# Digits and Rounding

> Establish a relationship between '*accurate digits*' and rounding error.

$$\pi = 3.1415...$$

$$\bar{\pi} \approx 3.\underbrace{142}_{n=4}$$

$$\text{relative error} = \frac{\pi - \bar{\pi}}{\pi} \leq 5 \cdot 10^{-4}$$

$$\leq 10^{-3}$$

## Condition Numbers

Methods $f$ take input $x$ and produce output $y = f(x)$.
Input has (relative) error $|\Delta x| / |x|$.
Output has (relative) error $|\Delta y| / |y|$.
**Q:** Did the method make the relative error bigger? If so, by how much?

$$k = \max_{x} \left( \frac{\text{relative change in } f(x)}{\text{rel. perturbation to } x} \right)$$

$$= \max_{x} \left( \frac{|f(x + \Delta x) - f(x)|}{|f(x)|} \bigg/ \frac{|\Delta x|}{|x|} \right)$$

# absolute condition number

$$\kappa_{abs} = \max_{x} \left( \frac{|f(x + \Delta x) - f(x)|}{|\Delta x|} \right)$$

$$= \max_{input} \left( \frac{\text{absolute change of output}}{\text{size of perturbation to input}} \right)$$

$$k(x) = \frac{|x||f'(x)|}{|f(x)|}$$

$$f(x) = a + bx + cx^2 + \ldots + y \cdot x^n$$

as $x \to \infty$, $f(x) \approx y x^n$

$$k(x) \approx \frac{|x||ny x^{n-1}|}{|y x^n|} = n$$

# $n$th-Order Accuracy

Often, *truncation error* is controlled by a parameter $h$.

Examples:

- distance from expansion center in Taylor expansions
- length of the interval in interpolation

A numerical method is called '$n$th-order accurate' if its truncation error $E(h)$ obeys

$$E(h) = O(h^n).$$

# Outline

# Wanted: Real Numbers... in a computer

Computers can represent *integers*, using bits:

$$23 = 1 \cdot 2^4 + 0 \cdot 2^3 + 1 \cdot 2^2 + 1 \cdot 2^1 + 1 \cdot 2^0 = (10111)_2$$

How would we represent fractions, e.g. 23.625?

$$23.625 = 10111.101$$

$$1 \cdot 2^{-1} + 0 \cdot 2^{-2} + 1 \cdot 2^{-3}$$

$$10111.101 \qquad \text{fixed-point}$$

$$\underbrace{\qquad}_{32} . \underbrace{\qquad}_{32}$$

# Fixed-Point Numbers

> Suppose we use units of 64 bits, with 32 bits for exponents $\geqslant 0$ and 32 bits for exponents $< 0$. What numbers can we represent?

$$2^{31} + 2^{30} + \ldots \underline{2^{-32}} \approx 2^{32} - 2^{-32} \approx 10^9$$

$$2^{-32}$$

> How many 'digits' of relative accuracy (think relative rounding error) are available for the smallest vs. the largest number?

9 decimal digits

$$\underbrace{2^{31} \quad 2^{30} \quad \ldots \ldots \quad 2^{-32}}_{1 \text{ accurate digit for } 2^{-32}}$$

# Floating Point numbers

> Convert $13 = (1101)_2$ into floating point representation.

$$1.2345\ldots \times 10^6$$

$$1.101 \times 2^3$$

$$\underbrace{(1101}_{\text{Significand}} \quad \underbrace{11}_{\text{exponent}}$$

> What pieces do you need to store an FP number?

Significand - binary digits
$$s_1 \ldots s_n$$

exponent - magnitude
$$e_1 \ldots e_k$$
$$n + k \approx \text{bits}$$

$$s_0 . s_1 s_2 s_3 \ldots s_n \cdot 2^{e_1 \ldots e_k}$$

$$s_0 = 1 \qquad s_1 \ldots s_n = 0$$

$$e_1, e_2 \ldots e_k = 1$$

$$1.0 \cdot 2^{\overbrace{11111}^{}}_{\underbrace{\quad}_{fraction}} = 2^{2^k - 1}$$

$$k = 10 \qquad 2^{1024 - 1} = 2^{1023}$$

$$1.\underbrace{01011}_{significand}$$

$$1 . s_0 s_1 \ldots s_n \cdot 2^{e_1 \ldots e_r}$$

$$= 0 . 1 s_0 \ldots s_n \cdot 2^{e_1 \ldots e_r + 1}$$

$$1 0 1 1 1 \; 0 1 . 1 1 0 0 1$$

$\longrightarrow \quad 1 . 0 1 1 1 0 1 1 1 0 0 1 \times 2^6$

$. 0 0 0 1 1 0 1$

$\longleftarrow$

$1 . 1 0 1 \cdot 2^{-4}$

64-bit

$$\underset{\text{sign bit}}{\underset{\uparrow}{1}} \quad \underset{\substack{\uparrow \\ \text{exponent}}}{13} \quad \underset{\substack{\uparrow \\ \text{significand}}}{50}$$

if exponent $= 0\ 0\ 0\ 0\ 0\ 0$

$\longrightarrow 2^{-2^{13}}$

roundoff error based machine
$\varepsilon$ is smallest such         epsilon
that $fp(1 + \varepsilon) \neq fp(1)$

$$exp = 000 \implies 2^{-4}$$
$$100 \implies 2^{0}$$
$$111 \implies 2^{3}$$

$$fp(1) = \underbrace{0}_{sign} \ \underbrace{100}_{exp} \ \underbrace{0 \ldots 0}_{s.g.}$$

$$exp = 10D \implies \cancel{2^{0}} \quad 2^{-4}$$

$$5 = 1.01 \cdot 2^2$$

$$s_1 \, s_2$$

$$6 = 1.10 \cdot 2^R$$

$$s_1$$

$$7 = 1.11$$

$$8 = 1.00 \cdot 2^3 \qquad 9$$

$\in$

$\overline{h}$

$df?$

$$fp(1) = \underbrace{0}_{sign} \quad \underbrace{0\,0\,0}_{exponent} \;,\; \underbrace{0\,0\,0\,0\,0}_{Significand}$$

$$fp(1+\xi) = 0 \quad 0\,0\,0 \quad 0\,0\,0\,0\,1$$

$$fp(exp) = 0\,0\,0 \implies 2^{-4} \implies 2^{-11}$$
$$= 1\,0\,0 \implies 2^{0}$$
$$= 1\,1\,1 \implies 2^{3}$$

**In-class activity:** Floating Point