

MATH50003

Numerical Analysis

II.3 Floating Point Arithmetic

Dr Sheehan Olver

Part II

Representing Numbers

1. **Reals** via floating point
2. **Floating point arithmetic** and bounding errors
3. **Interval arithmetic** for rigorous computations

Rounding

How does a computer round a real to a float?

$$x = \pm 2^{q-g} (1 \text{ } b_1 b_2 \text{ } \text{---} \text{ } b_s b_{s+1} b_{s+2} \text{ } \text{---})_2$$

Definition 7 (rounding). $\text{fl}_{\sigma, Q, S}^{\text{up}} : \mathbb{R} \rightarrow F_{\sigma, Q, S}$

$$\text{fl}^{\text{up}}(x) := \min \{ y \in F \quad \text{s.t.} \quad y \geq x \}$$

$\in \mathcal{F}$ $\text{fl}_{16}^{\text{up}}(1/3) = \text{fl}^{\text{up}}(2^{-2} \times (1.0101010101 \text{ } 0101 \text{ } \text{---})_2)$
 15, 5, 10 $= 2^{-2} \times (1.0101010110)_2$

$\text{fl}_{\sigma, Q, S}^{\text{down}} : \mathbb{R} \rightarrow F_{\sigma, Q, S}$

$$\text{fl}^{\text{down}}(x) := \max \{ y \in F \quad \text{s.t.} \quad y \leq x \}$$

$\in \mathcal{F}$ $\text{fl}^{\text{down}}(1/3) = 2^{-2} (1.0101010101)_2$

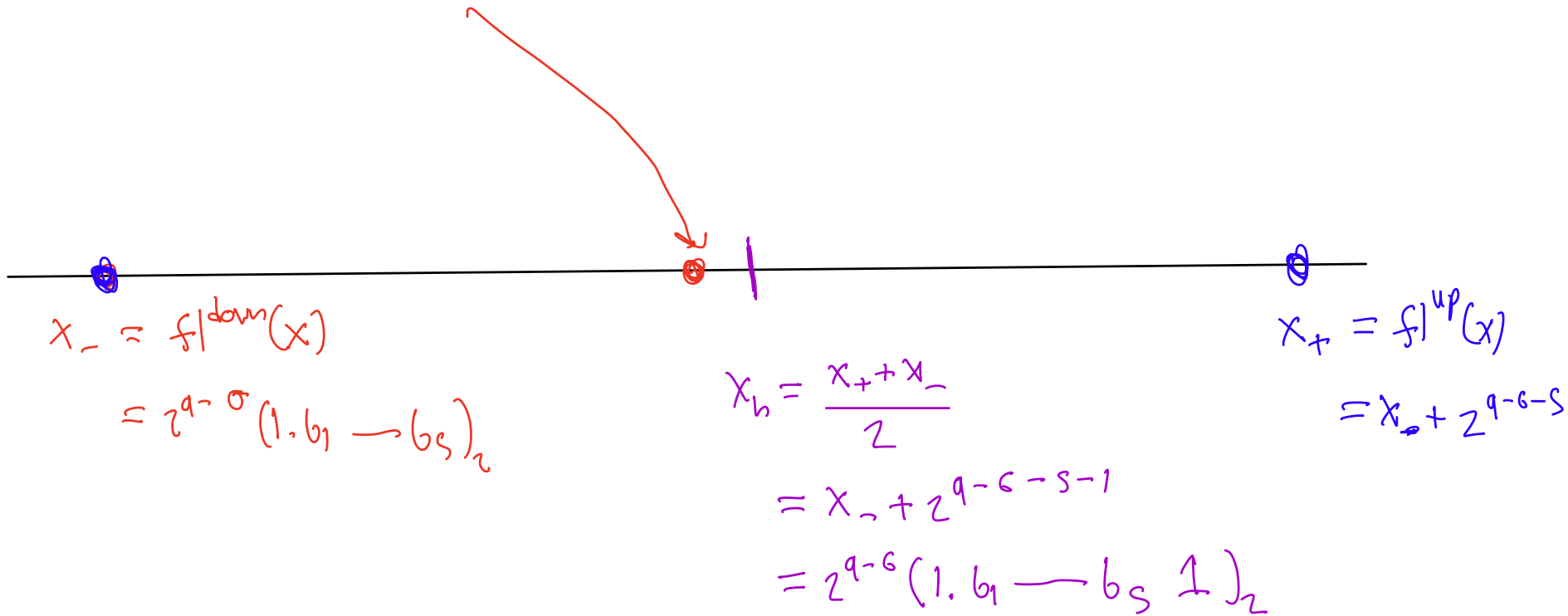
Note $\text{fl}^{\text{up}}(-x) = - \text{fl}^{\text{down}}(x)$

Default is round-to-nearest.

$$\underbrace{f_{\sigma, Q, S}^{\text{nearest}}}_{=: f|} : \mathbb{R} \rightarrow F_{\sigma, Q, S}$$

Consider

$$x = 2^{q-g} (1. b_1 b_1 \dots b_s b_{s+1} b_{s+2} \dots)_2 > 0$$



Different cases:

down $x_- \leq x < x_h \Rightarrow f|^{near}(x) = f|^{down}(x) = x_-$

up $x_b < x \leq x_+ \Rightarrow f|^{near}(x) = f|^{up}(x) = x_+$

half $x = x_b \Rightarrow f|^{near}(x) = \begin{cases} x_- & \text{if } b_s = 0 \\ x_+ & \text{otherwise} \end{cases}$

so that $f|^{near}(x)$ has last bit 0,

For $x < 0$

$$f|_G(x) = -f|(-x)$$

Arithmetic

Operations are exact up to rounding

$$x \oplus y := \text{fl}(x + y)$$

$$x \ominus y := \text{fl}(x - y)$$

$$x \otimes y := \text{fl}(x * y)$$

$$x \oslash y := \text{fl}(x / y)$$

$$x \uparrow y := \text{fl}(x^y)$$

← "↑" on a computer adds
exactly, then rounds

Eg: in F_{16} ,

$$\begin{aligned} 1 \uplus 2^{-11} &= \text{fl}(1 + 2^{-11}) \\ &= \text{fl}((1.000000000001)_2) \\ &= 1 \end{aligned}$$

Here, $x, y \in F$, not general reals!

Example 8 (decimal is not exact).

What's $1.1 + 0.1$ on a computer?

When we type "1.1" it makes $fl(1.1)$.

So this is

$$fl(1.1) \oplus fl(0.1)$$

We have

$$\begin{aligned} fl(1.1) &= fl((1.0001100110 \overset{\text{red arrow}}{\underset{b_{s+1}=0 \Rightarrow \text{round down}}{01100}})_2) \\ &= (1.000110011)_2 \end{aligned}$$

$$\begin{aligned} fl(0.1) &= fl(2^{-4} (1.1001100110 \overset{\text{red arrow}}{\underset{b_{s+1}=0 \Rightarrow \text{round down}}{0110}})_2) \\ &= 2^{-4} (1.100110011)_2 \end{aligned}$$

\Rightarrow

$$f_1(1.1) \oplus f_1(0.1) = f_1(f_1(1.1) + f_1(0.1))$$

$$= f_1((1.0011001100 \text{ } 011)_2)$$

round
down

↘

$$= (1.00110011)_2$$

But

$$f_1(1.2) = f_1(1.0011001100 \text{ } 110011 \text{ } \text{---})_2$$

↑
 $b_{s+1} = 1$

round
up

→

$$= (1.0011001101)_2$$

II.2.1 Bounding errors

Analysis on rounding errors

Definition 8 (machine epsilon/smallest positive normal number/largest normal number).

Machine epsilon is denoted

$$\epsilon_{m,S} := 2^{-S}.$$

For F_{64} , $S = 52 \Rightarrow \epsilon_m = 2^{-52}$

$$\approx 2.22 \times 10^{-16}$$

Definition 9 (normalised range). The *normalised range* $\mathcal{N}_{\sigma,Q,S} \subset \mathbb{R}$ is the subset of real numbers that lies between the smallest and largest normal floating-point number:

$$\mathcal{N}_{\sigma,Q,S} := \{x : \min |F_{\sigma,Q,S}^{\text{normal}}| \leq |x| \leq \max F_{\sigma,Q,S}^{\text{normal}}\}$$

smallest
positive
normal

largest
normal

$$x \in \mathcal{N}_{\sigma,Q,S} \Rightarrow fl^{\text{mode}}(x) \in F^{\text{normal}}$$

Proposition 2 (round bound). If $x \in \mathcal{N}$ then

$$\text{fl}^{\text{mode}}(x) = x(1 + \delta_x^{\text{mode}}) = x + \overbrace{x \delta_x^{\text{mode}}}^{\text{relative error}}$$

where the relative error is bounded by:

$$|\delta_x^{\text{nearest}}| \leq \frac{\epsilon_m}{2}$$

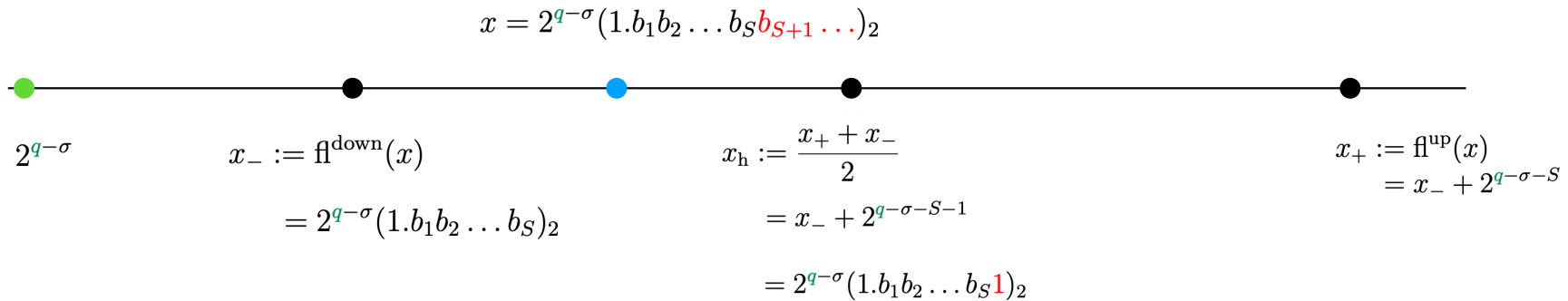
$$|\delta_x^{\text{up/down}}| < \epsilon_m.$$

← machine $\epsilon \approx 2.22 \times 10^{-16}$

Proof Only show mode = near, Assume $x > 0$.

Suppose $fl(x) = fl^{\text{down}}(x)$, i.e. $x_- \leq x \leq x_h$.

(Round Down)



We have

$$fl(x) = x_- = x \left(1 + \underbrace{\frac{x_- - x}{x}}_{\delta_x} \right)$$

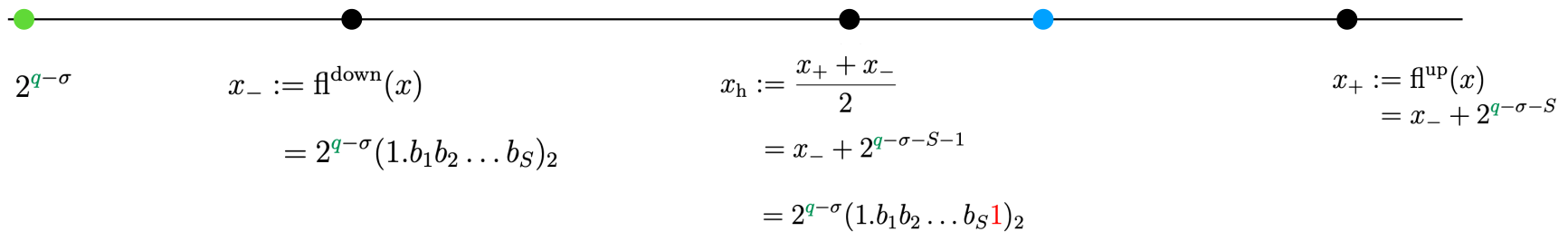
where

$$|\delta_x| \leq \frac{x - x_-}{x} \stackrel{\substack{\text{since } x \leq x_h \\ \downarrow}}{\leq} \frac{x_h - x_-}{x} \stackrel{\substack{x \geq x_- \\ \downarrow}}{\leq} \frac{x_h - x_-}{x_-}$$

$$\stackrel{\text{Round Up}}{\leq} \frac{2^{q-\sigma-S-1}}{2^{q-\sigma}} = 2^{-S-1} = \frac{\epsilon_m}{2}$$

$x_- \geq 2^{q-\sigma}$

$$x = 2^{q-\sigma}(1.b_1b_2\dots b_S \textcolor{red}{b}_{S+1}\dots)_2$$



Example 9 (bounding a simple computation).

$$(1.1 + 1.2) \times 1.3 = 2.99$$

What's a bound on error when done on computer?

$$\underbrace{[f_1(1.1) \oplus f_1(1.2)]}_{(*)} \otimes f_1(1.3) = 2.99 + \delta$$

absolute
error

Show

$$|\delta| \leq 23 \epsilon_m \leq 50 \times 10^{-16}$$

pessimistic bound

$$f_1(1.1) = 1.1 (1 + \delta_1)$$

$\epsilon_{N,Q,S}$

round
bound

$$|\delta_1| \leq \frac{\epsilon_m}{2}$$

$$f(1.2) = 1.2(1 + \delta_2)$$

$|\delta_2| \leq \frac{\epsilon_m}{2}$

\Rightarrow

$$(*) = f(\underbrace{1.1(1 + \delta_1)}_{f(1.1)} + 1.2(1 + \delta_2))$$

from round bound
↙

$$= (2.3 + 1.1\delta_1 + 1.2\delta_2)(1 + \delta_3)$$

$|\delta_3| \leq \frac{\epsilon_m}{2}$

$$= 2.3 + 1.1\delta_1 + 1.2\delta_2 + 2.3\delta_3 + 1.1\delta_1\delta_3 + 1.2\delta_2\delta_3$$

true
1.1 + 1.2

$=: \epsilon_1$

where

$$|e_1| \leq \underbrace{(1.1)}_{\leq 1} + \underbrace{1.2}_{\leq 2} + \underbrace{2.3}_{\leq 3} \frac{\epsilon_m}{2} + \underbrace{(1.1)}_{\leq 1} + \underbrace{1.2}_{\leq 2} \frac{\epsilon_m^2}{4}$$

$$\leq \frac{\epsilon_m}{4} 2^{-5} \leq \frac{\epsilon_m}{4}$$

$$\leq (2+2+3+1+1) \frac{\epsilon_m}{2} \leq 5 \epsilon_m.$$

error from
 \otimes
 \downarrow

\Rightarrow

(*) \otimes $f(1.3) = \underbrace{(2.3 + \epsilon_1)}_{1.3(1+\delta_4)} 1.3(1+\delta_4)(1+\delta_5)$

$$= 1.99 + \underbrace{1.3}_{\leq \frac{3}{2}} (\underbrace{\epsilon_1}_{\leq 3} + \underbrace{2.3}_{\leq 3} \delta_4 + \underbrace{2.3}_{\leq 3} \delta_5 + \underbrace{\epsilon_1 \delta_4}_{\leq \frac{5\epsilon_m}{2} \leq \frac{5}{2} \epsilon_m} + \underbrace{\epsilon_1 \delta_5}_{\leq \frac{5\epsilon_m}{2}} + \underbrace{2.3}_{\leq 3} \underbrace{\delta_4 \delta_5}_{\leq \frac{\epsilon_m}{4}} + \underbrace{\epsilon_1 \delta_4 \delta_5}_{\leq \frac{5\epsilon_m}{4}})$$

S

where $|s| \leq 23 \epsilon_m$



II.2.2 Idealised floating point

A simplified model for analysis

Definition 10 (idealised floating point). An idealised mathematical model of floating point numbers for which the only subnormal number is zero can be defined as:

$$F_{\infty, S} := \{\pm 2^q \times (1.b_1b_2b_3 \dots b_S)_2 : q \in \mathbb{Z}\} \cup \{0\}$$

Handwritten annotations:
- A red bracket under 2^q is labeled $Q=W$.
- A red arrow points from q to the word "anything".
- A red bracket under $b_1b_2b_3 \dots b_S$ is labeled S digits.

Bound bound applies $\forall x \in F_{\infty, S}$

II.2.3 Divided differences floating point error bound

Explain the unexplained error in divided differences

General model of a function implemented in floating point:

$\forall x,$

$$f(x) = f^{\text{FP}}(x) + \delta_x^f$$

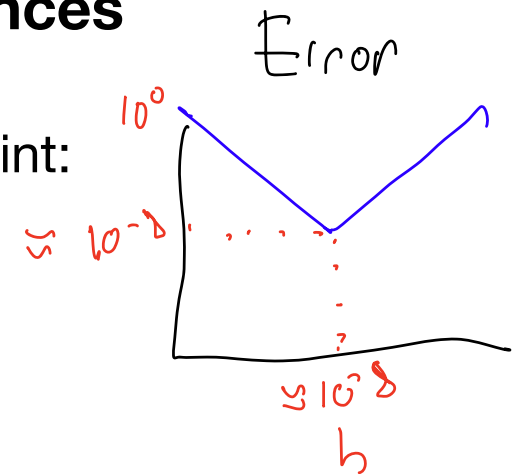
absolute error

such that

$$|\delta_x^f| \leq c\epsilon_m$$

\uparrow
some

reasonable c



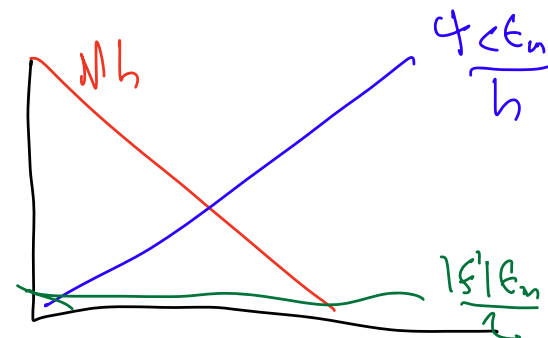
Theorem 4 (divided difference error bound). For $x \in F_{\infty, S}$,

$$\frac{f^{\text{FP}}(x+h) \ominus f^{\text{FP}}(x)}{h} = f'(x) + \delta_{x,h}^{\text{FD}}$$

where

$$|\delta_{x,h}^{\text{FD}}| \leq \underbrace{\frac{|f'(x)|}{2} \epsilon_m}_{\substack{\text{small} \\ \approx 10^{-16}}} + \underbrace{Mh}_{\substack{\rightarrow 0 \\ h \rightarrow 0}} + \underbrace{\frac{4c\epsilon_m}{h}}_{\substack{\rightarrow \infty \\ h \rightarrow 0}}$$

for $M = \sup_{x \leq t \leq x+h} |f''(t)|$.



Of course, bound is just upper bound so
haven't shown it will blow up.

Proof

$$\frac{f^{\text{FP}}(x+h) \ominus f^{\text{FP}}(x)}{h} = \frac{(f(x+h) - \int_{x+h}^x f - f(x) + \int_x^x f)}{h} (1 + \delta_1)$$

$|\delta_1| \leq \frac{6\epsilon}{2}$

$$= \underbrace{\frac{f(x+h) - f(x)}{h} (1 + \delta_1)}_{\text{Prop 1}} + \frac{f'_x - f'_{x+h}}{h} (1 + \delta_1)$$

Prop 1

$$= f'(x) - \delta^{dd}$$

where

$$|\delta^{dd}| \leq \frac{Mh}{2}$$

$$= f'(x) + \underbrace{f'(x)\delta_1}_{\substack{|\cdot| \leq \frac{|f'(x)|}{2} \epsilon_m}} - \underbrace{\delta^{dd} (1 + \delta_1)}_{\substack{\leq |\delta^{dd}| (1 + |\delta_1|) \\ \leq Mh}}$$

$$+ \frac{f'_x - f'_{x+h}}{h} (1 + \delta_1)$$

$$1.1 \leq \frac{|f_x^f| + |f_{x+h}^f|}{h} (1 + |f_1|) \leq \frac{4\epsilon_m}{h}$$

Corollary 2 (divided differences in practice). We have (Non examinable)

$$(f^{\text{FP}}(x \oplus h) \ominus f^{\text{FP}}(x)) \oslash h = \frac{f^{\text{FP}}(x+h) \ominus f^{\text{FP}}(x)}{h}$$

whenever $h = 2^{j-n}$ for $0 \leq n \leq S$ and the last binary place of $x \in F_{\infty, S}$ is zero, that is $x = \pm 2^j(1.b_1 \dots b_{S-1}0)_2$.

No error in \oslash since
$$\frac{\pm 2^q (1.b_1 \dots b_S)_2}{2^{j-n}}$$

$$= \pm 2^{q+j-n} (1.b_1 \dots b_S)_2 \in F_{\infty, S}$$

Also, since $b_S = 0$ we have

$$x \oplus h = x + h.$$

Heuristic (divided difference with floating-point step)

What's the best choice of h ?

Balance

$$Mh \lesssim \frac{4c \epsilon_m}{h} \Rightarrow$$

$$h^2 \lesssim \frac{4c}{M} \epsilon_m$$

\uparrow might not know

$$\Rightarrow h = \text{const.} \sqrt{\epsilon_m}$$

$$\approx 1.49 \times 10^{-4}$$



Now to Lab 3
To see rounding modes.