
Subject: Re: machine precision

Posted by [Wout De Nolf](#) on Wed, 20 May 2009 13:08:47 GMT

[View Forum Message](#) <> [Reply to Message](#)

Ok, so I was reading the Sky Is Falling paper and the Goldberg paper again. I learned some things I thought I'd share (since this is a recurring issue, despite the Sky Is Falling paper).

A floating point number is stored like this:

$f(\text{binary}) = \text{sign} \mid \text{exponent} \mid \text{mantissa without leading 1}$

sign: 1bit

exponent: 8bits (11bits when double)

mantissa: 23bits (52bits when double)

The real number it represents can be found like this

$f = \text{sign} \cdot \text{mantissa} \cdot \text{base}^{(\text{exponent} - \text{bias} - n_{\text{mantissa}})}$

sign: -1 when sign-bit=1, +1 when sign-bit=0

base: 2 (ibeta from MACHAR)

exponent: 8bit number

bias: 127 (1023 when double)

n_{mantissa} : number of mantissa bits (23, 52 when double)

We will rewrite this as

$f = \text{sign} \cdot \text{mantissa} \cdot \text{eps} \cdot \text{base}^{\text{exp}}$

eps: $\text{base}^{(-n_{\text{mantissa}})}$ (eps from MACHAR)

exp: exponent-bias

For example: $f = 470$.

$f(\text{binary}) = 0 \mid 10000111 \mid 110101100000000000000000$

sign = +1

exp = $135 - 127 = 8$

mantissa = 15400960

eps = $2.^{-23}$

$f(\text{stored}) = 15400960 \cdot 2.^{-15}$

The difference between a stored floating point number $f1$ and its closest neighbour $f2$:

$\text{abs}(f1 - f2) = \text{eps} \cdot (\text{mantissa1} \cdot \text{base}^{\text{exp1}} - \text{mantissa2} \cdot \text{base}^{\text{exp2}})$

smallest possible difference when:

$\text{exp1} = \text{exp2} = \text{exp}$

$\text{mantissa1} = \text{mantissa2} + 1$

$= \text{eps} \cdot \text{base}^{\text{exp}} = 1 \text{ ulp (unit in last place)}$

The absolute error made when storing a real number is therefore

$\text{abserr} = \text{abs}(f_{\text{real}} - f) \leq c \text{ ulp}$

where $c=1$ for truncation and $c=0.5$ for rounding

The relative error made is

$\text{relerror} = \text{abs}(\text{freal}-f)/\text{abs}(\text{freal})$

$\leq c.\text{eps}.\text{base}^{\text{exp}}/\text{abs}(\text{freal})$

$\leq c.\text{eps}$ (not sure about this last step....)

Finally, two numbers are considered equal if

$\text{relerr} = \text{abs}(f1-f2)/(\text{abs}(f1)>\text{abs}(f2)) \leq \text{eps}$

I'm not really sure about this one either (e.g. what should be in the denominator, what about c,...)

All this doesn't deal with accumulated errors in floating point arithmetic, only with errors introduced by storing a real number.
