

1) The three components are

- 1) sign bit S
- 2) Exponent (E)
- 3) fractional part (M)

Number is $F = (-1)^S M \times 2^E$

Two kinds of encoding done

- i) single precision (32 bits, E 8 bits, M 23 bits)
- ii) double precision (64 bits, E 11 bits, M 52 bits)

1 sign bit for each type

⇒ The precision depends on fractional part.

⇒ The number of significant digits depends on the number of bits in M

when $M = 24$ bits

we have x significant digits

where $2^{24} = 10^x \Rightarrow x = 7.2 \dots$

⇒ $x = 7 \Rightarrow$ up to seven significant

decimal places we can represent in binary floating point representation.

⇒ The range of number depends on the number of bits in exponent that is $2^8 - 1 = 255$

⇒ The exponent is encoded as a biased value

$E = \text{exp} + \text{bias}$ where $\text{bias} = 127 = (2^{8-1} - 1)$ for single precision

$\text{bias} = 1023 (2^{11-1} - 1)$ for double precision

⇒ when the number's ^{decimal part} lies between 0 to $\frac{1}{2^{52}}$ it can be represented accurately using double precision representation.

⇒ when number is

$1.x$

if $0 < x < \frac{1}{2^{23}}$
and it's represented in single precision,
it has 100% accuracy

⇒ when number is

$1.x$

if $0 < x < \frac{1}{2^{52}}$

and it's represented in double precision,
it has 100% accuracy.

2) When the content of the exponent $(e) \neq 0$ & significant (m) is $\neq 0$ then the subnormal number is

$$(-1)^s \times 0.m \times 2^{-126}$$

normalized number

subnormal

Smallest 2^{-126}

$$2^{-149} \approx 0$$

Largest 3.4×10^{38}

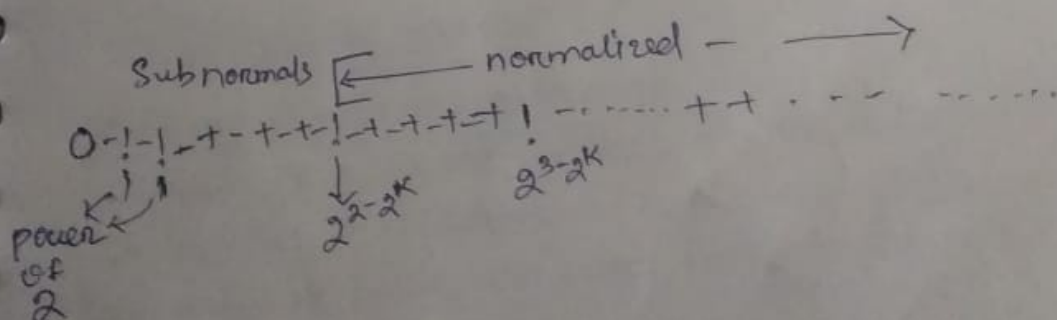
$$0.99999988 \times 2^{-126}$$

\approx smallest normalized number

\Rightarrow So there are total $[0.0, 2^{-126}]$ that means 2^{23} numbers within the range.

\Rightarrow The smallest difference between 2 normalized number is 2^{-149} which is equal to difference between any two consecutive subnormal numbers.

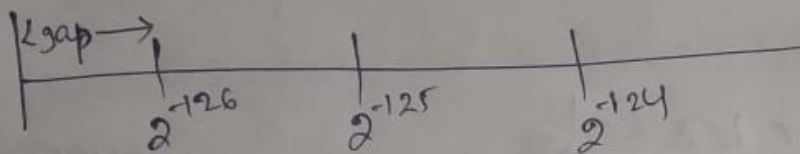
\Rightarrow Meanwhile the largest difference b/w 2 consecutive numbers is 2^{104}



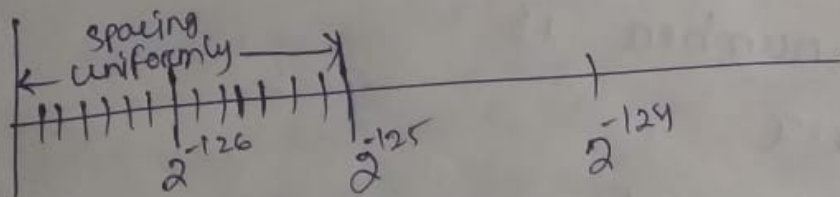
⇒ Subnormals extend range of magnitudes ^{page-5} representable but have less precision than normalised numbers.

⇒ for a 32 bit precision type the number line distinguishes between normal & subnormal values in the figure below

i) without subnormal

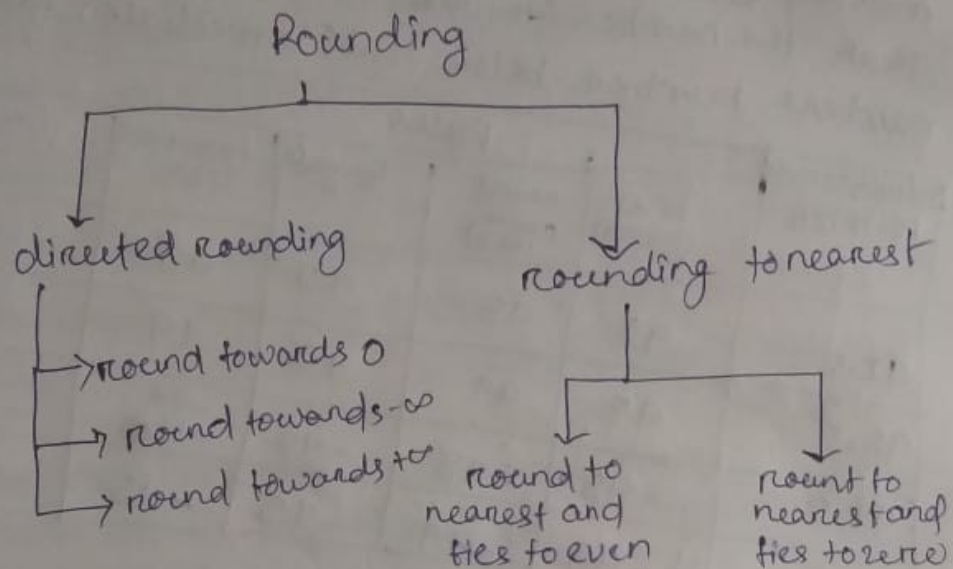


ii) with subnormal



Q3)

page-6
⇒ There are basically five rounding rules.
And they are categorized into two subsets.



a) Directed rounding (truncation)

i) Round towards 0

Directed rounding towards zero

ii) Round towards $+\infty$ (ceiling)

Directed rounding towards $+\infty$

iii) Round towards $-\infty$ (floor)

Directed rounding towards $-\infty$

b) Rounding to nearest

i) Round to nearest, ties to even

rounds to nearest value if the value falls in middle, it is rounded to the nearest number i.e it must have an even least significant digit.

ii) Round to nearest, ties away from zero | page - 7

Rounds to the nearest value, if the number falls midway then for positive number number greater than the number (^{nearest} ~~integer~~) taken and for negative numbers number below that (immediate) taken.

Values to be rounded Number	Rules				
	nearest (to even)	nearest (away from 0)	towards 0	towards ∞	towards $-\infty$
97.5	98	98	97	98	97
98.5	98	99	98	99	98
-97.5	-98	-98	-97	-97	-98
-98.5	-98	-99	-98	-98	-99