

chapter 1 : finite precision arithmetic

number systems

$$\text{decimal} : (2013)_{10} = 2 \cdot 10^3 + 0 \cdot 10^2 + 1 \cdot 10^1 + 3 \cdot 10^0$$

$$\text{binary} : (101.01)_2 = 1 \cdot 2^2 + 0 \cdot 2^1 + 1 \cdot 2^0 + 0 \cdot 2^{-1} + 1 \cdot 2^{-2} = (5.25)_{10}$$

question : how are numbers stored in a computer?

answer : floating point representation

$$x = \pm(0.d_1d_2 \cdots d_n)_\beta \cdot \beta^e, \quad n \text{ significant digits}, \quad d_1 \neq 0 : \text{normalization}$$

$$\beta : \text{base}, \quad d_i : \text{digits (or bits if } \beta = 2), \quad 0 \leq d_i \leq \beta - 1$$

$$(0.d_1d_2 \cdots d_n) : \text{mantissa}, \quad e : \text{exponent}, \quad -M \leq e \leq M$$

ex : consider a computer with  $\beta = 2$ ,  $n = 4$ ,  $M = 3$

$$x_{\max} = (0.1111)_2 \cdot 2^3 = \left(\frac{1}{2} + \frac{1}{4} + \frac{1}{8} + \frac{1}{16}\right) \cdot 8 = 4 + 2 + 1 + \frac{1}{2} = 7.5$$

$$x_{\min} = (0.1000)_2 \cdot 2^{-3} = \frac{1}{2} \cdot \frac{1}{8} = 0.0625$$

question : how many different numbers can be represented in this system?

answer : hw

note

1. In IEEE double precision format, each number is stored as a string of 64 bits.

±	±	mantissa = 52 bits	exponent = 10 bits
---	---	--------------------	--------------------

The first 2 bits are for the sign of the mantissa and exponent. Hence we have  $\beta = 2, n = 52, M = (111111111)_2 = 1 + 2 + 4 + 8 + \cdots + 2^9 = 2^{10} - 1 = 1023$ .

2. If  $x$  is a number and  $\text{fl}(x)$  is its floating point representation, then  $x - \text{fl}(x)$  is the roundoff error.

ex 1

$$\pi = 3.14159\ 26535\ 89793\ 23846 \dots : \text{exact value to 20 decimal digits}$$

$$\text{fl}(\pi) = 3.14159\ 26535\ 89793 : \text{Matlab gives 15 decimal digits}$$

$$\Rightarrow \text{roundoff error} = 0.23846 \dots \cdot 10^{-15} \approx 2^{-52}$$

ex 2

$$0.1234 - 0.1233 = 0.0001 = (0.1000)_{10} \cdot 10^{-3}$$

$\Rightarrow$  the result has only 1 significant digit

Hence if 2 floating point numbers with  $n$  significant digits are subtracted, the result may have less than  $n$  significant digits; this is called loss of significance due to cancellation of digits.

ex 3 : quadratic formula

$$ax^2 + bx + c = 0 \Rightarrow x = \frac{-b \pm \sqrt{b^2 - 4ac}}{2a}, \quad \text{pf : ...}$$

$$0.2x^2 - 47.91x + 6 = 0 \Rightarrow x = 239.4247, 0.1253 : \text{Matlab}$$

Now use 4-digit arithmetic (i.e. each step is rounded to 4 digits, as if we're using a computer with  $n = 4$ ).

$$\begin{aligned} x &= \frac{47.91 \pm \sqrt{47.91^2 - 4(0.2)6}}{2(0.2)} = \frac{47.91 \pm \sqrt{2295 - 4.8}}{0.4} = \frac{47.91 \pm \sqrt{2290}}{0.4} \\ &= \frac{47.91 \pm 47.85}{0.4} = \begin{cases} \frac{47.91 + 47.85}{0.4} = \frac{95.76}{0.4} = 239.4 : \text{all 4 digits are correct} \\ \frac{47.91 - 47.85}{0.4} = \frac{0.06}{0.4} = 0.15 : \text{only 1 digit is correct} \end{cases} \end{aligned}$$

question : what caused the problem?

answer : loss of significance in the subtraction  $47.91 - 47.85$

remedy 1 : Matlab (higher precision arithmetic)

remedy 2 : reformulate the arithmetic

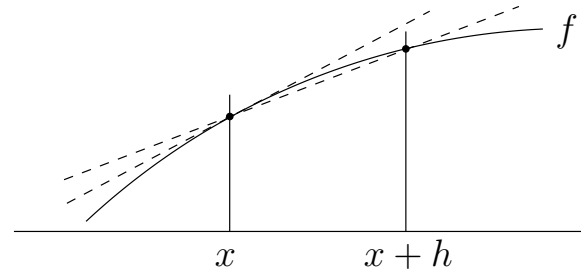
$$\begin{aligned} x &= \frac{-b - \sqrt{b^2 - 4ac}}{2a} \cdot \frac{-b + \sqrt{b^2 - 4ac}}{-b + \sqrt{b^2 - 4ac}} = \frac{b^2 - (b^2 - 4ac)}{2a(-b + \sqrt{b^2 - 4ac})} = \frac{2c}{-b + \sqrt{b^2 - 4ac}} \\ &= \frac{2 \cdot 6}{47.91 + 47.85} = \frac{12}{95.76} = 0.1253 : \text{now all 4 digits are correct} \end{aligned}$$

remedy 3 : solve for  $x$  using an iterative method (more later)

ex 4 : finite-difference approximation of a derivative

forward difference

$$f'(x) \approx \frac{f(x+h) - f(x)}{h} = D_+f(x)$$



2
Tues
1/15

question : how large is the error?

Taylor series :  $f(x) = f(a) + f'(a)(x-a) + \frac{1}{2}f''(a)(x-a)^2 + \dots$

equivalent form :

$$\left. \begin{array}{l} x \rightarrow x+h \\ a \rightarrow x \end{array} \right\} \Rightarrow f(x+h) = f(x) + f'(x)h + \frac{1}{2}f''(x)h^2 + \dots$$

$$\Rightarrow \frac{f(x+h) - f(x)}{h} = f'(x) + \underbrace{\frac{1}{2}f''(x)h + \dots}_{\text{truncation error}}$$

$\uparrow$                        $\uparrow$                        $\uparrow$   
 approximation      exact value      truncation error

Hence the error is proportional to  $h$ ; we write this as  $D_+f(x) = f'(x) + O(h)$ , where the symbol  $O(h)$  means “order  $h$ ”.

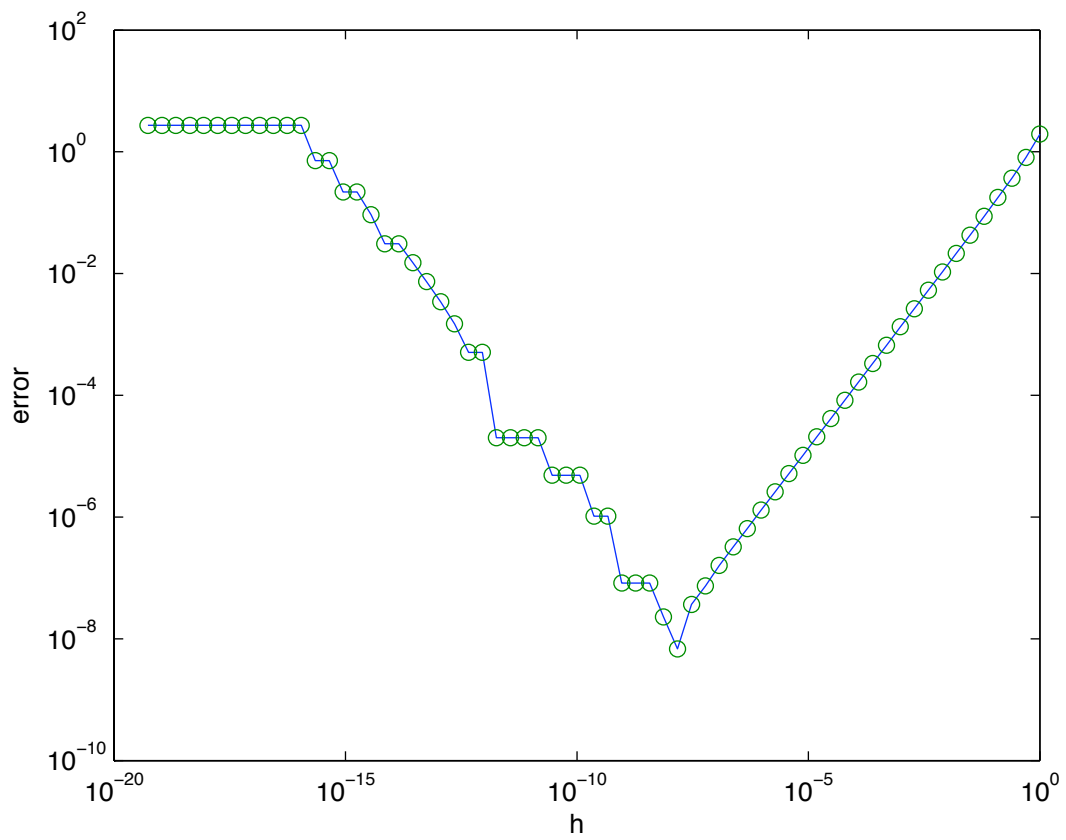
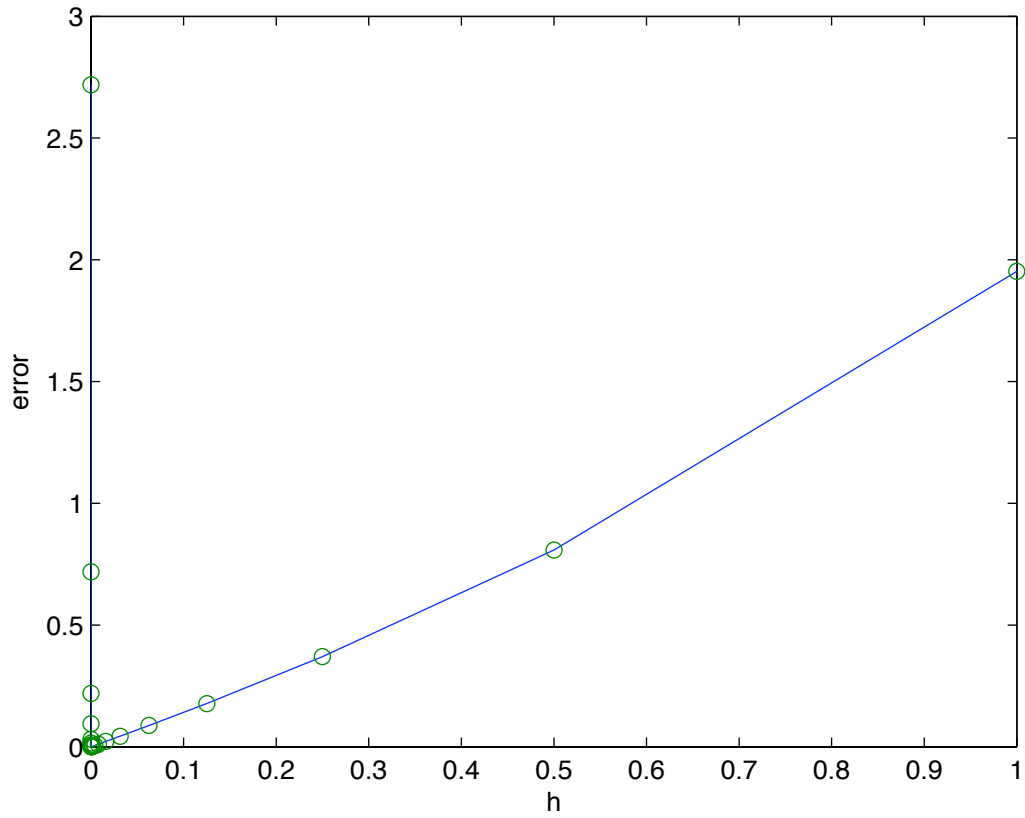
For example, if  $f(x) = e^x$ ,  $x = 1$ , then  $f'(1) = e = 2.71828\dots$  is the exact value.

$h$	$D_+f(1)$	$ D_+f(1) - f'(1) $	$ D_+f(1) - f'(1) /h$
0.1	2.8588	0.1406	1.4056
0.05	2.7874	0.0691	1.3821
0.025	2.7525	0.0343	1.3705
↓	↓	↓	↓
0	$e$	0	$\frac{e}{2} = \frac{1}{2}f''(1)$

% Matlab

```

exact_value = exp(1);
for j=1:65
    h(j) = 1/2^(j-1);
    computed_value = (exp(1+h(j)) - exp(1))/h(j);
    error(j) = abs(computed_value - exact_value);
end
plot(h,error,h,error,'o'); xlabel('h'); ylabel('error')
loglog(h,error,h,error,'o'); ...
  
```

finite-difference approximation of a derivative

note

If error  $\approx ch^p$ , then  $p$  is called the order of accuracy of the approximation. In this case we have  $\log(\text{error}) \approx \log(ch^p) = \log c + p \log h$ ; this implies that the slope of the data on the log-log plot gives the value of  $p$ . We see that  $p = 1$  for large  $h$  (expected) and  $p = -1$  for small  $h$  (unexpected).

question : why does the error increase for small  $h$ ?

1. The computed value has two sources of error: truncation error is due to replacing the exact derivative  $f'(x)$  by the finite-difference approximation  $D_+f(x)$ , and roundoff error is due to using finite precision arithmetic.
2. The truncation error is  $O(h)$  and the roundoff error is  $O(\epsilon/h)$ , where  $\epsilon \approx 10^{-15}$  in Matlab.
3. The total error is  $O(h) + O(\epsilon/h)$ . Hence, for large  $h$  the truncation error dominates the roundoff error, but for small  $h$  the roundoff error dominates the truncation error.

note

backward difference : 
$$D_-f(x) = \frac{f(x) - f(x-h)}{h}$$

centered difference : 
$$D_0f(x) = \frac{f(x+h) - f(x-h)}{2h} \dots \text{hw}$$