

## Understanding Errors in Computing

Storage space in a computer or any device is finite. They have the following effects :-

- a new number line called
  - floating point number line.
- a floating point arithmetic which features
  - largest and smallest number
  - round off mechanism
  - machine precision.
- need to take care (during computations)
  - loss of significance in +
  - close enough versus =
  - cancellation occurring in -

- introduction of Truncation Error when transferring Symbolic Analysis for Computation.
- need to keep track of number of Computations; i.e Order of Computations

### Computer storage :-

Bit - a short form for binary digit.  
smallest unit of storage. These  
can be thought of as on/off switches  
to store information.

Byte - 8 bits make one byte.  
original intent was that 8 bits  
was enough to assign a unique  
storage value for language characters.

26 - Capital A-Z	}	. Standardise
26 lower case a-z		American
10 0-9		Standard
32 Punctuation		Code
...		for

American

Standard

Code  
for

Information

Interchange

Around 127 were standardised.

Word - "A machine is 32-bit or 64-bit"  
 # of bits a particular Computer's CPU  
 Can deal with in one go.

Storing and Working with numbers.

Integers : are stored in R using  
 32-bit integer format.

$\overline{\overline{\overline{\overline{\quad}}}}$   
 ↑  
 one bit for sign       $\overbrace{\quad \quad \quad \quad}$   
 31 bits for number in basis binary

$\Rightarrow$  The largest positive integer =  $2^3 - 1$   
that can be stored (integer storage)

E.g. largest 3-bit integer

$$\begin{array}{c} 111 \\ \swarrow \quad \searrow \\ \text{most significant bit} \quad \text{least significant bit} \end{array} \quad (\Rightarrow) \quad 2^2 + 2^1 + 2^0 = 2^3 - 1$$

This the number of bits limit the size of integer that can be represented.

- There are some advantages:-
  - Convenient format of entering binary #.
  - Hexadecimal code [Colours]

Floating point numbers (default R uses to store numbers)

I - notation for

E lexical

E lectronic

E ngineer

double precision storage.

Any number say 12.3456 is  
written as

$$\text{Sign} \underset{\text{mantissa}}{\sim} + 0.123456 \times 10^2 \quad \text{Exponent}$$

Sign - 1 bit

mantissa - 52 bits

Exponent -  $\frac{11}{64}$  bits [ 1 bit for sign ]

Again finite storage space will result  
in range - largest and smallest [ Exponent ]

Precision - limits on # of significant  
digits [ mantissa ]

Range :-

largest number  $\sim 1.797693 \times 10^{308}$

Smallest number  $\sim 2.225074 \times 10^{-308}$   
[ Positive ]

$$\begin{array}{r}
 0.1 \\
 0.25 \\
 0.0625 \\
 \hline
 0.8125
 \end{array}
 = 1 \cdot 2^0 + 1 \cdot 2^1 + 0 \cdot 2^{-2} + 1 \cdot 2^{-4}$$

Precision :  $0.8125 \Leftrightarrow 1101$

[Exact] — no loss in storage  
 mantissa binary

$0.1 \Leftrightarrow$  loose some precision.

- cannot be represented by a finite number of binary digits.
- we will need to approximate such numbers by a round off mechanism
- most real numbers cannot be stored exactly.

### Machine precision :- ( $\epsilon$ )

- captures the spacing of the floating point line.

machine  $\epsilon_m$  :- is the smallest  $\epsilon > 0$  such that  $1 + \epsilon > 1$

- this is not as small as the "smallest number"

$$\epsilon_m = 2.220446 \times 10^{-12}$$

This is a real limit. Also thus

Changes with  $x$ . If smallest  $\epsilon > 0$

$$x + \epsilon_x > x$$

directly related to and proportional to value of  $x$ . It grows with  $x$

Thus many numbers cannot be represented by machine. Such numbers are represented by approximations. Such induced errors are called round-off errors.

Accuracy :-

floating point comparison "close enough" instead of equality

Absolute difference :  $|x - y|$

Relative difference :  $\frac{|x-y|}{|x|}$

We must address this in any calculation

Q:

if  $x == y$  [not to be done]

if  $\text{abs}(x-y) < \text{tol}$  [✓]

### Absolute Error

$\alpha$  - exact value

$\hat{\alpha}$  - Some Computed value

$$E(\hat{\alpha}) = |\alpha - \hat{\alpha}|$$

### Relative error

$$E_{rel}(\hat{\alpha}) = \frac{|\alpha - \hat{\alpha}|}{|\alpha|}$$

$$E_{rel}(\hat{\alpha}) = \frac{|\alpha - \hat{\alpha}|}{\hat{\alpha}_{\text{act}}}$$

Could be used to determine convergence

(Absolute)

$x = \dots$

$x_{old} = \dots$

while  $\text{abs}(z - z_{old}) > tol$

$x_{old} = z$

update  $z$

end

(Relative)

$x = \dots$

$x_{old} = \dots$

while  $\frac{\text{abs}(z - z_{old})}{x_{old}} > tol$

$x_{old} = z$

update  $z$

end

Cancellation due to Round off  
Errors

---

### Quadratic Equation

$$ax^2 + bx + c = 0$$
$$x = \frac{-b \pm \sqrt{b^2 - 4ac}}{2a}$$

$$a=1, \quad b=54.32, \quad c=0.1$$

Solution (11-digits)

$$x_1 = 54.3218158995$$

$$x_2 = 0.0018416049576$$

Note:  $b^2 = 2950.7$   
 $\ggg 4ac = 0.4$ .

Floating point 4-digit [Decimal]  
Arithmetic :-

$$\sqrt{b^2 - 4ac} = \sqrt{2951 - 0.4000}$$
$$= \sqrt{2951} = 54.32$$

$$x_{1,4} = \frac{108.6}{2.000} = 54.30$$

$$x_{2,4} = \frac{54.32 - 54.32}{2.000} = 0.$$

Relative Error :  $1 \rightarrow 0.4\%$   
 $2 \rightarrow 100\%$ .

Alternate: see due to cancellation of  
 2 numbers close in 4 digit

Rationalize!

$$\frac{x}{z_{1,4}} = \frac{2c}{-b + \sqrt{b^2 - 4ac}} = \frac{0.2}{108.6} = 0.0018$$

Relative error  $\equiv$  0.05 percent  
 reduced.

- Alas alteration if used for  $x_{1,4}$  would imply division by zero.  
 - Catastrophic Error

Robust Solution :

Define:  $a_r = -\frac{1}{2} [b + \text{Sign}(b) \sqrt{b^2 - 4ac}]$   
 $(b \neq 0)$

$$\text{Sign}(b) = \begin{cases} 1 & b > 0 \\ -1 & b < 0 \end{cases}$$

Then the two solutions are :-

$$x_1 = \frac{av}{a} \quad x_2 = \frac{c}{av}$$

- Finite precision causes round off in individual calculations
- Effect of round off accumulates
- Subtracting nearly equal numbers or very different in magnitude numbers leads to large loss of precision.
- $\alpha \gg \beta$  or  $\beta \gg \alpha$  then  $\alpha + \beta$  or  $\alpha - \beta$  will have cancellation error
- $\alpha \approx \beta$  then  $\alpha - \beta$
- Error caused in one operation itself lead to a huge loss.



