out by stevenj.

First, the FMA implemented by PowerPC and by Intel differ in (at least) one subtle way. PowerPC computes $+/(ab \ +/- \ c)$, whereas Intel implements $+/(ab) \ +/- \ c$. The two differ when the result is -0 ; for example, $-(+0+0 - \ +0) = -0$, whereas $+0 - \ +0+0 = +0$. So a pedantic (I dare not say correct) compiler cannot translate the PowerPC expression into the x86 expression.

Second, common-subexpression elimination makes things interesting. For example, FFTs tend to compute pairs of expressions of the form $a \ +/- \ bc$, which reduces to two fma's. However, a sufficiently smart compiler will tend to compute $t=bc$; $a \ +/- \ t$, which is three flops. Historically in fftw, we were able to defeat gcc-3.x with some hackish workaround, but gcc-4.[0-2] had a more formidable optimizer that was always isolating the common subexpression no matter what. I haven't tried more recent gcc's or llvm in this respect, but that's something to watch for.

Σφάλματα και επανόρθωση στον πολλαπλασιασμό? I

- Οι παραπάνω ιδέες μπορούν να επεκταθούν και στον πολλαπλασιασμό α.κ.υ.
- Δείτε το άρθρο του Dekker που αναρτήθηκε νωρίτερα και θυμηθείτε και τη συζήτηση στην τάξη.
- Επιγραμματικά μόνον αναφέρουμε ορισμένες από τις ιδέες.
- Περισσότερα στις αναφορές.

Ερώτημα Μπορούμε να υπολογίσουμε το ακριβές σφάλμα στον πολλαπλασιασμό α.κ.υ.?

Αναζητούμε τα εξής: Δοθέντων 2 α.κ.υ. x, y να υπολογίσουμε το γινόμενο σε α.κ.υ. και α.κ.υ. b τ.ω. $xy = fl(xy) + b$ (σε αριθμητική άπειρης ακρίβειας).

Απαιτείται πρώτα να διαχωρίσουμε την είσοδο σε 2 μέρη: Αν $eps = 2^{-p}$ τότε $s = \lceil p/2 \rceil$. Για παράδειγμα αν $p = 53$ τότε $s = 27$. Ο αλγόριθμος του Dekker διαχωρίζει τον α.κ.υ. a σε δύο τμήματα x, y ώστε $a = x_h + x_l$ αλλά $|x_l| \leq |x_h|$ και τα x_h, x_l δεν έχουν επικάλυψη.

Listing 8: Διαχωρισμός του Veltkamp

```
function [x_h,x_l]=SplitV(x,s)
% split precision-p radix-B fp into two floats x_h, x_l
% so that for given s<p, the significand of x_h fits
% in p-s digits and the significand of x_l fits into s digits
% and x=x_h+x_l exactly
F = B^s+1;
g = F*x;
d = x-g;
x_h = g / F; x_l = x-x_h;
```

Παράδειγμα (M^+10)

Για βάση 10, ακρίβεια $p = 8$ ψηφία, διαχωρισμός σε 2 τμήματα $s = 4$ ψηφίων.

Καλούμε $\text{SplitV}(1.2345678, 4)$

$$F = 10^s + 1; \Rightarrow F = 10001$$

$$g = F * x; \Rightarrow g = 12346.9125678 \Rightarrow g = 12346.913$$

$$d = x - g; \Rightarrow x - g = -12345.6784322 \Rightarrow d = -12345.678$$

$$x_h = g + d; \Rightarrow g + d = 1.235 \Rightarrow x_h = 1.235$$

$$x_l = x - x_h; \Rightarrow x - x_h = -0.0004322 \Rightarrow x_l = -0.0004322$$

Προσέξτε ισχύει ακριβώς ότι

$$x_h + x_l = 1.235 - 0.0004322 = 1.2345678$$



W. Gautschi.

Numerical Analysis: An Introduction.

Birkhauser, Boston, 1997.



N.J. Higham.

Accuracy and Stability of Numerical Algorithms.

SIAM, Philadelphia, 2nd edition, 2002.



S I Kabanikhin.

Definitions and examples of inverse and ill-posed problems.

J. Inverse and Ill-posed Problems, 16(4), 2008.



Ming-Yang Kao and Jie Wang.

Linear-time approximation algorithms for computing numerical summation with provably small errors.

SIAM J. Comput., 29(5):1568–1576, 2000.



A.N. Langville.

Catastrophic cancellation on the high seas.

The Pi Mu Epsilon Journal, 11(4):205–208, 2001.

<http://www4.ncsu.edu/~ipsen/ma798I/langville.pdf>.



J.-M. Muller et al.

Handbook of Floating-Point Arithmetic.

Birkhäuser Boston, 2010.



A. Quarteroni, R. Sacco, and F. Saleri.

Numerical Mathematics.

Springer, New York, 2000.



J.H. Wilkinson.

The perfidious polynomial.

