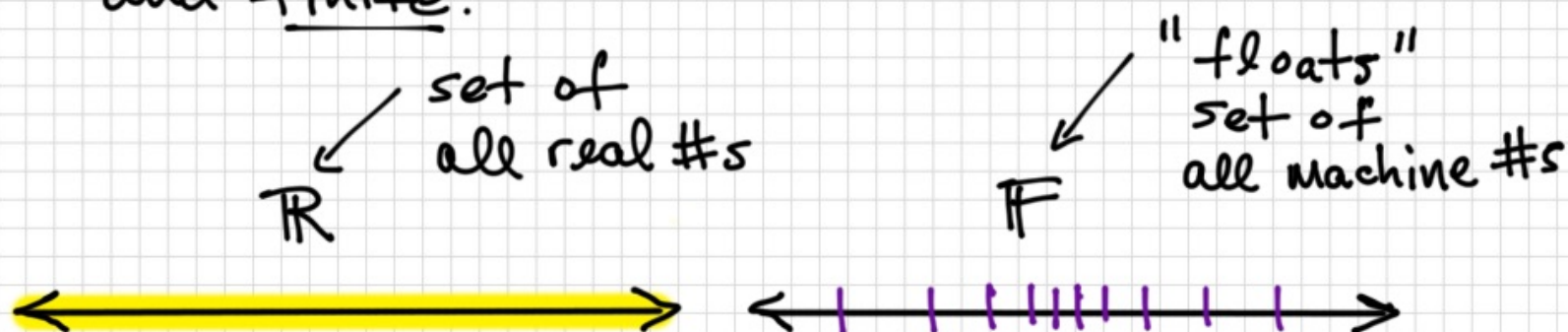


Machine Numbers

float

- * The real line \mathbb{R} is a continuum
- * The number system on a computer is discrete and finite.



Created with Doceri



Mostly, given $x \in \mathbb{R}$ of interest, there is an $fl(x) \in \mathbb{F}$ which represents x on the computer. How do we find $fl(x)$?

$fl(x)$ is closely related to binary representation of x .

Problem Find $fl(9.4)$

→ Step ① : find binary rep of 9.4
focus on this step now
Step ② : express in scientific notation
Step ③ : truncate/round if necessary

$$\begin{aligned}\text{Recall: } x &= (b_m \dots b_2 b_1 b_0 . b_{-1} b_{-2} \dots)_2 \\ &= b_m 2^m + \dots + b_2 2^2 + b_1 2^1 + b_0 2^0 \\ &\quad + b_{-1} 2^{-1} + b_{-2} 2^{-2} + \dots\end{aligned}$$

Where each integer bit is either 0 or 1.

$$9.4 = 9 + 0.4$$

integer part fractional part

Convert these to binary separately.

Created with Doceri



Integer part $q = b_m 2^m + \dots + b_1 2^1 + b_0 2^0$

$$9/2 = 4 \text{ remainder } 1 \quad (\text{get } b_0)$$

$$4/2 = 2 \text{ remainder } 0 \quad (\text{get } b_1)$$

$$2/2 = 1 \text{ remainder } 0 \quad (\text{get } b_2)$$

$$1/2 = 0 \text{ remainder } 1 \quad (\text{get } b_3 \text{ and stop})$$

$$q = (1001)_2, \quad 1 \cdot 2^3 + 0 \cdot 2^2 + 0 \cdot 2^1 + 1 \cdot 2^0$$

Created with Doceri

$$= 8 + 0 + 0 + 1 = 9$$

Fractional part

$$0.4 = b_{-1} \frac{1}{2} + b_{-2} \frac{1}{4} + \dots$$

$$0.4 * 2 = 0.8 \text{ plus } 0$$

(get b_{-1})

$$0.8 * 2 = 0.6 \text{ plus } 1$$

(get b_{-2})

$$0.6 * 2 = 0.2 \text{ plus } 1$$

(get b_{-3})

$$0.2 * 2 = 0.4 \text{ plus } 0$$

(get b_{-4})

$$0.4 * 2 \text{ ————— pattern repeats}$$

use
later

$$0.4 = (0.\overline{0110})_2 \Rightarrow 2^4 = (1001.\overline{0110})_2$$



Floating point #s (normalized)

$$\text{fl}(x) \rightarrow V = (-1)^s (1.b_1b_2\cdots b_{52})_2 \times 2^{F-1023}$$

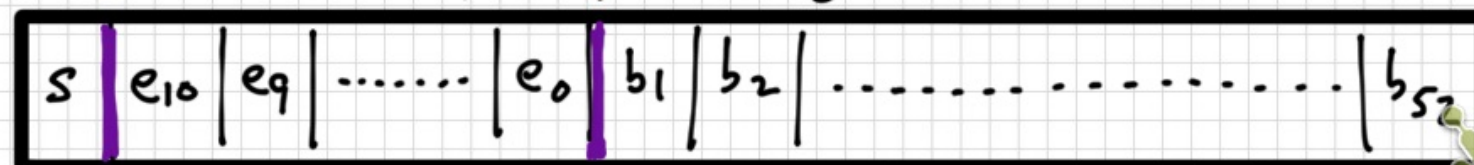
careful
before $(b_1b_0.b_{-1}b_{-2}\cdots)_2$
note:
 $(b_1b_0.b_{-1})_2 = (b_1.b_0b_{-1})_2 \times 2^1$

$$1 \leq F \leq 2046$$

↑
exponent

↖ $F-1023$
true exponent

64-bit storage



sign ① exponent ②

mantissa ⑤②



$$V = (-1)^s (1.b_1 b_2 \dots b_{52})_2 \times 2^{F-1023}$$

$1 \leq F \leq 2046$ normalized #s

$F=0$ unnormalized #s

$F=2047$ NaNs and $\pm \text{Inf}$

$$F = (e_{10} e_9 e_8 e_7 e_6 e_5 e_4 e_3 e_2 e_1 e_0)_2$$

$$\text{note} = (\underbrace{1111111111}_{11 \text{ ones}})_2 = 2^{11} - 1 = 2048 - 1 = 2047 \quad \checkmark$$

Created with Doceri



$$V = (-1)^S (1.b_1b_2 \dots b_{52})_2 \times 2^{F-1023}$$

Normalized $1 \leq F < 2046$

$$V = (-1)^S (0.b_1b_2 \dots b_{52})_2 \times 2^{-1022}$$

Unnormalized $F = 0$

$$V = \text{NaN or } \pm \text{Inf} \quad F = 2047$$

only if all mantissa bits are 0.

$$V_{\max, \text{norm}} = (1.111\dots 1)_2 \times 2^{1023} \simeq 1.7977 \times 10^{308}$$

$$V_{\min, \text{unnorm}} = (0.0000\dots 1)_2 \times 2^{-1022} \simeq 4.9407 \times 10^{-324}$$

Notice

| 0 | 1000000010 | 001 0110 0110 0110 0110 0110 0110 0110 0110 0110 0110 0110 0110 0110 |

is

| 0100 0000 0010 0010 1100 1100 1100 1100 1100 1100 1100 1100 1100 1100 1100 1101 |

4022cccccccccd

$(0)_{16} = (0000)_2$ $(4)_{16} = (0100)_2$ $(8)_{16} = (1000)_2$ $(c)_{16} = (1100)_2$

$(1)_{16} = (0001)_2$ $(5)_{16} = (0101)_2$ $(9)_{16} = (1001)_2$ $(d)_{16} = (1101)_2$

$(2)_{16} = (0010)_2$ $(6)_{16} = (0110)_2$ $(a)_{16} = (1010)_2$ $(e)_{16} = (1110)_2$

$(3)_{16} = (0011)_2$ $(7)_{16} = (0111)_2$ $(b)_{16} = (1011)_2$ $(f)_{16} = (1111)_2$