

## Week 1: Wednesday, Jan 25

### Binary floating point encodings

Binary floating point arithmetic is essentially scientific notation. Where in decimal scientific notation we write

$$\frac{1}{3} = 3.333\dots \times 10^{-1},$$

in floating point, we write

$$\frac{(1)_2}{(11)_2} = (1.010101\dots)_2 \times 2^{-2}.$$

Because computers are finite, however, we can only keep a finite number of bits after the binary point.

In general, a *normal floating point number* has the form

$$(-1)^s \times (1.b_1b_2\dots b_p)_2 \times 2^E,$$

where  $s \in \{0, 1\}$  is the *sign bit*,  $E$  is the *exponent*, and  $(1.b_2\dots b_p)_2$  is the *significand*. In the 64-bit double precision format,  $p = 52$  bits are used to store the significand, 11 bits are used for the exponent, and one bit is used for the sign. The valid exponent range for normal floating point numbers is  $-1023 < E < 1024$ ; this leaves two exponent encodings left over for special purpose. One of these special exponents is used to encode *subnormal numbers* of the form

$$(-1)^s \times (0.b_1b_2\dots b_p)_2 \times 2^{-1022};$$

the other special exponent is used to encode  $\pm\infty$  and NaN (Not a Number).

For a general real number  $x$ , we will write

$$\text{fl}(x) = \text{correctly rounded floating point representation of } x.$$

By default, “correctly rounded” means that we find the closest floating point number to  $x$ , breaking any ties by rounding to the number with a zero in the last bit<sup>1</sup>. If  $x$  exceeds the largest normal floating point number, then  $\text{fl}(x) = \infty$ .

---

<sup>1</sup>There are other rounding modes beside the default, but we will not discuss them in this class

## Basic floating point arithmetic

For basic operations (addition, subtraction, multiplication, division, and square root), the floating point standard specifies that the computer should produce the *true result, correctly rounded*. So the MATLAB statement

```
% Compute the sum of x and y (assuming they are exact)
z = x + y;
```

actually computes the quantity  $\hat{z} = \text{fl}(x+y)$ . If  $\hat{z}$  is a normal double-precision floating point number, it will agree with the true  $z$  to 52 bits after the binary point. That is, the relative error will be smaller in magnitude than the *machine epsilon*  $\epsilon_{\text{mach}} = 2^{-53} \approx 1.1 \times 10^{-16}$ :

$$\hat{z} = z(1 + \delta), \quad |\delta| < \epsilon_{\text{mach}}.$$

More generally, basic operations that produce normalized numbers are correct to within a relative error of  $\epsilon_{\text{mach}}$ . The floating point standard also recommends that common transcendental functions, such as exponential and trig functions, should be correctly rounded, though compliant implementations that do not comply with this recommendation may produce results with a relative error just slightly larger than  $\epsilon_{\text{mach}}$ .

The fact that normal floating point results have a relative error less than  $\epsilon_{\text{mach}}$  gives us a useful *model* for reasoning about floating point error. We will refer to this as the “ $1 + \delta$ ” model. For example, suppose  $x$  is an exactly-represented input to the MATLAB statement

```
z = 1-x*x;
```

We can reason about the error in the computed  $\hat{z}$  as follows:

$$\begin{aligned} t_1 &= \text{fl}(x^2) = x^2(1 + \delta_1) \\ t_2 &= 1 - t_1 = (1 - x^2) \left( 1 + \frac{\delta_1 x^2}{1 - x^2} \right) \\ \hat{z} &= \text{fl}(1 - t_1) = z \left( 1 + \frac{\delta_1 x^2}{1 - x^2} \right) (1 + \delta_2) \\ &\approx z \left( 1 + \frac{\delta_1 x^2}{1 - x^2} + \delta_2 \right), \end{aligned}$$

where  $|\delta_1|, |\delta_2| \leq \epsilon_{\text{mach}}$ . As before, we throw away the (tiny) term involving  $\delta_1 \delta_2$ . Note that if  $z$  is close to zero (i.e. if there is *cancellation* in the subtraction), then the model shows the result may have a large relative error.

## Exceptions

We say there is an *exception* when the floating point result is not an ordinary value that represents the exact result. The most common exception is *inexact* (i.e. some rounding was needed). Other exceptions occur when we fail to produce a normalized floating point number. These exceptions are:

**Underflow:** An expression is too small to be represented as a normalized floating point value. The default behavior is to return a subnormal.

**Overflow:** An expression is too large to be represented as a floating point number. The default behavior is to return `inf`.

**Invalid:** An expression evaluates to Not-a-Number (such as  $0/0$ )

**Divide by zero:** An expression evaluates “exactly” to an infinite value (such as  $1/0$  or  $\log(0)$ ).

When exceptions other than *inexact* occur, the usual “ $1 + \delta$ ” model used for most rounding error analysis is not valid.

## Limits of the “ $1 + \delta$ ” model

Apart from the fact that it fails when there are exceptions other than *inexact*, the “ $1 + \delta$ ” model of floating point does not reflect the fact that some computations involve no rounding error. For example:

- If  $x$  and  $y$  are floating point numbers within a factor of two of each other,  $\text{fl}(x - y)$  is computed without rounding error.
- Barring overflow or underflow to zero,  $\text{fl}(2x) = 2x$  and  $\text{fl}(x/2) = x/2$ .
- Integers between  $\pm(2^{53} - 1)$  are represented exactly.

These properties of floating point allow us to do some clever things, such as using ordinary double precision arithmetic to simulate arithmetic with about twice the number of digits. You should be aware that these tricks exist, even if you never need to implement them – otherwise, I may find myself cursing a compiler you wrote for rearranging the computations in a floating point code that I wrote!

## Finding and fixing floating point problems

Floating point arithmetic is not the same as real arithmetic. Even simple properties like associativity or distributivity of addition and multiplication only hold approximately. Thus, some computations that look fine in exact arithmetic can produce bad answers in floating point. What follows is a (very incomplete) list of some of the ways in which programmers can go awry with careless floating point programming.

### Cancellation

If  $\hat{x} = x(1 + \delta_1)$  and  $\hat{y} = y(1 + \delta_2)$  are floating point approximations to  $x$  and  $y$  that are very close, then  $\text{fl}(\hat{x} - \hat{y})$  may be a poor approximation to  $x - y$  due to *cancellation*. In some ways, the subtraction is blameless in this tail: if  $x$  and  $y$  are close, then  $\text{fl}(\hat{x} - \hat{y}) = \hat{x} - \hat{y}$ , and the subtraction causes no additional rounding error. Rather, the problem is with the approximation error already present in  $\hat{x}$  and  $\hat{y}$ .

The standard example of loss of accuracy revealed through cancellation is in the computation of the smaller root of a quadratic using the quadratic formula, e.g.

$$x = 1 - \sqrt{1 - z}$$

for  $z$  small. Fortunately, some algebraic manipulation gives an equivalent formula that does not suffer cancellation:

$$x = (1 - \sqrt{1 - z}) \left( \frac{1 + \sqrt{1 - z}}{1 + \sqrt{1 - z}} \right) = \frac{z}{1 + \sqrt{1 - z}}.$$

### Sensitive subproblems

We often solve problems by breaking them into simpler subproblems. Unfortunately, it is easy to produce badly-conditioned subproblems as steps to solving a well-conditioned problem. As a simple (if contrived) example, try running the following MATLAB code:

---

```
x = 2;
for k = 1:60, x = sqrt(x); end
for k = 1:60, x = x^2;     end
disp(x);
```

---

In exact arithmetic, this should produce 2, but what does it produce in floating point? In fact, the first loop produces a correctly rounded result, but the second loop represents the function  $x^{2^{60}}$ , which has a condition number far greater than  $10^{16}$  — and so all accuracy is lost.

## Unstable recurrences

We gave an example of this problem in the last lecture notes when we looked at the recurrence

$$\begin{aligned} E_0 &= 1 - 1/e \\ E_n &= 1 - nE_{n-1}, \quad n \geq 1. \end{aligned}$$

No single step of this recurrence causes the error to explode, but each step amplifies the error somewhat, resulting in an exponential growth in error.

## Undetected underflow

In Bayesian statistics, one sometimes computes ratios of long products. These products may underflow individually, even when the final ratio is not far from one. In the best case, the products will grow so tiny that they underflow to zero, and the user may notice an infinity or NaN in the final result. In the worst case, the underflowed results will produce nonzero subnormal numbers with unexpectedly poor relative accuracy, and the final result will be wildly inaccurate with no warning except for the (often ignored) underflow flag.

## Bad branches

A NaN result is often a blessing in disguise: if you see an unexpected NaN, at least you *know* something has gone wrong! But all comparisons involving NaN are false, and so when a floating point result is used to compute a branch condition and an unexpected NaN appears, the result can wreak havoc. As an example, try out the following code in MATLAB.

---

```
x = 0/0;
if x < 0 then      disp('x_is_negative');
elseif x >= 0 then disp('x_is_non-negative');
else              disp('Uh... ');
end
```

---

## Problems to ponder

1. In double precision, is  $\text{fl}(0.2)$  larger, smaller, or equal to 0.2?
2. How do we accurately evaluate  $\sqrt{1+x} - \sqrt{1-x}$  when  $x \ll 1$ ?
3. How do we accurately evaluate  $\ln \sqrt{x+1} - \ln \sqrt{x}$  when  $x \gg 1$ ?
4. How do we accurately evaluate  $(1 - \cos(x))/\sin(x)$  when  $x \ll 1$ ?
5. How would we compute  $\cos(x) - 1$  accurately when  $x \ll 1$ ?
6. The *Lamb-Oseen vortex* is a solution to the 2D Navier-Stokes equation that plays a key role in some methods for computational fluid dynamics. It has the form

$$v_{\theta}(r, t) = \frac{\Gamma}{2\pi r} \left( 1 - \exp\left(\frac{-r^2}{4\nu t}\right) \right)$$

How would one evaluate  $v(r, t)$  to high relative accuracy for all values of  $r$  and  $t$  (barring overflow or underflow)?

7. For  $x > 1$ , the equation  $x = \cosh(y)$  can be solved as

$$y = -\ln\left(x - \sqrt{x^2 - 1}\right).$$

What happens when  $x = 10^8$ ? Can we fix it?

8. The difference equation

$$x_{k+1} = 2.25x_k - 0.5x_{k-1}$$

with starting values

$$x_1 = \frac{1}{3}, \quad x_2 = \frac{1}{12}$$

has solution

$$x_k = \frac{4^{1-k}}{3}.$$

Is this what you actually see if you compute? What goes wrong?

9. Considering the following two MATLAB fragments:

---

```
% Version 1
f = (exp(x)-1)/x;

% Version 2
y = exp(x);
f = (1-y)/log(y);
```

---

In exact arithmetic, the two fragments are equivalent. In floating point, the first formulation is inaccurate for  $x \ll 1$ , while the second formulation remains accurate. Why?

- Running the recurrence  $E_n = 1 - nE_{n-1}$  *forward* is an unstable way to compute  $\int_0^1 x^n e^{x-1} dx$ . However, we can get good results by running the recurrence *backward* from the estimate  $E_n \approx 1/(N+1)$  starting at large enough  $N$ . Explain why. How large must  $N$  be to compute  $E_{20}$  to near machine precision?
- How might you accurately compute this function for  $|x| < 1$ ?

$$f(x) = \sum_{j=0}^{\infty} (\cos(x^j) - 1)$$