

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 $\left\{ \begin{array}{l} \text{number} \\ \text{location of the radix point.} \end{array} \right.$

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 EEEEEEE 1. FFFFFFF FFFFFFF FFFFFFF FFFFFFF FFFFFFF
 Sign 8-bit 23-bit
 bit biased exponent 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)

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

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

Rounded :

Truncation = Round to zero.

Convert to decimal:

$$\begin{array}{l} S = 0 \quad S = + \\ S = 1 \quad S = - \end{array}$$

$$E = \begin{array}{c} (255) \\ 11111111 \end{array} \left\{ \begin{array}{l} F = 0 \quad \infty \\ F \neq 0 \quad \text{NaN (not a number)} \end{array} \right.$$

$$\begin{array}{c} (0) \\ = 00000000 : 2^{-126} (1-127) \\ 15 \leq 254 : 2^{E-127} \end{array}$$

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

Convert to 32-bit IEEE-754 FP

$$S = \begin{array}{c} + \quad 0 \\ - \quad 1 \end{array} (-1)^S \quad \text{underflow.}$$

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

$$\Rightarrow \text{exponent} = -126 \Rightarrow E = 00000000$$

Round the number to 23 bits nearest

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

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 $(n+127)_2$

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

$> -126 : 1. \dots_2$