

Floating point Arithmetic

①

Consider -321.416

↓

$$- (3 \cdot 10^2 + 2 \cdot 10 + 1 \cdot 10^0$$

$$+ 4 \cdot 10^{-1}$$

$$+ 1 \cdot 10^{-2} + 6 \cdot 10^{-3})$$

Decimal representation

$$= -3.21416 \times 10^2$$

$$= -0.321416 \times 10^3$$

A similar representation is
used in a computer,

"floating-point arithmetic"

- . 3 2 1 4 1 6 \times (10) ⁽³⁾ \rightarrow exponent (2)
 \downarrow sign fraction \downarrow base

In general

$$\pm f \times \beta^e$$

\downarrow
sign

$\beta = 2$ binary number

10 decimal :

16 hexadecimal :

f : fraction,

digits from 0, 1, ... $\beta-1$

e : exponent

digit, from 0 ... $\beta-1$

(3)

Binary numbers:

$$\underbrace{b_m \cdots b_2 b_1 b_0}_{\text{integer}} \cdot \underbrace{a_1 a_2 \cdots a_n}$$

each digit b_i, a_j takes 0 or 1

This number in base 10 is

$$\begin{aligned} & b_m \cdot 2^m + b_{m-1} 2^{m-1} + \cdots + b_1 \cdot 2^1 \\ & \quad \uparrow \\ & + b_0 2^0 + a_1 2^{-1} + a_2 2^{-2} \\ & \quad \quad \quad + \cdots a_n 2^{-n} \\ & \quad \quad \quad \text{"B"} \end{aligned}$$

Note: $(0.1101)_2$

$$= (1.101)_2 \times 2^{-1}$$

$$= (0.001101)_2 \times 2^2$$

↓ B

② To convert between binary ($b=2$) and decimal ($b=10$)

Example:

1) $x = (1.1011)_2$, convert x to a decimal

$$\begin{aligned} x &= 1 \cdot 2^0 + 1 \cdot 2^{-1} + 0 \cdot 2^{-2} + 1 \cdot 2^{-3} \\ &= 1 + \frac{1}{2} + 0 + \frac{1}{8} + \frac{1}{16} \\ &= \frac{27}{16} \end{aligned}$$

$$\begin{aligned} 2) \quad x &= (1.101010 \dots 10 \dots)_2 \\ &= (1.\overline{10})_2 \end{aligned}$$

$$\begin{aligned} x &= 1 \cdot 2^0 + 1 \cdot 2^{-1} + 1 \cdot 2^{-3} + \cancel{1 \cdot 2^{-5}} \\ &\quad + 1 \cdot 2^{-7} + \dots \end{aligned}$$

Recall geometric series if
 $1 + r + r^2 + r^3 + \dots = \frac{1}{1-r}, |r| < 1$

$$x = 1 + \frac{1}{2} + \left(\frac{1}{2}\right)^3 + \left(\frac{1}{2}\right)^5 + \left(\frac{1}{2}\right)^7 + \dots \quad (5)$$

$$= 1 + \frac{1}{2} \left(1 + \frac{1}{4} + \left(\frac{1}{4}\right)^2 + \left(\frac{1}{4}\right)^3 + \dots \right)$$

$$= 1 + \frac{1}{2} \frac{1}{1 - \frac{1}{4}} \quad \downarrow r = \frac{1}{4}$$

$$= 1 + \frac{1}{2} \frac{1}{\frac{3}{4}} = \frac{5}{3}$$

Alternatively:

$$x = (1.\overline{10})_2$$

$$= 1 \cdot 2^0 + (\overline{0.\overline{10}})_2 \dots (*)$$

$$= 1 + (10.\overline{10})_2 \cdot 2^{-2}$$

$$y = (\overline{0.\overline{10}})_2 = (\overline{10.\overline{10}})_2 \cdot 2^{-2}$$

$$= \left\{ \underbrace{(10)_2}_{=2} + \underbrace{(\overline{0.\overline{10}})_2}_y \right\} 2^{-2}$$

$$\Rightarrow y = (2 + y) 2^{-2}$$

$$\Rightarrow 4y = 2 + y \Rightarrow y = \frac{2}{3} \cdot \frac{1}{4} \quad (**)$$

$$\Rightarrow x = \frac{5}{3}$$

Example: Convert 14.8125 to (6)
a binary number.

we are looking for

$$14.\boxed{8125} = (\underbrace{b_m b_{m-1} \dots b_1 b_0}_{\text{integer part}} \underbrace{. a_1 a_2 \dots a_n}_{\text{fractional part}})_2$$

fractional part

$$0.8125 = (. a_1 a_2 \dots a_n)_2$$

$$= a_1 \cdot 2^{-1} + a_2 2^{-2} + \dots + a_n 2^{-n}$$

- $\times 2$

$$\boxed{1} 6250 = \cancel{2} \boxed{a_1} + a_2 2^{-1} + \dots + a_n 2^{-(n-1)}$$
$$\boxed{a_1 = 1}$$

$$\Rightarrow 0.6250 = a_2 2^{-1} + \dots + a_n 2^{-(n-1)}$$

- $\times 2$

$$\boxed{1} 2500 = \boxed{a_2} + a_3 2^{-1} + \dots + a_n 2^{-(n-2)}$$
$$\boxed{a_2 = 1}$$

$$0.25 \times 2 = \boxed{0.50}$$

$$0 = a_3$$

⑦

$$0.50 \times 2 = \boxed{1.0}$$

$$\boxed{1 = a_4}$$

$$\Rightarrow .8125 = (.a_1 a_2 a_3 a_4)_2$$

$$\quad \quad \quad \begin{matrix} ' & '' & '' & '' \\ & 1 & 0 & 1 \end{matrix}$$

$$= (.1101)_2$$

< collect integer part, ordered from
radix point >

Integer part

$$14 = (b_m \dots b_2 b_1 b_0)_2$$

$$= b_m \cdot 2^m + b_{m-1} 2^{m-1} + \dots + b_1 2^1 + b_0$$

$$+ \dots + b_1 2^1 + b_0$$

divided by 2

$$14/2 = \overline{) R \boxed{0}} = (b_m 2^{m-1} + \dots + b_1) R \boxed{b_0}$$

$$\boxed{b_0 = 0}$$

$$7/2 = 3 \text{ R } \boxed{1}$$

$$b_1 = 1$$

(8)

$$3/2 = 1 \text{ R } \boxed{1}$$

$$b_2 = 1$$

$$1/2 = 0 \text{ R } \boxed{1}$$

$$b_3 = 1$$

$$\Rightarrow 14 = (1110)_2$$

$$\Rightarrow 14.8125 = (1110.1101)_2$$

floating point number

$$\pm f \cdot \beta^e \quad (\beta = 2)$$

(9)

f : (fraction) : the number of digits in f determines the precision

e : (exponent) : the number of digits in e determine the range of representable numbers

we follow

[IEEE 754 floating point standard]

normalized form

$$f = 1. b_m b_{m-1} + \dots b_1 b_0$$

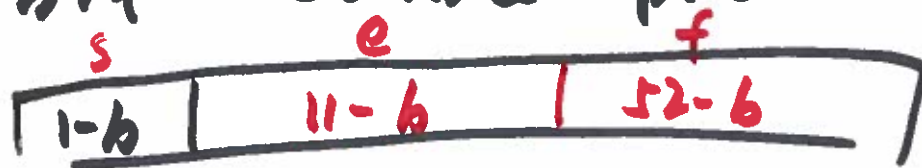
\uparrow leading digit is 1.
(0.010101110...)

Advantage: leading 1 need not be stored.

32-bit single precision (16)



64-bit double precision



The represented number is
 $(-1)^s \cdot (1 + f) \cdot \underbrace{2^{e - e_0}}$

e : un-signed, e_0 : exponent bias

$e - e_0$: can be either positive and negative.

- Let's focus on
'e', or equivalently
 $2^{e - e_0}$

↓
to represent
small num.

Single-precision

(11)

$$e \in [e_{\min}, e_{\max}]$$

e with all 0s or all 1s
needs to be interpreted
in a special way

$$e_{\min} = (\underbrace{0 \cdots 0}_{8\text{-bit}} 1)_2 = 1$$

$$e_{\max} = (\underbrace{11 \cdots 10}_{8\text{-bit}})_2 = 1 \cdot 2 + 2^2 + 2^3 + \cdots + 2^7$$
$$= 2 \left(\frac{1-2^8}{1-2} \right) = 254$$

$$\underline{a + a^2 + a^3 + \cdots + a^n = \frac{1-a^{n+1}}{1-a} \quad a \neq 1}$$

$$e_0 = 127$$

$$\Rightarrow 2^{e-e_0} \in [2^{-126}, 2^{127}]$$
$$\approx [10^{-38}, 10^{38}]$$

Double-precision

$$e \in [e_{\min}, e_{\max}]$$

"1"

$$2^1 + 2^2 + \cdots + 2^{10} = 2046$$

$$e_0 = 1023$$

$$2^{e-e_0} \in [2^{-1022}, 2^{1023}] \approx [10^{-308}, 10^{308}]$$

(12)

Next : fraction f and precision
using double-precision as
an example:

- How to store a number
- How to do calculation

Consider

$$x_1 = \frac{27}{16} = (1.1011)_2$$

$$x_2 = \frac{5}{3} = (1.\overline{10})_2$$

$$x_3 = \frac{2}{3} = (1.\overline{10})_2 = (1.0\overline{1})_2 \times 2^{-1}$$

$$x_4 = 1 = (1.0)_2$$

$$x_5 = 1 \times 2^{-52}$$

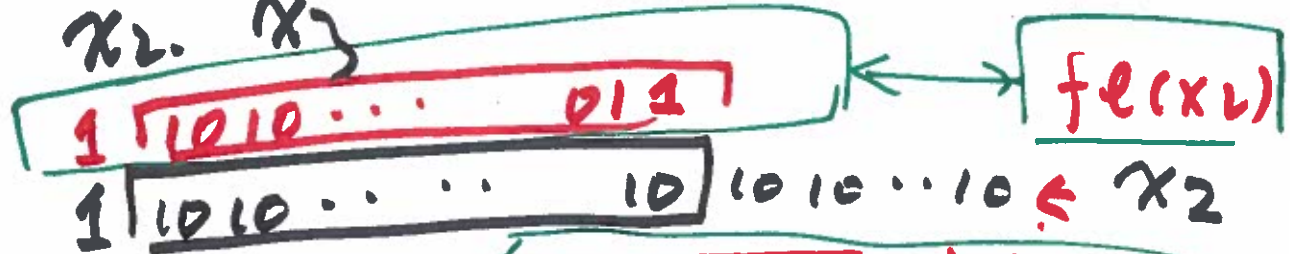
$$x_6 = 1 \times 2^{-53}$$

52-bit

Exactly

- 1. $1.01100\dots 0$ $\dots x_1$
- 1. $1.0\dots 0$ $\dots x_4$
- 1. $1.0\dots 0$ $2^{-52} x_5$
- 1. $1.0\dots 0$ $2^{-53} e-e_0 x_6$

Now x_2, x_3



$x_3 = (1.\overline{01})_2 \times 10^{-2}$

$1.\overline{01010101} \cdot 2^{-1}, fl(x_3)$

$1.0101\dots 01 \cdot 2^{-1}$

$\cdot 2^{-1}$

we follow

IEEE: rounding to nearest rule.



following the rule.

the computer represented number for x : $fl(x)$

IEEE Rounding to Nearest Rule

<Double Precision>

check the 53rd bit to the right of binary point.

If it is 0, then truncate after 52nd bit (round down)

If it is 1, then

- add 1 to 52nd bit (round up)

UNLESS all digits to the right of the 1 are 0's,

in which case 1 is added to bit 52 if and only if bit 52 is 1