

How are numbers stored on the computer? (§1.2) 1/20/2011 LL

First, remember what we call "scientific notation"

For any decimal number x (assume finite nonzero digits)

we can write

$$x = a \times 10^b$$

where $1 \leq |a| < 10$

Exception: When $x=0$, we simply set $a=b=0$

Examples

x (decimal notation)	x (scientific notation)
2011	2.011×10^3
412	4.12×10^2
3.14	3.14×10^0
0.000789	7.89×10^{-4}
0.2091	2.091×10^{-1}

Every decimal (or base-10) number can be written

$$a_n \cdots a_3 a_2 a_1 a_0 \cdot a_{-1} a_{-2} a_{-3} \cdots a_{-l} = \sum_{i=-l}^{+k} a_i 10^i$$

↑
decimal point

$a_i \in \{0, 1, 2, \dots, 9\}$

For example

	a_4	a_3	a_2	a_1	a_0	a_{-1}	a_{-2}	a_{-3}	a_{-4}
3.14			0	3	1	4	0	...	
$= 3 \times 10^0 + 1 \times 10^{-1} + 4 \times 10^{-2}$									
0.037				0	3	7	0	...	
$= 3 \times 10^{-2} + 7 \times 10^{-3}$									
2001			0	2	0	1	1	0	...
$= 2 \times 10^3 + 1 \times 10^1 + 1 \times 10^0$									

Binary numbers (base -2) are written

$$\sum_{i=k}^n b_i 2^i \quad \text{but equal} \quad \sum_{i=-l}^m b_i 2^i \quad (2)$$

(compare with $\sum_{i=-l}^k b_i 10^i$) and every $b_i = 0$ or ± 1 .

example :

$$5.75 = 4 + 1 + 0.5 + 0.25$$

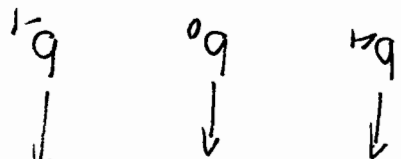
$$= 1 \times 2^2 + 1 \times 2^0 + 1 \times 2^{-1} + 1 \times 2^{-2}$$



(2)

$$17.5 = 16 + 1 + 0.5$$

$$= 1 \times 2^4 + 1 \times 2^0 + 1 \times 2^{-1}$$



(2)

Note that certain numbers which are finite decimals, actually are periodic in binary

$$\text{e.g. } 0.4_{(10)} = 0.01100110011\dots_{(2)} = 0.0110011_{(2)}$$

Machine numbers (binary floating point numbers)

The numbers stored on the computer are, essentially "binary numbers" in scientific notation

$$X = \pm a \cdot 2^b$$

a is called the mantissa
 b is called the exponent

The convention here is that $\frac{1}{2} \leq a < 1$. The idea

is that, for any number x , we can always divide it by an appropriate power of 2, s.t. the result will be

within $[\frac{1}{2}, 1)$.

e.g. $x = 5$

Try $b = 0$

$b = 2$

$b = 3$

$S = x = a \cdot 2^0 \Rightarrow a = 5$ too much

$S = x = a \cdot 2^2 \Rightarrow a = 1.25$ still too much

$S = a \cdot 2^3 \Rightarrow a = 0.625$ ok!

lastly, we have to write $a = 0.625$ in binary

$0.625 = 0.5 + 0.125 = 1 \times 2^{-1} + 1 \times 2^{-3} = 0.101^{(2)}$
Thus $S = x = + 0.101^{(2)} \times 2^3$

How about a different way?

Remember, in decimal:

$314 = 314 \times 10^0 = 0.314 \times 10^1 = \dots$ etc.

In binary

$S_{(10)} = 101^{(2)} = 10.1 \times 2^1 = 1.01 \times 2^2 = 0.101 \times 2^3$

In general, a machine number is stored as:

$$X = \pm 0.1a_1a_2 \dots a_{k-1}a_k \cdot 2^b$$

Single Precision

$$k = 23$$

$$-126 \leq b \leq 127$$

(max number $\approx \pm 3.4 \times 10^{38}$)

Double Precision

$$k = 52$$

$$-1022 \leq b \leq 1023$$

(max number $\approx \pm 1.8 \times 10^{308}$)

S.P has "23 binary significant digits". How much is that in decimal?

Rule of thumb $2^{10} \approx 10^3$

i.e. 10 binary significant digits ≈ 3 decimal digits

Thus: single precision \Rightarrow about 7 decimal significant digits

double precision \Rightarrow a bit more than 15!

Absolute and relative error

All computations are approximate, due to the limited precision on computer.

Exact quantity: q

Quantity on the computer: \hat{q}

2 measures of Error:

Absolute error = $|q - \hat{q}|$

important when we want to examine accuracy within a certain interval

Relative error = $\frac{|q - \hat{q}|}{|q|}$

important when we care about error as % of actual value.

Rounding Round to closest When only n digits ≈ 3.1

3.142
3.1415

(2 s.d)
(3 s.d)

truncate
3.1
3.14
3.1415

Rounding also occurs in binary

# sig. digit	Rounded	Truncated
1	1.0	0.1
2	0.11	0.10
3	0.110	0.101
4	0.1011	0.1011
5	0.10111	0.10110

Machine ϵ (epsilon)

There are a number of (slightly different) definitions in literature.

We will define the machine ϵ as the smallest positive machine representable number, such that

$$1 + \epsilon \neq 1 \text{ (on the computer)}$$

Why isn't the above inequality always true for any $\epsilon > 0$?

e.g.

$$1 = 0.10000 \dots 0 \cdot 2^0$$

$$+ 2^{-25} = 0.00 \dots 00 \dots 01 \cdot 2^1$$

$$1 + 2^{-25} = 0.1000 \dots 01 \cdot 2^0$$

$$1 = 0.10 \dots 01 \cdot 2^0 \text{ which is rounded to } 1$$

However

$$1 + 2^{-24} = 0.1000\dots 01 \times 2^1$$

23 digits

$$\text{(rounded to)} = 0.1000\dots 01 \times 2^1$$

23

> 1 (just barely)

Thus, $\epsilon = 2^{-24}$ ($\approx 6 \times 10^{-8}$)

Roundoff error

The importance of ϵ is that it corresponds to the maximum relative error incurred when converting an exact value to a machine-representable number

$$\frac{|q|}{|q - q'|} \leq \epsilon$$

q : actual number
 q' : machine representation

Another way to express this: let $\epsilon_1 = \frac{q'}{q}$

Then $q' = (1 + \epsilon_1)q$, where $|\epsilon_1| \leq \epsilon$.

For example, assume the real result of a computation

is $1,009,009,000$. ($=q$)

$$q = q \pm \epsilon q = 1,009,000,000 \pm 60$$

$\approx 6 \times 10^{-8} \times 10^9$

We have no way of telling whether the actual

solution was $1,009,009,005$ or $999,999,970$ \Rightarrow all look the same to the computer.

(Partial) Solution: Use higher precision

(Better) Solution: Make sure that

(1) You can tolerate this error

(2) Other algorithms that use this result can tolerate this error.

Square roots & quadratic equ revisited

$$ax^2 + bx + c = 0$$

In theory, just use the formula $x_{1,2} = \frac{-b \pm \sqrt{b^2 - 4ac}}{2a}$

In practice, we saw it's dangerous w/ errors.

Alternative?

First try simpler equation:

$$x^2 + a = 0$$

(we know the solution, $x = \sqrt{-a}$!)

Even now computing $\sqrt{\dots}$ can be expensive.

How about this:

→ Start with some initial guess $x^{(0)} = x^*$

→ Iterate the sequence

$$x^{(k+1)} = \frac{x^{(k)} (x^{(k)} + a)}{x^{(k)} + a}$$

x^k

Let us assume that this sequence converges.
 What will the limit be?

$$\text{let } \lim x_n = A$$

$$\lim (x_{n+1}) = \lim \left(\frac{x_n^2 + a}{2x_n} \right)$$

$$A = \frac{A^2 + a}{2A}$$

$$\Rightarrow 2A^2 = A^2 + a \Rightarrow A^2 = a \Rightarrow A = \sqrt{a} \text{ (or } -\sqrt{a})$$

Thus, if x_n has a limit, it will be the desired

square root. Note that only fundamental operations (addition/multiplication/division) were used.