

HW2 Elijah Williams

Problem 1.

Recall that double-precision format can be used to represent any real number $x \in \mathbb{R}$ of the form:

$$x = (-1)^s \times 2^{c-1023} \times (1.f)_2$$

where $\begin{cases} S \in \{0, 1\} \\ 0 < C < 2^{17} \rightarrow C - 1023 \in [-1022, 1023] \\ 1 \leq f.f \leq 2 - 2^{-52} \end{cases}$

(a) The largest machine number we can represent is: 1023 52) - 1024

$$(1) \times 2^{1023} \times (2 - 2^{52}) < 2^{1024}$$

$$\text{So } m_1 = 1023$$

~~(b) 25.3~~ \equiv ~~1023~~ $\begin{array}{|c|c|c|c|c|c|c|c|c|c|c|c|c|} \hline 0 & 1 & 1 & 1 & 1 & 1 & 1 & 1 & 1 & 1 & 1 & 0 & 0 & 0 & \dots & 0 & 0 & 0 \\ \hline \end{array}$

Sign, 1 Exponent, 11 Mantissa, 52

(b) Similarly smallest sized machine number is

$$(1) \times 2^{-1022} \times (1.0)_2 = 2^{-1022}$$

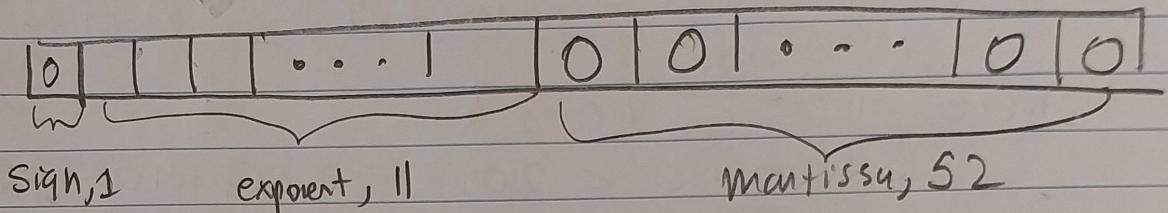
$$2^{-1022} \equiv \underbrace{[0}_\text{Sign} \underbrace{b}_\text{exponent} \underbrace{000000000001}_\text{mantissa} \underbrace{00}_{\dots} \underbrace{000000}_\text{mantissa}$$

$$so \quad \boxed{m_s = -1022}$$

(c) Suppose $m \in \mathbb{Z}$ s.t. $2 \leq 2^m \leq 10$
 $-1022 \leq m \leq 1023$

Then 2^m is a machine number: that

looks like:



- If $m \neq m_2$, then the machine number immediately to the right of 2^m is

$$r = \boxed{(1 + 2^{-52})2^m} = \boxed{2^m(2 + 2^{-52})}$$

- If $m \neq m_s$, then the machine number to the left of 2^m is:

$$l = \boxed{2^{m-1}(2 - 2^{-52})} = 2^m - 2^{-53+m}$$

$$\left\{ \begin{array}{l} 2^m - l = 2^{-53+m} \\ r - 2^m = 2^{-52+m} \end{array} \right.$$

- The number immediately to the left is closer than the number immediately to the right by a factor of 2

②

Idea: When we order the numbers in a list in terms of magnitude, we can decrease the roundoff error of the sum.

(a) See Pictures Below
(b) List =

(b) This is not always the case.

Suppose we are working on a decimal machine with two digits allocated to the mantissa.

$$\text{List} = [0.25 \times 10^{-2}, 0.34 \times 10^{-4}, 0.51 \times 10^{-5}, 0.61 \times 10^{-2}]$$
$$(b) \text{sum}(\text{list}) = [0.31491, 0.0034, 0.0051, 0.061]$$

Suppose addition take place in a sorted digit accumulator. Then, [0.11, 0.25, 0.31, 0.61]

Unordered Sum:

$$\begin{aligned} & 25.2500 \times 10^{-2} \\ & + 0.0034 \times 10^{-2} \\ & + 0.2534 \times 10^{-2} \\ & + 0.0005 \times 10^{-2} \\ & \hline \text{foc} \Rightarrow 25.05 \times 10^{-2} \end{aligned}$$

$$\begin{aligned} \text{sum_list} = \text{foc} \Rightarrow & 25.00 \times 10^{-2} \\ & + 0.0610 \times 10^{-2} \\ & \hline \text{sum_sort} = 25.0610 \times 10^{-2} \\ \text{foc} & : 31 \times 10^{-2} \end{aligned}$$

Rounding Error $.31491 - .31 = \boxed{.491}$

Ordered Sum: $[.51 \times 10^{-5}, .34 \times 10^{-4}, .61 \times 10^{-3}, .25 \times 10^{-2}]$

$$\begin{aligned} & 0.510 \times 10^{-4} \\ + & \underline{3400 \times 10^{-4}} \\ & 3910 \times 10^{-4} \\ f(1) \Rightarrow & 0.390 \times 10^{-3} \\ + & \underline{6100 \times 10^{-3}} \\ & 6490 \times 10^{-3} \\ f(2) \Rightarrow & 0.650 \times 10^{-2} \\ + & \underline{2500 \times 10^{-2}} \\ & 3150 \times 10^{-2} \\ \Rightarrow & 32 \times 10^{-2} \end{aligned}$$

Rounding Error:

$$.31491 - .32 = \boxed{-.509}$$

The roundoff error for the ordered sum
is larger in magnitude than the
unordered sum.

3.) How many machine numbers are there in single-precision arithmetic?

Suppose $x \in \mathbb{R}$ is a single-precision machine number (non zero)

$$\text{Then } x = (-1)^s \times 2^{c-127} \times (1.f)_2$$

where $s \in \{0, 1\}$

$$c \in \{1, 2, \dots, 254\}$$

$$f = (b_1 b_2 b_3 \dots b_{52})$$

Counting the possibilities:

$$\underline{2} \cdot \underline{254} \quad \underline{2}^{52}$$

$$\boxed{\underline{\underline{= 2^{53} \cdot 254}} \text{ single-precision numbers}}$$

excluding $\pm 0, \pm \infty$

④ In mathematics,

$$(a+b)+c = a+(b+c)$$

However this is not always true

when adding machine numbers

ex)

- In my solution for problem

2, ~~you~~ I already found

3 floating point numbers

for which this is true:

- Unsorted-list = $[a, b, c]$

- Sorted-list = $[c, b, a]$

$$\boxed{(a+b)+c \neq (c+b)+a}$$

$$\textcircled{5} \quad \sin(x) + \cos(x) - 1, \quad x \text{ near zero}$$

• When x is near zero,

$$\sin(x) \approx 0$$

$$\cos(x) \approx 1$$

and we get subtractive cancellation.

We can try rationalization:

$$(\sin(x) + \cos(x) - 1) \cdot \frac{(\sin(x) + \cos(x) + 1)}{\sin(x) + \cos(x) + 1}$$

$$\begin{aligned} &= \cancel{\sin^2(x)} + \sin(x)\cos(x) + \sin(x) \\ &\quad \cancel{\cos^2(x)} + \sin(x)\cos(x) + \cos(x) - 1 \\ &\quad - \cancel{\sin(x)} - \cancel{\cos(x)} - \cancel{1} \\ &\quad \cancel{\sin(x) + \cos(x) + 1} \end{aligned}$$

$$\boxed{-\frac{\sin(2x)}{\sin(x) + \cos(x) + 1}}$$

No

Subtractive
cancellation

Problem 6

- The purpose of this code is to find out what the smallest positive number we can add to 1 that causes the computer to round up to another machine number.

$$y = \underbrace{\text{Machine epsilon}}_{\text{Once}} = \frac{\text{unit roundoff error}}{\text{error}}$$

- I calculated

$$y = 2.2 \times 10^{-16}$$

This is the same machine epsilon we calculated for double precision numbers

$$\epsilon = 2^{-52} \approx 2.2 \times 10^{-16}$$

Problem 7

(a) Evaluate $f(x) = (x - \sin x) / x^3$, x close to zero

- See Program

(b) We can eliminate subtractive cancellation
by dividing by x :

$$f(x) = \frac{x}{x^3} - \frac{\sin(x)}{x^3} = \frac{1}{x^2} - \frac{\sin(x)}{x^3}$$

>> HW2_P2

List: [678.7351548577735230, 0.7577401305783334, 0.0007431324681249]

Exact Sum: 679.49363812082003732939483597874641

Sorted Sum: 679.49363812082003732939483597874641

Sorted Error: 0.0000000000005600196955474201488

Unsorted Sum: 679.49363812081992364255711436271667

Unsorted Error: -0.0000000000005768486816687401486

>>

