

CH-1

Floating point arithmetic

Floating point arithmetic is more accurate and range is more.

In a 2-bit system, without a sign bit there are 4 representations

Comb... sbib		β^{-1}
0	0	00
0	1	01
1	0	10
1	1	11

∴ In a 64-bit system there are 2^{64} representations.

Fixed point numbers:

$$n = \pm (d_1 d_2 \dots d_{k-1} d_k d_{k+1} \dots d_n)_{\beta}$$

where $d_1, \dots, d_n \in \{0, 1, \dots, \beta-1\}$

Floating Point Numbers:

L-1D

$$F \subset R$$

Floating point numbers belong to subset of real numbers

$$F = \pm \underbrace{(0.d_1d_2 \dots d_m)}_{\text{fraction / mantissa / significand}} \times \beta^e$$

where $\beta, d_i, e \in \mathbb{Z}$; $0 \leq d_i \leq \beta-1$ and

we can now represent a large set of numbers but the numbers are not equally spaced. We can only represent a finite set of numbers.

$$\{1.0, \dots, 1.0\} \cup \{1.0, \dots, 1.0\}$$

Conventions:

General / Standard convention

$$\text{1) } \pm (0.d_1 d_2 d_3 \dots d_m)_\beta \times \beta^e$$

In all the conventions the first number has to be non-zero so that the number can never be 0. This will omit the confusion between +0 and -0.

2) Normalized Form:

$$\pm (0.d_1 d_2 d_3 \dots d_m)_\beta \times \beta^e$$

3) Denormalized Form:

$$\pm (1.d_1 d_2 d_3 \dots d_m)_\beta \times \beta^e$$

Example: If given, $m=3$, $e_{\min}=-1$ and $e_{\max}=2$

- Find the highest possible number
- Find the lowest possible number
- Find the lowest possible non-negative number
- Find the number of numbers that can be represented using the system
- If the numbers are sequential or non-sequential.

a) Highest possible value = $+ (0.111)_2 \times 2^2$

$$= (2^{-1} + 2^{-2} + 2^{-3}) \times 4 = 11$$

Lowest possible value = $- (0.111)_2 \times 2^2$

$$= -(2^{-1} + 2^{-2} + 2^{-3}) \times 4 = -3.5$$

c) Lowest possible non-negative value = $(0.100)_2 \times 2^{-1}$

$$= (2^{-1}) \times 2^{-1} = 0.25$$

d) Total representation

$$2 \times 2^2 \times 4 = 32$$



positive, 4 numbers possible
negative in 2 slots

e) $0.111 \times 2^0 \rightarrow 0.875$ } difference 0.125

$0.100 \times 2^1 \rightarrow 0.5$ } difference 0.25

$0.101 \times 2^1 \rightarrow 0.25$

So the numbers are non-sequential numbers

IEEE standard (1985) for double precision 64 bit

64bit distribution : 52 bits for fraction
11 bits for exponent
1 bit for signbit.

$$\rightarrow (0.) d_1 d_2 d_3 \dots d_{52})_2 \times 2^e$$

Now $2^{11} = 2048$. So possible values for e is

from 0 to ... 2047.

But there is a problem - to it has to be 11 bits.

In this state, lowest possible non-negative number is

$$(0.1000 \dots 0)_2 \times 2^0 = (0.1)_2 = (0.5)_{10}$$

But we see we need numbers lower than this quite often. Like 0.001 or 0.005 etc.

To achieve this the e_{\min} has to go lower. But that makes e_{\max} to go down. Which decreases the range.

So IEEE found a suitable position where numbers can be of good range as well as of good precision

The process to achieve this is called exponent biasing.
And for IEEE biasing is done by 1022.

$$\therefore e_{\min} = 0 - 1022 = -1022$$

$$\text{and } e_{\max} = 2047 - 1022 = 1025$$

But we need to reserve some special numbers.

like highest number $+ (0.111\dots1)_2 \times 2^{1025}$ is +infinity

again lowest number $- (0.111\dots1)_2 \times 2^{1025}$ is -infinity

again non-negative lowest $+ (0.1000\dots0)_2 \times 2^{-1022}$ is zero.

So, actual highest = $+ (0.111\dots1)_2 \times 2^{1024}$

and actual non-negative lowest = $+ (0.1000\dots0)_2 \times 2^{-1021}$

numbers greater than +infinity or lower than -infinity gives overflow.

number less than hypothetical zero and actual zero gives underflow.

$(19.25)_{10}$ convert to floating point number where the general convention will be followed for $m=4$

$$\begin{array}{r} 2 \longdiv{19.25} \\ 2 \longdiv{9-1} \\ 2 \longdiv{1-1} \\ 2 \longdiv{0-0} \\ 2 \longdiv{0-0} \\ \hline & 0.25 \\ & \times 2 \\ & 0.50 \\ & \times 2 \\ & 1.00 \end{array}$$

$$19.25 \rightarrow 1001.01$$

1st convention : $\underline{0.1001101} \times 2^5$

so here we need rounding.

$$\frac{1}{2^1} + \frac{0}{2^2} = \frac{0}{2^1} + \frac{0}{2^2} = \frac{0+0.0}{2^1} = 0.0$$

$$(1001.0) =$$

$$(00100010)_2 \approx 1001.0100100010$$

Rounding:

let we have $(0.1000100)_2$,
we have to store this in a system following standard convention for $m=3$. Where $e_{\min} = -3$, $e_{\max} = 3$.

$$\begin{aligned} & \text{Diagram: A number line from } 0.100 \times 2^0 \text{ to } 0.101 \times 2^0. \\ & \text{Left side: } (2^{-1}) \times 1 \\ & \text{Middle: mid} = 0.5 \\ & \text{Right side: } (2^1 + 2^{-3}) \times 2^0 \\ & \text{Calculation: } \frac{0.5 + 5/8}{2} = \frac{9}{16} = \frac{8}{16} + \frac{1}{16} \\ & \quad = 2^{-1} + 2^{-4} \\ & \quad = (0.1001)_2 \end{aligned}$$

Since $0.1000100 < 0.1001$ so fl(0.1000100)

$$= 0.100 \times 2^0$$

if the number is lower than the mid then we round to the left.

if the number is greater than the mid then we round to the right.

If at the center then we go to the even number.

(the number that ends with 0, is an even number).

If we are given to calculate $fl(0.5 \times 0.7)$ then we actually need to calculate $fl(0.5) \times fl(0.7)$

Rounding error:

let say we want to calculate the length of an object of 20cm. But we calculated 18cm.

\therefore error = 2cm.

But this doesn't give the whole idea of error.

For example, if we were measuring a 2cm object but we calculated 0.0m. Error is still 0cm but it is actually a massive error. This error is more significant than prior found in the previous scenario.

For this we need to find scale invariant error S ,

$$= \frac{\text{actual value} - \text{calculated value}}{\text{actual value}}$$

So for rounding error we will consider scale invariant error.

$$S_1 = \frac{|2|}{|20|} \quad S_2 = \frac{|12|}{|12|}$$
$$= 0.1 \quad = 1$$

$S_2 > S_1$. So Error₂ was higher

Machine Epsilon: is not super accurate, but

Highest possible δ is the Machine Epsilon, ϵ_M

For Standard Form,

$$\delta = \frac{|\text{actual value} - \text{calculated value}|}{\text{bin}(\text{actual value})}$$

to maximise δ we need to increase the numerator and decrease denominator.

$$\therefore x_{\min}, \text{ for } m=3,$$

$$= 0.100 \times 2^e$$

$$\text{for } m=4; x_{\min} = 0.1000 \times 2^e$$

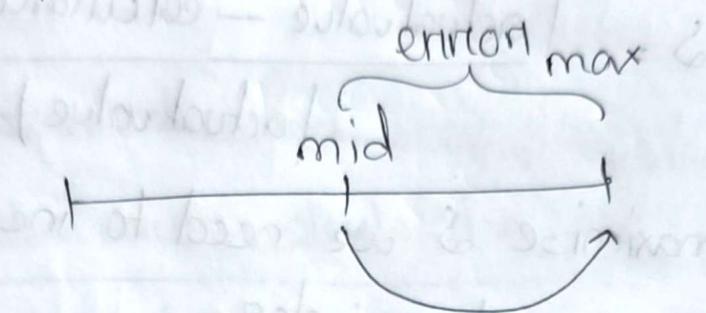
$$\therefore \text{for any } m, x_{\min} = 0.1 \times 2^e$$

$$= \emptyset 2^{-1} \times 2^e$$

$$= \beta^{-1} \beta^e$$

Now maximum value for numerator comes when input is at the center and rounded to one of the ends

\therefore we need to find the distance between midpoint and one of the ends.



for $m=3$

left can be 0.100×2^e

and right " " 0.101×2^e

$$\therefore \text{difference} = 0.001 \times 2^e = 2^{-3} \times 2^e$$

for $m=4$

left can be 0.1000×2^e

and right can be 0.1001×2^e

$$\therefore \text{difference} = 0.0001 \times 2^e = 2^{-4} \times 2^e$$

$$\therefore \text{difference} = 2^{-m} \times 2^e = \beta^{-m} \times \beta^e$$

\therefore numerator: $\frac{1}{2} \times \beta^{-m} \times \beta^e$ [\because we actually need half the distance]

$$G_N = S_{\max} = \frac{\frac{1}{2} \beta^{-m} \beta^e}{\beta^{-1} \beta^e} = \frac{\frac{1}{2} \beta^{1-m}}{\beta^0}$$

Normalized form,

$$\text{For } m=3; x_{\min} = 0.1000 \times 2^e$$

$$\text{For } m=4; x_{\min} = 0.10000 \times 2^e$$

$$\therefore \text{for any } m; x_{\min} = 0.1 \times 2^e$$

$$= 2^{-1} \times 2^e$$

$$= \beta^{-1} \times \beta^e$$

$$\text{for } m=3$$

$$\text{left can be } 0.1000 \times 2^e$$

$$\text{right can be } 0.1001 \times 2^e$$

$$\therefore \text{difference} = 0.0001 \times 2^e$$

$$= 2^{-4} \times 2^e$$

$$\therefore \text{numerator} = \frac{1}{2} \times 2^{-m-1} \times 2^e = \frac{1}{2} \times \beta^{-m-1} \times \beta^e$$

$$\text{for } m=4$$

$$\text{left can be } 0.10000 \times 2^e$$

$$\text{right can be } 0.10001 \times 2^e$$

$$\therefore \text{difference} = 0.00001 \times 2^e$$

$$= 2^{-5} \times 2^e$$

$$\therefore G_M = S_{\max} = \frac{\frac{1}{2} \beta^{-m-1} \times \beta^e}{\beta^{-1} \times \beta^e} = \frac{1}{2} \beta^{-m}$$

For denormalized form,

$$\text{For } m=3, x_{\min} = 1.000 \times 2^e$$

$$\text{for } m=4, x_{\min} = 1.0000 \times 2^e$$

$$\therefore \text{for any } m, x_{\min} = 2^0 \times 2^e$$

$$\text{for } m=3$$

$$\text{left can be } 1.000 \times 2^e$$

$$\text{right can be } 1.001 \times 2^e$$

$$\therefore \text{difference} > 0.001 \times 2^e$$

$$= 2^{-3} \times 2^e$$

$$\text{for } m=4$$

$$\text{left can be } 1.0000 \times 2^e$$

$$\text{right } " = 1.0001 \times 2^e$$

$$\therefore \text{difference} = 0.0001 \times 2^e$$

$$= 2^{-4} \times 2^e$$

$$\therefore \text{for any } m, \text{ numerator} = \frac{1}{2} \times 2^{-m} \times 2^e$$

$$\frac{1}{2} \beta^{-m} \beta^e$$

$$\therefore G_M = S_{\max} = \frac{\frac{1}{2} \beta^{-m} \beta^e}{\beta^e} = \frac{1}{2} \beta^{-m}$$

loss of significance,

we know,

$$fl(n) = n + s_1 n$$

$$\therefore fl(y) = y + s_2 y$$

Now, $n+y = fl(n) + fl(y)$

$$= (n + s_1 n) + (y + s_2 y)$$

$$= (n+y) \left(1 + \frac{n s_1 + y s_2}{n+y} \right)$$

scale invariant error

for any two numbers n & y if we try n & y then the error becomes massive

because, for $n-y$; error $\left(1 + \frac{x s_1 - y s_2}{n-y} \right)$

and if $n-y \approx 0$; then error $\approx \infty$

So for subtraction of two very close numbers we get
loss of significance

example:

$$x^2 - 56x + 1 = 0.$$

$$x = \frac{-b \pm \sqrt{b^2 - 4ac}}{2a}$$

$$x_1 = 28 + \sqrt{783} = 55.98$$

$$x_2 = 28 - \sqrt{783} = 0.01786$$

Let take a system calculating upto 4 significant place,

$$\therefore x_1 = 28 + \sqrt{783} = 28 + 27.98 = 55.98$$

$$x_2 = 28 - \sqrt{783} = 28 - 27.98 = 0.0200$$

$$\therefore 0.02 \approx 0.01786$$

This happens because of loss of significance as we subtracted two very close numbers 28 and 27.98.

Solution:

$$\text{We know } n^2 - (\alpha + \beta)n + \alpha\beta$$

$$\alpha = n_1,$$

$$\beta = \frac{1}{\alpha} = \frac{1}{n_1} = n_2$$

again, (Alternative solution)

equating $an^2 + bn + c = 0$,

$$n_1 \times n_2 = \frac{c}{a}$$

$$\therefore n_1 n_2 = 1$$

$$\therefore n_2 = \frac{1}{n_1}$$

$$\text{Now } n_1 = 59.98$$

$$\therefore n_2 = 0.01667.$$