

Fixed point \rightarrow decimal position of representing a number is fixed

$\Rightarrow 12345.67$
 $\underbrace{\hspace{1.5cm}}_{5 \text{ positions}} \quad \underbrace{\hspace{1.5cm}}_{\text{only 2 positions}}$

$$\begin{aligned}
 -123456 \times 10^{-1} &= -12345.6 \times 10^0 \\
 &= -1.23456 \times 10^4 \text{ (normalised)} \\
 \text{Error} - \left[\begin{aligned} &\approx -0.12345 \times 10^5 \\ &\approx -0.01234 \times 10^6 \end{aligned} \right]
 \end{aligned}$$

Sign Exponent significand

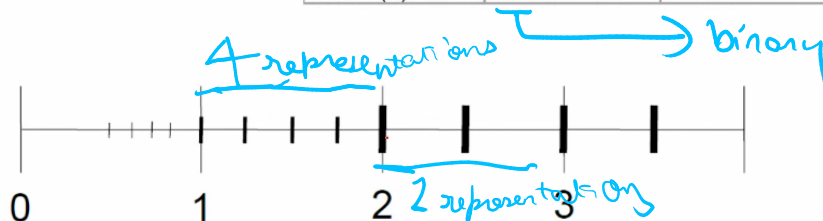
$$(-1)^{\text{sign}} \times \text{significand} \times \text{radix}^{\text{exponent}}$$

Distribution of Floating Point Numbers

- 3-bit mantissa
- Exponent: $\{-1, 0, +1\}$

unequally distributed

e = -1	e = 0	e = 1
$1.00 \times 2^{(-1)} = 1/2$	$1.00 \times 2^0 = 1$	$1.00 \times 2^1 = 2$
$1.01 \times 2^{(-1)} = 5/8$	$1.01 \times 2^0 = 5/4$	$1.01 \times 2^1 = 5/2$
$1.10 \times 2^{(-1)} = 3/4$	$1.10 \times 2^0 = 3/2$	$1.10 \times 2^1 = 3$
$1.11 \times 2^{(-1)} = 7/8$	$1.11 \times 2^0 = 7/4$	$1.11 \times 2^1 = 7/2$



Fixed point: Represent less, Floating point: Represent more but unbalanced

$\begin{aligned} 1.00 &\rightarrow 1 \\ 1.01 &\rightarrow 1.25 \\ 1.10 &\rightarrow 1.5 \\ 1.11 &\rightarrow 1.75 \end{aligned}$

more bits for significand \Rightarrow more accuracy

More bits for exponent increases range

$$-101.001101 \times 2^{111001} = -1.01001101 \times 2^{111011} \quad (\text{normalized})$$

if significand $\in (0, 1]$ (without sign)
then it's normalized

For decimal

Significand $\in (0, 1]$

$$0.001 \times 10^0 \text{ to } 9.999 \times 10^{99}$$

max $\times 10$, exp 2, Significand 4 places

fixed point is much less even with all 6 ~~bits~~ places
for integer and 0 for decimal
max number possible = ~~1000~~ 999999

IEEE 754 FPR

approx of Real numbers - float double
approx of Integers - int

float = 32 bit single precision

double: 64 bit double precision

GC treats [float] as double precision
[constant]

Sign magnitude ~~X~~

not 2's complement

31	30 - 23	22 - 0
sign	exp	significand

31	30 - 20	19 - 0 upper
sign		31 - 0 lower

→ sign 0 - positive 1 - negative

→ Mantissa is normalized, always has a leading
binary point 1 bit $(1.0 \leq \text{significand} < 2)$
'0' is a special case

Exponent biased exponent = actual + Bias
Ensures exponent is unsigned
bias of 127 for single precision
1023 for double p

Why not sign bit?

(1) You need 2's complement for signed bit
(unsigned is easier to handle)

(2) Special numbers cannot be represented
eg. 0, NaN etc.

Normalized

Types of Data

- Data represented in this format are classified in five groups.

- Normalized numbers, ✓✓
- Zeros, ✓
- Subnormal(denormal) numbers, ✓✓
- Infinity and not-a-number (nan), ✓✓

- Single Precision data interpretation

Single Precision		Data Type
Exponent	Significand	
0 ✓✓	0 ✓✓	± 0
0	nonzero	\pm subnormal number
1 - 254	anything	\pm normalized number
255	0	$\pm \infty$
255	nonzero	NaN (not a number)

if we used sign:

-128 to 127

we may use -128 for '0' } or other way
 127 for '0' }

our range is -127 to 126

biased: -126 to 127 } this representation is way
more significant than

2^{1-127}