

Scientific Computing

Is an important part of of Computational Science that covers:

- Development of mathematical models (equations),
- **Development of algorithms to solve equations numerically,**
- **Implementation of algorithms in software,**
- **Numerical simulation of physical phenomena using computer software,**
- Interpretation, Validation and Visualisation of results.
- ← Redesign needed?

History of SC

- Most of the concepts formulated 200 years ago by Newton, Gauss, Euler, Jacobi and many others.
- The motivation was obtaining approximate solutions for mathematical problems that arose in physics, astronomy and other fields of science.
- Efficient use of computational resources (pencil, paper, brain power),
- with the advent of computers the problem sizes are increasing,
- rounding errors are becoming critical, because the precision is not under human control,
- computation (simulation of reality) is becoming as important as measurements and theory.

Propagation Error

- Let e express the relative error in representing a nonzero floating point number.
- The sum: $a(1 \pm e) + b(1 \pm e) = (a \pm ae) + (b \pm be) = (a + b) \pm e(a + b) = (a+b)(1 \pm e)$
The sum error is in the same range as the error of factors. Similar is valid for difference, but suppose that a and b are of similar values then $a-b \approx 0$ and e may become as large as result!
- The product: $a(1 \pm e) \cdot b(1 \pm e) = (a \pm ae) \cdot (b \pm be) = ab \pm abe \pm bae \pm abe^2 \approx ab(1 \pm 2e)$
We get double error what leads to the pessimistic estimation of propagation error. Similar is valid for division $((1 \pm e)^{-1} = 1 \pm e + e^2 \pm \dots \approx 1 \pm e)$.

Propagation Error (Cont.)

- For an exponent of n the relative error of result is n -times greater than initial error:
 $a(1 \pm e)^n = a^n [(1 \pm ne \pm n(n-1)e^2/2! \pm \dots) \approx a^n (1 \pm ne)$
- What happens if the exponent is smaller than 1? Is the error in result smaller? No.
- Similar is valid for m consecutive multiplication or division operations. Error can increase in the worst case for a factor of m .
- But in practical consecutive calculations we usually desire that the final result should have the similar absolute error as the less accurate input data.

Floating-Point Numbers

- Floating-point number system represents approximately the real number system.
- Floating point numbers are used in a similar way as scientific notation, with an exponent.
- Examples:
 $2347 = 2.347 \cdot 10^3$,
 $0.0007396 = 7.396 \cdot 10^{-4}$.
- The name floating-point is used because the decimal point floats as the power of 10 changes.

Floating-Point Numbers IEEE SP-normalized

- radix or base $b = 2$,
- precision $p=24$ (23 bits for mantissa, 1bit for sign)
- exponent range $[L=-126, U=127]$ (8 bits for exponent)
- 32 bits used for representation,
 $(1 \ 00110001101001001001101 \ 10100101)_2 =$
 $= - (2^{-3} + 2^{-4} + 2^{-8} + 2^{-9} + 2^{-11} + 2^{-14} + 2^{-17} + 2^{-20} + 2^{-21} + 2^{-23}) \cdot 2^{-37} =$
 $(- 1.4113822957573241012596554355696 \cdot 10^{-12})_{10}$
- the least significant bit in mantissa $1 \rightarrow 0$, the gap between two consecutive numbers:
 $2^{-23} \cdot 2^{-37} = 2^{-60} =$
 $= 8.6736173798840354720596224069 \cdot 10^{-19}$
- gaps are equally spaced between powers of b , but become smaller and smaller if we approaching to zero.

