

VALIDATED NUMERICS #1 by W.T.

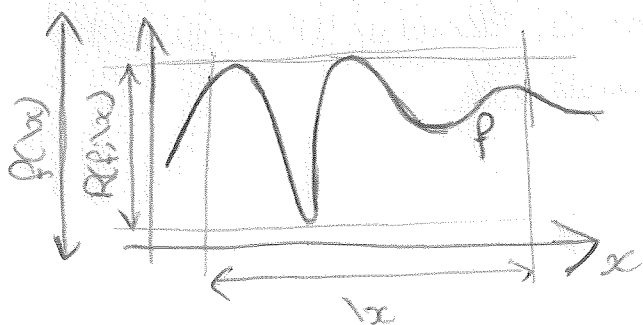
(mathematics) \cap (Scientific Computing)

Goals: Computational proofs aimed at continuous problems.

- We will use "amp. repr." sets (intervals) to explore the continuum of \mathbb{R} .
- Operations on intervals are described by interval analysis.
- Gives us the ability to enclose the range of a function: $R(f; \alpha) = \{f(x); x \in \alpha\}$

$$R(f; \alpha) \subseteq f(\alpha)$$

← enclosure, computable.



Applications: evaluation, equation solving, global optimization, quadrature, ode/bvp/pde(s)

Strategy: Given a "purely" mathematical problem, break it up into a mathematical part and a finite, numerical part.

Example: compute $S = \sum_{k=1}^{\infty} 1/k^2$ to 12 correctly rounded digits. Then we write:

$$S = \underbrace{\sum_{k=1}^N 1/k^2}_{S_n} + \underbrace{\sum_{k=N+1}^{\infty} 1/k^2}_{S_n^*} \quad \text{where } S_n^* \in \left[\frac{1}{N+1}, \frac{1}{N} \right]$$


Any real number x can be expressed in base β as $x = \sigma \sum_{k=0}^n b_k \beta^k = \sigma (b_n b_{n-1} \dots b_0 b_{-1} \dots)_{\beta}$. This expression is not unique. First, we use scientific notation $x = \sigma (\underbrace{b_0 b_1 b_2 \dots}_{\text{mantissa}})_{\beta} \beta^e$. There were props:

① $\forall i$, we have $0 \leq b_i \leq \beta - 1$; ② For $x \neq 0$, $1 \leq b_0 \leq \beta - 1$.

③ For as many i , $0 \leq b_i \leq \beta - 2$; These props. gives a unique repr. of x .

More conditions (4) $\check{e} \leq e \leq \hat{e}$; (5) $m = (b, b, \dots, b_{p-1})$ where p is precision.

Properties (1)-(5) gives $\mathbb{F} = \mathbb{F}(\beta, p, \check{e}, \hat{e})$, a finite set of "floating point numbers"

$\mathbb{F}(2, 3, -1, +1) =$  There is a lot of space around zero.

(3') for $e > \check{e}$ $1 \leq b_0 \leq \beta - 1$. When $e = \check{e}$, b_0 is allowed to be zero (subnormal number)

With subnormal numbers, $x = y \Leftrightarrow x \ominus y = 0$
 $x, y \in \mathbb{F}$

Problem: \mathbb{F} is not closed under arithmetic operations: for $x, y \in \mathbb{F}$, $x * y \notin \mathbb{F}$, $x \in \mathbb{F}, y \notin \mathbb{F}$

Rounding: operation $\circ: \mathbb{R} \rightarrow \mathbb{F}$ should satisfy: (R1) $x \in \mathbb{F} \Rightarrow \circ(x) = x$; (R2) $x, y \in \mathbb{R}, x \leq y \Rightarrow \circ(x) \leq \circ(y)$

Examples: "Truncate" (Round towards 0), "Round up" (Round towards $+\infty$), "Round down" (Round towards $-\infty$), Round towards nearest even.

Theorem ~~If~~ If x and y are normal fp numbers with $x * y \neq 0$, then

$$\left| \frac{x * y - x \circledast y}{x * y} \right| < \epsilon_M$$

where ϵ_M is machine epsilon, with $\epsilon_M = \beta^{-(p-1)}$.

There is no information when two operations are done successively. Because the operands of the 2nd operation are not fp numbers a priori.