

Number Systems

Representation of Integers - We use decimal number system which consists of digits 0, 1, 2, ..., 9.

$$\text{e.g. } 1256 = 1 \times 10^3 + 2 \times 10^2 + 5 \times 10^1 + 6 \times 10^0$$

Hence 10 is the base of this system.

∴ 1256 is expressible as a polynomial in the base 10 with integral coefficients lies between 0 and 9.

2. Decimals

$$\bullet 1256 = 1 \times 10^3 + 2 \times 10^2 + 5 \times 10^1 + 6 \times 10^0$$

$$\text{Now } 42.965 = 4 \times 10^1 + 2 \times 10^0 + 9 \times 10^{-1} + 6 \times 10^{-2} + 5 \times 10^{-3}$$

Floating-Point Representation of Numbers

$$\frac{8}{3} = 2.666\ldots = 2 \times 10^0 + 6 \times 10^{-1} + 6 \times 10^{-2} + 6 \times 10^{-3} + \dots$$

Now Number Systems are of the following forms.

Name	Base	Digits
1. Decimal Number System	10	0, 1, 2, ..., 9
2. Binary System	2	0, 1
3. Octal	8	0, 1, 2, ..., 7
4. Hexadecimal	16	0, 1, 2, ..., 9 and A, B, C, D, E, F

We represent base by β .

$$\therefore \text{From above, } 42.965 = 4 \times \beta^1 + 2 \times \beta^0 + 9 \times \beta^{-1} + 6 \times \beta^{-2} + 5 \times \beta^{-3}$$

$$(111.011)_2 = 1 \times 2^2 + 1 \times 2^1 + 1 \times 2^0 + 0 \times 2^{-1} + 1 \times 2^{-2} + 1 \times 2^{-3}.$$

Examp.

$$(127)_8 = 1 \times 8^2 + 2 \times 8^1 + 7 \times 8^0$$

$$\text{Ans} \leftarrow \text{d}_1 \text{d}_2 \dots \text{d}_n \times \beta^e \quad \text{Now } 2.666\ldots = 0.2666\ldots \times 10^1 \\ = 0.02666\ldots \times 10^2 \\ = 0.0002666\ldots \times 10^4$$

Defn: A floating point number is a number represented in the form $\bullet d_1 d_2 \dots d_n \times \beta^e$ where d_1, d_2, \dots, d_n are integers and satisfy $0 \leq d_i < \beta$ and the exponent.

is such that $m \leq e \leq M$.

The fractional part $0.a_1a_2\ldots a_m$ is called mantissa, and it lies between +1 and -1.

Defn $\therefore fl(x) = \pm (0.a_1a_2\ldots a_m)\beta^e$

For each real number x , we associate a floating point representation denoted by $fl(x)$, given by

$$fl(x) = \pm (0.a_1a_2\ldots a_m)\beta^e$$

where β based fraction, i.e., $(0.a_1a_2\ldots a_m)_\beta$ is called mantissa with all a_i 's integers and e is known as exponent. This representation is called β -based floating point representation of x .

$$42.965 = 0.42965 \times 10^2$$

$$0 = 0.000\ldots 0 \times \beta^e$$

$$-0.00234 = -0.234 \times 10^{-2} \quad \leftarrow (\text{Not unique})$$

$$(\because = -0.0234 \times 10^{-1})$$

Defn: (Normal form) - A non-zero floating-point number is in normal form if the value of mantissa lies in $(-1, -\frac{1}{\beta}] \cup [\frac{1}{\beta}, 1)$

$$\text{e.g. } 0.4296 \in [0.1, 1)$$

$$-0.234 \in (-1, -0.1]$$

If $\beta = 10$ then $\frac{1}{\beta} = 0$.

On A floating-point number $(0.a_1a_2\ldots a_m)_\beta \times \beta^e$ is said to be normal if $a_1 \neq 0$

Also the range of exponent is also restricted. There are integers m and M such that $-m \leq e \leq M$.

Rounding and Chopping: Let x be any real number, and $fl(x)$ be its machine approximation. There are two ways to do cutting to store a real number.

$$x = \pm (0.a_1a_2\ldots a_m a_{m+1}a_{m+2}\ldots)_\beta \times \beta^e, a_1 \neq 0$$

Take $fl(x)$ as the approximation upto only m -digits.

(i.e. normal form)

(i) Chopping: $fl(x) = \pm(0.a_1 a_2 \dots a_m) \beta \times \beta^e$

(ii) Rounding:

$$fl(x) = \begin{cases} \pm(0.a_1 a_2 \dots a_m) \beta \times \beta^e, & 0 \leq a_{m+1} < \beta/2 \text{ (rounding down)} \\ \pm(0.a_1 a_2 \dots a_{m+1}) \beta \times \beta^e, & \frac{\beta}{2} \leq a_{m+1} < \beta \text{ (rounding up)} \end{cases}$$

\downarrow

$$\pm[(0.a_1 a_2 \dots a_m) \beta + (0.000\dots 01) \beta] \times \beta^e$$

Ex $x = 6/7 = 0.85714 \times 10^0$

Chopping and rounding off upto two decimal places,

$$fl(x) = \begin{cases} 0.85 \times 10^0 \text{ (chopping)} \\ 0.86 \times 10^0 \text{ (rounding)} \end{cases}$$

Rules for rounding off numbers

- (1) If digit to be dropped is greater than 5, the last retained digit is increased by one. For example, 12.6 is rounded to 13.
- (2) If the digit to be dropped is less than 5, the last retained digit is left as it is. For example, 12.4 is rounded to 12.
- (3) If the digit to be dropped is 5 and any digit following it is non-zero, the last remaining digit is increased by one. For example, 12.51 is rounded to 13.
- (4) If the digit to be dropped is 5 and is followed by only zeros, the last remaining digit is increased by one if it is odd, but left as it is if even. For example, 11.5 is rounded to 12 and 12.5 is rounded to 12.

Errors: Let x be any real number and $fl(x)$ is its approximation

(1) Absolute Error = $|x - fl(x)|$

(2) Relative Error = $\frac{|x - fl(x)|}{|x|}$

$x = 0.33, fl(x) = 0.3$
$AE = x - fl(x) = 0.03$
$RE = \frac{0.03}{0.33} = 0.0909 \approx 0.1$
$x = 0.33 \times 10^{-5}, fl(x) = 0.3 \times 10^{-5}$
$AE = 0.03 \times 10^{-5}, RE = 0.1 \times 10^{-5}$

As a measure of accuracy, the absolute error may be misleading and relative error is more meaningful.

Theorem (Error bound theorems in case of chopping and Rounding off)

Let x be any real number in decimal system. and $fl(x)$ be its approximation in computer. Then find the error bounds for Absolute and relative errors in case of (i) Chopping and (ii) Rounding off.

Proof (i) Chopping

$$\text{let } x = (a_0.a_1 a_2 \dots a_m a_{m+1} \dots) \times 10^e, a_1 \neq 0 \\ = \left(\sum_{i=1}^{\infty} \frac{a_i}{10^i} \right) \times 10^e, a_1 \neq 0$$

$$fl(x) = (a_0.a_1 a_2 \dots a_m) \times 10^e \\ = \left(\sum_{i=1}^m \frac{a_i}{10^i} \right) \times 10^e$$

\therefore Absolute Error is given by,

$$A.E = |x - fl(x)| = \left(\sum_{i=m+1}^{\infty} \frac{a_i}{10^i} \right) \times 10^e$$

Since each $a_i \leq 9$ in chopping $a_i < 10$ or $a_i \leq (10-1) = 9$

$$\therefore |x - fl(x)| \leq \left(\sum_{i=m+1}^{\infty} \frac{9(10-1)}{10^i} \right) \times 10^e \\ = \left(\sum_{i=m+1}^{\infty} \frac{9}{10^i} \right) \times 10^e = 9 \left(\sum_{i=m+1}^{\infty} \frac{1}{10^i} \right) \times 10^e \\ = 9 \times \cancel{9} \times \cancel{9} \times \left(\frac{1}{10^{m+1}} + \frac{1}{10^{m+2}} + \dots \right) \times 10^e \\ = 9 \times \frac{\frac{1}{10^{m+1}}}{1 - \frac{1}{10}} \times 10^e \quad (\text{Geometric Series with } A = \frac{1}{10^{m+1}} \text{ and } r = \frac{1}{10}, \text{ since } |r| = |\frac{1}{10}| < 1) \\ = 9 \times \frac{\frac{1}{10^{m+1}}}{\frac{9}{10}} \times 10^e = \frac{1}{10^m} \times 10^e \\ = 10^{e-m}$$

$$\text{Now } |x| = (0.a_1 a_2 \dots a_m a_{m+1} \dots)_{10} \times 10^e$$

$$\geq (0.1) \times 10^e = \frac{1}{10} \times 10^e = 10^{e-1}$$

$$\Rightarrow \frac{1}{|x|} \leq \frac{1}{10^{e-1}}$$

Relative error,

$$R.E = \frac{|x - fl(x)|}{|x|} \leq \frac{10^{e-n}}{10^{e-1}} = 10^{1-n}$$

(ii) Rounding off

$$fl(x) = \begin{cases} (0.a_1 a_2 \dots a_m)_{10} \times 10^e & a_{m+1} < 5 \\ (0.a_1 a_2 \dots [a_{m+1}])_{10} \times 10^e & a_{m+1} \geq 5 \end{cases}$$

Case I For $a_{m+1} < 5$ i.e. $a_{m+1} \leq 4$

$$fl(x) = (0.a_1 a_2 \dots a_m)_{10} \times 10^e = \sum_{i=1}^m \frac{a_i}{10^i} \times 10^e$$

$$\text{Now } |x - fl(x)| = \left(\sum_{i=n+1}^{\infty} \frac{a_i}{10^i} \right) \times 10^e = \left[\frac{a_{n+1}}{10^{n+1}} + \sum_{i=n+2}^{\infty} \frac{a_i}{10^i} \right] \times 10^e$$

$$\leq \left(\frac{4}{10^{n+1}} + \sum_{i=n+2}^{\infty} \frac{9}{10^i} \right) \times 10^e$$

$$= \left(\frac{4}{10^{n+1}} + 9 \left[\frac{1}{10^{n+2}} + \frac{1}{10^{n+3}} + \dots \right] \right) \times 10^e$$

$$= \frac{4}{10^{n+1}} + 9 \cdot \left(\frac{\frac{1}{10^{n+2}}}{1 - \frac{1}{10}} \right) \times 10^e$$

$$= \left(\frac{4}{10^{n+1}} + \frac{1}{10^{n+1}} \right) \times 10^e = \left(\frac{5}{10^{n+1}} \right) \times 10^e$$

$$= \frac{5}{10} \times 10^{e-n} = \frac{1}{2} \times 10^{e-n}$$

Case II For $a_{m+1} \geq 5$

$$fl(x) = (0.a_1 a_2 \dots [a_{m+1}])_{10} \times 10^e$$

$$= \left[(0.a_1 a_2 \dots a_m)_{10} + (0.0000 \dots 01)_{10} \right] \times 10^e$$

$$= \left(\sum_{i=1}^n \frac{a_i}{10^i} + \frac{1}{10^n} \right) \times 10^e$$

(2)

$$\text{Now } |x - f(x)| = \left| \left(\sum_{i=n+1}^{\infty} \frac{a_i}{10^i} - \frac{1}{10^n} \right) \times 10^e \right|$$

$$= \left| \frac{1}{10^n} - \sum_{i=n+1}^{\infty} \frac{a_i}{10^i} \right| \times 10^e$$

$$= \left| \frac{1}{10^n} - \frac{a_{n+1}}{10^{n+1}} - \sum_{i=n+2}^{\infty} \frac{a_i}{10^i} \right| \times 10^e$$

$$\leq \left| \frac{1}{10^n} - \frac{a_{n+1}}{10^{n+1}} \right| \times 10^e$$

$$\leq \left| \frac{1}{10^n} - \frac{5}{10^{n+1}} \right| \times 10^e$$

$$= \left| \frac{1}{10^n} - \frac{1}{2} \frac{1}{10^n} \right| \times 10^e$$

$$= \frac{1}{2} \cdot \frac{1}{10^n} \times 10^e = \frac{1}{2} \times 10^{e-n}$$

$$\begin{cases} a_{n+1} \geq 5 \\ -a_{n+1} \leq -5 \end{cases}$$

Therefore for both cases I & II,

$$|x - f(x)| \leq \frac{1}{2} \times 10^{e-n}$$

$$\text{Now } \frac{|x - f(x)|}{|x|} \leq \frac{1}{2} \frac{10^{e-n}}{10^{e-1}} = \frac{1}{2} 10^{1-n}$$

Significant figures: Rules for deciding the

Number of significant figures in measured quantity.

- (1) All nonzero digits are significant. e.g. 1.234 has 4 significant figures, 1.2 has 2 significant figures.
- (2) Zeros between nonzero digits are significant.
e.g. 1002 has 4 significant figures.
- (3) Leading zeros to the left of the first nonzero digit are not significant. e.g. 0.001 has only 1 significant figure.
- (4) Trailing zeros that are also to the right of a decimal point in a number are significant.
e.g.: 0.0230 has 3 significant figures. $\cdot 230 \times 10^{-3}$ or 23.0×10^{-4}
- (5) When a number ends in zeros that are not to the right of a decimal point, the zeros are not necessarily significant. e.g. 190 may be 2 or 3 significant figures,
 50600 may be 3 or 4 or 5 significant figures.
or 50600×10^3 or 0.5068×10^5 or 0.506×10^5

Exact numbers have an infinite number of significant figures. e.g. 23 students, 5 km, 500m

* Significant figures. All measurements are approximations. No measuring device can give perfect measurements without experimental uncertainty. The number of significant figures in a result/calculation is simply the number of figures that are known with some degree of reliability. (3) Here leading zeros only help to fix the position of the decimal point. Hence these are not significant.

Significant figures

Take x in normalized floating point number. Then

$$x = \underbrace{(0.a_1 a_2 \dots a_n)}_{\text{No. of digits in mantissa}} \times 10^e$$

No. of digits in mantissa of normal form are the no. of significant digits in that number.

$$209.3050 \\ = 0.\underline{2093050} \times 10^3 \\ 7 \text{ digits are significant.}$$

Example:

$$25.923$$

$$= 0.\underline{25923} \times 10^2$$

5 digits are significant.

From exact no.

$$256 \\ = 0.256 \times 10^3 \\ = 0.2560 \times 10^3$$

At least three significant digits.

$x = (0.a_1 a_2 \dots a_n a_{n+1} \dots) \times 10^e$

$x \rightarrow$ exact value, $\text{fl}(x) \rightarrow$ Approximation.
 $x^* \rightarrow$ Approximation.

$$\frac{|x - x^*|}{|x|} \leq \frac{1}{2} 10^{1-t} \text{ or } \frac{1}{2} 10 \times 10^{-t} = 5 \times 10^{-t}$$

The number x^* is said to approximate x to t -significant digits if t is the largest non-negative integer for which

$$\frac{|x - x^*|}{|x|} \leq 5 \times 10^{-t}$$

Example $x = 25.923, x^* = 25.92$

$$\frac{|x - x^*|}{|x|} = \frac{0.003}{25.923} = 0.000115727 \leq 0.0005$$

$$p = 0.54617, q = 0.54601 \\ \text{Use 4-digit Arithmetic} \\ \text{find } x = p - q \text{ and find} \\ \text{significant digit in approximation} \\ x = p - q = 0.00016 \text{ (exact value)} \\ \text{fl}^* = p^* - q^* = 0.5462 - 0.5460 \\ = 0.0002 \\ R.E. = \frac{|x - x^*|}{|x|} = \frac{0.0002}{0.0005} = 0.4 \leq 0.5 \\ = 5 \times 10^{-4} \quad \boxed{1-\text{significant digit.}}$$

$\Rightarrow x^* = 25.92$ is approximating x to 4-significant digits.

The number x^* is said to approximate x correctly upto n decimal places if

$$|x - x^*| \leq \frac{1}{2} 10^{-n}$$

e.g. $x = 25.9235 \quad x^* = 25.92$

$$|x - x^*| = 0.0035 \leq 0.005 = 5 \times 10^{-3} = \frac{5 \times 10^{-2}}{10} = \frac{1}{2} 10^{-2}$$

$\Rightarrow x^*$ is approximating x correctly upto 2 decimal places.

$$x = 25.9263, \quad x^* = 25.93$$

$$|x - x^*| = 0.0037 \leq 0.005 = \frac{1}{2} 10^{-2}$$

Arithmetic Operations

(Finite-digit Arithmetic)

① Addition: Let x and y be two numbers.

$$x \oplus y = \text{fl}(\text{fl}(x) + \text{fl}(y))$$

Do all operations

② Subtraction, $x \ominus y = \text{fl}(\text{fl}(x) - \text{fl}(y))$

③ Multiplication, $x \otimes y = \text{fl}(\text{fl}(x) \times \text{fl}(y))$

④ Division, $x \oslash y = \text{fl}(\text{fl}(x) \div \text{fl}(y))$

Use 4-digit rounding arithmetic to find

$$x+y, x-y, x \times y, x \div y$$

$$x = \frac{5}{7} = 0.7142857142857, \text{fl}(x) = 0.714\overline{28}$$

$$y = \frac{1}{3} = 0.333333333, \text{fl}(y) = 0.333\overline{33}$$

$$x+y = \text{fl}(0.714\overline{28} + 0.333\overline{33}) = \text{fl}(1.0476\overline{6}) \\ = \underline{1.0476} \quad 1.048$$

$$x-y = \text{fl}(0.714\overline{28} - 0.333\overline{33}) = \text{fl}(0.381\overline{6}) = \underline{0.381\overline{6}} \quad 0.3810$$

$$x \times y = \text{fl}(0.714\overline{28} \times 0.333\overline{33}) = \text{fl}(0.23807679\overline{9}) \\ = \underline{0.23807679} \quad 0.2381$$

$$x \div y = \text{fl}(0.714\overline{28} \div 0.333\overline{33}) = \text{fl}(2.1431143) \\ = 2.143$$