

Nearby Machine Numbers

mantissa = significand

$$x = q \cdot 2^m \quad (1 \leq q < 2 \quad -126 \leq m \leq 127)$$

Say $x = (1.a_1 a_2 \dots a_{23} a_{24} a_{25} \dots)_2 \cdot 2^m \quad a_i = 0 \text{ or } 1$

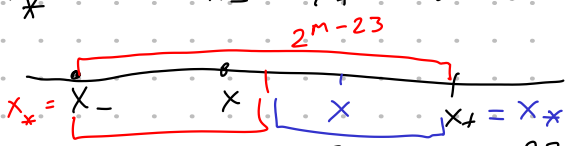
Let $x_- = (1.a_1 a_2 \dots a_{23})_2 \cdot 2^m$ (chopping) $(1.a_1 a_2 \dots)_2 \cdot 2^{-2}$
 $0.00\dots001 = 2^{-23}$ $0.01000 = 2^{-2}$

$x_+ = ((1.a_1 a_2 \dots a_{23})_2 + 2^{-23}) \cdot 2^m$ (rounding up)
 $x_+ = 5.2319700000$
 $x_- = 5.23196824$
 $x = 5.2319600000$
 $x_- \leq x \leq x_+$

Either $\left| \frac{x-x_-}{x} \right| \leq 2^{-24}$ or $\left| \frac{x-x_+}{x} \right| \leq 2^{-24}$

Q) why?

A) Let x_* denote x_- or x_+ whichever is closer to x .



Clearly $x_+ - x_- = 2^{-23+m} = 2^{m-23}$

So $|\text{abs. error}| = |x - x_*| \leq \frac{1}{2} 2^{m-23} = 2^{m-24}$ $\begin{matrix} q \geq 1 \\ \frac{1}{q} \leq 1 \\ \downarrow \end{matrix}$

$|\text{Relative error}| = \left| \frac{\text{abs. error}}{x} \right| \leq \frac{2^{m-24}}{q \cdot 2^m} \leq \frac{1}{q} 2^{-24} \leq 2^{-24}$

$|\text{Relative error}| = \left| \frac{x-x_*}{x} \right| \leq 2^{-24}$: Unit Roundoff Error

Machine Epsilon = 2 * Unit Roundoff Error $\begin{matrix} \text{64 bits} \\ \text{mach. eps} = 2^{-52} \end{matrix}$

Exercise: On a modern machine using double precision, 52 bits are reserved for the mantissa. Try the following on your Python interpreter

```
u = 2*x - 53          for n in range(1, 60):
e = 2**x - 52        u = 2**x - (n+1)
1+e == 1             e = 2**x - n
1+u == 1             if not (1+e==1) and
                    (1+u==1):
                    print(n)
                    break
```

Notation $fl(x) = x_*$ the floating point machine number closest to x .

Example For $x = \frac{2}{3}$, compute the nearby machine numbers x_- and x_+ (in the architecture we described last time). Find $fl(x)$. What are the absolute round off error and the relative round off error in representing x by $fl(x)$? $52.879 \times 10^{-52} = 5.2879 \times 10^{-53}$

$\frac{2}{3} < 1$ so $\frac{2}{3} = (0.a_1 a_2 a_3 \dots)_2$

$2 \cdot \frac{2}{3} = \frac{4}{3} = a_1.a_2 a_3 \dots$ Since $\frac{4}{3} > 1$, $a_1 = 1$
 $1 \quad 1.000000$

Subtract 1 from both sides

$\frac{1}{3} = 0.a_2 a_3 a_4 \dots$

$2 \cdot \frac{1}{3} = \frac{2}{3} = a_2.a_3 a_4 \dots < 1$ so $a_2 = 0$

$2 \cdot \frac{2}{3} = \frac{4}{3} = a_3.a_4 a_5 \dots$ so $a_3 = 1$

Thus, $\frac{2}{3} = (0.101010 \dots)_2 = (0.\overline{10})_2 \quad \begin{matrix} x = q \cdot 2^m \\ 1 \leq q < 2 \end{matrix}$

$\frac{2}{3} = (0.\overline{10})_2 = (1.0101 \dots)_2 \cdot 2^{-1} \quad m = -1$

$q = (1.0101 \dots)_2 = (1.f)_2$ so $f = (0.010101 \dots)_2$

1st, 3rd, 5th ... bits (after binary point)

are 0. So 23rd bit is also 0 (followed by a 1).

Thus, $x_- = (1.0101 \dots 010)_2 \cdot 2^{-1}$
 $x_+ = (1.0101 \dots 011)_2 \cdot 2^{-1}$
└──────────┘
23 bits

To find $fl(x)$ we need to compute $x_+ - x$ and $x - x_-$ (and choose the smaller one).

Here is another way: the number that is exactly in the middle of x_+ and x_- is $\frac{x_+ + x_-}{2} = \text{ave}$

$x_- = (1.0101 \dots 010000 \dots)_2 \cdot 2^{-1}$ $\begin{matrix} \text{ave} \\ \text{---} \bullet \text{---} \\ x_- \quad x_+ \end{matrix}$

$\text{ave} = (1.0101 \dots 010100 \dots)_2 \cdot 2^{-1}$

$x_+ = (1.0101 \dots 011000 \dots)_2 \cdot 2^{-1}$

since $x > \text{ave}$, $fl(x) = x_+$

$\frac{2}{3} = 1.0101 \dots 010101 \dots > 1.0101 \dots 01010000$
└──────────┘
ave
 x_+ is closer to x .

