

Example

x = exact number (often don't know this)

\bar{x} = approx number

$e_x = x - \bar{x}$ = error in \bar{x}

$E_x = \frac{x - \bar{x}}{x}$ = relative error in \bar{x} .

Let $x = 92.8711149$ (in base 10).

- a) Find \bar{x} = binary # used by computer to approximate x with 16 bit precision, with standard rounding [f | c | s | B]
10 | 5 | 1 | 15
IEEE 754

Solution: Because our significand will have 11 digits, so we need to calculate x in binary up to 11 signif. digits.

$$x = 92.8711149$$

$$92 = \underbrace{64 + 16 + 8 + 4}_{2^6} = 1011100_2$$

7 signif. digits.
Just need 4 more.

$$\begin{array}{r} .8711149 \\ \times .5 \\ \hline .3711149 \\ + .01_2 \\ \hline .1211149 \\ + .0001_2 \\ \hline .0625 \end{array}$$

$$x = 1011100.\overline{1101} \dots$$
$$\bar{x} = 1011100.1110_2$$

rounded.

$$\begin{array}{r}
 \overline{x} = 92.875 \\
 = 1011100,1110_2
 \end{array}
 \quad
 \left.
 \begin{array}{l}
 \text{• } 0586149 \\
 \text{• } 0312500 + .000011_2 \\
 \vdots \\
 \cdot \frac{5}{250} \\
 \cdot \frac{125}{.875}
 \end{array}
 \right\} .111_2$$

b) Write the 16-bit binary expression (in float format) that is used to store \bar{x} .

$f(10)$	$e(5)$	$s(1)$
01110001110	10101	0 ↑

$$\tilde{x} = 1.0111001110 \times 2^6 \quad (1 \text{ for sign})$$

6 = e - B

6 = e - |S| \Rightarrow e = 21

$$e = 21 = 16 + 4 + 1 = 10101$$

Answer:

c) Find e_x & E_x :

$$e_x = x - \bar{x} = 92.87 - 92.875 = -0.0038851$$

$$\begin{array}{r} 928750000 \\ - 928711149 \\ \hline .0039851 \end{array}$$

$$E_x = \frac{e_x}{\bar{x}} = \frac{-0.0038851}{92.875} \boxed{\approx -0.0000418315}$$

(d) Suppose that x is instead truncated to binary number with a 400-digit significand (including the leading 1). To how many significant decimal digits is this accurate?

(Note: Accurate to d digits means the number can at most be off by 1 in the d^{th} slot, i.e.

$$|e_x| \leq \frac{1}{2} \cdot b^{\text{last slot}}$$

Recall $x = 92.8711149$

binary: $\bar{x} = 1011100.11011\dots \dots \#$
 $\underbrace{}_{\text{3 digits}}, \underbrace{}_{\text{393 digits}}$

$$|e_x| \leq 0.0\dots \underbrace{9}_{392} \cdot 2^{-393}$$

ie $|4.957 \times 10^{-119}| \approx 4.957 \times 10^{-119} \leftarrow \begin{matrix} \text{accurate} \\ \text{upto } 10^{-118} \text{ in decimal} \end{matrix}$

$$\bar{x} = 92.871114$$

$\Rightarrow 120$ (significant decimal places of accuracy)

↑
 10^{-118} slot.

Propagation of error -

Error analysis.

Idea: Using approximate numbers in calculations results in more errors. How bad is it?

Can we estimate or bound the error?

Recall x = exact number
 \bar{x} = approximation to x

$$e_x = x - \bar{x}$$

$$\epsilon_x = \frac{x - \bar{x}}{\bar{x}} = \frac{e_x}{\bar{x}}$$

$$\Rightarrow x = \bar{x} + e_x$$

$$e_x = \bar{x} \epsilon_x$$

Example: Use floating point representation

Note:-
assume
 $e \neq 0$
 $e \neq \text{null}$

$$\bar{x} = (-1)^f b^{-t} b^{e-B} \quad (f | e | s)$$

↑
(including leading 1)
↓
adding ready.

$$|e_x| \leq \underbrace{1 \cdot b^{-t} b^{e-B}}_{1 \text{ in the last significant}}$$

$$|e_x| = \left| \frac{e_x}{\bar{x}} \right| \leq \left| \frac{b^{-t} b^{e-B}}{(-1)^f b^{-t} b^{e-B}} \right| = \frac{|b^{-t} b^{e-B}|}{|(-1)^f b^{-t} b^{e-B}|}$$

$$= \frac{b^{-t} b^{e-B}}{f b^{-t} b^{e-B}} = \frac{1}{f}$$

$$|e_x| \leq \frac{1}{f} \leq \frac{1}{10_6^{m-1}} = \frac{1}{b^{m-1}} = b^{1-m}$$

lowest possible denom.
↑
sign digits

Example: approximating a decimal # to 15 sign digits $|e_x| \leq 10^{-15} = 10^{-14}$

② approx - a binary # to 39 sign
digits $\Rightarrow |E_x| \leq 2^{1-39} = 2^{-38}$

think of as $(100\% - 2^{-38})$
accurate.

These give error estimates & relative
error estimate for floating point
representation.

Operations

Addition: Given \bar{x}, \bar{y} as
approximations to x, y (so we know
 e_x, E_x, e_y, E_y , or we know
bounds on $|e_x|, |e_y|, |\varepsilon_x|, |\varepsilon_y|$).

How close is $\bar{x} + \bar{y}$ to $x + y$?

i.e. What are $e_{x+y}, \varepsilon_{x+y}$ in terms
of $e_x, e_y, \varepsilon_x, \varepsilon_y$?

We use $\bar{x} + \bar{y} = \bar{x} + \bar{y}$ as our
estimate of $x + y$.

$$\begin{aligned}
 e_{x+y} &= x + y - (\bar{x} + \bar{y}) \\
 &= (\bar{x} + e_x) + (\bar{y} + e_y) - (\bar{x} + \bar{y})
 \end{aligned}$$

\Rightarrow

$e_{x+y} = e_x + e_y$

TRIANGLE INEQUALITY

$$\begin{aligned}
 \text{i)} |A+B| &\leq |A| + |B| \\
 \text{ii)} |A+B| &\geq |A| - |B| \\
 |A+B| &\geq |B| - |A| \\
 |A-B| &\geq |A| - |B| \\
 |A-B| &\geq |B| - |A|.
 \end{aligned}$$

\Rightarrow

$|e_{x+y}| = |e_x + e_y|$

$|e_{x+y}| \leq |e_x| + |e_y|$

Relative error:

$$e_x = \bar{x} \epsilon_x$$

$$e_y = \bar{y} \epsilon_y$$

$$e_{x+y} = (\bar{x} + \bar{y}) \epsilon_{x+y}$$

Plug in: $|(x+y)\epsilon_{x+y}| \leq |x|\epsilon_x + |y|\epsilon_y$

$$|AB| = |A||B| \Rightarrow |x+y||\epsilon_{x+y}| \leq |x||\epsilon_x| + |y||\epsilon_y|$$

If $x, y > 0$: $|x| = x$, $|y| = y$ etc.

$$(x+y) |\epsilon_{x+y}| \leq x |\epsilon_x| + y |\epsilon_y|$$

If $A \leq B$

- C positive $\Rightarrow \frac{A}{C} \leq \frac{B}{C}$
- C negative $\Rightarrow \frac{A}{C} \geq \frac{B}{C}$

$$|\epsilon_{x+y}| \leq \frac{x |\epsilon_x| + y |\epsilon_y|}{x+y}$$

$\boxed{\begin{aligned} & \text{If } x, y > 0 \\ & \Rightarrow |\epsilon_{x+y}| \leq \frac{x |\epsilon_x| + y |\epsilon_y|}{x+y} \\ & \Rightarrow |\epsilon_{x+y}| \leq \frac{x}{x+y} |\epsilon_x| + \frac{y}{x+y} |\epsilon_y| \\ & \Rightarrow |\epsilon_{x+y}| \leq |\epsilon_x| + |\epsilon_y| \end{aligned}}$

\star better estimate.