

Sep-6

Note Title

06-09-2012

Floating point representation.

What is a FP representation?

3.18924 {FP number}

-10.192

2×10^{-15}
 -4.23×10^{31}

Fixed Point numbers : Currency

↓[₹]
7.22
2500.99

How to represent floating point numbers in binary?

$$\begin{aligned}(11) &= 8 + 0 + 2 + 1 \\ &= 1 \times 2^3 + 0 \times 2^2 + 1 \times 2^1 + 1 \times 2^0 \\ &= (1011)_2\end{aligned}$$

$$\begin{aligned}1.75 &= 1 + 0.5 + 0.25 \\ &= 1 \times 2^0 + 1 \times 2^{-1} + 1 \times 2^{-2} \\ &= (1.11)_2\end{aligned}$$

$$\begin{aligned}1.5625 &= 1 + 0.5 + 0 \times 0.25 + 0 \times 0.125 + 1 \times 0.0625 \\ &= (1.1001)_2\end{aligned}$$

decimal \rightarrow binary

binary \rightarrow decimal

$$(1.0011)_2 \rightarrow 1 \times 2^0 + 0 \times 2^{-1} + 0 \times 2^{-2} + 1 \times 2^{-3} + 1 \times 2^{-4}$$

$$= 1 + \frac{1}{8} + \frac{1}{16}$$

$$= \left[1 \frac{3}{16}\right]_{\text{decimal}}$$

floating point \rightarrow $\left\{ \begin{array}{l} +/- \\ \cdot \text{ point} \\ \text{exponent } (1.37 \times 10^{-19}) \end{array} \right.$

IEEE 754 Floating Point Format

Standard form representation.

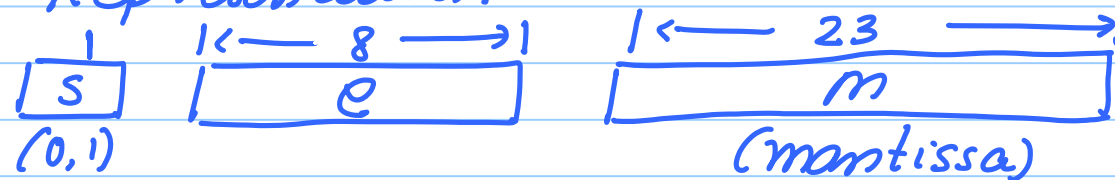
$$N = (-1)^s (1.m) \times 2^{\text{exp}}$$

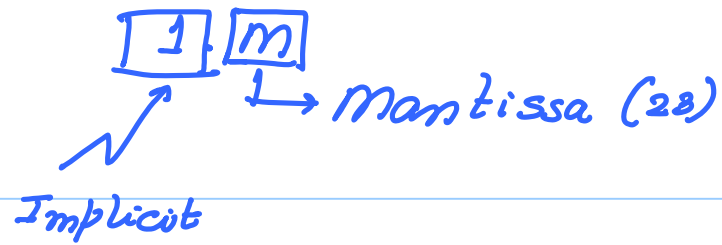
$s \in \{0, 1\}$

$$2.5 = (-1)^0 (1.25) \times 2^1$$

$$4.5 = (-1)^0 (1.125) \times 2^2$$

IEEE 754 Representation



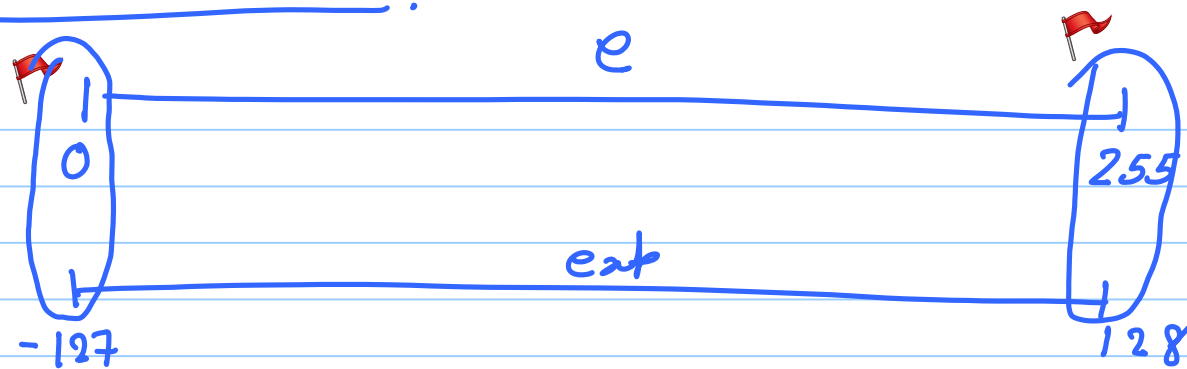


have both (+)ve and (-)ve exponents

Instead of representing (e) in 2s complement form
→ biased notation

$$\begin{array}{l} \text{exp} = e - \text{BIAS} \\ \begin{array}{l} \text{---}^{128} \\ | \\ \text{---}^{-127} \end{array} \quad \begin{array}{l} \text{---}^{255} \text{ (127)} \\ | \\ \text{---}^0 \end{array} \end{array}$$

Special Cases.



The range of (e, exp) for normal FP numbers

$$\begin{aligned} e &\rightarrow 1 \dots 254 \\ \text{exp} &\rightarrow -126 \dots 127 \end{aligned}$$

Limits of (+)ve FP (normal) numbers:

Smallest (+)ve normal FP number:

$$1.\underbrace{0 \dots 0}_{23} \times 2^{-126} = 2^{-126}$$

Largest (+)ve normal FP number:

$$\begin{aligned} & 1.\underbrace{1 \dots 1}_{23} \times 2^{127} \\ &= 2^{127} \times (2^0 + 2^{-1} \dots \dots 2^{-23}) \\ &= 2^{127} \times (2^1 - 2^{-23}) \\ &= [2^{128} - 2^{104}] \end{aligned}$$

Special Cases

$e = 0, e = 255$

$$e = 255, m = 0, s = 0$$
$$s = 1$$

~~∞~~
 ∞
 $-\infty$

$$e = 255, m \neq 0$$

NAN
(Not a Number)
{ $\sin^{-1}(5)$ }
{ $\log(-4)$ }

$$2 + \text{NAN} = \text{NAN}$$
$$\text{NAN} \times \pi = \text{NAN}$$

$$e = 0, m = 0$$

~~$\frac{\infty}{0}$~~

$$e = 0, m \neq 0$$

[denormal numbers]

$$\left. \begin{array}{l} x = 2^{-126}; \\ \text{if } (x/2 == 0) \\ \quad \text{printf} ("hi \n"); \end{array} \right\} \quad \begin{array}{l} \text{Violates Axiom} \\ x > 0 \Rightarrow x/2 > 0 \end{array}$$

To protect against such potential contradictions
we have some buffer in the form of
denormal numbers

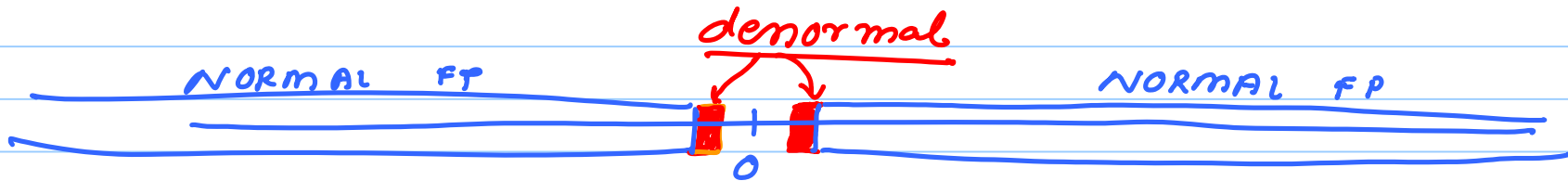
$$\text{Denormal number: } [N = (-1)^s 0.m \times 2^{-126}]$$

Limits of (+)ve denormal numbers:

$$\begin{aligned} \text{smallest} &: 0.0 \overbrace{\dots 0}^{22} 1 \times 2^{-126} \\ &= 2^{-23} \times 2^{-126} = 2^{-149} \end{aligned}$$

$$\begin{aligned}
 \text{largest: } & 0.\underbrace{11\dots1}_{23} \times 2^{-126} \\
 &= 2^{-126} \times (2^{-1} + \dots + 2^{-23}) \\
 &= 2^{-126} \times (2^0 - 2^{-23}) \\
 &= 2^{-126} \times (1 - 2^{-23}) \\
 &= [2^{-126} - 2^{-149}]
 \end{aligned}$$

NUMBER LINE



Now, this piece of code will work as expected

```
 $x = 2^{-126};$   
 $\text{if } (x/2 == 0)$   
     $\text{printf}("hi");$ 
```

REASON: $x/2 = 2^{-127}$ can be represented as a valid denormal number

```
 $x = 2^{-149}$   
 $\text{if } (x/2 == 0)$   
     $\text{printf}("hi");$ 
```

Output: hi

I want to solve this problem.

Instead of a float, I will use a double.

float \rightarrow floating point, single precision
(32)

double \rightarrow double precision
(64)



double

$$\text{exp} = e - \text{BIAS}$$

$$= e - 1023$$

largest double precision number:

$$\{ 2^{e_{\max}/2} - 2^{e_{\max}/2 - 1 - (m)} \}$$

$$\begin{array}{r} 1023 \\ - 52 \\ \hline 971 \end{array}$$

$$2^{1024} - 2^{1024 - 1 - 52}$$

$$= \{ 2^{1024} - 2^{971} \} \approx \textcircled{10^{300}}$$

Range of double⁽⁶⁴⁾ numbers $\approx (-10^{300}$ to $10^{300})$
of float⁽³²⁾ numbers $\approx (-10^{40}$ to $10^{40})$

Intel machines have a format called
extended precision (80 bits)

You can still have mathematical contradictions.

$$\Delta x > 0$$

$$x + \Delta x > x$$

$$\left\{ \begin{array}{l} x = 2^{100} \\ \Delta x = 2^{-100} \\ y = 2^{100} \\ z = x + \Delta x - y \end{array} \right.$$

$$\left\{ \begin{array}{l} \text{if } (z > 0) \\ \quad \text{printf ("good");} \end{array} \right.$$

Expected Output: good

$$\left[\begin{array}{l} z = (x + \Delta x) - y \\ = x - y = 0 \end{array} \right]$$

comp 1

$$\left[\begin{array}{l} z = (x - y) + \Delta x \\ = 0 + \Delta x \\ = \Delta x > 0 \end{array} \right]$$

comp 2

The result clearly depends on the order of computations.

Conclude:

1) FP arithmetic is approximate

2) Can lead to the violation of basic mathematical axioms.