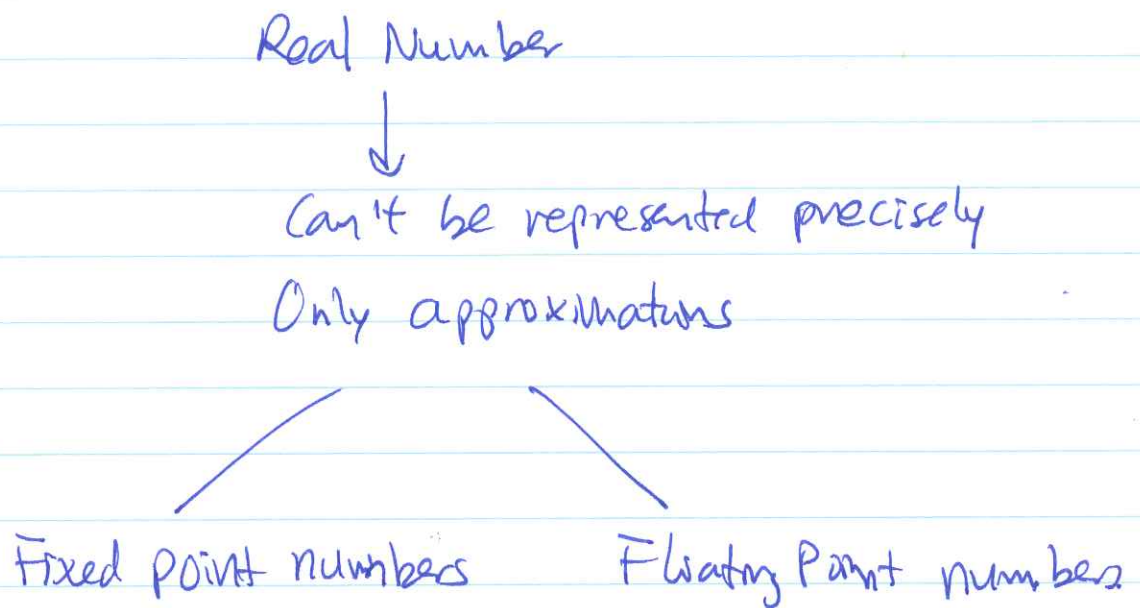


Real Numbers

4



- Before we discuss the representation of fixed and floating point numbers, we look at how fractions are converted between decimal & binary number systems.

Converting fractions

• Binary → Decimal

$$(0.01101)_2 = \begin{matrix} 0 & -1 & -2 & -3 & -4 & -5 \\ 0 & 0 & 1 & 1 & 0 & 1 \\ & & 2^{-2} & 2^{-3} & & 2^{-5} \end{matrix}$$

$$\frac{1}{4} + \frac{1}{8} + \frac{1}{32} = \frac{8+4+1}{32} = \frac{13}{32} = (0.40625)_{10}$$

• Decimal → Binary : multiply by 2 & note integer part.

$$(0.6875)_{10} = (\dots)_2$$

	x2	<u>0.6875</u>	
take remainder		<u>1.3750</u>	1
	x2	<u>0.3750</u>	
take remainder		<u>0.750</u>	0
	x2	<u>0.750</u>	
take remainder		<u>1.5</u>	1
	x2	<u>0.5</u>	
take rem.		<u>1.0</u>	1
		0	

read.

↓

answer.

$$(0.1011)_2$$

done.

$$(0.1)_{10} = (\dots)_2$$

	x2	<u>0.4</u>	0
	x2	<u>0.8</u>	0
	x2	<u>1.6</u>	1
	x2	<u>0.6</u>	
	x2	<u>1.2</u>	1
	x2	<u>0.2</u>	
	x2	<u>0.4</u>	0

repeat

ans: 0.011001100110... repeats.

- Decimal \rightarrow octal multiply by 8, note integer part.

$$(0.12)_{10} = (\dots)_8 = (0.0753\dots)_8$$

$$\begin{array}{r} 0.12 \\ \times 8 \\ \hline 0.96 \end{array} \quad 0$$

$$\begin{array}{r} 0.96 \\ 8 \\ \hline 7.68 \end{array} \quad 7.$$

$$\begin{array}{r} 0.68 \\ 8 \\ \hline 5.44 \end{array} \quad 5.$$

$$\begin{array}{r} 0.44 \\ 8 \\ \hline 3.52 \end{array} \quad 3$$

\vdots
 \vdots etc.

Fixed point ~~int~~ numbers

• integer :
$$\begin{array}{cccccccc} & 7 & 6 & 5 & 4 & 3 & 2 & 1 & 0 \\ (10101010)_2 & = & 2^7 & + & 2^5 & + & 2^3 & + & 2^1 \\ & = & 128 & + & 32 & + & 8 & + & 2 & = & \underline{\underline{170}} \end{array}$$

It is assumed that the decimal point lies at the right most digit :

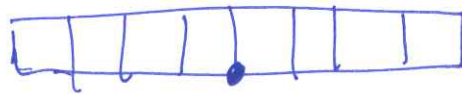
10101010.

• Fixed point number :

is very similar to integer representation, except the decimal point is assumed to lie elsewhere.

Note: computer CANNOT place a decimal point explicitly. It must assume ~~the local case~~ a fix place for it.

eg: Assume decimal point lies at :



↑
imaginary point !!!

Then the ~~number~~ bit pattern

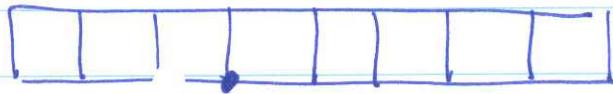
10101010

means:
$$\begin{array}{cccccccc} & 3 & 2 & 1 & 0 & -1 & -2 & -3 & -4 \\ & 1 & 0 & 1 & 0 & . & 1 & 0 & 1 & 0 \end{array}$$

$$8 + 2 \cdot \frac{1}{2} + \frac{1}{8} = 10 \cdot \frac{5}{8} = \underline{\underline{10.625}}$$

Converting Fixed Point Binary numbers.

Assume binary fixed point representation is:



Binary to decimal conversion:

Pattern: $\boxed{01101100}$

This implies (assume decimal point):

011.01100

$2^1 2^0 2^{-1} 2^{-2} 2^{-3}$

$$\text{Value} = 2 + 1 + \frac{1}{4} + \frac{1}{8}.$$

$$= 3 \frac{3}{8}$$

$$= 3.375$$

Decimal to binary (fixed point) conversion:

$$(3.14)_{10} = ?$$

Split 3.14 into integer part $\rightarrow 3$
and fraction part $\rightarrow 0.14$.

Integer part converted as follows:

$$\begin{array}{r} 2 \overline{) 3} \\ \underline{2} \\ 1 \\ 2 \overline{) 1} \\ \underline{0} \\ 0 \end{array}$$

$$(3)_{10} = (11)_2$$

Fraction part converted as follows:

$$\begin{array}{r} \times 2 \quad \underline{0.14} \\ \quad \underline{0.28} \end{array} \rightarrow 0$$

$$\begin{array}{r} \times 2 \quad \underline{0.28} \\ \quad \underline{0.54} \end{array} \rightarrow 0$$

$$\begin{array}{r} \times 2 \quad \underline{0.54} \\ \quad \underline{1.12} \end{array} \rightarrow 1$$

$$\begin{array}{r} \times 2 \quad \underline{0.12} \\ \quad \underline{0.24} \end{array} \rightarrow 0$$

$$\begin{array}{r} \times 2 \quad \underline{0.24} \\ \quad \underline{0.48} \end{array} \rightarrow 0$$

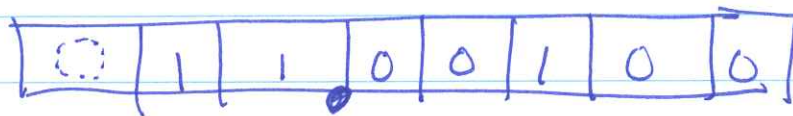
$$\begin{array}{r} \times 2 \quad \underline{0.48} \\ \quad \underline{0.96} \end{array} \rightarrow 0$$

$$(0.14)_{10} = (0.001000\dots)_2$$

So:

$$(3.4)_{10} = (11.001000\dots)$$

Fit in 8 bits and line up the decimal point:



↑
answer.

Conclusion: ~~Real numbers are NOT represented precisely.~~

- Integer numbers can be represented precisely by computers.
- Non-integer numbers CANNOT.
There will be a round off / truncation error for non-integer numbers.
The magnitude of the error depends on the number of bits used to represent the fraction.

Representing Floating Point numbers

- Floating point \neq Real !

Real number can have infinite number of digits.

- Floating point number n :

$$n = m \times 10^e \quad (\text{for humans}).$$

m ← mantissa ↑ exponent.

eg: $19.47 = 0.1947 \times 10^2$
 $= 1947.0 \times 10^{-2}$ etc.

~~• Not every possible number~~

- Mantissa & exponent are of finite length.
- Not every possible numerical value can be expressed.

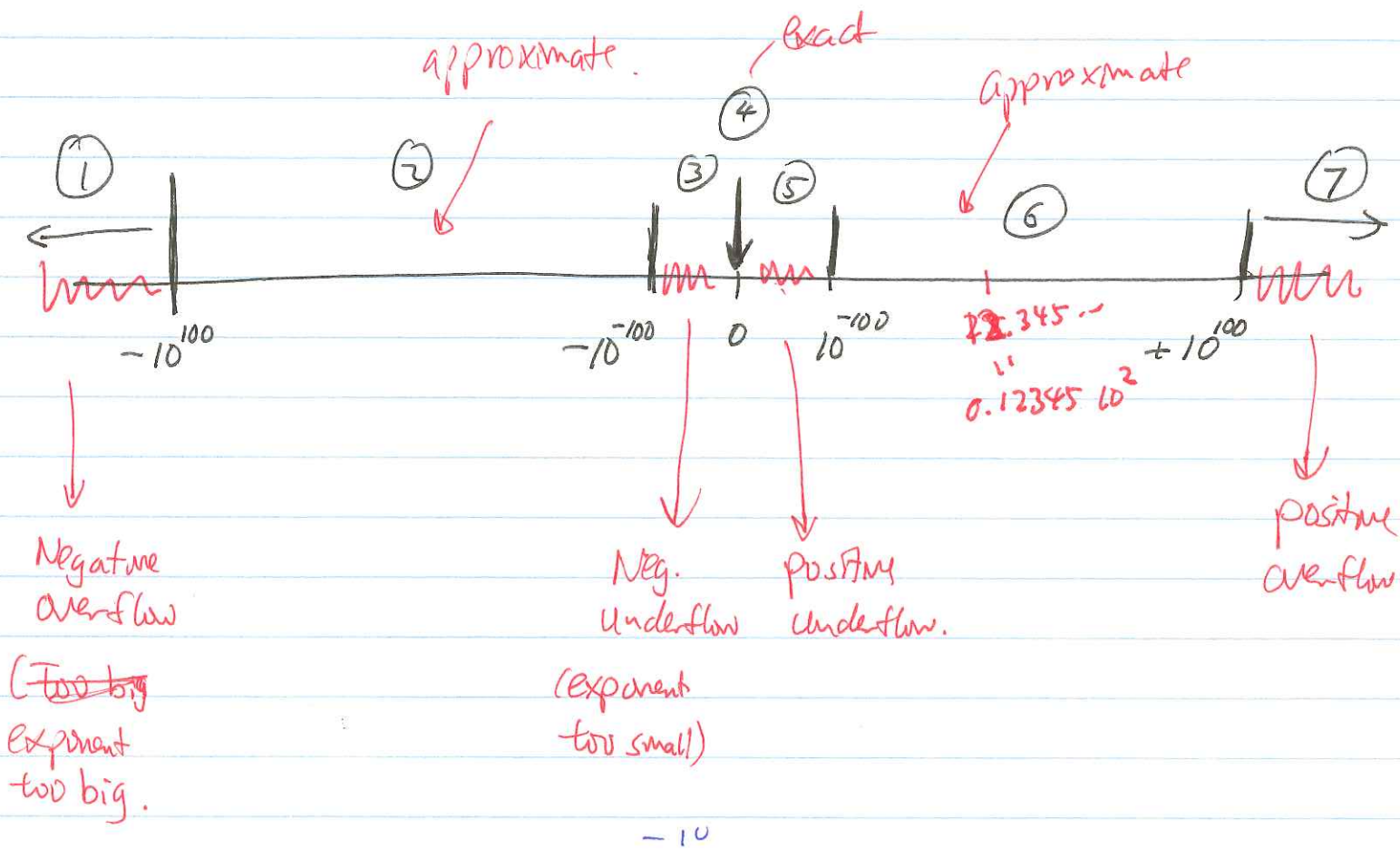
Eg: S_1

--	--	--	--	--

 3 digit

Example: 3 decimal digit mantissa & decimal |exponent| ≤ 99

The real numbers ~~then~~ can be divided into 7 regions:



Any number between $(-10^{100}, -10^{-10})$.

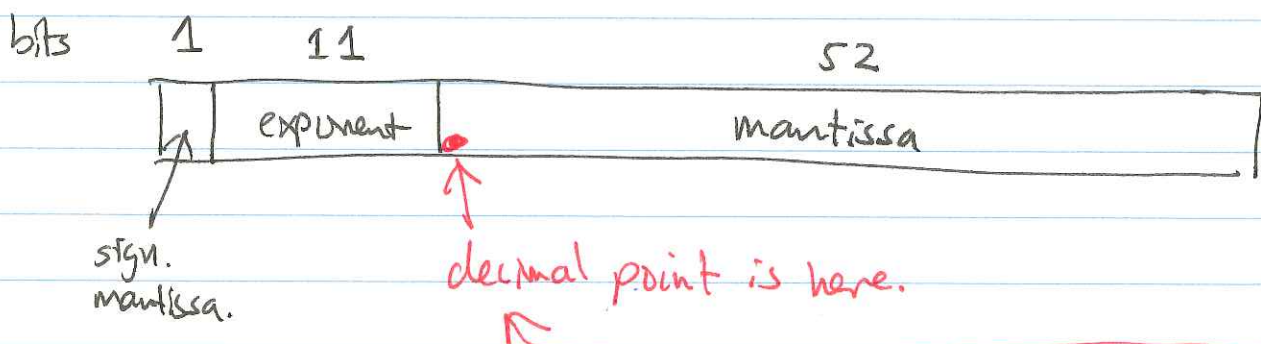
IEEE floating-point standard 754

- IEEE 754 defines 3 standard formats for floating point numbers.

(Most computers now use these formats).

- single precision. (32 bits)
- double precision (64 bits)
- Extended precision (80 bits)

- Double precision format: (others are similar)



- the mantissa is between 1 & 2. :

$$1 \leq m < 2. \quad (\text{Decimal point assumed HERE})$$

- exponent is (2^e) in excess 2^{10} .

(Thus: the mantissa is signed magnitude.

the exponent is excess 2^{10}).

Operations with floating point numbers:

- multiply & divide : (1) multiply / divide mantissa
(2) add / subtract exponent

(3) Normalize !

↓

bring result back to a form that $1 \leq \text{mantissa} < 2$

- add & subtract : (1) bring both operands to same exponent value.

(2) add / subtract mantissa
(exponent = exponent of one operand)

(3) Normalize.

Floating Point Number

- Short coming of fixed point number : bound by range
- A better way to represent real number is the "scientific" notation :

$$\text{real number} = \text{fraction} \times 10^{\text{exponent}}$$

↑
fraction is also called mantissa

- The computer version of this way of representing real numbers is called "floating point".
- Example :

$$\begin{aligned} 314 &= 3.14 \times 10^2 &= 0.314 \times 10^3 \\ 0.000012 &= 1.2 \times 10^{-5} &= 0.12 \times 10^{-4} \end{aligned}$$

- ~~There~~ There are many ways to represent the same number :

$$3.14 = 3.14 \times 10^0 = 0.314 \times 10^1 \text{ etc.}$$

↑
form 1.

One form is usually chosen as the standard.
(ie: where to put the decimal point).

- Standard forms typically chosen are:

(1) mantissa ~~is~~ consists of 1 digit ($\neq 0$)

number stored as:

eg:	$314 = 3.14 \times 10^2$	31400 +02
	$31.4 = 3.14 \times 10^1$	31400 +01
	$1255 = 1.255 \times 10^3$	12550 +03
	$0.00012 = 1.2 \times 10^{-4}$	12000 -04

(2) mantissa is 0 and leading digit of fraction must not be 0.

number stored as

eg:	$314 = 0.314 \times 10^3$	31400 +03
	$31.4 = 0.314 \times 10^2$	31400 +02
	$1255 = 0.1255 \times 10^4$	12550 +04
	$0.00012 = 0.12 \times 10^{-3}$	12000 -03

- Now do this in binary:

~~$$101 = 10.1 \times 2^1 = 1.01 \times 2^2 = 0.101 \times 2^3$$~~

Note that:

$$101 = 10.1 \times 2^1$$

$$\uparrow = 1.01 \times 2^2$$

$$\text{"5"} = 0.101 \times 2^3$$

Here also, there are many way to represent the same floating point number.

One way is picked as the standard.

- Eg: mantissa consists of 1 digit ($\neq 0$)

Then: number	standard repr	internally stored as:
101	$= 1.01 \times 2^2$	10100000 0010
0.101	$= 1.01 \times 2^{-1}$	10100000 "0001"
-101	$= -1.01 \times 2^2$	"-10100000" 0010"
-0.101	$= -1.01 \times 2^{-1}$	"-10100000" "0001"

SU: mantissa AND exponent can be pos. & neg. numbers !!!

- Single precision floating point number: 4 bytes.

Standard form: 1 digit mantissa ($\neq 0$).

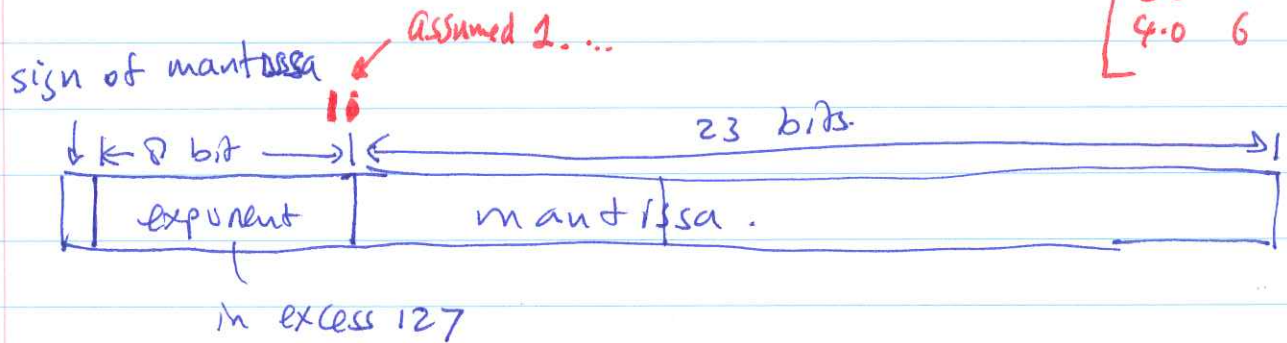
Example:

$$101 = 1.01000 \dots \times 2^2$$

`printfloat.c`

use

1.0	1.5	1.75
2.0	3	3.5
4.0	6	7



- mantissa is stored in sign-magnitude representation
 unsigned value
 or absolute value.

where the assumed decimal point is immediately after the leading digit.

(mantissa is always 1.)

Also: the leading 1 in mantissa is NOT stored, but it is assumed (since it's always there!)

- Exponent is stored in 8 bits with excess 127 encoding (i.e.: value is stored as absolute (value+127))

o Example:

$$(1.0)_{10} = (1.0)_2 \times 2^0$$

$$\text{Mantissa} = \underbrace{1.00000 \dots 0}_{23 \text{ bits.}} \quad \text{sign} = \emptyset.$$

$$\text{exponent} = 0$$

representation in excess 127

value x is repr. by absol. value $x+127$

0 is repr. by ~~abs~~: unsigned value

$$0+127 = 127 = 0111.1111$$

(8 bits).

Put it in the format specified:



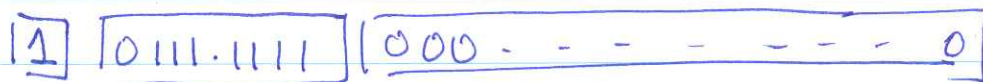
Example: $(-1.0)_{10}$

$$(-1.0)_{10} = (-1.0)_2 \times 2^0$$

mantissa = sign = 1
magnitude = $1.000\dots 0$
23 bits.

exponent = 0 \rightarrow 127 = 0111.111

Repr:



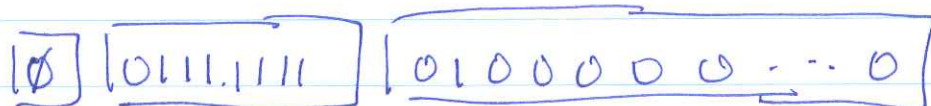
Example: $(1.25)_{10}$

$$(1.25)_{10} = (1.01)_2$$
$$= (1.01)_2 \times 2^0$$

mantissa : sign = 0.
magnitude = $1.010000\dots$
23 bits.

exponent = 0.

Repr:



Example:

$$(16.25)_{10} = (10000.01)_2 \\ = (1.000001)_2 \times 2^4$$

mantissa: sign = 0
magnitude = $\underbrace{1.0000010000000}_{23 \text{ bits}}$

Exponent = 4

in excess 127 repr. by unsigned $4 + 127$
 $= 131$

$$131 = 128 + 3 \\ = 1000.0000 + 11 \\ = 1000.0011$$

Representation is:

$$\boxed{0} \boxed{1000.0011} \boxed{00001000 \dots 0}$$

Given: repr:

$\boxed{1} \boxed{0111.1110} \boxed{110000 \dots 0}$

What is the ~~number~~ value?

Mantissa: sign = 1 \rightarrow neg.

magnitude = 1.110000

$$= \cancel{1 + \frac{1}{2} + \frac{1}{4}} = \cancel{1.75}.$$

so:

$$\cancel{-1.75}$$

Exponent: 0111.1110 = 126

repr. value

$$x + 127 = 126$$

$$x = 126 - 127$$

$$= -1.$$

so:

$$\text{value in binary} = -1.11 \times 2^{-1}$$

$$= -0.111$$

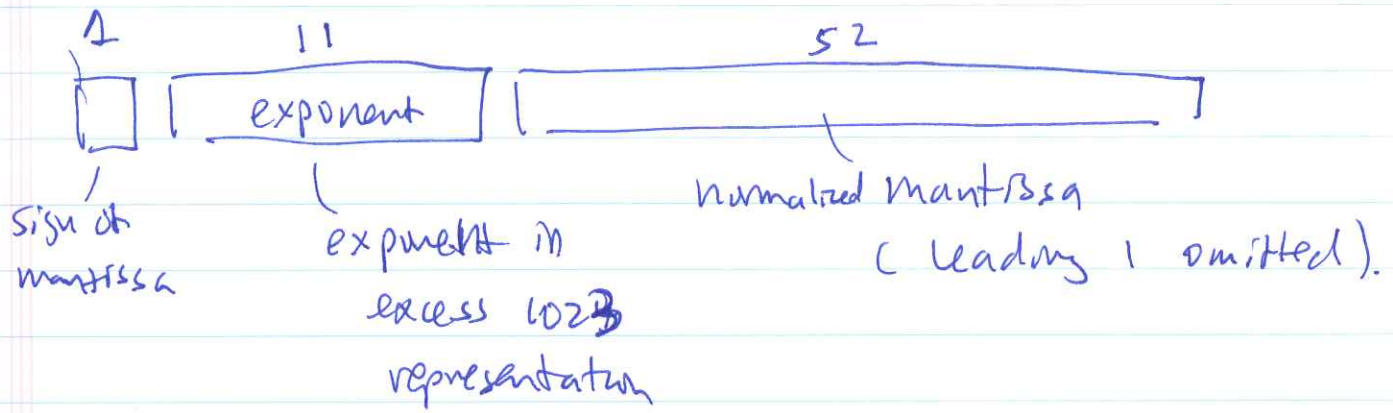
$$= -\left(\frac{1}{2} + \frac{1}{4} + \frac{1}{8}\right)$$

$$= -\frac{7}{8}$$

$$= -0.875$$

==

• Double precision IEEE format :



The IEEE Floating point standard

- Upto about 1980, each computer manufacturer has its own way to represent floating point numbers.
- To ease exchange of data (information) between computers, a universal standard for floating point numbers was established by IEEE (Institute of Electrical & Electronics Engineers) in 1985 (The IEEE standard 754)
- Most modern CPUs (Motorola, SPARC, MIPS) include a floating point processor and all of them conform to the IEEE standard.
- The standard defines 3 types of floating point numbers:

- | | |
|----------------------|--------------------|
| (1) Single precision | (float in C) |
| (2) double precision | (double in C) |
| (3) Quad precision | (long double in C) |

They are all similar.