
Subject: Re: Floats

Posted by [Paul Van Delst\[1\]](#) on Fri, 10 Mar 2006 16:45:59 GMT

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

Sheldon wrote:

- > Boy, you guys really went off on this issue. But it is good.
- > Much is revealed about IDL, the art of programming, and, an added
- > extra, the personalities of the writers :)
- > To answer an earlier question of why I asked. Well, maybe I am just too
- > curious about things :)
- > I used Matlab before and it had the possibility to limit the precision.
- > But I could not control it too much.
- > I am working with a large number of arrays and averaging millions of
- > pixel values. As a result I am getting rounding errors.
- > I stop the program here and there and check the data at different
- > points. As such I only want a precision to the nearest hundreth.
- > It makes for quick assessment. Silly but, I am still learning :)
- > I need to understand more fundemental things about ROUND, FIX, FLOAT,
- > and LONG so as to eliminate some of these annoying errors. (Thanks
- > David for the info)
- > Using double precision seem to be the next step for me.

Ah. You mentioning rounding errors as "annoying" sends up a red flag for me. :o) You should /expect/ this to occur if you're adding a whole bunch of numbers, and your algorithm should handle it - especially if the end result is affected (sometimes rounding errors only show up in intermediate results and the end result is fine.) I don't think use of ROUND, FIX, FLOAT, and LONG should be considered until the problem is better understood.

You might want to consider using a compensated summation routine to sum your millions of pixel values to minimise rounding errors. A popular (or, at least, better known) method is also called Kahan summation. There is also doubly compenated summation which requires the data being summed to be sorted in ascending order.

Depending on your problem (e.g. sorting the data first may be too onerous), one or the other should do. (Although maybe the TOTAL function in IDL already does this?)

Anyway, check out chapter 4 of "Accuracy and Stability of Numerical Algorithms" by Nicholas Higham. Accuarate summation of floating point numbers is exhaustively dealt with.

paulv

--

Paul van Delst
CIMSS @ NOAA/NCEP/EMC
