

Storing numbers in the computer ("real numbers")

Basic idea: scientific notation.

Note: Computers store numbers in binary form, and the only numbers that can be stored exactly are rational numbers of the form $\frac{p}{2^n}$, where $p \in \mathbb{Z}$.

Otherwise - we just get approximations to other real numbers by numbers of the form $\frac{p}{2^n}$.

Floating Point Numbers:

Ingredients: b = base of the number system
 $(b=10 \text{ for decimals}, b=2 \text{ for binary } \#s)$

f = fraction / significand / mantissa

e = exponent

s = sign

t = radix point placement
 $(\text{eg decimal point})$

B = bias.

Example -4537.28

$$= -4.5372800 \times 10^3$$

\uparrow
 $s = \text{Sign} = 1$
 $\text{Sign} = (-1)^1$

$f = 45372800$
 \uparrow
 $t = \# \text{ of places} = 7$

Decimal
 $b = 10$
 base
 $e - B = 3$
 exponent.
 Sp. B = 7
 $\Rightarrow e = 10$

To store in the computer \rightarrow put fes into one number. (for a particular system, b, B, t are fixed ahead of time.)

$$\text{fes} = 45372800 \begin{matrix} [0010] \\ f \quad e \quad s \end{matrix}$$

\Rightarrow

\uparrow negative
 0 would mean positive

The corresponding number
that this means is $x = -4.5372800 \times 10^{10-7}$

In general

$$\underline{x = (-1)^s f \cdot b^{-t} b^{e-B}}$$

In reality, we use binary #'s in the computer, & numbers are displayed as decimals using a conversion program.

How does this work with binary #'s?

For reference: IEEE754 Standard

precision type	f	e	s	bias
half-precision (16 bit)	10	5	1	15
single precision (32 bit)	23	8	1	127
double precision (64 bit)	52	11	1	1023
quadruple precision (128 bit)	112	15	1	16383

16 bit example fes

1011101001011101 $B=15$

TRICK for binary - for most #s
assume the first digit is 1, and don't store it.

$f = 1011101001$ $t = 10$ for most #s.

$$\Rightarrow X = 1.1011101001 \times 2^{\frac{14-15}{2-1}} = [0.11011101001_2]$$

When $e = 00000$, then we take
 $f = (\text{exactly the } 10 \text{ digits}) \quad t = 9$

That way,
 $B = 000000000000|00000|e$.

$$x = 01\ 11000000|000000|0$$

$$= 0.111 \times 2^{0-15} = 0.111 \times 2^{-15}$$

$$= 1.11 \times 2^{-16}$$

~~101101001~~

Another special exponent:
when $e=11111$, $f=\underbrace{1111111111}_{0000000000}, s=0$
reserved for $+\infty$
 $e=11111, f=\underbrace{1111111111}_{0000000000}, s=1$
 $- \infty$.

$e=11111, f=\text{anything else}$
 $\rightarrow \text{NaN}$ "Not a number"