# Floating point numbers.

The fact that infinitely many real nos. have to be stored in finite amount of space gives rise to 2 limitations:

    ① the represented numbers cannot be arbitrarily large or small

    ② there will have to be gaps between them.

Since any real number has to rounded off to the closest represented number, this introduces rounding errors.

Several diff. representations have been proposed but by far the most widely used is the floating point representation.

The floating point number system is a sub-set $\mathbb{F} \subseteq \mathbb{R}$ determined by a base $\beta$ (an integer $\geq 2$) and an integer $p \geq 1$, known as precision.

The elements of $\mathbb{F}$ are $0$ along with all numbers of the form $\pm \dfrac{m}{\beta^p} \times \beta^e$, where $m$ is an integer

$$1 \leq m \leq \beta^p,$$

$e$ is an arbitrary integer, called exponent.

A floating point number is represented as −

$$\pm \underbrace{d.d \cdots d}_{\text{significand}} \times \underset{\text{base}}{\beta}^{\overset{e \;\rightarrow\; \text{exponent}}{}}$$

$\underset{\text{sign}}{}$    significand (has $p$ digits)

More precisely, the number $\pm \left( d_0 + d_1 \beta^{-1} + \cdots + d_{p-1} \beta^{-(p-1)} \right) \times \beta^e$

                           (each $0 \leq d_{\pm} < \beta$)

is stored as −

$$d_0 . d_1 \cdots d_{p-1} \times \beta^e .$$

eg ① $\beta = 10$, $p = 3$ :   $0.1$ is represented as   $1\ 0\ 0 \times 10^{-1}$

② if $\beta = 2$, $p = 24$, then $0.1$ cannot be represented exactly. It is approximated by the nearest floating point number

$$1.100110011001100110011001101 \times 2^{-4}.$$

$$\left( \begin{array}{l} 0.1 = \dfrac{1}{10} = \dfrac{1}{1010} \quad, \quad \text{long division} \quad 1010\overline{)1} \\ \text{(decimal)} \quad \text{(binary)} \qquad\qquad \text{in binary} \end{array} \right)$$

③ Let $\beta = 2$, $p = 3$, $e_{min} = -1$, $e_{max} = 2$.

There are 16 normalized floating point numbers.

⌣ ⌣ ⌣ , allowed digits: $0, 1$ $\left( \text{since } 0 \le d_t < \beta \right)$.

$$\left( d_0 . d_1 d_2 \times 2^e \right) \qquad \left( d_0 + d_1 \times 2^{-1} + d_2 \times 2^{-2} \right) \times 2^e$$

Floating pt. representation · · · · · · · · · real number being represented.

| Floating pt. representation | real number |
|---|---|
| $1.00 \times 2^{-1}$ | $\longrightarrow$ 0.5 |
| $1.01 \times 2^{-1}$ | 0.625 |
| $1.10 \times 2^{-1}$ | 0.75 |
| $1.11 \times 2^{-1}$ | 0.875 |
| $1.00 \times 2^{0}$ | $\longrightarrow$ 1 |
| $1.01 \times 2^{0}$ | 1.25 |
| $1.10 \times 2^{0}$ | 1.5 |
| $1.11 \times 2^{0}$ | 1.75 |
| $1.00 \times 2^{1}$ | $\longrightarrow$ 2 |
| $1.01 \times 2^{1}$ | 2.5 |
| $1.10 \times 2^{1}$ | 3 |
| $1.11 \times 2^{1}$ | 3 5 |
| $1.00 \times 2^{2}$ | $\longrightarrow$ 4 |
| $1.01 \times 2^{2}$ | 5 |
| $1.10 \times 2^{2}$ | 6 |
| $1.11 \times 2^{2}$ | 7 |

$$1.00 \leftrightarrow 1 \times 2^{-1} = \dfrac{1}{2} = 0.5$$

$$1.01 \leftrightarrow \left( 1 + 1 \times 2^{-2} \right) \times 2^{-1} = \left( 1 + \dfrac{1}{4} \right) \times 2^{-1}$$

$$= \dfrac{5}{4} \times \dfrac{1}{2} = \dfrac{5}{8} = 0.625.$$

& so on.

Notice that the floating point numbers are not equally spaced.

The numbers represented are: 



and their negatives.

So in this case, $|F| = 33$ $\left( 0, \text{positives & negatives} \right)$.