

Models for arithmetic by computer.

Last time, we defined interval arithmetic and showed that calculations with slightly fuzzy numbers are... weird.

We're now going to introduce two more models which get closer to what a computer does.

Definition. (Base-10, n-digit ~~floating~~ fixed point)

We define fixed point numbers to have at most n digits to the right of the decimal point.

$f_x(x, n)$ = the closest* number to x in the fixed point #s with n-digits.
round up!

(2)

Note. $f_x(x, n)$ is always rational,
in fact always an ^{integer}
multiple of 10^{-n} .

We define the ^{fixed point} arithmetic operations
by

$$x \odot y = f_x(f_x(x, n) \odot f_x(y, n))$$

where \odot is one of $+, -, \times, \div$.

~~Examples.~~ We observe that the results
of the fixed point operations are not
always correct! *

Example. Let $n=5$.

$$100 \times 10^{-6} = 10^{-4}$$

$$f_x(f_x(100) \times f_x(10^{-6}, 5)) = f_x(100 \times 0) = f_x(0) = 0.$$

(3)

Example. Let $n=4$.

$$\begin{array}{r} 0.1036 \\ \times 0.2122 \\ \hline \end{array}$$

$$0.2122 \times 0.2081 = 0.044158 \xrightarrow{\text{fx}} 0.04415$$

$$0.1036 \times 0.4247 = 0.04399 \xrightarrow{\text{fx}} 0.0440$$

We can carry this one step further

$$(0.1036 \times 0.4247 - 0.2122 \times 0.2081)$$

\downarrow multiply
and round

$$(0.0440 - 0.0442)$$

\downarrow subtract
and round

$$-0.0002$$

The correct answer is -0.0001599 .

The difference might not look like much, but it's almost 30% of the correct answer!

The difference between x and $f_x(x, n)$ is called roundoff error.

Definition. If a is correct and b is an approximation to a , then

absolute error in b is $|b-a|$

relative error in b is $\frac{|b-a|}{|a|}$

Observation. The absolute roundoff error $|x - f_x(x, n)| \leq \frac{1}{2} \times 10^{-n+1} = 5 \times 10^{-n-1}$.

The relative roundoff error is unbounded! (Just make x very small.)

(precision-linear-equations-slideshow.nb)

(5)

Fixed point arithmetic seems to be something you could program, but there's still a problem - we haven't bounded the number of digits left of the decimal point. This leads to a natural idea - bound the total number of digits by writing everything in scientific notation.

Definition. (Base 10 n-digit ~~floating point~~ floating point)

A number in n-digit floating point is written $\pm d_1.d_2\cdots d_n \times 10^e$ where e is an integer between bounds $m \leq e \leq M$, and $d_1 \neq 0$ or all $d_i = 0$.

(6)

Note that all these are rational numbers and that the set is weird.

Example. 1-digit with $0 \leq e \leq 1$.

The numbers are ±

1 2 3 4 5 6 7 8 9 ~~10 20 30 40 50 60 70 80 90~~
 10 20 30 40 50 60 70 80 90

We define $f1(x)$ to be the nearest floating point number to x in the current system, and *(round up for ties)

$$x \odot y = f1(f1(x) \odot f1(y))$$

where \odot is one of $+, -, *, \frac{*}{*}$.

Example.

$$2+5=f1(f1(2)+f1(5))=f1(7)=7.$$

$$2+30=f1(f1(2)+f1(30))=f1(32)=30. (?!)$$

7

$$3 \times 6 = f_1(f_1(3) \times f_1(6)) = f_1(18) = 20. (?!!)$$

$$24 \div 5 = f_1(f_1(24) \div f_1(5)) = f_1(20 \div 5) = f_1(4) = 4.$$

(This last one is particularly disturbing because it's not even $f_1\left(\frac{24}{5}\right) = f_1(4.8) = 5$.)

Even in this number system, we can see several problems are going to happen:

$$30 \times 5 = f_1(30 \times 5) = f_1(150) = ?$$

We could take the definition of f_1 literally and round to 100, but computers instead signal an overflow error (and stop computing). when the result

$$2 \div 30 = f_1\left(\frac{1}{16}\right) = ?$$

is larger than the largest representable number.

(8)

Again, we could round (to 0) but computers instead signal underflow when the result is between 0 and the smallest nonzero representable number.

Floating point is better than fixed point, but there are still error problems.

Proposition. The absolute error $|x - f_1(x)|$ is bounded by $\frac{1}{2} \times 10^{M-n}$. The relative error

$$\frac{|x - f_1(x)|}{|x|} \leq \frac{1}{2} \times 10^{m-n}$$

Proof. We know that any x between $-9.9\dots9 \times 10^M$ and $9.9\dots9 \times 10^M$

(9)

We say that x is in range if

$$-\underbrace{9.9\dots 9}_{n \text{ digits}} \times 10^M \leq x \leq \underbrace{-1.000\dots 0}_{n \text{ digits}} \times 10^m$$

or

$$\underbrace{1.0\dots 0}_{n \text{ digits}} \times 10^m \leq x \leq \underbrace{9.9\dots 9}_{n \text{ digits}} \times 10^M.$$

(That is, if x does not trigger an underflow or overflow error.)

Proposition. If x is in range, then

$$|x - f_1(x)| \leq \frac{1}{2} 10^{M-n+1}$$

Proof. Since x is in range, it can be written as a (possibly infinite decimal)

$$x = \pm d_1.d_2\dots d_n d_{n+1}\dots \times 10^e$$

where $d_1 \neq 0$ and $m \leq e \leq M$. Then

(10)

Rounding $\star d_1.d_2 \dots d_n d_{n+1} \dots$ to n digits changes its value by at most $\frac{1}{2} 10^{-n+1}$. This is multiplied by 10^e , so

$$\begin{aligned}|x - f(x)| &\leq \frac{1}{2} 10^{-n+1} \times 10^e \\ &\leq \frac{1}{2} 10^{e-n+1} \\ &\leq \frac{1}{2} 10^{M-n+1}\end{aligned}$$

□

Exercise. What's $f(9.9\dots 96 \times 10^e)$? Does it still obey the proposition?

Proposition. If x is in range, then

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

Proof. Again,

$$x = \pm d_1.d_2 \dots d_n \times 10^e$$

and

$$|x - f_1(x)| \leq \frac{1}{2} 10^{e-n+1}$$

Now $|x| \geq 1.0 \dots 0 \times 10^e$, so

$$\frac{|x - f_1(x)|}{|x|} \leq \frac{1}{2} 10^{-n+1}$$

□

Note: The absolute error depends on M (and may be large!), while the relative error depends only on n (and is quite small).

Corollary. If x is in range,

$$f_1(x) = x(1+\delta) \quad (|\delta| \leq \frac{1}{2} 10^{-n+1})$$

Proof. Let

$$\delta = \frac{f_1(x) - x}{x}$$

(12)

(Simplifying, we have

$$x \left(1 + \frac{f_1(x) - x}{x} \right) = f_1(x).$$

We just proved $\left| \frac{f_1(x) - x}{x} \right| \leq \frac{1}{2} 10^{-n+1}$.