

## Floating Point Form of Numbers:

In decimal notation, every real number is represented by a finite or an infinite sequence of decimal digits. Most of computers have two ways of representing numbers called Fixed Point and Floating Point.

In fixed point system, all numbers are given with a fixed number of digits after decimal (and sometimes also before decimal).

Standard Form:

I I I I . F F F F (Decimal System)

9 9 9 9 . 9 9 9 9 (largest number)

0 0 0 0 . 0 0 0 1 (smallest number)

For examples, 62.3248, 0.0142, 0.1000 etc.

Fixed point system (representation) is impractical.

In floating point system, we write a number in the form

$$x = \pm m \cdot 10^e \rightarrow \text{exponent}$$

↑  
mantissa

where  $0.1 \leq |m| < 1$ , and  $e \in \mathbb{I}$ .

For examples,  $0.6247 \cdot 10^3$ ,  $0.1735 \cdot 10^{-13}$ ,  $-0.2000 \cdot 10^{-1}$ , etc.

In floating point system, the number of significant digits is kept fixed whereas the decimal point is floating.

### Significant Digits (Rules)

- ① Non-zero digits are always significant.
- ② Any zero between two significant digits are significant.
- ③ A final zero or trailing zeros in the decimal portion ONLY are significant.

Example:

<u>Number</u>	<u>No. of Significant Digits</u>
<u>406</u>	3
0.00 <u>500</u>	3
0.0 <u>3040</u>	4
<u>136000</u>	6
0.00 <u>1360</u>	4

## Round-Off

An error caused by Chopping or rounding is called rounding error or round-off error.

Let  $x = 0.d_1d_2, \dots, d_K d_{K+1} d_{K+2}, \dots \times 10^e$  be a real number. The floating point form of  $x$ , denoted by  $fl(x)$ , is obtained by terminating the mantissa of  $x$  at  $K$  decimal digits either by chopping or rounding.

→ The chopping produces

$$fl(x) = 0.d_1d_2 \dots d_K \times 10^e.$$

→ The rounding adds  $5 \times 10^{e-(K+1)}$  to  $x$  and then chops the result to obtain a number of the form

$$fl(x) = 0.d_1d_2d_3 \dots d_K \times 10^e.$$

→ For rounding, when  $d_{K+1} \geq 5$ , we add 1 to  $d_K$ , and discards the digits after  $K^{\text{th}}$  place to obtain  $fl(x)$ , i.e., we round up.

→ When  $d_{K+1} < 5$ , we simply chop off all but the first  $K$  digits, so we round down.

If we round down, then  $\delta_i = d_i$  for  $i = 1, 2, \dots, k$ .  
However, if we round up, the digits (and even the exponent) might change.

Example:  $\pi = 3.14159265\dots$

Normalized decimal form:  $\pi = 0.314159265\dots \times 10^1$

(a) Floating point form of  $\pi$  using 5-digit Chopping is

$$fl(x) = 0.31415 \times 10^1 = 3.1415$$

(b) Floating Point form of  $\pi$  using 5-digit rounding is

$$fl(x) = 0.31416 \times 10^1 = 3.1416$$