

Advanced Algorithms II

Todd Davies

February 13, 2016

Overview

This course will explore certain classes of algorithms for modelling and analysing complex systems, as arising in nature and engineering. These examples include: flocking algorithms - e.g., how schools of fish or flocks of birds synchronised; optimisation algorithms; stability and accuracy in numerical algorithms.

Aims

By the end of the course, students should:

- Appreciate the role of using nature-inspired algorithms in computationally hard problems.
- Be able to apply what they learnt across different disciplines.
- Appreciate the emergence of complex behaviours in networks not present in the individual network elements.

Syllabus

Attribution

These notes are based off of both the course notes (found on Blackboard). Thanks to the course staff (*{names}*) for such a good course! If you find any errors, then I'd love to hear about them!

Contribution

Pull requests are very welcome: <https://github.com/Todd-Davies/third-year-notes>

Contents

1	Finite precision computation	3	1.1.2	Error in floating point numbers	4
1.1	Floating point numbers	3	1.2	Finding the number of significant digits	6

1 Finite precision computation

Unfortunately, the world is not solely restricted to integers, and computers often need to work with real numbers \mathbb{R} . With integers, the main problem we have in computer terms is overflow, and since there is a finite distance from one to the next, they are easy to encode in a computer.

On the other hand, between any two real numbers, there are infinitely many more real numbers. Since computers are discrete, we need to sample the real numbers so that we can find a representation for them in the computer. However, this introduces errors, since we can't represent every value exactly and therefore most approximate.

1.1 Floating point numbers

One problem we have with computation is that we don't know what the error is with computations; how 'good' is the result of an algorithm or computation? We would like to know the error bounds of a solution, and have the output be reliable.

In the 70's, it was realised that different floating point implementations produced different results. This had significant concerns for reproducibility, and as a result the ANSI IEEE standard for binary floating point arithmetic was created.

Each floating point number is represented as four integers; the base, the precision, the exponent and the mantissa.

$$x = \pm m \times b^{e-n}$$

Where:

m		The Mantissa (the bit before the decimal place)
b		The base (or radix), usually two or ten
e		The exponent (the power of the radix)
n		The precision (the number of digits in the mantissa)

We can represent different numbers in different ways, for example:

$$0.121e10^3 = 0.0121e10^4$$

In this case, we can normalise the way in which we represent numbers and at least all computers will get the same errors.

The amount of numbers we can represent with the floating points depends on the values permissible for b, n and e . When $b = 2, n = 2$ and $e = [-2 \dots 2]$:

2	×	2^{-2}	=	0.500000
3	×	2^{-2}	=	0.750000
2	×	2^{-1}	=	1.000000
3	×	2^{-1}	=	1.500000
2	×	2^0	=	2.000000
3	×	2^0	=	3.000000
2	×	2^1	=	4.000000
3	×	2^1	=	6.000000
2	×	2^2	=	8.000000
3	×	2^2	=	12.000000

The Python code used to generate this is found in the /COMP36212/programs folder of the source for these notes.

The mantissa is always 2 or 3 since we're using an explicit one, so the binary values are either $[1, 0]$ or $[1, 1]$.

Floating point numbers are relatively spaced; even though they might not be the same distance apart, the ratio between them is the same. The unit round off (basically the last digit) is called the *relative machine precision*.

The **Relative Machine Precision** is given by $u = 0.5 \times b^{1-n}$, and is the largest possible difference between a real number and its floating point representation. In the above example, $u = 0.5 \times 2^{-1} = \frac{1}{4}$. The value $2u$ is called the **Machine Precision**.

In the exam, assume explicit storage of leading bit of mantissa.

1.1.1 Real world floats

For the sign bit, 0 is for positive numbers and 1 is for negative ones. The exponent must also be able to represent negative numbers (in the case of 24×2^{-2} for example), and thus in single precision floats, a bias of +127 is added to the exponent and that value is stored. The exponent values -127 and +128 are reserved for special numbers.

The first bit of the mantissa is implicitly 1 in the IEEE base two floating point representation. This is because normalised numbers always have 1 as the first digit of their mantissa, and then we can get another digit of precision.

Sixty-four bit floating point numbers have one sign bit, 11 exponent bits and 52 mantissa bits. This means their bias will be $2^{11} = 2048$ and the range will be $2 - 2^{52} \times 2^{2^{11}}$.

The standard also has special values built in:

Zero: When the exponent is all zeros and the mantissa equal to zero.

Denormalised number: If the exponent is all zero, but the mantissa is non-zero, then the number is $-1^{sign} \times 0.m \times 2^{-126}$.

Infinity: Exponent is all 1's and mantissa is all 0's. The sign dictates between positive and negative infinity.

NaN: Not your grandma, this is when the value isn't a real number, such as when a division by zero occurs. The exponent is all 1's and the mantissa is non-zero.

1.1.2 Error in floating point numbers

When a real number is converted to floating point number, it may lose precision. If the real number is x and the floating point representation is \bar{x} , then the error is:

$$e = \bar{x} - x$$

You can find how many significant digits a floating point number approximates a real number to by doing:

$$|\bar{x} - x| = 10^{-\text{significant digits}}$$

However, the absolute error e does not give us a very good description of the accuracy (if the error is 10^{-6} but the value is 10^{-7} then we're very, very inaccurate)! To rectify this, we have relative error:

$$r = \frac{e}{x} = \frac{\bar{x} - x}{x}$$

When we're getting the floating point number from the real number, we can truncate the (possibly infinite) digits so that it fits in the mantissa. Simply chopping the number so it fits in m bits is called **simple truncation**.

If we round numbers instead of using simple truncation, then we can reduce the error. We need some rules though:

- If the part of the mantissa to be chopped off is less than 0.5, use simple truncation.

- If it's greater, then increment the last digit of the mantissa, and then truncate.
- If it's equal to 0.5, then we can do either (though IEEE says to round up).

Now the relative error is:

$$|r| = \frac{|e|}{|x|} = \frac{0.5 \times b^{e-n}}{|m| \times b} = \frac{1}{2 \times m \times b^n}$$

There are three types of errors that computers can make:

- Essential errors are ones that cannot be avoided (e.g. from erroneous input).
- Rounding errors are when we have to approximate real numbers with floating point ones. This error can be measured and controlled.
- Methodology errors come into play when replacing one problem by another similar, easier but less accurate problem is done. The solution is close, but not exact.

Unfortunately, errors can propagate through a computation. We must know the errors introduced by every operation a computer performs on floating point numbers. If we know $e_x = \bar{x} - x$ and $e_y = \bar{y} - y$, what is $e_{x \cdot y}$?

The error introduced by addition, subtraction, multiplication and division is:

Addition:

$$\bar{x} + \bar{y} = (x + e_x) + (y + e_y) = (e_x + e_y) + (x + y)$$

$$e_{x+y} = e_x + e_y$$

Subtraction:

$$e_{x-y} = e_x - e_y$$

Multiplication:

$$\bar{x} \times \bar{y} = (x + e_x) \times (y + e_y) = xy + xe_y + ye_x + e_x e_y$$

$$e_{x \times y} \approx xe_y + ye_x$$

Division:

$$e_{\frac{x}{y}} \approx \frac{1}{y}e_x - \frac{x}{y^2}e_y$$

Relative error can also be calculated:

Addition:

$$r_{x+y} = \frac{e_{x+y}}{x+y} = \frac{x}{x+y}r_x + \frac{y}{x+y}r_y$$

Subtraction:

$$r_{x-y} = \frac{e_{x-y}}{x-y} = \frac{x}{x-y}r_x + \frac{y}{x-y}r_y$$

Multiplication:

$$r_{x \times y} = \frac{e_{x \times y}}{x \times y} \approx r_x + r_y$$

Division:

$$r_{x/y} = \frac{e_{x/y}}{x/y} \approx r_x - r_y$$

In general, $\bar{x} \circ \bar{y} = (x \circ y)(1 + r_{x \circ y})$.

While it is useful to know the error of one operation, we also need to be able to work out the error of consecutive operations. That is to say given $e_x = \bar{x} - x$ and $e_y = \bar{y} - y$, determine $e_{x \circ y}$.

Remember, error propagation is not associative. The error from a multiplication and then an add is probably not the same as doing the add then the multiplication.

$$\begin{aligned}\bar{x} &= FP(\bar{x}_1 \circ \bar{x}_2) \\ &= \bar{x}_1 \circ \bar{x}_2 \\ &= (\bar{x}_1 \circ \bar{x}_2)(1 + u)\end{aligned}$$

Where $|u| \leq 0.5 \times b^{-n+1}$.

The total relative error is given by:

$$r_z^t = \frac{e_z^t}{x} = r_{x \circ y} + u = a_x r_x + a_y r_y + u$$

a_x and a_y is the error introduced by the \circ operation.

Lets put all that into an example. Given the numbers x, y and z with their relative round off errors r_x, r_y and r_z , determine the relative error in $u = (x + y)z$:

$$\begin{aligned}r_{x+y}^t &= \frac{x}{x+y} r_x + \frac{y}{x+y} r_y + r_+ \\ r_u^t &= \frac{x}{x+y} r_x + \frac{y}{x+y} r_y + r_+ + r_z + r_*\end{aligned}$$

1.2 Finding the number of significant digits

We want to find an integer that represents how many digits in our number are non-nonsense (i.e. how many significant digits we have). The number of significant digits in the floating point number \bar{x} where its real equivalent is x is:

\mathbb{Z} is the set of integers.

$$l = \mathbb{Z}(\log_b \frac{|x|}{|\bar{x} - x|})$$

If we rearrange this, the relative error is:

$$r_x \approx b^{-l}$$

If we have a computation that takes m real numbers as arguments and outputs a real number, if the arguments are floating point numbers with l_i significant digits then we can estimate:

$$|e| \approx \left| \sum_{i=1}^m x_i \frac{\delta f}{\delta x_i} b^{-l_i} \right|$$

Also:

$$\left| \sum_{i=1}^m x_i \frac{\delta f}{\delta x_i} b^{-l_i} \right| \leq b^{-l_{min}} \left| \sum_{i=1}^m x_i \frac{\delta f}{\delta x_i} \right|$$

The number of significant digits in the answer is:

$$l = l_{min} - \delta$$

Where δ is the loss of significant digits:

$$\delta = Z(\log_b(\frac{\sum_{i=1}^m |x_i \frac{\delta f}{\delta x_i}|}{|f(x_1, \dots, x_m)|}))$$

If we try and subtract numbers that are close in magnitude, then we will lose lots of significant digits. If we do $\sqrt{2.01} - \sqrt{2}$ (where both numbers are known to 9 significant digits), then we get:

$$\delta = Z(\log_{10}(\frac{|\sqrt{2.01}| + |\sqrt{2}|}{|\sqrt{2.01} - \sqrt{2}|})) = 3$$

Our answer would be to six significant figures. In order to get all of the significant figures, we need to use a different method:

$$\begin{aligned} z &= \sqrt{2.01} - \sqrt{2} \\ &= (\sqrt{2.01} - \sqrt{2}) \frac{\sqrt{2.01} + \sqrt{2}}{\sqrt{2.01} + \sqrt{2}} \\ &= \frac{\sqrt{2.01}^2 - \sqrt{2}^2}{\sqrt{2.01} + \sqrt{2}} \\ &= \frac{0.01}{\sqrt{2.01} + \sqrt{2}} \end{aligned}$$