

## Errors And Their Analysis

Following are the broad sources of errors in numerical analysis

### (i) Input errors

The input information is rarely exact since it comes from the experiments and any experiment can give ~~exact~~ results of only limited accuracy. Moreover, the quantity used can be represented in a computer for only a limited number of digits.

### (ii) Algorithmic errors

If direct algorithms based on a finite ~~steps~~ sequence of operations are used, errors due to limited steps does not amplify the existing errors but if infinite algs are used, ideally exact results are expected only after an infinite number of steps.

(iii) ~~error~~. As this cannot be done in practice, the algorithm has to be stopped after a finite number of steps and as a consequence the result are not exact.

### (iii) Computational errors

Even when elementary operations, such as multiplication and division are used, the number of digits increases greatly so the results cannot be held fully in register available in a given computer. In such cases, a certain number of digits must be discarded. Furthermore, the errors here accumulates one after another from operation to operations, changing during the process and producing new errors.

Rounding off  
There are  
22

## Accuracy of Number

### 1. Approximate Numbers →

There are two types of numbers: Exact and Approximate

Exact numbers are  $2, 4, 9, \frac{1}{2}, 6.45, \dots$  etc.

But there are numbers such that

$$\frac{4}{3} = 1.333\dots, \sqrt{2} = 1.414213\dots, \pi = 3.141592\dots$$

Which can not be expressed by a finite number of digits. These may be approximated by numbers 1.333, 1.4141 and 3.1416 respectively. Such numbers which represent the given numbers to a certain degree of accuracy are called approximate numbers.

### 2. Significant digits.

The digits used to express a number are called significant digits.

The digits 1, 2, 3, 4, 5, 6, 7, 8, 9, are significant digits. '0' is also a significant digit except when it is used to fix the decimal point or to fill the places of unknown or discarded digits.

For example number 7845, 3.589, .4758 contains 4 significant figures, while the number .00386, .000587, .00000296 contain only three significant figures. Since zeros only help to fix the position of decimal point.

Similarly, in the number .0003090, the first four '0's are not significant digits since they serve only to fix the position of the decimal point and indicate the place values of the other digits.

The other two '0's are significant.

## Accuracy of Number

### 1. Approximate numbers →

There are two types of numbers: Exact and Approximate

Exact numbers are  $2, 4, 9, \frac{1}{2}, 6.45, \dots$  etc.

But there are numbers such that

$$\frac{4}{3} = 1.3333\dots, \sqrt{2} = 1.414213\dots, \pi = 3.141592\dots$$

which can not be expressed by a finite number of digits. There may be approximated by numbers  $1.333, 1.4141$  and  $3.1416$  respectively such numbers which represent the given numbers to a certain degree of accuracy are called approximate numbers.

### 2. Significant digits.

The digits used to express a number are called significant digits.

The digits  $1, 2, 3, 4, 5, 6, 7, 8, 9$ , are significant digits. '0' is also a significant digit except when it is used to fix the decimal point or to fill the places of unknown or discarded digits.

For example number  $7845, 3.589, .4758$  contains 4 significant figures, while the number  $.00396, .000587, .0000296$  contain only three significant figures. Since zeros only help to fix the position of decimal point.

Similarly, in the number  $.0003090$ , the first four '0's are not significant digits since they serve only to fix the position of the decimal point and indicate the place values of the other digits.

The other two '0's are significant.

Note:

1. The significant figure in a number in positional notation consists of
  - (i) All non-zero digits
  - (ii) Zero digits which
    - (a) Lie between significant digits
    - (b) Lie to the right of decimal point and at the same time, to the right of non-zero digits.
    - (c) Are specifically indicated to be significant.
2. The significant figure in a number written in scientific notation ( $\text{eg } \pm M \times 10^K$ ) consists of all the digits explicitly in  $M$ . Here  $\frac{1}{10} \leq M < 1$

Number	significant digits	No. of significant digits
3969	3, 9, 6, 9	0 4
3060	3, 0, 6	0 3
3900	3, 9	0 2
39.69	3, 9, 6, 9	0 4
•3969	3, 9, 6, 9	0 4
39.00	3, 9, 0, 0	0 4
•00039	3, 9	0 2
•00390	3, 9, 0	0 3
3.0069	3, 0, 0, 6, 9	0 5
$3.9 \times 10^6$	3, 9	0 2
$3.909 \times 10^5$	3, 9, 0, 9	0 4
$6 \times 10^{-2}$	6	0 1

### Rounding off

There are numbers with large number of digits e.g.  
 $\frac{22}{7} = 3.142857143$ . In practice it is desirable to limit such numbers to a manageable ~~number~~ number of digits such as 3.14 or 3.143. This process of dropping unwanted digits is called rounding-off.

Numbers are rounded-off according to following rule:

To round off a number to  $n$  significant digits, discard all digits to the right of  $n$ th digit and if this discarded number is

(i) Less than 5 in  $(n+1)$ th place, leave the  $n$ th digit unaltered.  
e.g. 7.893 to 7.89

(ii) Greater than 5 in  $(n+1)$ th place, increase the  $n$ th digit by unity e.g. 6.3456 to 6.346.

(iii) Exactly 5 in  $(n+1)$ th place, increase the  $n$ th digit by unity if is odd otherwise leave it unchanged

$$\text{e.g. } \begin{aligned} 12.675 &\approx 12.68 \\ 12.685 &\approx 12.69 \end{aligned}$$

Number	Rounded-off to		
	Three digits	Four digits	Five digits
.543241	.543	.5432	.54324
39.5255	39.5	39.52	39.526
69.4155	69.4	69.42	69.416
.667676	.668	.6677	.66768

### Errors

$$\text{Errors} = \text{True Value} - \text{Approximate Value}$$

In Any numerical ~~approximate~~ Computation, We come across following type of errors :

1. Inherent errors
2. Rounding errors
3. Truncation errors
4. Absolute errors
5. Relative errors
6. Percentage errors.

#### 1. Inherent Errors :

Errors which are already present in the statement of a problem before its solution are called inherent errors. Such errors arise either due to the given data being approximate or due to limitations of mathematical tables, calculators or the digital computer.

Inherent errors can be minimized by taking better data or using high precision.

Accuracy refers to the no. of significant digit in a value e.g. 53.965 is accurate to 5 significant digits. Precision Refers no. of decimal position or order of magnitude of the last digit in value. e.g. in 53.965, Precision is  $10^{-3}$ .

#### 2. Rounding Errors :

They arise from the process of rounding off the numbers during the computation. It is also called procedural error or numerical error. Such errors are unavoidable in most of the calculations due to limitations of computing aids.

These errors can be reduced however by

- (i) Changing the calculation procedure so as to avoid subtraction of nearly numbers or division by a small number
- (ii) Retaining atleast one more significant digit at each step and rounding off at last step.

Rounding off may be executed in two ways

- (a) Chopping: In it extra digits are dropped by truncation of numbers.
- (b) Symmetric round off: In it, the last retained significant digit is rounded up by unity if the first discarded digit is  $\geq 5$  otherwise the last retained digit is unchanged.

### (3) Truncation Errors

They are caused by using approximate results or on replacing an infinite process by a finite one.

For Example If  $s = \sum_{i=1}^{\infty} a_i x_i$  is replaced by or truncated to  $s = \sum_{i=1}^n a_i x_i$  the error developed is called Truncation error.

Truncation error is a type of algorithm error. Also

$$\text{If } e^x = 1 + x + \frac{x^2}{2} + \frac{x^3}{3} + \frac{x^4}{4} + \dots + \infty = x \text{ is truncated to}$$

$$1 + x + \frac{x^2}{2} + \frac{x^3}{3} = x' \text{ (say)}, \text{ then Truncation error} = x - x'$$

### (4) Absolute Error

Absolute Error is the numerical difference between the true value of a quantity and its approximate value.

Thus if  $x$  is the true value of a quantity and  $x'$  is the approximate value then  $|x - x'|$  is called absolute error  $e_a$ .

$$e_a = |x - x'| = |\text{Error}|$$

### (5) Relative Error :

Relative error is denoted by  $e_r$  and defined as

$$e_r = \left| \frac{x - x'}{x} \right|$$

Where  $x$  is true value and  $x - x'$  is error.

### (6) Percentage Error :

Percentage error is denoted by  $e_p$  and defined as

$$e_p = \left| \frac{x - x'}{x} \right| \times 100$$

Note:

1. If a number is ~~not~~ correct to  $n$  decimal places, then

$$\text{Error} = \frac{1}{2} (10^{-n})$$

$$\text{Error} = \frac{1}{2} (10^{-n})$$

For Example If the number 3.1416 is correct to 4 decimal

$$\text{places, then Error} = \frac{1}{2} (10^{-4}) = 0.0005$$

2. If the first significant digit of a number is  $k$  and the number is correct to  $n$  significant digits then,

$$\text{Relative error} < \frac{1}{(k \times 10^{n-k})}$$

Note: 2

To estimate the errors which creep in when the numbers in a calculation are truncated or rounded off to certain number of digits, the following rules are useful

If the approximate value of a number  $x$  having  $n$  decimal digits is  $x'$

Then Absolute error due to truncation to  $n$  digits =  $|x - x'| < 10^{n-k}$

2. Absolute error due to rounding off to  $n$  digit =  $|x - x'| < \frac{1}{2} 10^{n-k}$

3. Relative error due to truncation to  $n$  digits =  $\left| \frac{x - x'}{x} \right| < 10^{1-k}$

It . 333 is approximate  
relative and a

(4) Relative error due to Rounding off to  $k$  digits.

$$= \left| \frac{x-x'}{x} \right| < \frac{1}{2} 10^{1-k}$$

Note:

① If a number is correct to  $n$  significant digits,  
then maximum relative error  $\leq \frac{1}{2} (10^{-n})$

If a number is correct to  $d$  decimal places, then absolute  
error  $\leq \frac{1}{2} 10^{-d}$

② If the first significant figure of a number is  $k$  and the number is  
correct to  $n$  significant figures, then relative error  $< \frac{1}{k(10^{n-1})}$

① Suppose  $1.414$  is used as an approximation to  $\sqrt{2}$ . Find the  
absolute and relative errors.

Sol<sup>3</sup> True Value =  $\sqrt{2} = 1.41421356$

Approximate Value =  $1.414$

$$\begin{aligned} \text{Error} &= \text{True Value} - \text{Approximate value} \\ &= 1.41421356 - 1.414 \\ &= 0.00021356 \end{aligned}$$

$$\therefore \text{Absolute Error} = e_a = |\text{Error}| = |0.00021356| \\ = 2.1356 \times 10^{-5}$$

$$\text{Relative Error} = e_r = \frac{e_a}{\text{True Value}}$$

$$= \frac{2.1356 \times 10^{-5}}{1.41421356} = 1.51 \times 10^{-5}$$

- ② If .333 is approximate value of  $\frac{1}{3}$ , find absolute, relative and percentage errors.

Sol: Given True Value  $x = \frac{1}{3}$

Approximate value  $x' = .333$

$$\text{Absolute error } e_a = |x - x'|$$

$$= \left| \frac{1}{3} - .333 \right| =$$

$$= \left| .33333 - .333 \right| = .000333$$

$$\text{Relative error } e_r = \frac{e_a}{x} = \frac{.000333}{\underline{\underline{\frac{1}{3}}}}$$

$$\text{Absolute error } e_n = \frac{e_a}{x} = \frac{.000333}{\text{True value}}$$

$$= \frac{.000333}{\frac{1}{3}} = .000999$$

$$\text{Percentage Error } e_p = e_p \times 100$$

$$= \text{Relative error} \times 10^0$$

$$= .000999 \times 10^0$$

$$= .099\%$$

- ③ An approximate value of  $\pi$  is given by 3.1428571 and its true value is 3.1415926. Find absolute and relative errors.

$$\text{True value} = 3.1415926$$

$$\text{Approximate value} = 3.1428571$$

$$\text{Error} = \text{True value} - \text{Approximate value}$$

$$= 3.1415926 - 3.1428571$$

$$= -.0012645$$

$$\text{Absolute error} = |\text{Error}| = |-0.0012645| = .0012645$$

$$\text{Relative Error} = \frac{\text{Absolute Error}}{\text{True value}} = \frac{.0012645}{3.1415926} = .000402502 \approx$$

Q.4. Three approximate values of number  $\frac{1}{3}$  are given as .30, .33 and .34 which of these three is the best approximation?

Sol<sup>b</sup>) The best Approximation will be the one which has least absolute error.

$$\text{True value} = \frac{1}{3} = .3333$$

Case-I Approximate value = .30

$$\begin{aligned}\text{Absolute Error} &= |\text{True Value} - \text{Approximate Value}| \\ &= |.3333 - .30| = .0333\end{aligned}$$

Case-II Approximate value = .33

$$\begin{aligned}\text{Absolute error} &= |\text{True value} - \text{Approximate value}| \\ &= |.3333 - .33| = .0033\end{aligned}$$

Case-III Approximate value = .34

$$\begin{aligned}\text{Absolute error} &= |\text{True value} - \text{Approximate value}| \\ &= |.3333 - .34| = .00667\end{aligned}$$

Since absolute error is least in case II. Hence .33 is the best approximation.

Q. Find the relative error if  $\frac{2}{3}$  is approximated to .667

Sol<sup>b</sup>) True value =  $\frac{2}{3} = .66666$

Approximate value = .667

$$\begin{aligned}\therefore \text{Absolute error } e_a &= |\text{True value} - \text{Approximate value}| \\ &= |.66666 - .667| = .000334\end{aligned}$$

$$\text{Relative Error } e_r = \frac{.000334}{.66666} = .0005.$$

Q.1 Suppose 1.414 is used as an approximation to  $\sqrt{2}$ . Find the absolute and relative errors.

$$\text{True value} = \sqrt{2} = 1.41421356$$

$$\text{Approximate value} = 1.414$$

$$\text{Error} = \text{True Value} - \text{Approximate Value}$$

$$= \sqrt{2} - 1.414 \\ = 1.41421356 - 1.414 = .00021356$$

$$\text{Absolute error} = |\text{Error}| \\ = |.00021356| = 21356 \times 10^{-3}$$

$$\therefore \text{Relative error} = \frac{e_a}{\text{True value}} = \frac{21356 \times 10^{-3}}{\sqrt{2}} = 151 \times 10^{-3}$$

Q.2 If 0.333 is the approximate value of  $\frac{1}{3}$ , find absolute, relative and percentage errors.

$$\text{True value } x = \frac{1}{3}$$

$$\text{Approximate value} = x' = .333$$

$$\therefore \text{Absolute error} \cdot e_a = |x - x'| = \left| \frac{1}{3} - .333 \right| = .000333$$

$$\text{Relative Error} \cdot e_r = \frac{e_a}{x} = \frac{.000333}{.333} = .000999$$

$$\text{Percentage error} e_p = e_r \times 100 = .000999 \times 100 = .099\%$$

Q3. An Approximate value of  $\pi$  is given by 3.1428571 and its true value is 3.1415926. Find Absolute and relative errors.

$$\text{Given True value} = 3.1415926$$

$$\text{Approximate value} = 3.1428571$$

$$\therefore \text{Error} = \text{True value} - \text{Approximation value} \\ = 3.1415926 - 3.1428571 = -.0012645$$

$$\therefore \text{Absolute Error} = e_a = |\text{Error}| = .0012645$$

$$\text{Relative Error} = \frac{e_a}{\text{True value}} = \frac{.0012645}{3.1415926} = .000402502$$

a. Round off  
significant  
figures

Find the

Q. Round off the number 865250 and 37.46235 to four significant figures and compute  $e_a$ ,  $e_r$ ,  $e_p$  in each case.

(i) Number rounded off to four significant digits = 865200

$$\therefore x = 865250$$

$$x' = 865200$$

$$\therefore \text{Error} = x - x' = 865250 - 865200 = 50$$

$$\therefore \text{Absolute error} = e_a = |\text{error}| = 50$$

$$\text{Relative error} = e_r = \frac{e_a}{x} = \frac{50}{865250} = 5.77 \times 10^{-5}$$

$$\text{Percentage error} = e_p = e_r \times 100 = 5.77 \times 10^{-5} \times 10^2 = 5.77 \times 10^{-3}$$

(ii) Number rounded off to four significant digits = 37.46

$$\text{Given True Value} = 37.46235 = x$$

$$\text{Approximate Value} = x' = 37.46$$

$$\text{Error} = \text{True value} - \text{Approximate value}$$

$$= 37.46235 - 37.46 = .00235$$

$$\therefore \text{Absolute Error} = e_a = |\text{Error}|$$

$$e_a = .00235$$

$$\text{Relative error} = \frac{\text{Absolute Error}}{\text{True value}} = \frac{.00235}{37.46235} = 6.2729 \times 10^{-5}$$

$$\text{Percentage Error} = \text{Relative Error} \times 100$$

$$= 6.2729 \times 10^{-5} \times 10^2 = 6.2729 \times 10^{-3}$$

Q. Round off the number 75462 to four significant digits and then calculate the absolute error and percentage error.

Sol) Number Rounded off to Four significant digits = 75460

$$\text{Absolute error} = |\text{True Value} - \text{App. Value}|$$

$$= |75462 - 75460| = 2$$

$$\text{Relative error} = e_r = \frac{\text{Absolute error}}{\text{True value}}$$

$$e_r = \frac{2}{75462} = .0000265$$

$$\text{Percentage error} = e_p = e_r \times 100$$

$$= .0000265 \times 100$$

$$= 0.0265 \text{ Ans}$$

Q. Find the absolute, relative and Percentage errors if  $x$  is rounded off to three decimal digits. Given  $x = 1.005998$ .

Sol: The number rounded off to three decimal ~~places~~ digits  
= 1.006

$$\text{Error} = 1.005998 - 1.006 = -.000002$$

$$\text{Absolute error} \cdot e_a = |\text{Error}| = .000002$$

$$\text{Relative Error} = \frac{\text{Absolute error}}{\text{True value}} = \frac{.000002}{1.005998} = .0033344$$

$$\text{Percentage Error} = e_p = e_r \times 100$$

$$= .0033344 \times 100$$

$$= 33.344 \text{ Ans}$$

Q. Evaluate the sum  $s = \sqrt{3} + \sqrt{5} + \sqrt{7}$  to 4 significant digits and find its absolute and relative errors.

Sol:  ~~$\sqrt{3} = 1.732$~~ ,  ~~$\sqrt{5} = 2.236$~~ ,  ~~$\sqrt{7} = 2.6456$~~

$$s = 1.732 + 2.236 + 2.646 = 6.614$$

$$\therefore \sqrt{3} = 1.732050808 \text{ (True value)}$$

$$\sqrt{3} \text{ to 4 significant digit} = 1.732 \text{ (Approximate value)}$$

$$\therefore \text{error} = |1.732050 - 1.732| = 0.0005$$

Similarly  $\sqrt{5} = 2.235067977$   
 $\sqrt{5} \text{ to 4 significant digit} = 2.235$

$$\text{error} = 0.0005$$

Similarly error in  $\sqrt{3}$  after 4 significant digits

$$\therefore \text{error} = 0.0005$$

$$\therefore \text{Total error} = 0.0005 + 0.0005 + 0.0005 \\ = 0.0015$$

The total absolute errors shows that the sum is correct to 3  
significant figure only

$$\text{Now } S = \sqrt{3} + \sqrt{5} + \sqrt{7} \\ = 6.61$$

$$e_n = \frac{0.0015}{6.61} = 0.0002 \text{ Ans}$$

Q1. Find the absolute error if the number  $x = 100545828$  is

- (i) Truncated to three decimal ~~places~~ digits  
(ii) Rounded off to three decimal digits.

Subs

$$x = 100545828 \\ = 1.545828 \times 10^8$$

(i) After Truncation to three decimal places its approximate value

$$x' = 1.545 \times 10^8$$

$$\therefore \text{Absolute error} = |x - x'|$$

$$= |1.545828 \times 10^8 - 1.545 \times 10^8| \\ = 0.000828 \times 10^8$$

$$\therefore \text{Absolute error due to truncation to } k \text{ digit} = |x - x'| < 10^{n-k}$$

$$\therefore 0.000828 \times 10^8 < 10^{-2-3}$$

$$0.2 \times 10^{-5} < 10^{-5}$$

This proves Rule-I

(ii) After rounding off to three decimal places, its approximate value  $x' = 1.546 \times 10^8$

$$\therefore \text{Absolute error} = |x - x'|$$

$$= |1.545828 - 1.546| \times 10^8 \\ = 0.000172 \times 10^8 = 0.172 \times 10^{-5}$$

Which is  $< 0.5 \times 10^{-2-3}$ . This proves Rule-II

### Errors in Numerical Computations

① If  $u = \frac{4\pi^2 y^3}{3^4}$  and error in  $x, y, z$  be 0.001, compute the relative max. error in  $u$  when  $x=y=z=1$

Sol

$$u = \frac{4\pi^2 y^3}{3^4}$$

$$\delta u = \frac{\partial u}{\partial x} \delta x + \frac{\partial u}{\partial y} \delta y + \frac{\partial u}{\partial z} \delta z$$

$$= \frac{8\pi y^3}{3^4} \delta x + \frac{12\pi^2 y^2}{3^4} \delta y + \frac{16\pi^2 y^3}{3^5} \delta z$$

$$(\delta u)_{\text{max}} = \left| \frac{8\pi y^3}{3^4} \delta x \right| + \left| \frac{12\pi^2 y^2}{3^4} \delta y \right| + \left| \frac{16\pi^2 y^3}{3^5} \delta z \right|$$

$$= \frac{8(1)(1)}{1} \times 0.001 + \frac{12(1)(1)}{1} (0.001) + \frac{16(1)(1)}{1} (0.001)$$

$$= 0.008 + 0.012 + 0.016$$

$$\delta u = 0.036$$

$$\therefore \text{Max Relative error} = \frac{\delta u}{u}$$

$$= \frac{0.036}{4}$$

$$u = \frac{4\pi^2 y^3}{3^4}$$

$$n=3, z=2=1$$

$$= 4$$

$$= 0.009 \text{ Ans}$$

② Compute the Percentage error in the time period  $T = 2\pi\sqrt{\frac{l}{g}}$  for  $l=1 \text{ m}$  if the error in the ~~length~~ measurement of  $l$  is ~~±~~ ± 0.1. Here  $g$  is constant

$$T = 2\pi\sqrt{l/g}$$

$$\log T = \log 2\pi + \frac{1}{2} [\log l - \log g]$$

$$\frac{1}{T} \delta T = 0 + \frac{1}{2} \cdot \frac{1}{l} \delta l - \cancel{\frac{1}{2} \cdot \frac{1}{g} \delta g} 0$$

$$= \frac{1}{2} \cdot \frac{1}{1} (0.1) - = \frac{0.1}{2} = 0.005$$

$$\therefore \delta T \times 100 = 0.005 \times 100 = 1.5\%, \text{ Ans}$$

Q. If  $u = 2v^6 - 5v$ , find the percentage error in  $u$ ,  
 $v=1$  if error in  $v$  is  $0.05$ .

$$u = 2v^6 - 5v$$

diff. it

$$\delta u = 12v \delta v - 5 \delta v$$

$$\delta u = 12v \delta v - 5 \delta v$$

$$\begin{aligned}\delta u &= 12 \times 1(0.05) - 5(0.05) \\ &= (12-5)0.05 = 0.35\end{aligned}$$

$$\therefore \frac{\delta u}{u} = \frac{0.35}{2v^6 - 5v} = \frac{0.35}{2-5} = -0.1166$$

$$\frac{\delta u}{u} \times 100 = -0.1166 \times 100 = -11.66\%.$$

Percentage error =  $-11.66\%$

Q. If  $r = 3h(h^6 - 2)$ , find the percentage error in  $r$  at  $h=1$ ,  
if the percentage error in  $h$  is  $5\%$ .

$$r = 3h(h^6 - 2)$$

$$r = 3h^7 - 6h$$

$$\delta r = 21h^6 \delta h - 6 \delta h$$

$$\left| \begin{array}{l} \delta r = 21(1)(5) - 6(5) = 105 \\ \frac{\delta r}{r} = \frac{75}{r} = \frac{75}{3(1^6 - 2)} = \frac{75}{3(-1)} = -25 \end{array} \right.$$

$$\delta r = (21h^6 - 6) \delta h$$

$$\frac{\delta r}{r} = \frac{(21h^6 - 6)}{3h(h^6 - 2)} \cdot \frac{\delta h}{h}$$

$$\frac{\delta r \times 100}{r} = \frac{21h^6 - 6}{3(h^6 - 2)} \cdot \frac{15}{5} \times 100$$

$$= \frac{21-6}{3(-1)} (5) = \frac{15 \times 5}{-3} = -25\%.$$

Given Percentage error in  $h$   
i.e.  $\frac{\delta h}{h} \times 100 = 5$

Q. The discharge  $Q$  over a notch for head  $H$  is calculated by the formula  $Q = k H^{5/2}$ , where  $k$  is a given constant. If the head is 75 cm and an error of 15 cm is possible in its measurement, estimate the percentage error in computing the discharge.

Sol's

$$Q = k H^{5/2}, \text{ Here } k \text{ is constant}$$

$$\log Q = \frac{5}{2} [\log k + \log H]$$

$$\frac{1}{Q} \delta Q = 0 + \frac{5}{2} \cdot \frac{1}{H} \delta H$$

$$\frac{\delta Q}{Q} \times 100 = \frac{5}{2} \frac{1}{H} \delta H \times 100$$

$$= \frac{5}{2} \cdot \frac{1}{75} (0.15) \times 100 = \frac{0.5}{75} \times 100 = \frac{5}{75} \%$$

## Floating Point Representation of Number

There are two type of arithmetic operations available in Computer

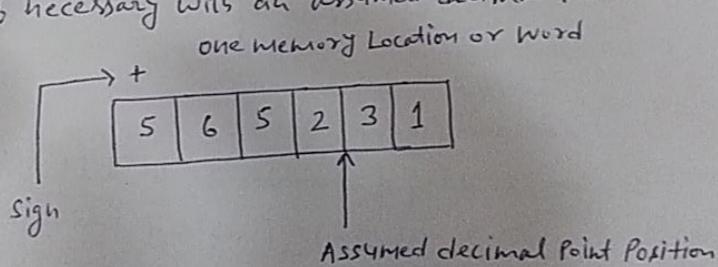
### (i) Integer Arithmetic

Integer arithmetic deals with integer operands and used mainly in counting and as subscripts.

### (ii) Real or ~~Frac~~ Floating Point Arithmetic →

Real arithmetic uses numbers with fractional parts as operands and is used in most computations. ~~Computers~~  
Computers are usually designed such that each location, called Word, in memory stores only a finite number of digits.

Let us assume a hypothetical computer having memory in which each location can store 6 digits and having provision to store one or more signs. One method of representing real numbers in that computer would be to assume a fixed position for the decimal point and store all numbers after appropriate shifting if necessary with an assumed decimal point.



(A memory location storing number 5652.31)

In such ~~convention~~ convention, the max. and minimum possible numbers to be stored are 9999.99 and 0000.01 respectively in magnitude.

For this a new convention is adopted which aims to preserve the maximum number of significant digits in a real number and also increase the range of values of real numbers stored. This representation is called the normalized floating point mode representation.

In this mode, a real number is expressed as combination of a mantissa and an exponent. The mantissa is made less than 1 or  $\geq 1$  and the exponent is the power of 10 which multiplies the mantissa.

for example : The number  $43.76 \times 10^6$  is represented in this notation as  $.4376 E 8$ , where E8 is used to represent  $10^8$ . The mantissa is .4376 and Exponent is 8.

The number is stored in memory location as :

Sign of mantissa							Sign of exponent
+	4	3	7	6	0	8	+
↑	mantissa						

Implied decimal pt.

Moreover, the shifting of mantissa to the left till its most significant digit is non-zero is called normalization.

e.g : The number 1006031 may be stored as ~~.1006031~~  
 $.6031 E -2$ .

Q1 Add the following floating point number

(i) ~~.5433E5~~, .4546 E5 and .5433 E5

Here Exponents are equal

∴ Mantissa are added

$$\therefore \text{sum} = .9979 \text{ E}5 \text{ Ans}$$

(ii) .4546 E5 and .5433 E7

Here Exponents are not equal.

The operand with the larger exponent is kept as it is

$$\begin{array}{r} .5433 \text{ E}7 \\ .0045 \text{ E}7 \\ \hline .5478 \text{ E}7 \end{array}$$

(iii) .4546 E3 and .5433 E7

$$\begin{array}{r} .5433 \text{ E}7 \\ .0000 \text{ E}7 \\ \hline .5433 \text{ E}7 \end{array}$$

(iv) .6434 E3 and .4845 E3

$$\begin{aligned} .6434 \text{ E}3 + .4845 \text{ E}3 &= 1.1279 \text{ E}3 \\ &= \underline{\underline{1.127}} \text{ E}4 \\ &= .1127 \text{ E}4 \text{ Ans} \end{aligned}$$

(v) .6434 E99 and .4845 E99

$$\text{Here } .6434 + .4845 = 1.1279 \text{ E}99 = \underline{\underline{.1127}} \text{ E}100$$

Here Exponent part cannot store more than two digits.

This condition is called an overflow condition and the arithmetic unit will intimate an error condition.

(ii)

Q2. Subtract the following floating point numbers:

(i)  $0.9432 \times 10^{-4}$  from  $0.5452 \times 10^{-3}$

$$\begin{array}{r} 0.5452 \times 10^{-3} \\ - 0.9432 \times 10^{-4} \\ \hline 0.4509 \times 10^{-3} \end{array}$$

(ii)  $0.5424 \times 10^{-3}$  from  $0.5452 \times 10^{-3}$

$$\begin{array}{r} 0.5452 \times 10^{-3} \\ - 0.5424 \times 10^{-3} \\ \hline 0.028 \times 10^{-3} \end{array}$$

In normalized floating point, the mantissa is  $\geq 1$ . So result is

(iii)  $0.5424 \times 10^{-99}$  from  $0.5452 \times 10^{-99}$   $= 0.28 \times 10^{-100}$

$$\begin{array}{r} 0.5452 \times 10^{-99} \\ - 0.5424 \times 10^{-99} \\ \hline 0.0028 \times 10^{-99} \text{ or } 0.28 \times 10^{-100} \end{array}$$

But exponent becomes -100

But exponent part cannot store more than two digits

Such condition is called underflow.

Q.3

(i)  $0.4546 \times 10^3 + 0.5454 \times 10^8$

Sol:

$$\begin{array}{r} 0.4546 \times 10^3 \\ + 0.5454 \times 10^8 \\ \hline 0.5454 \times 10^8 \text{ Ans} \end{array}$$

$$(ii) \quad .9432 E-4 - .6363 E-5$$

$$\begin{array}{r} \text{Sub} \\ \hline = \\ \begin{array}{r} .9432 E-4 \\ - .6363 E-5 \\ \hline .3069 E-4 \end{array} \end{array}$$

Multiplication : —

Two numbers are multiplied in the normalized floating point mode by multiplying the mantissas and adding the exponents. After the multiplication of the mantissas, the result mantissa is normalized as in addition or subtraction operation and the exponent appropriately adjusted.

Q.1 Multiply the following point number

$$(i) \quad .5543 E12 \text{ and } .4111 E-15$$

$$\begin{array}{r} \text{Sol} \\ \hline = \\ .5543 E+12 * .4111 E-15 = .2278 E-3 \text{ Ans} \end{array}$$

$$(ii) \quad .1111 E10 \text{ And } .1234 E15$$

$$\begin{array}{r} \text{Sol} \\ \hline = \\ .1111 E10 * .1234 E15 = .013709 E25 \\ = .1370 E24 \text{ Ans} \end{array}$$

$$(iii) \quad .1111 E51 \text{ and } .4444 E50$$

$$\begin{array}{r} \text{Sol} \\ \hline = \\ .1111 E51 * .4444 E50 = .04937284 E101 \\ = .49372 E100 \end{array}$$

overflow

$$(iv) \quad .1234 E-49 \text{ And } .1111 E-54$$

$$\begin{array}{r} \text{Ans} \\ \hline = \\ .1234 E-49 * .1111 E-54 = .01370974 E-103 \\ = .1370 E-1024 \end{array}$$

The result Underflow

### Division

In division, the mantissa of the numerator is divided by that of the denominator. The denominator exponent is subtracted from the numerator exponent. The quotient mantissa is normalized to make the most significant digit non-zero and the exponent appropriately adjusted. The mantissa of the result is chopped down to occupy 4 digits.

$$Q1 \text{ (i) } .9998 E1 \div .1000 E-99 \quad ??$$

$$\text{Sol}^{\text{s}} \quad \frac{.9998 E1}{.1000 E-99} = 9998 E100$$

$$\text{(ii) } .1000 E5 \div .9999 E3$$

$$\text{Sol}^{\text{s}} \quad \frac{.1000 E5}{.9999 E3} = .1000 E2 \quad \text{Ans}$$

$$\text{(iii) } .9998 E-5 \div .1000 E98$$

$$\text{Sol} \quad \frac{.9998 E-5}{.1000 E98} = .9998 E-104$$

Hence result is underflow Ans

$$Q2 \text{ For } x = .4845 \text{ and } y = .4800, \text{ calculate the value of } \frac{x^2 - y^2}{x+y} \text{ using normalized floating point arithmetic.}$$

Compare with the value of  $x-y$ .

$$\text{Sol}^{\text{s}}$$

$$x^2 = .4845 E0 \times .4845 E0 = .2347 E0$$

$$y^2 = .4800 E0 \times .4800 E0 = .2304 E0$$

$$\therefore x^2 - y^2 = .2347 E0 - .2304 E0 = .0043 E0 \approx .0043$$

$$x+y = .4845 E0 + .4800 E0 = .9645 E0$$

Now  
x-y is divided by 1000  
x-y is 3600.

$$\therefore \frac{x^2 - y^2}{x+y} = \frac{.0043 E_0}{.9645 E_0} = .004458 E_0 \\ = .4458 E^{-2}$$

Now  
 $x-y = .4845 E_0 - .4800 E_0$   
 $= .0045 E_0 = .4500 E^{-2}$

$$\text{Relative Error.} = \frac{.4500 - .4458}{.4500} = .0093 \text{ or} \\ = .93\%.$$

(6)

## Loss of significance

Q Let  $p = .54617$  and  $q = .54601$ . Use 4 digit arithmetic to approximate  $p-q$  and determine the absolute and relative error using rounding.

Sol

$$p = .54617, q = .5460$$

After 4 digit rounding

$$\bar{p} = .5462, \bar{q} = .5460$$

$$\text{Error} = \frac{(p-q) - (\bar{p}-\bar{q})}{p-q} = \frac{(.54617-.54601) - (.5462-.5460)}{.54617-.5460}$$

$$ex = \left| \frac{.00016 - .0002}{.00016} \right| = .25$$

$$ea = .00016 - .0002 = .00004$$

② Compute  $\sqrt{n+1} - 1$ , when  $n = .12345678 \times 10^{-5}$  use 8 digit rounding. Rewrite the expression to avoid the subtraction.

$$\underline{\text{Sol}}' f(n) = \sqrt{n+1} - 1$$

$$\text{Here } \sqrt{n+1} = \sqrt{1 + .12345678 \times 10^{-5}} = \sqrt{1.00001234} \\ = 1.0000062$$

(Rounding 8 significant digits)

$$\therefore \sqrt{n+1} - 1 = 1.0000062 - 1 = .0000062$$

Here only 2 significant digits  
i.e. Loss of significant digit

so to fix this problem rationalize the expression

$$\frac{\sqrt{1+x}-1}{x} \times \frac{\sqrt{1+x+1}}{\sqrt{1+x+1}} = \frac{1}{\sqrt{1+x+1}} = 6.17281915 \times 10^{-7}$$

i.e. 8 significant digit

minimum height =  $\sqrt{P} - \sqrt{P-9}$  m  
 maximum height =  $\sqrt{P+9} - \sqrt{P-9}$  m  
 average height =  $\frac{1}{2}(\sqrt{P+9} + \sqrt{P-9})$  m

$$\text{minimum height} = \sqrt{P} - \sqrt{P-9} = 9$$

(approximate value)

$$\text{maximum height} = \sqrt{P+9} - \sqrt{P-9} = 9$$

$$\text{average height} = \frac{(\sqrt{P+9} + \sqrt{P-9})}{2} = 9.147$$

$$\text{error} = \left| \frac{9.147 - 9}{9} \right| = 0.052$$

$$\text{percentage error} = \frac{0.052}{9} \times 100\% = 0.578\%$$

$\sqrt{P+9} - \sqrt{P-9} = \sqrt{P+9} - \sqrt{P} + \sqrt{P} - \sqrt{P-9}$  (approx.)

$$\sqrt{P+9} - \sqrt{P} + \sqrt{P} - \sqrt{P-9} = \sqrt{P+9} - \sqrt{P-9}$$