

Επιστημονικός Υπολογισμός I

ΗΥ 343: ΔΙΑΛΕΞΗ 11

Ε. Γαλλόπουλος

Τμήμα Η/Υ & Πληροφορικής
Πανεπιστήμιο Πατρών



Πανεπιστήμιο Πατρών



Αντιστάθμιση (compensation) - βασική ιδέα

Πόρισμα Αν οι αριθμοί $x \geq y$ είναι στο σύστημα α.κ.υ. IEEE και χρησιμοποιούμε στρόγγυλη προς τον πλησιέστερο και αν

$$s = \text{fl}(x + y), w = \text{fl}(s - x), c = \text{fl}(y - w)$$

τότε

$$s + c = x + y$$

ΔΗΛΑΔΗ για κάθε ζευγάρι προσθετέων (x, y) μπορούμε να υπολογίσουμε ακριβώς την προσέγγιση του αθροίσματος $s := \text{fl}(x + y)$ και έναν (γενικά πολύ μικρό) α.κ.υ. c που σε αριθμητική άπειρης ακρίβειας μπορεί να διορθώσει το s και να το επαναφέρει στην τιμή του ακριβούς αθροίσματος $x + y$.

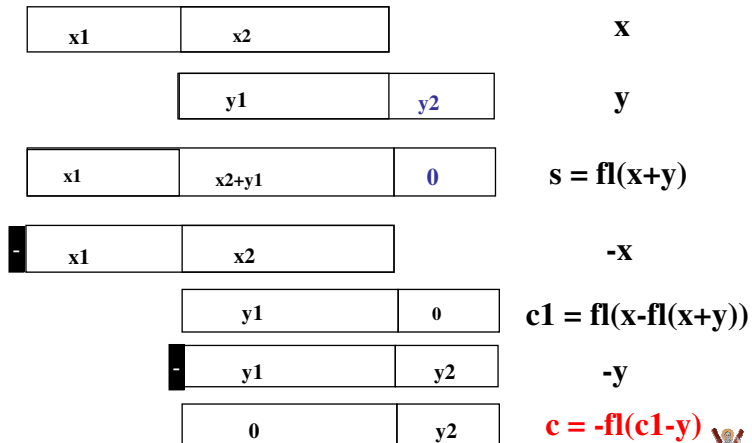
```
function [s,e] = kahan0(x,y);  
if (abs(x)< abs(y))  
    temp =a; a=b; b=temp;  
end  
s = x+y;  
e = y-(s-x);
```



Πανεπιστήμιο Πατρών



Η διαδικασία σχηματικά (λίγο παραλλαγμένη η σειρά)



TMHYU
Τμήμα Ηλεκτρονικών Υπολογιστών & Τεχνολογίας

Πανεπιστήμιο Πατρών



Πως χρησιμοποιείται το c (σαν carry);

- Το **fl(s)** είναι το καλύτερο δυνατό (ακριβή στρογγύλευση)
- άρα η διόρθωση **fl(fl(c)+fl(x+y))** δεν οδηγεί σε ακριβέστερο άθροισμα
- Αλλά αν έχουμε να προσθέσουμε και άλλον αριθμό, **z**, τότε μπορούμε να «διορθώσουμε» με το **c**

```
function [s] = kahan1(x)
s = 0; e = 0;
for j = 1:length(x)
    temp = s
    y = x(j) + e
    s = s + y
    e = (temp-s)+y
end
```



TMHYU
Τμήμα Ηλεκτρονικών Υπολογιστών & Τεχνολογίας

Πανεπιστήμιο Πατρών



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

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

We shall derive some formulas for calculating zz which use only singlelength floating-point addition and subtraction.

First we consider the formula

$$(4.3) \quad w = fl(z - x), \quad zz = fl(y - w)$$

and prove some theorems stating sufficient conditions for the validity of this formula. For practical computation, formulas (4.1) and (4.3) can be written as the following sequence of ALGOL 60 statements:

$$(4.4) \quad "z := x + y; \quad zz := y - (z - x)";$$

in which w remains anonymous.

Let $x, y \in R$ be representable such that their exponents satisfy

$$(4.5) \quad ex \geq ey.$$

In particular, this holds if $x, y \in R$ satisfy

(4.7) **Theorem.** If R has the form (2.2), where $\beta = 2$ or 3 and M is a multiple of β , and if, moreover, floating-point addition is optimal and subtraction faithful, then, for all x and y satisfying (4.5) and for z obtained according to (4.1), formula (4.3) yields the correction term zz defined by (4.2).



Αντιστάθμιση - βασική ιδέα

Πόρισμα Αν οι αριθμοί $x \geq y$ είναι στο σύστημα α.κ.υ. IEEE και χρησιμοποιούμε στρωγύλευση προς τον πλησιέστερο και αν

$$s = \text{fl}(x + y), w = \text{fl}(s - x), c = \text{fl}(y - w)$$

τότε

$$s + c = x + y$$

ΔΗΛΑΔΗ για κάθε ζευγάρι προσθετέων (x, y) μπορούμε να υπολογίσουμε ακριβώς την προσέγγιση του αθροίσματος $s := \text{fl}(x + y)$ και έναν (γενικά πολύ μικρό) α.κ.υ. c που σε αριθμητική άπειρης ακρίβειας μπορεί να διορθώσει το s και να το επαναφέρει στην τιμή του ακριβούς αθροίσματος $x + y$.

```
function [s,e] = kahan0(x,y);
if (abs(x)< abs(y))
    temp =a; a=b; b=temp;
end
s = x+y;
e = y-(s-x);
```



TMHYP
Τμήμα Ηλεκτρονικών Υπολογιστών & Τεχνολογίας

Πανεπιστήμιο Πατρών



Η διαδικασία σχηματικά (λίγο παραλλαγμένη η σειρά)

x1	x2
----	----

x

y1	y2
----	----

y

x1	x2+y1	0
----	-------	---

s = fl(x+y)

-	x1	x2
---	----	----

-x

y1	0
----	---

c1 = fl(x-fl(x+y))

-	y1	y2
---	----	----

-y

0	y2
---	----

c = -fl(c1-y)



TMHYP
Τμήμα Ηλεκτρονικών Υπολογιστών & Τεχνολογίας

Πανεπιστήμιο Πατρών



Πως χρησιμοποιείται το c (σαν carry);

- Το $\text{fl}(s)$ είναι το καλύτερο δυνατό (ακριβή στρογγύλευση)
- άρα η διόρθωση $\text{fl}(\text{fl}(c) + \text{fl}(x+y))$ δεν οδηγεί σε ακριβέστερο άθροισμα
- Αλλά αν έχουμε να προσθέσουμε και άλλον αριθμό, z , τότε μπορούμε να «διορθώσουμε» με το c

```
function [s] = kahan1(x)
s = 0; e = 0;
for j = 1:length(x)
    temp = s
    y = x(j) + e
    s = s + y
    e = (temp-s)+y
end
```



ΤΜΗΥΠ

Πανεπιστήμιο Πατρών



Ιδιότητα αντισταθμισμένης άθροισης

- Γενικά πιο ακριβής από κανονική άθροιση.
- Αποδεικνύεται ότι το άθροισμα n στοιχείων
$$\text{fl}(s) = \sum (1 + \mu_j) x_j$$

όπου $|\mu_j| \leq 2u + O(nu^2)$
- Δηλαδή
σε πρώτη προσέγγιση, το πίσω σφάλμα είναι ανεξάρτητο από το n .



ΤΜΗΥΠ

Πανεπιστήμιο Πατρών



Παράδειγμα

- Έστω $x = 1/n * \text{ones}(n,1)$
 - Τότε θεωρητικά, $\text{sum}(x) = 1$

- Εκτελούμε σε **MATLAB**

n	sum(x)	kahan1(x)	abs(1-sum(x))	abs(1-kahan1(x))
10	9.999999999999999e-01	1	1.1102e-016	0
100	1.0000000000000001e+00	1	6.6613e-016	0
10 ⁴	9.999999999980838e-01	1	1.9162e-012	0



TMHYP
ΤΕΧΝΙΚΗ ΜΗΧΑΝΙΚΗ ΥΠΟΛΟΓΙΣΤΩΝ & ΕΠΕΞΕΡΓΑΣΙΑΣ

Πανεπιστήμιο Πατρών



Journal of Economic Literature
Vol. XXXVII (June 1999), pp. 633–665

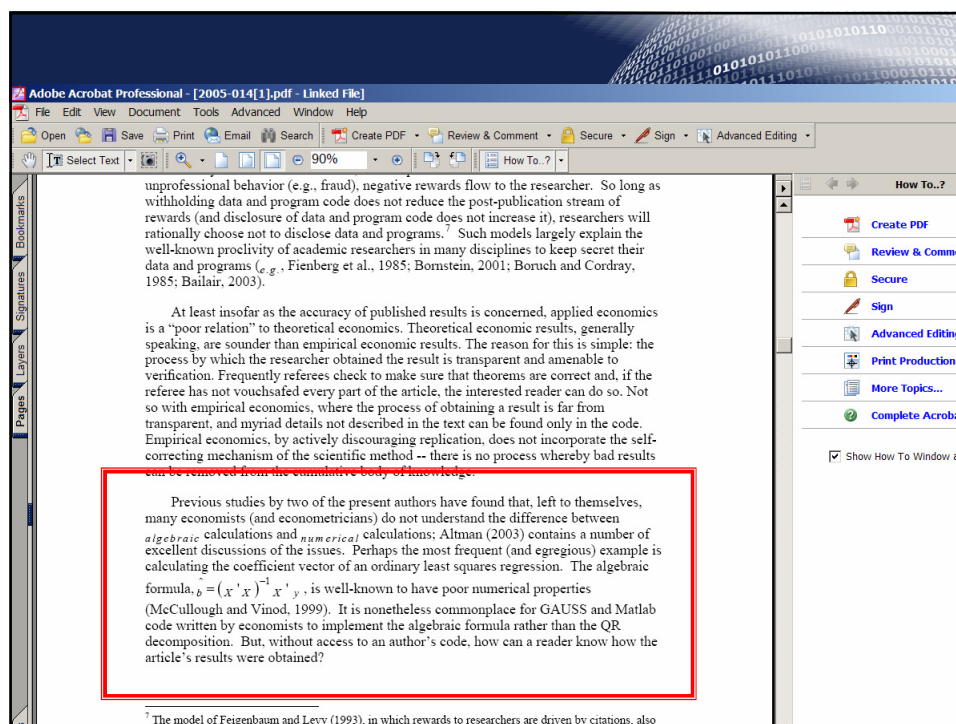
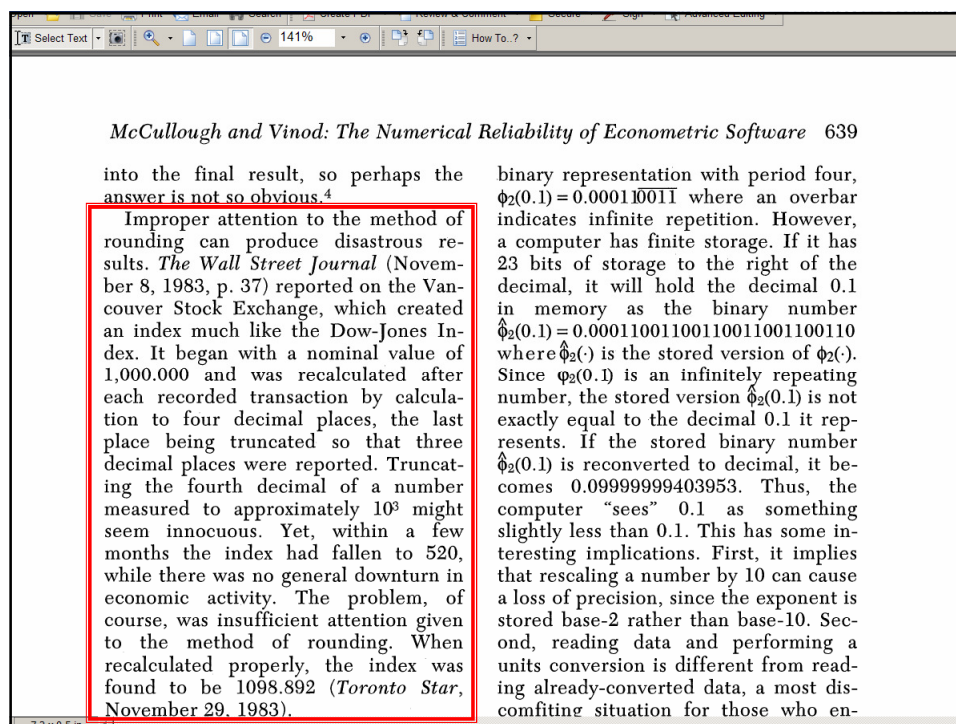
The Numerical Reliability of Econometric Software

B. D. McCULLOUGH
and
H. D. VINOD¹

1. Introduction

Numerical software is central to our computerized society; it is used . . . to analyze future options for financial markets and the economy. It is essential that it be of high quality; fast, accurate, reliable, easily moved to new machines, and easy to use. (Ford and Rice 1994)

APART FROM COST considerations, economists generally choose their 1990–97, over 120 reviews appeared. All but three paid no attention to numerical accuracy, and only two applied more than a single test of numerical accuracy (Michael Veall 1991, and McCullough 1997, but see also Vinod 1989 and Colin McKenzie 1998). Since computation is the *raison d'être* of an econometric package, this lacuna is all the more puzzling given the failure of many statistical pack-



High-Precision Floating-Point Arithmetic in Scientific Computation

David H. Bailey

28 January 2005

Abstract

At the present time, IEEE 64-bit floating-point arithmetic is sufficiently accurate for most scientific applications. However, for a rapidly growing body of important scientific computing applications, a higher level of numeric precision is required: some of these applications require roughly twice this level; others require four times; while still others require hundreds or more digits to obtain numerically meaningful results. Such calculations have been facilitated by new high-precision software packages that include high-level language translation modules to minimize the conversion effort. These activities have yielded a number of interesting new scientific results in fields as diverse as quantum theory, climate modeling and experimental mathematics, a few of which are described in this article. Such developments suggest that in the future, the numeric precision used for a scientific computation may be as important to the program design as are the algorithms and data structures.



Πανεπιστήμιο Πατρών



3. Climate Modeling

It is well-known that weather or climate simulations are fundamentally chaotic—if microscopic changes are made to the present state, within a certain period of simulated time the future state is completely different. Indeed, ensembles of these calculations are required to obtain statistical confidence in global climate trends produced from such calculations. As a result, computational scientists involved in climate modeling applications have long resigned themselves that their codes quickly diverge from any “baseline” calculation, even if they only change the number of processors used to run the code. As a result, it is not only difficult for researchers to compare results, but it is often problematic even to determine whether they have correctly deployed their code on a given system.

Recently Helen He and Chris Ding, two researchers at LBNL, investigated this non-reproducibility phenomenon in a widely-used climate modeling code. They found that almost all of the numerical variation occurred in one inner product loop in the atmospheric data assimilation step, and in a similar operation in a large conjugate gradient calculation. He and Ding found that a straightforward solution was to employ double-double arithmetic (using the DDFUN package above) for these loops. This single change dramatically reduced the numerical variability of the entire application, permitting computer runs to be compared for much longer run times than before. Details of their work can be read in [21].

In retrospect, it is not clear that handling these sums in this manner is the best solution—there may be other ways to preserve meaningful numerical accuracy, while at the same time yielding reproducible results. Such phenomena deserve further study and refinement and are being investigated. But in the meantime, He and Ding’s solution is a straightforward and effective means of dealing with this problem.



Πανεπιστήμιο Πατρών



9. Computational Geometry and Grid Generation

Grid generation, contour mapping and several other computational geometry applications crucially rely on highly accurate arithmetic. This is because small numerical errors can lead to geometrically inconsistent results. Such difficulties are inherent in the mathematics of the formulas commonly used in such calculations and cannot be easily remedied. For example, William Kahan and Joseph Darcy have shown that small numerical errors in the computation of the point nearest to a given point on a line of intersection of two planes can result in the computed point being so far from either plane as to rule out the solution being correct for a reasonable perturbation of the original problem [23].

Two commonly used computational geometry operations are the orientation test and the incircle test. The orientation test attempts to unambiguously determine whether a point lies to the left of, to the right of, or on a line or plane defined by other points. In a similar way, an incircle test determines whether a point lies inside, outside, or on a

circle defined by other points. Each of these tests is typically performed by evaluating the sign of a determinant that is expressed in terms of the coordinates of the points. If these coordinates are expressed as single or double precision floating-point numbers, roundoff error may lead to an incorrect result when the true determinant is near zero. In turn, this misinformation can lead an application to fail or produce incorrect results, as noted above.

To remedy such problems, Jonathan Shewchuk has produced a software package that performs "adaptive" floating-point arithmetic, which dynamically increases numeric precision until an unambiguous result is obtained. This software and some related information



TMHYP
ΤΕΧΝΗ ΜΗΧΑΝΙΚΗΣ ΥΠΕΡ ΤΕΧΝΟΛΟΓΙΑΣ ΚΑΙ ΣΥΜΦΩΝΗΣΕΩΣ

Πανεπιστήμιο Πατρών



$$\begin{aligned}
 \pi\sqrt{3} &= \frac{9}{32} \sum_{k=0}^{\infty} \frac{1}{64^k} \left(\frac{16}{6k+1} - \frac{8}{6k+2} - \frac{2}{6k+4} - \frac{1}{6k+5} \right) \\
 \pi^2 &= \frac{1}{8} \sum_{k=0}^{\infty} \frac{1}{64^k} \left[\frac{144}{(6k+1)^2} - \frac{216}{(6k+2)^2} - \frac{72}{(6k+3)^2} - \frac{54}{(6k+4)^2} + \frac{9}{(6k+5)^2} \right] \\
 \pi^2 &= \frac{2}{27} \sum_{k=0}^{\infty} \frac{1}{729^k} \left[\frac{243}{(12k+1)^2} - \frac{405}{(12k+2)^2} - \frac{81}{(12k+4)^2} - \frac{27}{(12k+5)^2} \right. \\
 &\quad \left. - \frac{(12k+6)^2}{9} - \frac{(12k+7)^2}{9} - \frac{(12k+8)^2}{5} - \frac{(12k+10)^2}{1} + \frac{(12k+11)^2}{1} \right] \\
 \log 3 &= \frac{1}{729} \sum_{k=0}^{\infty} \frac{1}{729^k} \left(\frac{729}{6k+1} + \frac{81}{6k+2} + \frac{81}{6k+3} + \frac{9}{6k+4} + \frac{9}{6k+5} + \frac{1}{6k+6} \right) \\
 \pi \log 2 &= \frac{1}{256} \sum_{k=0}^{\infty} \frac{1}{4096^k} \left[\frac{4096}{(24k+1)^2} - \frac{8192}{(24k+2)^2} - \frac{26112}{(24k+3)^2} + \frac{15360}{(24k+4)^2} \right. \\
 &\quad - \frac{(24k+5)^2}{768} + \frac{(24k+6)^2}{64} + \frac{(24k+8)^2}{408} + \frac{(24k+9)^2}{720} - \frac{(24k+10)^2}{512} \\
 &\quad + \frac{(24k+12)^2}{16} - \frac{(24k+13)^2}{196} + \frac{(24k+15)^2}{60} + \frac{(24k+16)^2}{37} \\
 &\quad \left. + \frac{(24k+17)^2}{196} + \frac{(24k+18)^2}{60} + \frac{(24k+20)^2}{37} - \frac{(24k+21)^2}{37} \right] \\
 0 &\stackrel{?}{=} -2J_2 - 2J_3 - 2J_4 + 2J_{10} + 2J_{11} + 3J_{12} + 3J_{13} + J_{14} - J_{15} - J_{16} \\
 &\quad - J_{17} - J_{18} - J_{19} + J_{20} + J_{21} - J_{22} - J_{23} + 2J_{25}
 \end{aligned}$$

Figure 2: Some new math identities found by high-precision computations



TMHYP
ΤΕΧΝΗ ΜΗΧΑΝΙΚΗΣ ΥΠΕΡ ΤΕΧΝΟΛΟΓΙΑΣ ΚΑΙ ΣΥΜΦΩΝΗΣΕΩΣ

Πανεπιστήμιο Πατρών



