

Why Numbers Lie: Floating Point Errors and How to Tame Them

July 2025

Pragya Goyal

In numerical computation, precision is both a goal and a challenge. Due to the finite nature of digital systems, real numbers are represented using *floating point arithmetic*, a format that inherently introduces approximation and rounding errors. This article delves into the foundations of *error analysis*—from the role of *machine epsilon* and the dangers of *catastrophic cancellation* to strategies like *pairwise summation* and *partial pivoting* in Gaussian elimination. Through clear examples and practical techniques, we explore how seemingly minor computational errors can accumulate, and how careful algorithm design can significantly enhance numerical stability and accuracy.

1. Floating Point Representation: Approximating Reality Digitally

Computers cannot store all real numbers exactly. Instead, they use floating point representation, a finite approximation of the real number set ($F \subseteq \mathbb{R}$). Every number is expressed as:

$$x = \pm m \times \beta^e$$

Where:

- β is the base (commonly 2 in binary systems),
- t is the precision (number of digits),
- e is the exponent, constrained between e_{\min} and e_{\max} ,
- m is the mantissa or significand, where $0 < m < \beta^{t-1}$.

This system supports a wide range of values but inevitably introduces rounding errors due to limited precision.

2. Machine Epsilon: Quantifying Precision Limits

A fundamental concept in error analysis is machine epsilon (ϵ)—the smallest number such that:

$$1 + \epsilon > 1$$

in the machine's floating point system.

It acts as an upper bound on relative rounding error in computations, and is essential for evaluating the precision and reliability of numerical results.

3. Addition and Subtraction: Precision Pitfalls

When adding or subtracting floating point numbers with different exponents, the smaller number is shifted to align exponents—potentially losing significant digits due to limited precision.

To mitigate this, a guard digit is used. It temporarily preserves an extra digit during intermediate steps to reduce the error. Ultimately, the result is rounded to fit the allowed precision.

4. Catastrophic Cancellation: Subtracting Dangerously

Catastrophic cancellation occurs when subtracting nearly equal numbers. Their significant digits cancel, leaving a result with poor precision.

For example, computing:

$$x^2 - y^2$$

vs.

$$(x - y)(x + y)$$

The factored form is more numerically stable, especially when $x \approx y$, since it avoids subtracting similar numbers and reduces error accumulation. Each step in a calculation introduces potential error, so reformulating expressions thoughtfully is crucial.

5. Error Propagation in Summation

When summing a sequence, rounding errors can compound at each step. If each operation has a small multiplicative error $(1 + \delta)$, summing x_1, x_2, x_3, \dots yields:

$$((x_1 + x_2)(1 + \delta) + x_3)(1 + \delta) \text{ and so on.}$$

Expanding this shows:

$$x_1(1 + \delta)^2 + x_2(1 + \delta)^2 + x_3(1 + \delta) + \dots$$

This means earlier terms accumulate more error. To minimize this, sum smaller numbers first, especially when all values are positive—so that the larger terms, which carry more weight, are affected by fewer compounded errors.

6. Pairwise Summation: Improving Accuracy

Pairwise summation is a recursive method that reduces error propagation:

- a) Split the list into pairs.
- b) Add each pair.
- c) Recursively sum the resulting values until one result remains.

This forms a balanced binary tree, minimizing the number of operations each number undergoes, and is significantly more accurate than naive left-to-right summation.

7. Matrix Multiplication: The Diver's Perspective

Matrix multiplication involves a complex series of arithmetic operations, each prone to rounding error. For matrices:

$$\begin{bmatrix} a & b \\ c & d \end{bmatrix} \begin{bmatrix} e & f \\ g & h \end{bmatrix}$$

The product is:

$$\begin{bmatrix} ae + bg & af + bh \\ ce + dg & cf + dh \end{bmatrix}$$

Understanding how each intermediate term contributes to overall error—especially with larger matrices—is vital for building numerically stable algorithms. Techniques like the diver's method (a systematic breakdown of computation steps) help track and limit such errors.

8. Gaussian Elimination: Solving Systems with Care

Gaussian Elimination solves systems of linear equations by transforming the coefficient matrix into upper triangular form using row operations:

- a) Row swapping,
- b) Scaling a row to make the pivot element 1,
- c) Row replacement to eliminate entries below the pivot.

Once in upper triangular form, the system is solved via back-substitution. However, the method is sensitive to rounding—especially when small pivot elements are used.

To reduce numerical error, partial pivoting is applied: at each step, the row with the largest available coefficient in the pivot column is swapped into place.

9. Conclusion: Why Error Analysis Matters

Whether you're simulating real-world systems, analyzing data, or developing numerical algorithms, understanding how errors originate and propagate is vital. Concepts like machine epsilon, catastrophic cancellation, and pairwise summation are not just academic—they're the tools that make computation accurate and trustworthy.

In numerical computing, the goal isn't just to compute a number—it's to compute the right number. That's the essence of error analysis.