Floating pt. numbers

scientific notation $x = \pm S \times 2^E$

$1 \leq S < 10$

ex: $365.25 = +3.6525 \times 10^2$

$S$ = significand , $E$ = exponent

decimal pt $\underset{\curvearrowright}{365.25}$ floats to $3.6525$

base 2 $\quad x = \pm S \times 2^E \quad , \quad 1 \leq S < 2,$

$\frac{11}{2} = (1.011)_2 \times 2^2 \qquad (x \neq 0)$

$S = (b_0 . b_1 b_2 \cdots)_2 \quad , \quad b_0 = 1 \quad \text{if } x \neq 0$

$\quad = (1. b_1 b_2 \cdots)_2$

$\qquad \qquad \underbrace{\qquad\qquad}$

$\qquad \qquad$ fractional part

normalized rep. if $b_0 = 1$

$\qquad \qquad \qquad$ 1 bit $\quad$ 8 bits $\quad$ 23 bits

computer word $\quad$ | sign | exp | fraction field |

$11/2 \quad = \quad$ | 0 | 8bits (2) | 011 $\cdots$ |

32 bit
word $\qquad\qquad$ sign $\quad 0 = + \quad , \quad 1 = -$

$\qquad\qquad \qquad \qquad \qquad \qquad -128$

$\qquad$ 8 bits for E $\therefore \leq exp \leq 127$

such numbers are floating pt. #'s

$71 = (1.000111)_2 \times 2^6$

| 0 | ebits(6) | 000111 0... |
|---|---|---|

ex $1 = (1.000\ldots)_2 \times 2^0$

$\longleftrightarrow$  | 0 | ebits(0) | 0.... |
|---|---|---|

$1024 = (1.000)_2 \times 2^{10}$

$\longleftrightarrow$  | 0 | ebits(10) | 0..... |
|---|---|---|

can even get $(1.011\cdots0)_2 \times 2^{127}$

23 place

$\longleftrightarrow$  | 0 | ebits(127) | 111 1... |  ← 23
|---|---|---|

$\uparrow_{127}$  all 1's
2

$=$

$$\text{largest \#} = +\,2^{127} \cdot (1.\underbrace{\phantom{----}}_{23 \text{ 1's}})$$

Double format   | ± | $a_1 a_2 \cdots a_{11}$ | $b_1 b_2 \cdots b_{52}$ |
|---|---|---|

1 + 11 + 52 = 64 bit word

notice

$(0)_{10} = $ | $a_1 \cdots a_{11} = 0 \cdots 0$ | $\longleftrightarrow \pm(0.b_1 \cdots b_{52})_2 \times 2^{-1022}$

$= (0\cdots 01)_2 \longleftrightarrow \pm(1.b_1 \cdots b_{52})_2 \times 2^{-1022}$

$= (0 \cdots 1\,0)_2 \longleftrightarrow \pm(1.b_1 \cdots b_{52})_2 \times 2^{-1021}$

$(1023)_{10} = (0111111)_2 \longleftrightarrow \pm(1.b_1 \cdots b_{52})_2 \times 2^0$

$(1024)_{10} = (10\cdots 0)_2 \longleftrightarrow \pm(1.b_1 \cdots b_{52})_2 \times 2^1$

$(2046)_{10} = (1111110)_2 \longleftrightarrow \pm(1.b_1 \cdots b_{52})_2 \times 2^{1023}$

$(2047)_{10} = (1111111)_2 \longleftrightarrow \pm\infty$ if $b_1 = \cdots = b_{52} = 0$, NaN else

|  | Emin | Emax | Nmin | Nmax |  | 38 |
|---|---|---|---|---|---|---|
| single | −126 | 127 | $2^{-126}$ | $\sim 2^{128} \sim 3.4\times10^{38}$ |
| double | −1022 | 1023 | $2^{-1022}$ | $\sim 2^{1024} \sim 1.8\times10^{308}$ |

machine precision $p$ = # bits in significand

machine epsilon $\epsilon$ = gap between 1 & next larger floating pt. #

↑

this differs from definition in book, where $\epsilon$ depends on method of rounding. (round up &

$$1+x > 1 \quad \text{if } x = \epsilon/2$$

double precision = ~~52~~ 53     ( 1 is included)

gap $\epsilon = 2^{-52}$     confusing ??

consider only base 2

floating pt. #'s $b_0 . b_1 \cdots b_m \times 2^E$

$E$ in some somewhat symmetric range $E_{min} \leq E \leq E_{max}$

in single $E$ is is det'd from (not equal to) $a_1 \cdots a_8$, $a_j = 0,1$

might assume $0 \leq E \leq 2^8 - 1 = 255$    256 numbers

instead $E \Leftrightarrow$

numbers between −126 & 127    254 numbers

what happened?

Well, $00000000$ and $0--01$ $_{-126}$
both correspond to $2$
+ $11111111$ corresponds to either ⬚⬚ $\pm\infty$
or $NAN$
depending on the result of a
calculation

then float number $0$ is

$$\pm (0.0 - 0) \times 2^{-126}$$
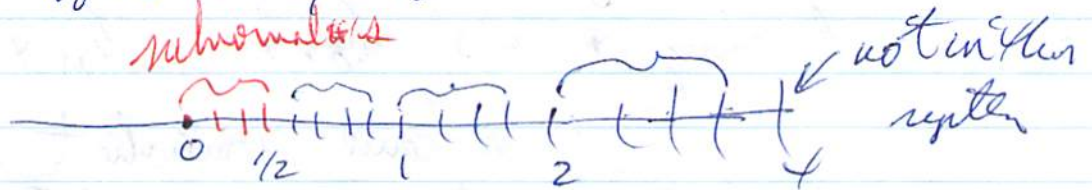$$\underbrace{\phantom{0.0-0}}_{\substack{23 \\ 0's}}$$

Toy number system not the float represent
most computers

$$E = -1, 0, 1, \qquad \pm(b_0 . b_1 b_2) \times 2^{E}$$

normalized if $b_0 \neq 0, (b_0 = 1)$ don't store it

+ $0 \longleftrightarrow b_0 = 0 \Rightarrow b_1 = b_2 = 0$



subnormals                                    not in the
                                              system

$E = -1 \quad (1.00, 1.01, 1.10, 1.11) \times 2^{-1} \qquad \epsilon = \frac{1}{4}$

$\left(\frac{4}{2}, \frac{5}{8}, \frac{6}{8}, \frac{7}{8}\right)$

$E = 0 \quad (1, \frac{5}{4}, \frac{6}{4}, 7/4) \qquad\qquad (\times 2^0)$

$E = +1 \quad (2, \frac{5}{2}, \frac{6}{2}, 7/2)$

add subnormal #'s $(0.00, 0.01, 0.10, 0.11) \times 2^{-1}$