

# CSC338 WINTER 2022

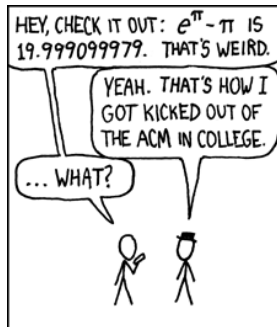
## WEEK 2 - FLOATING POINT NUMBERS

Ilir Dema

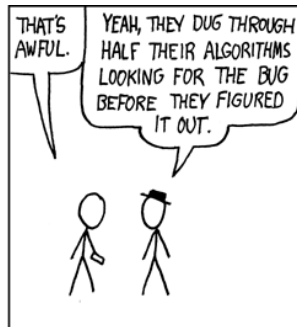
University of Toronto

Jan 21, 2022

# COMPUTE $e^\pi - \pi$



DURING A COMPETITION, I TOLD THE PROGRAMMERS ON OUR TEAM THAT  $e^{\pi} - \pi$  WAS A STANDARD TEST OF FLOATING-POINT HANDLERS -- IT WOULD COME OUT TO 20 UNLESS THEY HAD ROUNDING ERRORS.



<https://xkcd.com/217/>

## SCIENTIFIC NOTATION

## MOTIVATION

- [illegible]

# FLOATING-POINT NUMBERS

## DEFINITION

- Floating-point number system is the set of numbers defined by the tuple  $\mathcal{F}(\beta, p, L, U)$  where

 $\beta$  base, integer  $\geq 1$ 

$p$  precision, integer  $\geq 1$

$[L, U]$  exponent range, integers  $L \leq U$

- A number  $x \in \mathcal{F}$  if and only if

$$x = \pm \left( d_0 + \frac{d_1}{\beta} + \frac{d_2}{\beta^2} + \cdots + \frac{d_{p-1}}{\beta^{p-1}} \right) \beta^E$$

$$0 \leq d_j \leq \beta - 1, i = 0, \dots, p-1, \quad L \leq E \leq U$$

- $d_0 d_1 \dots d_{p-1}$  is called *mantissa*

- $d_1 \dots d_{p-1}$  is called *fractional part*

- $E$  is called *exponent*

- If  $\mathcal{F}$  is a fixed floating-point number system, its elements are called *machine numbers*.

## Typical Floating-Point Systems

Most computers use binary ( $\beta = 2$ ) arithmetic

Parameters for typical floating-point systems shown below

system	$\beta$	$p$	$L$	$U$
IEEE SP	2	24	-126	127
IEEE DP	2	53	-1022	1023
Cray	2	48	-16383	16384
HP calculator	10	12	-499	499
IBM mainframe	16	6	-64	63

IEEE standard floating-point systems almost universally adopted for personal computers and workstations

## EXAMPLE

$$\mathcal{F}(10, 3, -5, 5)$$

- $x = -20.5 \in \mathcal{F}$ , because  $x = -(2 + \frac{0}{10} + \frac{5}{10^2})10^1$ 
  - mantissa: 205
  - fractional part: 05
  - exponent: 1
- $0.012 = 1.2 \times 10^{-2} \in \mathcal{F}$ , however  $2.012 \notin \mathcal{F}$  (Why?)
- Notice  $0.012 = 0.12 \times 10^{-1}$ .
  - Same number may have more than one floating point representation.

# NORMALIZED FLOATING-POINT NUMBERS

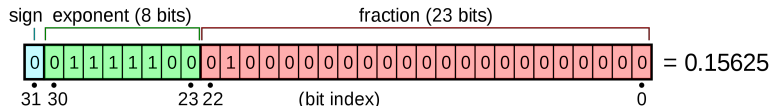
## DEFINITION

The floating point system  $\mathcal{F}(\beta, p, L, U)$  is called normalized if  $d_0 > 0$  for all machine numbers of the form  $x = \pm \left( \sum_{k=0}^{p-1} \frac{d_k}{\beta^k} \right) \beta^E$ .

NORMALIZED: IEEE-754:  $\mathcal{F}(2, 24, -127, 126)$

The 32 bits for single precision are divided as follows: 1 bit for the sign, 8 for the exponent, 23 for the fraction.  
(Where is the full mantissa?)

Here is an example (credit - Wikipedia):



# THE IEEE 754 FLOATING POINT SINGLE PRECISION STANDARD (32 BITS)

## NORMALIZED NUMBERS

- *sign* bit: 0 means positive, 1 means negative
- *exponent*: eeeeeeee is the 8 bits of the exponent

- written in excess 127, that is,
- real exponent = eeeeeeee-127
- for normalized numbers:

$$0[00000000]000000000000000000000000 = 0$$
$$0[11111111]000000000000000000000000 = \infty$$

previous slide:  $01111100 - 11111111 = -11 = (3)_{10}$

- ***mantissa***: 1.mmmmmmmmmmmmmmmmmmmmmmmmmmmmmmm  
  - 24 digits in 23 bits
  - prev. slide example: 1.010000000000000000000000



# THE IEEE 754 FLOATING POINT SINGLE PRECISION STANDARD (32 BITS)

## SPECIAL VALUES

- UFL:  $0[00000001]000000000000000000000000 = 2^{-126}$
- OFL:  $0[11111110]111111111111111111111111 = 2^{127}$
- Narrowing the gap between 0 and UFL, we have *denormalized numbers*:

$0[00000000]mmmmmmmmmmmmmmmmmmmmmmmmmmmmmm$

where not all digits  $m$  can equal zero

exponent switches to excess 126

so the smallest denormalized number is

$0[00000000]0000000000000000000000001 = 2^{-149}$

- NaN (Not a Number), not all  $m$  can equal 0:

$0[11111111]mmmmmmmmmmmmmmmmmmmmmmmmmmmmmm$

# THE IEEE 754 FLOATING POINT DOUBLE PRECISION STANDARD (64 BITS)

## NORMALIZED NUMBERS

- *sign* bit: 0 means positive, 1 means negative
- *exponent*: eeeeeeeeeeee is the 11 bits of the exponent
  - written in excess 1023, that is,
  - real exponent = eeeeeeeeeeee-1023
- *mantissa*: 1.m...m (53 digits in 52 bits)

# REPRESENTING REAL NUMBERS USING FLOATING POINT

## FLOATING-POINT NUMBERS

### THE REPRESENTATION OPERATOR $fl$

- It should be clear we are unable to represent irrational numbers precisely (i.e.  $\sqrt{2}, \pi, e$ ) (Why?)
- Also, many rational numbers have infinite decimal representations (i.e.  $1/3$ ).
- For a real number  $x \in \mathbb{R}$  and a floating-point system  $\mathcal{F}$  we would like  $fl(x)$  to be the element of  $\mathcal{F}$  that minimizes  $|x - y|$ ,  $\forall y \in \mathcal{F}$ .
- Compactness of  $\mathcal{F}$  implies  $fl(x)$  exists. In order to make it unique, we will impose a rounding condition.
- For IEEE-754 the rounding rule used is known as *round to even*.

## ROUND TO EVEN

- IEEE754 standard allows us to store 23 bits of our number within the mantissa.
- At the 23rd bit, we must round to nearest even.
- When rounding, take a look at what the next three bits (after the 23rd) would have been.
- Cases to consider:
  1. If the next (24th) bit is a 0, then you round down (do nothing)
  2. If the next bit is a 1, followed by either a 10, 01, or 11, you round up
  3. If the next three digits are "100" this is a tie (we are midway between two representable numbers). In this case:
    - A. If the last number in the mantissa (23rd bit) is a 1, then round up adding 1 to the mantissa's 0 least significant digit.
    - B. If the last number in the mantissa (23rd bit) is a 0, then do nothing

# ABSOLUTE REPRESENTATION ERROR

## DEFINITION AND CONSEQUENCES

- By definition, absolute representation error is  $fl(x) - x$ .
- If we chose a representation by chopping extra digits,  $fl(x) - x$  has the opposite sign of  $x$ , bounded in the interval  $[-\beta^{1-p}, 0]$  (Exercise: prove this!)
- If we use the round to even rule, the error is bounded in  $[-\beta^{-p}, \beta^{-p}]$
- If we apply chopping, generally the errors during computations tend to accumulate.
- As opposed to that, rounding gives better error propagation results.

$\epsilon_{mach}$ 

## DEFINITION

Smallest positive machine number such that  $fl(1 + \epsilon) > 1$  is called *epsilon of the machine* and denoted  $\epsilon_{mach}$ .

- **Exercise:**  $\epsilon_{mach} = \beta^{1-p}$  if rounding is done by chopping
- **Exercise:**  $\epsilon_{mach} = \frac{1}{2}\beta^{1-p}$  using rounding to the nearest
- **Exercise:** For any nonzero real number  $x$  within the normalized range

$$\left| \frac{fl(x) - x}{x} \right| \leq \frac{1}{2}\beta^{1-p}$$

. You may assume rounding to the nearest.

- Generally

$$0 < UFL < \epsilon_{mach} < OFL$$

# RELATIVE REPRESENTATION ERROR

## DEFINITION AND CONSEQUENCES

- By definition, for any nonzero  $x$ , relative representation error is  $\delta = \frac{fl(x) - x}{x}$ .
- Rearranging,  $fl(x) = (1 + \delta)x$ .
- The triangle inequality implies  $|fl(x)| = |(1 + \delta)x| \leq (1 + |\delta|)|x|$ . Since we aim for the best representation possible, it follows  $|\delta| \leq \epsilon_{mach}$ .

# PROPERTIES OF FLOATING-POINT NUMBERS

## $\mathcal{F}(\beta, p, L, U)$

- $\mathcal{F}$  is a finite set. In particular, a normalized floating-point system contains  $2(\beta - 1)\beta^{p-1}(U - L + 1) + 1$  elements (Why?).
- Largest positive number in  $\mathcal{F}$  is  $(\sum_{k=0}^{p-1} \frac{\beta^k}{\beta^k})\beta^U = \beta^{U+1}(1 - \frac{1}{\beta^p})$  (Why?) [Such number is called *overflow*]
- Smallest positive number in  $\mathcal{F}$  is  $\beta^L$ . [Such number is called *underflow*]



## ARITHMETIC OPERATIONS ON FLOATING-POINT VALUES

 $\mathcal{F}(\beta, p, L, U)$ 

- Usual operations  $(+, -, *, /)$  are not well defined on  $\mathcal{F}$ .
- Addition:  $\beta^{U+1}(1 - \frac{1}{\beta^p}) + 1 \notin \mathcal{F}$
- Division:  $\beta^L/\beta \notin \mathcal{F}$
- Subtraction and multiplication: exercise
- However, for two elements  $x, y \in \mathcal{F}$  we can compute their regular sum, product, etc. Then using the *fl* operator on the result we can obtain a number in  $\mathcal{F}$  which is the best approximation of the result.
- This idea can be extended for any two real number as follows: for any two  $x, y \in \mathbb{R}$  and any operator *op*, it is desirable to have:

$$fl(x) fl(op fl(y)) = fl(x op y).$$

# PROPAGATED ERROR

ASSUME VARIABLES  $x_1, \dots, x_n \in \mathcal{F}(\beta, p, L, U)$

- Consider the statement  $y = f(x_1, \dots, x_n)$ . Assuming  $x_1, \dots, x_n$  are floating point numbers and  $y$  is within representation range,  $fl(y) = (1 + \delta)f(x_1, \dots, x_n)$ .
- Now assume  $x_1, \dots, x_n$  are produced from prior computations. Then, in order to carry out computation of  $f$ , we apply  $f((1 + \delta_1)x_1, \dots, (1 + \delta_n)x_n)$ .
- The relative error into computing  $f$  is called *propagated error*.
- In the next slides, we will compute the propagated error for basic arithmetic operations and for functions of one variable. The ideas laid here certainly may apply to functions of many variables as well.

# FLOATING POINT ADDITION

$fl(x) +_f fl(y) = fl(x + y)$ ,  $x, y$  HAVE SAME SIGN

- As crazy as it sounds, this definition can be made to work, with a few caveats. (Note the plus operator suffixed with "f" indicates floating point addition. )
- The idea is as follows:
  - Step 1, is to rewrite both  $fl(x)$  and  $fl(y)$  using as exponent the maximal exponent of both exponents.
  - Step 2, is to add mantissas, round using the prescribed round rule, and normalize (if working with normalized numbers)
- Example:

$$\begin{array}{rcl}
 1.23450 \times 10^{12} & = & 1.23450 \times 10^{12} \\
 + 4.50000 \times 10^{10} & = & 0.04500 \times 10^{12} \\
 \hline
 & & = 1.27950 \times 10^{12}
 \end{array}$$

## PROPAGATED ERROR ANALYSIS

From before, if  $x, y$  are already in  $\mathcal{F}$ , the relative error produced by addition is  $\delta$ , for  $fl(x + y) = (x + y)(1 + \delta)$ .

As a result,  $fl(x + y) = (x + y)(1 + \delta) = x(1 + \delta) + y(1 + \delta)$ .

That means, letting  $\hat{x} = x(1 + \delta)$ ,  $\hat{y} = y(1 + \delta)$  we see result of floating point addition is exact for the operands  $x, y$  perturbed with a relative backward error  $|\delta|$ .

If two terms  $x, y$  originate from a prior computation, are both nonnegative, with at least one of them strictly positive, so  $fl(x) = x(1 + \delta_1)$ ,  $fl(y) = y(1 + \delta_2)$ , the propagated relative error is

$$\begin{aligned} & \frac{x(1 + \delta_1) + y(1 + \delta_2) - (x + y)}{x + y} \\ &= \frac{x}{x + y} \delta_1 + \frac{y}{x + y} \delta_2 \leq \delta_1 + \delta_2 \end{aligned}$$

# UNEXPCTED CONSEQUENCES OF THE FLOATING POINT OPERATIONS

## SOME (NORMAL ARITHMETIC) DIVERGENT SERIES MAY CONVERGE

Consider the harmonic series

$$\sum_{n=1}^{\infty} \frac{1}{n}.$$

**Exercise:** for any  $M > 0$ , find  $N$ , such that  $\sum_{n=1}^N \frac{1}{n} > M$ . This shows harmonic series diverge.

Now assume  $a_n = \frac{1}{n} \in \mathcal{F}$ . Let  $m$  be the first integer such that  $\frac{1}{m} < \epsilon_{mach}$ . It is not hard to show for all  $n > m$ , the equality  $\sum_{k=1}^n a_k = \sum_{k=1}^{n+1} a_k$  holds.

This implies **the harmonic series, given its general term is computed as a floating-point number, converges!**

**Exercise:** Sum this series in single precision IEEE-754.

# SOME FAMILIAR ALGEBRAIC LAWS NO LONGER HOLD (BUT SOME DO)

```
In [27]: eps = 1
         e = 1
         while 1+e>1:
             eps = e
             e=e/10
         print(eps)
```

1e-15

## ASSOCIATIVE PROPERTY NO LONGER HOLDS

```
In [48]: a = 0.1*eps
         print((((((((((1+a)+a)+a)+a)+a)+a)+a)+a)+a))
         print((1+(a+(a+(a+(a+(a+(a+(a+(a+(a+a)))))))))))
```

1.0

1.0000000000000001

Note, commutative property for the addition of two numbers holds. (Why?)

# REARRANGING TERMS IN AN EXPRESSION CAN HAVE SERIOUS CONSEQUENCES

```
In [49]: print(1+a+a+a+a+a+a+a+a+a)
         print(a+a+a+a+a+a+a+a+a+1)
```

```
1.0
1.0000000000000001
```

```
In [4]: from decimal import *
        getcontext().prec=6
        a,b = Decimal(1.2345e12), Decimal(1.25000e2)
        print(a+b-a, a-a+b)
```

```
0 125
```

# ERROR ANALYSIS FOR THE SUM OF $n$ NUMBERS

Let  $S = a_1 + a_2 + \dots + a_n$ . Then

$$fl(a_1 + a_2) = (a_1 + a_2)(1 + \delta_2)$$

$$fl(a_3 + S_2) = (a_3 + S_2)(1 + \delta_3)$$

...

$$fl(a_n + S_{n-1}) = (a_n + S_{n-1})(1 + \delta_n)$$

Therefore, denoting  $\Delta S$  the absolute error,

$$\Delta S = S_n - S$$

$$= (\delta_2 + \dots + \delta_n)a_1 + (\delta_2 + \dots + \delta_n)a_2 + (\delta_3 + \dots + \delta_n)a_3 + \dots + a_n\delta_n$$

It follows the bound for the relative error is

$|\Delta S/S| \leq (n-1)\varepsilon_{mach}$ . Also the error can be minimized if we multiply the larger sum of  $\delta_k$ , that is the coefficient to the first two factors, with smallest term. It follows we minimize the error if the terms of the sum  $S$  are sorted in increasing order.



# FLOATING POINT SUBTRACTION

The definition and implementation of subtraction is similar to the addition. However subtraction comes out with a serious problem.

## A SIMPLE EXAMPLE

Let  $0 < \varepsilon < \varepsilon_{mach}$ .

$$(1 + \varepsilon) - (1 - \varepsilon) = 0 \neq 2\varepsilon$$

Please note, both  $1 + \varepsilon$  and  $1 - \varepsilon$  values are close to one another. This phenomenon is known as *catastrophic cancellation*. Its effects can be mitigated by modifying the computational method in use.

**Exercise:** Do the propagated error analysis for subtraction.

# CATASTROFIC CANCELLATION

## EXAMPLE - QUADRATIC FORMULA

```
In [10]: from math import sqrt
         from decimal import *

         getcontext().prec = 4

         def solve_quad(a, b, c):
             a, b, c = Decimal(a), Decimal(b), Decimal(c)
             two, four = Decimal(2.0), Decimal(4.0)
             d = Decimal(sqrt(b*b-four*a*c))
             return (-b+d)/two/a, (-b-d)/two/a

         a, b, c = 0.05010, -98.78, 5.015
         #
         # Correct roots are 1971.605916, 0.05077069387
         #
         r1, r2 = solve_quad(a, b, c)
         print("Naively computed roots are: {}, {}".format(r1, r2))
```

Naively computed roots are: 1972, 0.07519

```
In [11]: #
         # Compare:
         #
         x1, x2 = 0.05077069387, 0.07519
         print(a*x1*x1+b*x1+c, a*x2*x2+b*x2+c)
```

4.5553694150157753e-10 -2.4119849578413906

# CATASTROFIC CANCELLATION

## IMPROVED QUADRATIC FORMULA

Clearly,  $-b$  and  $d$  have similar values. Let's take on this case and avoid subtraction:

$$\begin{aligned}\frac{-b - \sqrt{b^2 - 4ac}}{2a} &= \frac{(-b - \sqrt{b^2 - 4ac})(-b + \sqrt{b^2 - 4ac})}{2a(-b + \sqrt{b^2 - 4ac})} \\ &= \frac{2c}{-b + \sqrt{b^2 - 4ac}}\end{aligned}$$

The formula for the other root does not change.

**Exercise:** Write a Python function that computes the roots of quadratic equation avoiding catastrophic cancellation.

# FLOATING POINT MULTIPLICATION

## OVERFLOW

Although addition can result in overflow, multiplication is more prone to it. As an example, consider computation of Euclidean distance of a point from the origin:

$$d(x, y) = \sqrt{x^2 + y^2}$$

A straight forward implementation of it would suffer from unnecessary overflow. Consider instead, letting  $m = \max\{x, y\}$ :

$$d(x, y) = m \sqrt{\left(\frac{x}{m}\right)^2 + \left(\frac{y}{m}\right)^2}$$

Since both  $\frac{x}{m}, \frac{y}{m}$  are  $\leq 1$  we avoid the overflow as much as possible. Also note one of the terms is precisely 1. Exploit this to do an error analysis for this formula.

# FLOATING POINT MULTIPLICATION

## DEFINITION

Let the floating point representation for  $x, y$  be respectively  $fl(x), fl(y)$ . The floating point product of  $x$  and  $y$  is the floating point representation of the product  $fl(x)fl(y)$ . Then  $fl(fl(x)fl(y)) = fl(x)fl(y)(1 + \delta)$  where  $|\delta| \leq \epsilon_{mach}$ .

## RELATIVE ERROR ANALYSIS

$$\begin{aligned} r(xy) &= \frac{fl(fl(x)fl(y)) - xy}{xy} = \frac{fl(x)fl(y)(1 + \delta) - xy}{xy} \\ &= \frac{fl(x)fl(y) - xy}{xy} + \frac{fl(x)fl(y)}{xy}\delta \end{aligned}$$

Since  $fl(x)fl(y) \approx xy$ , we may conclude

$$r(xy) = \frac{fl(x)fl(y) - xy}{xy} + \delta$$

# FLOATING POINT MULTIPLICATION

## PROPAGATED ERROR

Let relative representation errors for  $x, y$  be

$r(x) = \frac{fl(x) - x}{x}$ ,  $r(y) = \frac{fl(y) - y}{y}$  and denote propagated error  $p(xy)$ .

Compute:

$$p(xy) = \frac{fl(x)fl(y) - xy}{xy} = \frac{(xr(x) + x)(yr(y) + y) - xy}{xy}$$

After expanding and cancelling, we find

$$p(xy) = r(x) + r(y) + r(x)r(y)$$

Ignoring the  $r(x)r(y)$  (what is its bound ?) we conclude

$$p(xy) = r(x) + r(y)$$

# FLOATING POINT DIVISION

## EXERCISE

We leave the error analysis for the floating point division as an exercise. In particular, you need to show that the propagated error for  $\frac{x}{y}$  is  $r(x) - r(y)$  where  $r(x), r(y)$  are the floating point relative errors of representations of  $x, y$ .

# ERROR ANALYSIS FOR THE COMPUTATION OF A SMOOTH FUNCTION $f(x)$

## PROPAGATED ERROR

Using Taylor's expansion,

$$fl(f(x)) = f(fl(x)) = f(x(1 + \delta)) = f(x) + x\delta f'(\xi)$$

where  $x \leq \xi \leq (1 + \delta)x$ . The relative error of computing  $f$  is bounded by

$$\epsilon_{mach} \left| \frac{xf'(x)}{f(x)} \right|$$



# EXERCISES

- ❶ Evaluate the relative error of computing  $f(x) = \frac{1}{x}$ .
- ❷ Assume we need to compute the list of numbers  $x_i = a + ih$ , where  $a, h$  are fixed, and  $i$  varies from 1 to  $n$ . Which one would you prefer to use given the computations are done in floating-point arithmetic?

(I) `x=[a+i*h for i in range(1,n+1)]`

(II) `x = []`  
`val = a`  
`i = 0`  
`while i < n:`  
    `val = val + h`  
    `x.append(val)`

# EXERCISES

- 1 Evaluate the relative error of computing dot product of two vectors:  $\vec{a} \cdot \vec{b} = \sum_{i=1}^n a_i b_i$ . Now suppose the coordinates are given in single precision and the answer is also expected to be in single precision. Give a solution with total relative error  $\varepsilon_{mach}$ .

# REFERENCES

Michael T. Heath  
Scientific Computing (Revised Second Edition)  
SIAM