

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?

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

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

$$\mathfrak{f}_{\sigma,Q,S}^{\text{nearest}} : \mathbb{R} \rightarrow F_{\sigma,Q,S}$$

Arithmetic

Operations are exact up to rounding

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

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

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

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

Example 8 (decimal is not exact).

II.2.1 Bounding errors

Analysis on rounding errors

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

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

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:

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

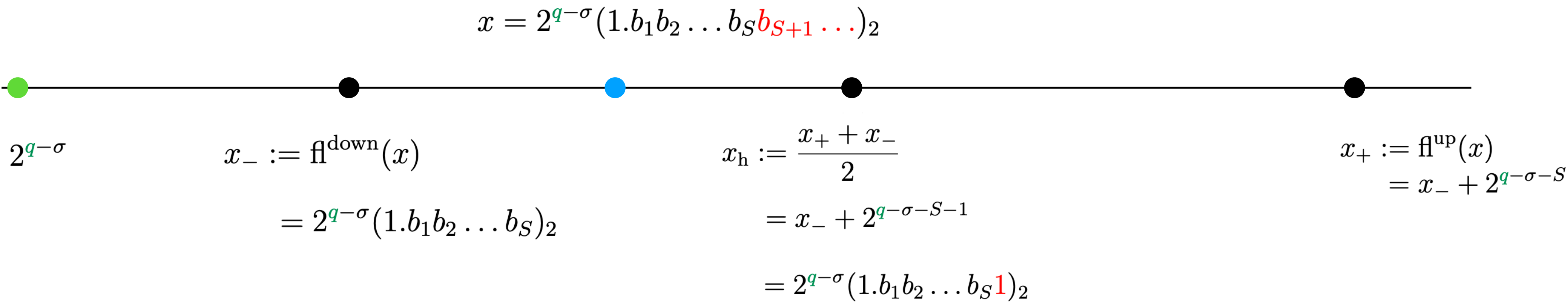
Proposition 2 (round bound). *If $x \in \mathcal{N}$ then*

$$\mathfrak{fl}^{\text{mode}}(x) = x(1 + \delta_x^{\text{mode}})$$

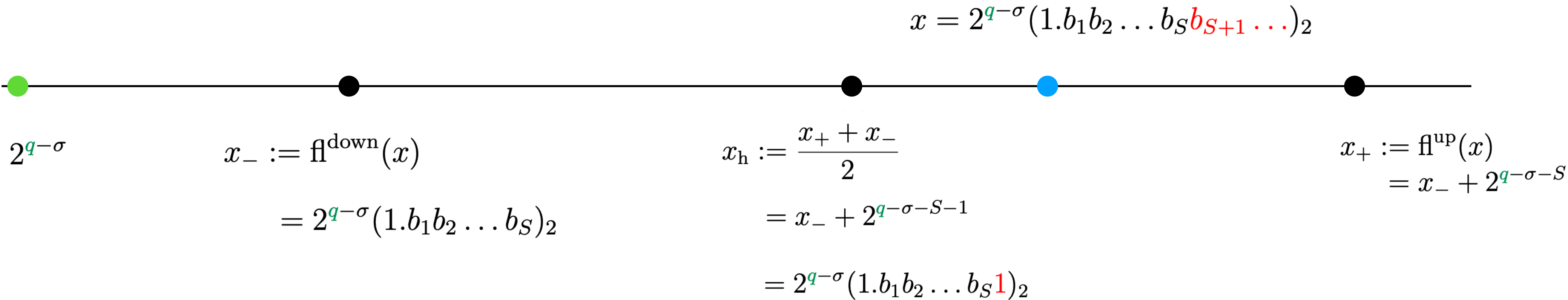
where the relative error is bounded by:

$$\begin{aligned} |\delta_x^{\text{nearest}}| &\leq \frac{\epsilon_{\text{m}}}{2} \\ |\delta_x^{\text{up/down}}| &< \epsilon_{\text{m}}. \end{aligned}$$

(Round Down)



(Round Up)



Example 9 (bounding a simple computation).

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:

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

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:

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

such that

$$|\delta_x^f| \leq c\epsilon_{\text{m}}$$

Theorem 4 (divided difference error bound).

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

where

$$|\delta_{x,h}^{\text{FD}}| \leq \frac{|f'(x)|}{2} \epsilon_{\text{m}} + Mh + \frac{4c\epsilon_{\text{m}}}{h}$$

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

Corollary 2 (divided differences in practice). *We have*

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

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 $x = \pm 2^j(1.b_1 \dots b_{S-1}0)_2$.

Heuristic (divided difference with floating-point step)

Now to Lab 3
To see rounding modes.