

In this class:

- *Understand the need for floating point arithmetic and some alternatives.*

- *Understand how the computer represents floating point numbers and how it manipulates them.*

September 2, 2016

HW DUE

Floating point

Next class

Floating point mathematics
Floating Point
G&C – Chapter 5

Next next class

QUIZ and Floating point Math
G&C – Chapter 5

Matlab demo of the following slides

Floating point

$$1/3 = 0.333333333333333333333333$$

$$2/3 = 0.666666666666666666666666$$

$$+ = 1.000000000000000000000000$$

$$3/10 = 0.299999999999999999999999$$

$$4/10 = 0.400000000000000000000002$$

$$+ = 0.699999999999999999999996$$

Computers can't subtract ????

```
>> x = 1e18;
```

```
>> x
```

```
x =
```

```
1.0000e+18
```

```
>> y = 1e18;
```

```
>> y = y + 1;
```

```
>> y - x
```

```
ans =
```

```
0
```

Shouldn't this be equal to 1?

Computers approximate!

```
>> x = 1e18;
>> y = 1e18;
>> x == y
ans =
    1
>> y = y + 1;
>> x == y
ans =
    1
>> explain_double(x)
IEEE 754 Double precision floating point representation (64-bits)
  sign = 0 (1 bit)
  exp  = 10000111010 (11 bits)
  frac = 101111000000101101101011001110100111011001000000000000
  val  = 1.00000000000000000000e+18
```

Alternatives to floating point

(From Nick Trefethen) Exact (rational) arithmetic, we want the roots of

$$p(x) = x^5 - 2x^4 - 3x^3 + 3x^2 - 2x - 1.$$

No “analytical” formula. (Galois 1820s)

Alternatives to floating point

Using Newton's Method to solve a polynomial

$$x^{(0)} = 0,$$

$$x^{(1)} = -\frac{1}{2},$$

$$x^{(2)} = -\frac{22}{95},$$

$$x^{(3)} = -\frac{11414146527}{36151783550},$$

Alternatives to floating point

Using Newton's Method to solve a polynomial

$$x^{(4)} = - \frac{43711566319307638440325676490949986758792998960085536}{138634332790087616118408127558389003321268966090918625} ,$$

$$x^{(5)} = - \frac{7243914791768201761290013818789259730350038836047543931178041194343579260105802744696299}{22974602373157587333399081666432003514775984720802108866006687478324948875098845198224797} \\ \frac{22882064184585670017703551996316651611596343634562735299921308664663139405767412052875538}{58228984471808467981536221568972260935865495325922571792991768547894449519518216876316931} \\ \frac{2012406424843006982123545361051987068947152231760687545690289851983765055043454529677921}{5683704659081440024954196748041166750181397522783471619066874148005355642107851077541250} .$$

20 iterations produce a 16 terabyte file

Exact arithmetic doesn't scale.

Fixed point

Use exact integer arithmetic. But keep a fixed decimal place.

e.g. compute with “cents” but report “dollars”

It's hard to control scale, but for some applications, it'll work great!

Uses of fixed point

Taxes!

Finance software (e.g. GNU Cash)

Many MP3 decoders.

Doom (the game)

What is and why floating point?

Floating point representations are like scientific notation.

number = significant digits \times base^{exponent}

$$12,342.1 = 1.23421 \times 10^4$$

It's a compromise!

A sense of floating point scale.

The floating point numbers we'll use have around 16 significant digits.

1 Astronomical Unit =

$1.49597871 \times 10^{17}$ microns

[More info](#)

That's enough of a dynamic range to (almost) represent a 10 micron particle between here and the sun.

Floating point arithmetic

Is $x^2 - y^2$ or $(x + y)(x - y)$ better?

How do I evaluate the roots of a quadratic?

$$ax^2 + bx + c = 0$$

$$f(x) = (x - t_1)(x - t_2) \cdots (x - t_k)$$

$$f'(x) = \sum_{i=1}^n \frac{f(x)}{x - t_i}$$

Floating point lectures

- 1) Mechanics of floating point (Section 5.3)
- 2) Mathematics of floating point (Section 5.5)
- 3) Study of floating point properties (Misc. ex)