

Week 3.

Rounding : Truncation (Rounding toward zero / Round down)
To the nearest
Toward infinity.

Normalization: $a.bcd \times 10^e$: with only a single digit before decimal point.

Floating-point value $\begin{cases} \text{number} \\ \text{location of the radix point.} \end{cases}$

Significand : the normalized digit part of the value.

In floating-points, the Significand is called Mantissa.

IEEE-754

Range : $1.000 \dots 0_2 \times 2^{-e} \sim 1.11 \dots 1_2 \times 2^e$.

The significand of an IEEE-754 Floating point number is represented in sign and magnitude form.

S.EEEEEEEEE 1.FFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFF

Sign 8-bit

biased exponent

23-bit
fractional significand.

$1 \leq E \leq 254$.

$$x_{10} = (-1)^S \times 2^{(E-B)} \times 1.F$$

$E > 0$: normalized

$E = 0$: not normalized (too small to represent)

$$\Rightarrow x = (-1)^S \times 2^{1-B} \times 0.F$$

$E = 0 \text{ \& \& } F \neq 0 \Rightarrow$ Denormalized underflow number.

Rounded :

Truncation = Round to zero.

Convert to decimal:

$$\begin{aligned} S = 0 & \quad S = + \\ S = 1 & \quad S = - \end{aligned}$$

$$E = \begin{matrix} & (255) \\ 1 & 1 & 1 & 1 & 1 & 1 & 1 & 1 \end{matrix} \begin{cases} F = 0 & \infty \\ F \neq 0 & \text{NaN (not a number)} \end{cases}$$

$$\begin{aligned} & \quad (0) \\ & = 0000 \ 0000 : 2^{-126} (1-127) \\ & \quad \leq 254 : 2^{E-127} \end{aligned}$$

$$\begin{aligned} F : E = 0 & : 0.\dots\dots_{10} \\ E \neq 0 & : 1.\dots\dots_{10} \end{aligned}$$

Convert to 32-bit IEEE-754 FP

$$S : \begin{matrix} + & 0 \\ - & 1 \end{matrix} (-1)^S \quad \text{underflow.}$$

$E: 2^n: n < -126$: too small to be represented as a normalized number \Rightarrow represent in an un-normalized form

$$\Rightarrow \text{exponent} = -126 \Rightarrow E = 0000 \ 0000$$

Round the number to 23 bits nearest

* if the rounded number is at the midway, keep the last digit 0

$$\text{e.g. } 000 \ 0000 \ 0000 \ 0000 \ 0000$$

$$0001 \ 1000 \ / \ 0000 \ 1000$$

$$\Rightarrow 0010 \ 1000 \ / \Rightarrow 0000 \ 0000$$

$n > 127$: too big to be represented, encoded as +inf,
 $\Rightarrow F = 000\ 0000\ 0000\ 0000\ 0000\ 0000$
 $E = 1111\ 1111_2$.

$-126 \leq n \leq 127$: convert to $(n+127)_2$

$F : E = -126 : 0. \dots_2$

$> -126 : 1. \dots_2$