

**Note** Floating-point numbers have limitations:

- limitations on significand

- limitations on exponent

- ↳ **overflow** <sup>→ set as Inf</sup> ⇒ exponent is too large (positive)

- ↳ **underflow** <sup>→ set as zero</sup> ⇒ exponent is too large (negative)

**Rmk.** With proper scaling, overflow / underflow can often be avoided.

## 1.2 Errors

Suppose  $p^* \in \mathbb{R}$  is an approximation of  $p \in \mathbb{R}$ .

- **absolute error:**  $e_a(p, p^*) = |p - p^*|$

- **relative error:**  $e_r(p, p^*) = \frac{|p - p^*|}{|p|}$  for  $p \neq 0$

Error bounds

- **absolute error bound:**  $e_a(p, p^*) \leq \varepsilon_a(p, p^*)$

- **relative error bound:**  $e_r(p, p^*) \leq \varepsilon_r(p, p^*)$

**Rmk.** Often we can only obtain a bound on error produced by algorithms.

Ways to reduce errors in finite digit precision:

I. Avoid subtraction of 2 nearly equal numbers

- reason: causes cancellation of significant digits (catastrophic cancellation)

**Ex. 1** Given 2 numbers  $x$  and  $y$ , with  $x > y$  and  $k$ -digit representation:

Then  $fl(x) = 0.d_1d_2 \dots d_p \alpha_{p+1} \alpha_{p+2} \dots \alpha_k \times 10^n$

$$fl(y) = 0.d_1d_2\dots d_p \beta_{p+1}\beta_{p+2}\dots \beta_n \times 10^n$$

$$\Rightarrow fl(fl(x) - fl(y)) = 0.\underbrace{00\dots 0}_{p \text{ times}} \sigma_{p+1}\sigma_{p+2}\dots \sigma_k \times 10^n$$

$$= 0.\sigma_{p+1}\sigma_{p+2}\dots \sigma_k \times 10^{n-p}$$

Then  $fl(fl(x) - fl(y))$  has  $k-p$  significant digits (loss of accuracy).

Ex. 2  $f(x) = \frac{1 - \cos x}{x^2}$

Fact  $0 < f(x) \leq \frac{1}{2}$ ,  $\lim_{x \rightarrow 0} f(x) = \lim_{x \rightarrow 0} \frac{\sin(x)}{2x}$  (L' Hospital)

$$= \lim_{x \rightarrow 0} \frac{\cos(x)}{2} = \frac{1}{2}$$

On MATLAB:  $f(1.2 \times 10^{-8}) \approx 0.77098$

$$f(1.2 \times 10^{-7}) = 0 \quad \text{NOT } \approx \frac{1}{2}!$$

Remedy: use trig identity  $\cos x = 1 - 2\sin^2(\frac{x}{2})$

$$\rightarrow f(x) = \frac{1 - \cos x}{x^2} = \frac{2\sin^2(\frac{x}{2})}{x^2}$$

(no longer requires subtraction in numerator)

## I. Avoid division by small numbers or multiplying by large numbers

-reason: causes overflow

Ex. Consider  $c = \sqrt{a^2 + b^2}$

If  $a = 10^{170}$ ,  $b = 1$ , then correctly rounded solution:  $c = 10^{170}$

However, with double-precision arithmetic:

$$a^2 = \text{Inf}, \quad a^2 + b^2 = \text{Inf} + 1 = \text{Inf}, \quad c = \sqrt{\text{Inf}} = \text{Inf}$$

Remedy: scale the data

$$c = s \sqrt{\left(\frac{a}{s}\right)^2 + \left(\frac{b}{s}\right)^2}, \quad \text{where } s = \max\{|a|, |b|\}$$

$$\text{Here, } s = 10^{170} \text{ and } c = 10^{170} \sqrt{(1)^2 + \left(\frac{1}{s}\right)^2} = 10^{170}$$

↳ underflow to 0 (ok)

## III. Reduce the number of arithmetic computations (+, -, \*, ÷)

reason: more computation  $\rightarrow$  more rounding errors

Ex. Evaluate  $f(x) = x^3 - 6.1x^2 + 3.2x + 1.5$  at  $x = 4.71$   
mult. 2 2 1  
add. 1 1 1 } 8

Now consider nested formulation:

$$f(x) = (x^3 - 6.1x^2 + 3.2x) + 1.5$$

$$= (x^2 - 6.1x + 3.2)x + 1.5$$

$$= ((x - 6.1)x + 3.2)x + 1.5 \quad 2 \text{ multiplications} + 3 \text{ additions}$$

$\rightarrow$  Horner's Method (Ch.2)

Truncation error: truncate infinite sum by finite sum

Thm. Taylor's Theorem

$$f(x) = \underbrace{f(x^*) + f'(x^*)(x - x^*) + \frac{f''(x^*)}{2!}(x - x^*)^2 + \dots + \frac{f^{(n)}(x^*)}{n!}(x - x^*)^n}_{p_n(x)} + \underbrace{\frac{f^{(n+1)}(\xi)}{(n+1)!}(x - x^*)^{n+1}}_{R_n(x)}$$

where  $\xi$  is between  $x$  and  $x^*$

Truncation error:  $R_n(x) = f(x) - p_n(x)$