# CS370 Lecture 1

Graham Cooper

January 4th, 2017

**Five Topics in the Course (Jan 4th)**

- Floating point numbers and Arithmetic

- Iterpolation, Splines, Parametric Curves

- Initial Value Problems - solve differencial equations

- Discrete Fourier Analysis

- Numerical Linear Algebra - solve equations - google pagerank

## Topic 1: Floating Point Arithmetic

Examples where problems come whehn using approximation

eg1. $e^{-5.5} =$

$$e^x = 1 + x + \frac{x^2}{2!} + \frac{x^3}{3!} + ... + \frac{x^n}{n!} + ...$$

$$e^{-x} = \frac{1}{x}$$

$$e^{-5.5} = \frac{1}{e^{5.5}} = \frac{1}{1 + 5.5 + \frac{5.5^2}{2} + ...}$$

Now do arithmetic keeping only 5 digits. In both cases infinite sums remain unchanged after 25 terms. There is no sense in going any further - we end up just truncating all of the terms after this as they are smaller than the 5th digit.

Method 1 gives $e^{-5.5} = 0.0026363$
Method 2 gives $e^{-5.5} = 0.0040868$

eg2. $ax^2 + bx + c = 0$
1

$$x_1 = \frac{-b + \sqrt{b^2 - 4ac}}{2a}$$

$$x_2 = \frac{-b - \sqrt{b^2 - 4ac}}{2a}$$

2

$$x_1 = \frac{2a}{-b + \sqrt{b^2 - 4ac}}$$

$$x_2 = \frac{2a}{-b - \sqrt{b^2 - 4ac}}$$

$$x^2 + 621x + 1 = 0$$

use 4 digit arithmetic

1: $x_1 = -0.200$ $x_2 = 62.1$
2: $x_1 = -0.0161$ $x_2 = -62.1$

# Floating point number systems (Jan 6th)

$$F(\beta, t, L, u)$$

Floating point number is either

$$0 +/- 0.x_1 x_2 ... x_t$$

Where $x_n$ is a digit

- $x_1 \neq 0$

- $0 \leq x_i < \beta$

- $L \leq d \leq u$

eg $F(10, 6, L, u)$ for $\pi = 0.314519 x 10^1$

Single Precision: [+/-][8 bits][23 digits] F(2, 24, -126, 127)
Double Precision: [+/-][11 bits][52 digits] F(2, 53, -1022, 1023)

Given any real number:

$$+/- 0.x_1 x_2 ... x_t x_{t+1} ... x \beta^d$$

you can truncate (rounding possible):

$$fl(x) = +/ - 0.x_1x_2...x_t$$

**Question:** How close is fl(x) to x? Relatively!

$$\delta = \frac{fl(x) - x}{x}$$

Claim: $|\delta| < \beta^{1-t}$ when truncating and $|\delta| < \frac{1}{2}\beta^{1-t}$ when rounding (Let this definition be $\epsilon$)

Trucating:

$$|\delta| = \frac{0.00...0x_{t+1}...x\beta^d}{0.x_1x_2...x\beta^d}$$

$$|\delta| = \frac{0.x_{t+1}x_{t+2}...x\beta^{-t}}{x_1.x_2x_3...x\beta^{-1}}$$

$$|\delta| < \beta^{1-t}$$

In general $|\delta| < \epsilon$

$$fl(x) = x(1 + \delta)$$

$$|\delta| < \epsilon$$

What about arithmetic an errors?

x,y real
$x \oplus y =$
eg. $x = 0.1111...$ F(2,4,L,u)
$y = 0.1110...$

$fl(x) = 0.1111$
$fl(y) = 0.1110$
when adding them: 1.1101 then truncate to $0.1110 \times 2^1$

x,y real
$x \oplus y = fl(fl(x) + fl(y))$

3

$$\left|\frac{x \oplus y - (x + y)}{x + y}\right|$$

$$= \left|\frac{fl(fl(x) + fl(y)) - (x + y)}{x + y}\right|$$

since $|\delta_1| < \epsilon$ and $|\delta_2| < \epsilon$ and $|\delta_3| < \epsilon$

$$= \left|\frac{(x(1 + \delta_1) + y(1 + \delta_2))(1 + \delta_3) - (x + y)}{x + y}\right|$$

$$\left|\frac{x + y + x\delta_1 + y\delta_2 + x\delta_3 + y\delta_3 + x\delta_1\delta_3 + y\delta_2\delta_3 - x - y}{x + y}\right|$$

$$\left|\frac{x\delta_1 + y\delta_2 + x\delta_3 + y\delta_3 + x\delta_1\delta_3 + y\delta_2\delta_3}{x + y}\right|$$

$$\leq \frac{|x||\delta_1| + |y||\delta_2| + |x||\delta_3| + |x||\delta_1||\delta_3| + |y||\delta_2\delta_3|}{|x + y|}$$

Since $|a + b| \leq |a| + |b|$ and $|a \times b| = |a||b|$

$$< \frac{(|x| + |y|)}{|x + y|}(2\epsilon + \epsilon^2)$$

$$\left|\frac{x \oplus y - (x + y)}{x + y}\right|$$

$$< \frac{|x| + |y|}{|x + y|}(2\epsilon + \epsilon^2)$$

If x and y have the same sign then $|x + y| = |x| + |y|$ and so addition has relative error bounded by $2\epsilon + \epsilon^2$

But if x and y are of opposite sign and perhaps of nearly the same size:

$$\frac{|x| + |y|}{|x + y|}$$

CATASTROPHIC CANCELLATION - Subtraction is deadly.
**Jan 9th:**
Recall: $F(\beta, t, L, u)$
For any non-zero real number

$x = 0.x_1 x_2 ... x \beta^d$

$fl(x) = 0.x_1 x_2 ... x_t$ then the exponent $x \beta^d$

We should the magnitude of the error:

$\delta = \frac{fl(x) - x}{x}$

showed: $|\delta| < \epsilon = \beta^{1-t}$ or $\epsilon = \frac{1}{2} \beta^{1-t}$

$\epsilon$ is called machine epsilon.

$fl(x) = x(1 + \delta)$ where $|\delta| < \epsilon$

## Operations

If x and y are two real numbers

$x \oplus y = fl(fl(x) + fl(x))$

$$\left| \frac{(x \oplus y) - (x + y)}{x + y} \right| < (2\epsilon + \epsilon^2)\left( \frac{|x| + |y|}{|x + y|} \right)$$

Notice: if x and y are of the same sign then relative error of addition $< 2\epsilon + \epsilon^2$

If x and y are of opposite sign and $x + y \approx 0$ then the upper bound can be significant. This is called catastrophic cancellation

$$x = 0.x_1 x_2 ... x_t ... x \beta^d$$

$$y = -0.x_1 x_2 ... \hat{x}_t ... x \beta^d$$

Subtracted $= 0.0.....0?$ with exponent $x \beta^d$

$$x = 0.1239...1231$$

$$y = -0.1229...1221$$

Subtracted $= 0.0010$

$$e^{-5.5} = 1 - 5.5 + \frac{5.5^2}{2} - \frac{5.5^3}{6} + ... = x - \hat{x}$$

5

$$e^{-5.5} = \frac{1}{e^{5.5}} = \frac{1}{1 + 5.5 + \frac{5.5^2}{2} + \ldots}$$

# Numerical Algorithm Issues

**Problem**: given $\alpha$ Solve numerically:

$$I_n = \int_0^1 \frac{x^n}{x + \alpha} dx$$

n = 0,1,2...

**Method**

$$I_0 = \int_0^1 \frac{1}{x + \alpha} dx = ln(x + \alpha)|_0^1$$

$$= ln(1 + \alpha) - ln(\alpha) = ln(\frac{1 + \alpha}{\alpha})$$

$$I_n = \frac{1}{n} - \alpha I_{n-1}$$

For $n \geq 1$:

$$I_n = \int_0^1 \frac{x^{n-1}(x + \alpha - \alpha)}{x + \alpha} dx$$

$$= \int_0^1 *1 x^{n-1} dx - \alpha \int_0^1 \frac{x^{n-1}}{x + \alpha} dx$$

$\alpha = 0.5$

$$I_0 = ln(3) = 1.09861228866..$$

$$I_0^a pp = 1.098612$$

$$I_1 00^a pp = 0.00664$$

$\alpha = 2.0$

$$I_0 = ln(1.5) = 0.405465108108..$$

$$I_0^a pp = 0.4054654$$

6

$$I_100^app = 2.1 \times 10^{22}$$

He used his own application and went through the solutions to the above. They are incorrect :S

Math: $I_0, I_n = \frac{1}{n} - \alpha I_{n-1}$
CS: $I_0^{app}, I_n^{app} = \frac{1}{n} - \alpha I_{n-1}^{app}$

$e_0 = I_0 - I_0^{app} = \left(\frac{1}{n} - \alpha I_{n-1}\right) - \left(\frac{1}{n} - \alpha I_{n-1}^{app}\right)$
$= -\alpha I_n + \alpha I_{n-1}^{app} = -\alpha e_{n-1}$

$$e_n = \alpha e_{n-1}$$
$$= \alpha^2 e_{n-2}$$
$$= -\alpha^3 e_{n-3}$$
$$...$$
$$= (-\alpha)^n e_0$$

IF $|\alpha| < 1$ then $|e_n| \to 0$ as $n \to \infty$ stable
$|\alpha| > 1$ then $|e_n| \to \infty$ is $n \to \infty$ unstable!!!

When $|\alpha| > 1$ then we can work backwards!

$$I_{200}, I_{199}, ..., I_{101} I_{100}$$
$$I_{n-1} = \frac{1}{\alpha n} - \frac{1}{\alpha} I_n$$
$$e_{n-1} = \frac{1}{2} e_n$$
$$|e_{199}| = \frac{1}{|\alpha|} |e_{200}|$$
$$|e_{100}| = \frac{1}{|\alpha|^{100}} |e_{200}|$$