

Rounding Errors and Stability of Algorithms

Rounding Errors

Let N_{\min} , N_{\max} be the floating point numbers of smallest and largest magnitude in a finite precision system.

Rounding Errors

Let N_{\min} , N_{\max} be the floating point numbers of smallest and largest magnitude in a finite precision system.

For any normalized number x (which is not necessarily a floating point number in the system) such that $N_{\min} \leq |x| \leq N_{\max}$,

$$fl(x) = x(1 + \delta), \quad |\delta| \leq u \quad (*)$$

This will be called the fundamental representation/rounding rule.

Rounding Errors

Let N_{\min} , N_{\max} be the floating point numbers of smallest and largest magnitude in a finite precision system.

For any normalized number x (which is not necessarily a floating point number in the system) such that $N_{\min} \leq |x| \leq N_{\max}$,

$$fl(x) = x(1 + \delta), \quad |\delta| \leq u \quad (*)$$

This will be called the fundamental representation/rounding rule.

For normalized numbers floating point numbers x, y such that $N_{\min} \leq |x \text{ op } y| \leq N_{\max}$,

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

where 'op' is any of the operations $+$, $-$, $/$, \times and the unit roundoff $u \approx 10^{-16}$ in IEEE double precision.

Rounding Errors

Let N_{\min} , N_{\max} be the floating point numbers of smallest and largest magnitude in a finite precision system.

For any normalized number x (which is not necessarily a floating point number in the system) such that $N_{\min} \leq |x| \leq N_{\max}$,

$$fl(x) = x(1 + \delta), \quad |\delta| \leq u \quad (*)$$

This will be called the fundamental representation/rounding rule.

For normalized numbers floating point numbers x, y such that $N_{\min} \leq |x \text{ op } y| \leq N_{\max}$,

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

where 'op' is any of the operations $+$, $-$, $/$, \times and the unit roundoff $u \approx 10^{-16}$ in IEEE double precision.

This is desirable as it implies that the maximum relative error in the operation is of the order of unit roundoff.

Rounding errors

Do similar statements hold for normalized numbers x and y that are not floating point numbers?

Rounding errors

Do similar statements hold for normalized numbers x and y that are not floating point numbers?

Suppose, $N_{\min} \leq |x|, |y|, |x \text{ op } y| \leq N_{\max}$. In such a case, x and y need to be rounded before performing the operation and

$$fl(x \text{ op } y) = fl(fl(x) \text{ op } fl(y)) = ((x(1+\epsilon_1)) \text{ op } (y(1+\epsilon_2)))(1+\epsilon_3),$$

where $|\epsilon_j| \leq u$, $j = 1, 2, 3$.

Rounding errors

Do similar statements hold for normalized numbers x and y that are not floating point numbers?

Suppose, $N_{\min} \leq |x|, |y|, |x \text{ op } y| \leq N_{\max}$. In such a case, x and y need to be rounded before performing the operation and

$$fl(x \text{ op } y) = fl(fl(x) \text{ op } fl(y)) = ((x(1+\epsilon_1)) \text{ op } (y(1+\epsilon_2)))(1+\epsilon_3),$$

where $|\epsilon_j| \leq u$, $j = 1, 2, 3$.

The answer depends on what 'op', x and y are.

Rounding errors

Do similar statements hold for normalized numbers x and y that are not floating point numbers?

Suppose, $N_{\min} \leq |x|, |y|, |x \text{ op } y| \leq N_{\max}$. In such a case, x and y need to be rounded before performing the operation and

$$fl(x \text{ op } y) = fl(fl(x) \text{ op } fl(y)) = ((x(1+\epsilon_1)) \text{ op } (y(1+\epsilon_2)))(1+\epsilon_3),$$

where $|\epsilon_j| \leq u$, $j = 1, 2, 3$.

The answer depends on what 'op', x and y are.

In the following

$O(u) \rightarrow$ quantities whose absolute values are small multiples of unit roundoff u ;

$O(u^2) \rightarrow$ quantities whose absolute values are small multiples of u^2 and can be ignored.

Rounding Errors

If $op = \times$, or $op = /$,

$$fl(x \ op \ y) = (x \ op \ y)(1 + \epsilon), \quad |\epsilon| \leq 3u + O(u^2) \quad (2)$$

Rounding Errors

If $op = \times$, or $op = /$,

$$fl(x \ op \ y) = (x \ op \ y)(1 + \epsilon), \quad |\epsilon| \leq 3u + O(u^2) \quad (2)$$

If $op = +$, or $op = -$,

$$fl(x \ op \ y) = (x \ op \ y) \left(1 + \left(\frac{x\epsilon_1}{x \ op \ y} \ op \ \frac{y\epsilon_2}{x \ op \ y} \right) \right), \quad (3)$$

where $|\epsilon_j| \leq 2u + O(u^2)$.

Rounding Errors

If $op = \times$, or $op = /$,

$$fl(x \ op \ y) = (x \ op \ y)(1 + \epsilon), \quad |\epsilon| \leq 3u + O(u^2) \quad (2)$$

If $op = +$, or $op = -$,

$$fl(x \ op \ y) = (x \ op \ y) \left(1 + \left(\frac{x\epsilon_1}{x \ op \ y} \ op \ \frac{y\epsilon_2}{x \ op \ y} \right) \right), \quad (3)$$

where $|\epsilon_j| \leq 2u + O(u^2)$.

However $\left| \frac{x\epsilon_1}{x \ op \ y} \ op \ \frac{y\epsilon_2}{x \ op \ y} \right| \not\approx O(u)$ if $|x \ op \ y| \ll |x|$ or $|x \ op \ y| \ll |y|$.

Rounding Errors

If $op = \times$, or $op = /$,

$$fl(x \ op \ y) = (x \ op \ y)(1 + \epsilon), \quad |\epsilon| \leq 3u + O(u^2) \quad (2)$$

If $op = +$, or $op = -$,

$$fl(x \ op \ y) = (x \ op \ y) \left(1 + \left(\frac{x\epsilon_1}{x \ op \ y} \ op \ \frac{y\epsilon_2}{x \ op \ y} \right) \right), \quad (3)$$

where $|\epsilon_j| \leq 2u + O(u^2)$.

However $\left| \frac{x\epsilon_1}{x \ op \ y} \ op \ \frac{y\epsilon_2}{x \ op \ y} \right| \not\approx O(u)$ if $|x \ op \ y| \ll |x|$ or $|x \ op \ y| \ll |y|$.

This is called *catastrophic cancellation* as it can result in sudden loss of accuracy.

Catastrophic Cancellation

Example 1: In the finite precision system $(10, 3, -1, 2)$ with round to nearest rounding, $\text{fl}(100 - 99.95) = 0$.

Example 2: The computation $z = 1 - \sqrt{1 - x^2}$ is prone to catastrophic cancellation for $x \approx 0$. The errors magnify when the computed value is multiplied by a large number.

This may be avoided by using the equivalent formulation $x^2/(1 + \sqrt{1 - x^2})$.

x	$10^{20}(1 - \sqrt{1 - x^2})$	$10^{20} \left(x^2/(1 + \sqrt{1 - x^2}) \right)$
$5.3105e - 007$	14100000	14101000
$4.7895e - 007$	11469000	11470000
$4.2684e - 007$	9114900	9109700
$3.7474e - 007$	7027700	7021400
$3.2263e - 007$	5206900	5204600
$2.7053e - 007$	3663700	3659200
$2.1842e - 007$	2387000	2385400
$1.6632e - 007$	1387800	1383000
$1.1421e - 007$	655030	652200
$6.2105e - 008$	199840	192850
$1e - 008$	11102	5000

Swamping \rightarrow Catastrophic Cancellation

If $0 < fl(a) \ll fl(b)$, or $fl(b) \ll 0 < fl(a)$, then $fl(a + b) \approx fl(b)$.

This is called *swamping*.

Swamping \rightarrow Catastrophic Cancellation

If $0 < fl(a) \ll fl(b)$, or $fl(b) \ll 0 < fl(a)$, then $fl(a + b) \approx fl(b)$.

This is called *swamping*.

Swamping can lead to *catastrophic cancellation*.

Swamping → Catastrophic Cancellation

If $0 < fl(a) \ll fl(b)$, or $fl(b) \ll 0 < fl(a)$, then $fl(a + b) \approx fl(b)$.
This is called *swamping*.

Swamping can lead to *catastrophic cancellation*.

Exercise: Perform GENP on

$$A = \begin{bmatrix} 0.1 & 10 & 9.985 \\ 1 & 0.05 & 0.15 \\ 1 & 0.04 & 0.19 \end{bmatrix}$$

in the finite precision system $(10, 3, -3, 2)$ with *round to nearest* as the rounding mode. Identify the computations where

- (a) swamping occurs;
- (b) catastrophic cancellation occurs;

Find L and U of the computed LU decomposition and compute $\|A - LU\|_1$. Repeat the calculations for GEPP.

Backward error Analysis and Stability of Algorithms

Some background

- ▶ In the early years of computing, the approach to understanding the accuracy of solution from algorithms involved bounding the rounding error at every stage of the computations.
- ▶ Apart from being practically very difficult in the presence of many computations, it was also not very successful as the risk of catastrophic cancellation which was not always possible to predict, loomed over the computations.
- ▶ Therefore, there was a general air of pessimism about the extent of errors in solutions computed via algorithms operating in a finite precision environment.
- ▶ In the early 1960s, a radically different approach to rounding error analysis was proposed by James Wilkinson.
- ▶ His idea was that instead of looking at errors at every stage of the computation, one should look at the computed answer from the algorithm as the exact answer of the same algorithm applied to perturbed data, the perturbations arising from a pushback of the errors in the computations back into data.

Backward Errors

- ▶ These relative perturbations to the data arising from the pushback are called the backward errors of the computed answer.
- ▶ If they are of the order of unit roundoff, then the algorithm is said to be backward stable or simply stable.
- ▶ Since working in finite precision environments and consequent rounding errors are inevitable for algorithms, the most that should be expected of them is that they are backward stable.
- ▶ After this, the accuracy of the computed solution depends upon the sensitivity of the problem to small changes in the data.
- ▶ The backward error analysis of the algorithm is combined with the sensitivity analysis of the problem to bound the relative errors in the solution.
- ▶ This approach separates the properties of the algorithm from those of the problem.
- ▶ To contrast with backward errors, the usual errors in computations are also referred to as forward errors.

Backward error analysis



**James Hardy Wilkinson, FRS
(1919-1986)**

Systematic analysis of backward errors and backward stability
was introduced by James Wilkinson.

Posteriori versus priori backward error analysis of GE

- ▶ The analysis of backward stability of algorithms is an analysis of the backward errors in the computations.
- ▶ For Gaussian Elimination one such (posteriori) analysis was already undertaken *after* solving a system of equations and using the residual vector associated with the computed solution to construct perturbed systems of equations of which the computed solution is an exact solution.
- ▶ This approach depends on the computed solution and can only say whether the algorithm is backward stable with respect to the given problem.
- ▶ To know about the backward stability property of an algorithm in general before using it to solve a problem a *priori backward error analysis* is required.

A priori backward error analysis

- ▶ This analysis is made without actually using the algorithm to solve a problem.
- ▶ It involves estimating maximum possible backward errors that can arise in standard arithmetic operations as an initial step.
- ▶ These are used to estimate maximum possible backward errors in the computed answer arising from the collective backward errors in the computations.
- ▶ The algorithm is declared to be backward stable if these maximum possible values of backward errors are $O(u)$.
- ▶ As the analysis should hold for any problem that the algorithm is designed to solve, it cannot make assumptions about the input data and computed solution.
- ▶ However, sometimes these backward errors are guaranteed to be of the order of unit roundoff only under certain conditions. In such cases the algorithm is *conditionally backward stable*.
- ▶ If no such conditions are required for backward stability, then the algorithm is *unconditionally backward stable*.

Backward errors in standard arithmetic operations

$$\frac{|fl(x \text{ op } y) - (x \text{ op } y)|}{|x \text{ op } y|} \rightarrow \text{relative (forward) error in } (x \text{ op } y).$$

Backward errors in standard arithmetic operations

$$\frac{|fl(x \text{ op } y) - (x \text{ op } y)|}{|x \text{ op } y|} \rightarrow \text{relative (forward) error in } (x \text{ op } y).$$

In contrast if

$$fl(x \text{ op } y) = \hat{x} \text{ op } \hat{y}$$

for some \hat{x}, \hat{y} , then,

$$\frac{|\hat{x} - x|}{|x|}, \frac{|\hat{y} - y|}{|y|} \rightarrow \text{relative backward errors in } (x \text{ op } y).$$

Backward errors in standard arithmetic operations

$$\frac{|fl(x \text{ op } y) - (x \text{ op } y)|}{|x \text{ op } y|} \rightarrow \text{relative (forward) error in } (x \text{ op } y).$$

In contrast if

$$fl(x \text{ op } y) = \hat{x} \text{ op } \hat{y}$$

for some \hat{x}, \hat{y} , then,

$$\frac{|\hat{x} - x|}{|x|}, \frac{|\hat{y} - y|}{|y|} \rightarrow \text{relative backward errors in } (x \text{ op } y).$$

The operation op is said to be *backward stable* if

$$\frac{|\hat{x} - x|}{|x|} \text{ and } \frac{|\hat{y} - y|}{|y|} \text{ are } O(u).$$

Backward errors in standard arithmetic operations

Since the relative representation error that arises when rounding normalized numbers with absolute values between N_{\min} and N_{\max} is very small, these may be ignored in the push back.

Backward errors in standard arithmetic operations

Since the relative representation error that arises when rounding normalized numbers with absolute values between N_{\min} and N_{\max} is very small, these may be ignored in the push back.

Therefore in backward error analysis it is assumed that all the standard arithmetic operations occur between floating point numbers.

Backward errors in standard arithmetic operations

Since the relative representation error that arises when rounding normalized numbers with absolute values between N_{\min} and N_{\max} is very small, these may be ignored in the push back.

Therefore in backward error analysis it is assumed that all the standard arithmetic operations occur between floating point numbers.

Theorem Let x, y be floating point numbers. Whenever op is any of the standard operations $+$, $-$, $/$, \times ,

$$fl(x \ op \ y) = \hat{x} \ op \ \hat{y}$$

where $\hat{x} = x(1 + \delta_1)$ and $\hat{y} = y(1 + \delta_2)$ with $|\delta_j| \leq u$, for $j = 1, 2$.

Backward errors in standard arithmetic operations

Since the relative representation error that arises when rounding normalized numbers with absolute values between N_{\min} and N_{\max} is very small, these may be ignored in the push back.

Therefore in backward error analysis it is assumed that all the standard arithmetic operations occur between floating point numbers.

Theorem Let x, y be floating point numbers. Whenever op is any of the standard operations $+, -, /, \times$,

$$fl(x \ op \ y) = \hat{x} \ op \ \hat{y}$$

where $\hat{x} = x(1 + \delta_1)$ and $\hat{y} = y(1 + \delta_2)$ with $|\delta_j| \leq u$, for $j = 1, 2$.

Hence all the standard arithmetic operations between two floating point numbers are backward stable operations.

A fundamental result in backward error analysis

Theorem Let $w_i, i = 1, \dots, n$, be floating point numbers. Then there exist $\gamma_i, i = 1, \dots, n$, satisfying $|\gamma_i| \leq (n-1)u + O(u^2)$, such that

$$\text{fl} \left(\sum_{i=1}^n w_i \right) = \sum_{i=1}^n w_i (1 + \gamma_i). \quad (4)$$

irrespective of the order of summation.

A fundamental result in backward error analysis

Theorem Let $w_i, i = 1, \dots, n$, be floating point numbers. Then there exist $\gamma_i, i = 1, \dots, n$, satisfying $|\gamma_i| \leq (n-1)u + O(u^2)$, such that

$$\text{fl} \left(\sum_{i=1}^n w_i \right) = \sum_{i=1}^n w_i (1 + \gamma_i). \quad (4)$$

irrespective of the order of summation.

Corollary Let $u_i, w_i, i = 1, \dots, n$, be floating point numbers. Then there exist $\gamma_i, i = 1, \dots, n$, satisfying $|\gamma_i| \leq nu + O(u^2)$, such that

$$\text{fl} \left(\sum_{i=1}^n u_i w_i \right) = \sum_{i=1}^n u_i w_i (1 + \gamma_i). \quad (5)$$

irrespective of the order of summation.